

# Maximum mutual information regularized classification

Jim Jing-Yan Wang<sup>a</sup>, Yi Wang<sup>b</sup>, Shiguang Zhao<sup>c</sup>, Xin Gao<sup>a,\*</sup>

<sup>a</sup>*Computer, Electrical and Mathematical Sciences and Engineering Division, King Abdullah University of Science and Technology (KAUST), Thuwal 23955-6900, Saudi Arabia*

<sup>b</sup>*Department of Computer Science and Engineering, The Ohio State University, Columbus OH 43210, USA*

<sup>c</sup>*Department of Neurosurgery, The First Affiliated Hospital of Harbin Medical University, Harbin, Heilongjiang 150001, P.R. China*

---

## Abstract

In this paper, a novel pattern classification approach is proposed by regularizing the classifier learning to maximize mutual information between the classification response and the true class label. We argue that, with the learned classifier, the uncertainty of the true class label of a data sample should be reduced by knowing its classification response as much as possible. The reduced uncertainty is measured by the mutual information between the classification response and the true class label. To this end, when learning a linear classifier, we propose to maximize the mutual information between classification responses and true class labels of training samples, besides minimizing the classification error and reducing the classifier complexity. An objective function is constructed by modeling mutual information with entropy estimation, and it is optimized by a gradient descend method in an iterative algorithm. Experiments on two real world pattern classification problems show the significant improvements achieved by maximum mutual information regularization.

*Keywords:* Pattern Classification, Maximum Mutual Information, Entropy, Gradient Descend

---

## 1. Introduction

The pattern classification problem is a problem of assigning a discrete class label to a given data sample represented by its feature vector [1, 2, 3, 4, 5]. It has many applications in various fields, including bioinformatics [6, 7, 8, 9, 10], biometrics verification [11, 12, 13], computer networks [14, 15, 16], and computer vision [17, 18, 19, 20]. For example, in the face recognition problem, given a

---

\*To whom all correspondence should be addressed. Tel: +966-12-8080323.

*Email addresses:* jimjywang@gmail.com (Jim Jing-Yan Wang),  
wayi@cse.ohio-state.edu (Yi Wang), guangsz@hotmail.com (Shiguang Zhao),  
xin.gao@kaust.edu.sa (Xin Gao)

face image, the target of pattern classification is to assign it to a person who has been registered in a database [21, 22]. This problem is usually composed of two different components — feature extraction [23, 24, 25, 26, 27, 28, 29, 30] and classification [31, 32]. Feature extraction refers to the procedure of extracting an effective and discriminant feature vector from a data sample, so that different samples of different classes could be separated easily. This procedure is usually highly domain-specific. For example, for the face recognition problem, the visual feature should be extracted using some image processing technologies, whereas for the problem of predicting zinc-binding sites from protein sequences, the biological features should be extracted using some biological knowledge [33]. In terms of feature extraction of this paper, it is highly inspired by a hierarchical Bayesian inference algorithm proposed in [24]. This new method created in [23] has advanced the ground-truth feature extraction field and has provided a more optimal method for this procedure. On the other hand, different from feature extraction, classification is a much more general problem. We usually design a class label prediction function as a classifier for this purpose. To learn the parameter of a classifier function, we usually try to minimize the classification error of the training samples in a training set and simultaneously reduce the complexity of the classifier. For example, the most popular classifier is support vector machine (SVM), which minimizes the hinge losses to reduce the classification error, and at the same time minimizes the  $\ell_2$  norm of the classifier parameters to reduce the complexity. In this paper, we focus on the classification aspect.

Mutual information [34, 35] is defined as the information shared between two sets of variables. It has been used as a criterion of feature extraction for pattern classification problems [36]. However, surprisingly, it has never been directly explored in the problem of classifier learning. Actually, mutual information has a strong relation to Kullback-Leibler divergence, and there are many works using KL-divergence for classifiers [37, 38]. Moreno et al. [37] proposed a novel kernel function for support vector classification based on Kullback-Leibler divergence, while Liu and Shum [38] proposed to learn the most discriminating feature that maximizes the Kullback-Leibler divergence for the Adaboost classifier. However, both these methods do not use the KL-divergence based criterion to learn parameters of linear classifiers. To bridge this gap, in this paper, for the first time, we try to investigate using mutual information as a criterion of classifier learning. We propose to learn a classifier by maximizing the mutual information  $I(f; y)$  between the classification response variable  $f$  and the true class label variable  $y$ . The classification response variable  $f$  is a function of classifier parameters and data samples. The insight is that mutual information is defined as the information shared between  $f$  and  $y$ . From the viewpoint of information theory, if the two variables are not mutually independent, and one variable is known, it usually reduces the uncertainty about the other one. Then mutual information is used to measure how much uncertainty is reduced in this case. To illuminate how the mutual information can be used to measure the classification accuracy, we consider the two extreme cases:

- On one hand, if the classification response variable  $f$  of a data sample

is randomly given, and it is independent of its true class label  $y$ , then knowing  $f$  does not give any information about  $y$  and vice versa, and the mutual information between them could be zero, i.e.,  $I(f; y) = 0$ .

- On the other hand, if  $f$  is given so that  $y$  and  $f$  are identical, knowing  $f$  can help determine the value of  $y$  exactly as well as reduce all the uncertainty about  $y$ . This is the ideal case of classification, and knowing  $f$  can reduce all the uncertainty about  $y$ . In this case, the mutual information is defined as the uncertainty contained in  $f$  (or  $y$ ) alone, which is measured by the entropy of  $f$  or  $y$ , denoted by  $H(f)$  or  $H(y)$  respectively, where  $H(\cdot)$  is the entropy of a variable. Since  $f$  and  $y$  are identical, we can have  $I(f; y) = H(f) = H(y)$ .

Naturally, we hope that the classification response  $f$  can predict the true class label  $y$  as accurately as possible, and knowing  $f$  can reduce the uncertainty about  $y$  as much as possible. Thus, we propose to maximize the mutual information between  $f$  and  $y$  with regard to the parameters of a classifier. To this end, we proposed a mutual information regularization term for the learning of classifier parameters. An objective function is constructed by combining the mutual information regularization term, a classification error term and a classifier complexity term. The classifier parameter is learned by optimizing the objective function with a gradient descend method in an iterative algorithm.

The rest parts of this paper are organized as follows: in Section 2, we introduce the proposed classifier learning method. The experiment results are presented in section 3. In section 4 the paper is concluded.

## 2. Proposed Method

In this section, we introduce the proposed classifier learning algorithm to maximize the mutual information between the classification response and the true class label.

### 2.1. Problem Formulation

We suppose that we have a training set denoted as  $X = \{\mathbf{x}_i\}_{i=1}^n$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  is the  $d$ -dimensional feature vector for the  $i$ -th training sample, and  $n$  is the number of training samples. The class label set for the training samples is denoted as  $Y = \{y_i\}_{i=1}^n$ , where  $y_i \in \{+1, -1\}$  is the class label of the  $i$ -th sample. To learn a classifier to predict the class label of a given sample with its feature vector  $\mathbf{x}$ , we design a linear function as a classifier,

$$g(\mathbf{x}; \mathbf{w}) = \text{sign}(f) = \text{sign}(\mathbf{w}^\top \mathbf{x}), \quad (1)$$

where  $\mathbf{w}$  is the classifier parameter vector,  $f = \mathbf{w}^\top \mathbf{x}$  is the classification response of  $\mathbf{x}$  given the classifier parameter  $\mathbf{w}$ , and  $\text{sign}(\cdot)$  is the signum function which transfers the classification response to the final binary classification result. We also denote the classification response set of the training samples as  $F = \{f_i\}_{i=1}^n$

where  $f_i = \mathbf{w}^\top \mathbf{x}_i \in \mathbb{R}$  is the classification response of the  $i$ -th training sample. To learn the optimal classification parameter  $\mathbf{w}$  for the classification problem, we consider the following three problems:

### 2.1.1. Classification Loss Minimization

To learn the optimal classification parameter  $\mathbf{w}$ , we hope the classification response  $f$  of a data sample  $\mathbf{x}$  obtained with the learned  $\mathbf{w}$  can predict its true class label  $y$  as accurately as possible. To measure the prediction error, we use a loss function to compare a classification response against its corresponding true class label. Given the classifier parameter  $\mathbf{w}$ , the loss function of the  $i$ -th training sample  $\mathbf{x}_i$  with its classification response  $f_i = \mathbf{w}^\top \mathbf{x}_i$  and true class label  $y_i$  is denoted as  $L(f_i, y_i; \mathbf{w})$ . There are a few different loss functions which could be considered.

**Hinge Loss** is used by the SVM classifier [39, 40, 41], and it is defined as

$$L(f_i, y_i; \mathbf{w}) = \max(0, 1 - y_i f_i) = \tau_i \times (1 - y_i \mathbf{w}^\top \mathbf{x}_i), \quad (2)$$

where  $\tau_i$  is defined as

$$\tau_i = \begin{cases} 1, & \text{if } y_i \mathbf{w}^\top \mathbf{x}_i \leq 1 \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

**Squared Loss** is usually used by regression problems [42, 43, 44], and it is defined as

$$L(f_i, y_i; \mathbf{w}) = (1 - y_i f_i)^2 = (1 - y_i \mathbf{w}^\top \mathbf{x}_i)^2. \quad (4)$$

**Logistic Loss** is defined as follows, and it is also popular in regression problems [45],

$$L(f_i, y_i; \mathbf{w}) = \log [1 + \exp(-y_i f_i)] = \log [1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)]. \quad (5)$$

**Exponential loss** is another popular loss function which could be used by both classification and regression problems [42], which is defined as

$$L(f_i, y_i; \mathbf{w}) = \exp(-y_i f_i) = \exp(-y_i \mathbf{w}^\top \mathbf{x}_i). \quad (6)$$

Obviously, to learn an optimal classifier, the average loss of all the training samples should be minimized with regard to  $\mathbf{w}$ . Thus the following optimization problem is obtained by applying a loss functions to all training samples,

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n L(f_i, y_i; \mathbf{w}). \quad (7)$$

2.1.2. *Classifier Complexity Reduction*

To reduce the complexity of the classifier to prevent the over-fitting problem, we also regularize the classifier parameter by a  $\ell_2$  norm term as

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2. \quad (8)$$

2.1.3. *Mutual Information Maximization*

We also propose to learn the classifier by maximizing the mutual information  $I(f; y)$  between the classification response variables  $f \in F$  and the true class label variables  $y \in Y$ . The mutual information between two variables  $f \in F$  and  $y \in Y$  is defined as

$$I(f; y) = H(f) - H(f|y), \quad (9)$$

where  $H(f)$  is the marginal entropy of  $f$ , which is used to measure the uncertainty about  $f$ , and  $H(f|y)$  is the entropy of  $f$  conditional on  $y$ , which is used as the measure of uncertainty of  $f$  when  $y$  is given. To use the mutual information as a criterion to learn the classifier parameters, we first need to estimate  $H(f)$  and  $H(f|y)$ .

**Estimation of  $H(f)$**  We use the training samples to estimate  $H(f)$ , and according to the definition of entropy, we have

$$\begin{aligned} H(f) &= - \int p(f) \log p(f) df \\ &\approx - \sum_{i=1}^n p(f_i) \log p(f_i), \end{aligned} \quad (10)$$

where  $p(f)$  is the probability density of  $f$ . It could be seen that the entropy of  $f$  is the expectation of  $-\log p(f)$  [46]. The non-parametric kernel density estimation (KDE) [47] is used to estimate the probability density function  $p(f)$ ,

$$p(f) = \frac{1}{n} \sum_{j=1}^n K(f - f_j; \sigma), \quad (11)$$

where  $K(z; \sigma) = \exp\left(-\frac{z^2}{2\sigma^2}\right)$  is the Gaussian kernel function [48] and  $\sigma$  is the bandwidth parameter [47].

**Estimation of  $H(f|y)$**  We also use the training samples to estimate  $H(f|y)$ , and according to its definition, we have

$$H(f|y) = \sum_{c \in \{+1, -1\}} p(c) H(f|y = c), \quad (12)$$

where  $p(c) = \frac{n_c}{n}$  is the probability density of class label  $c$ ,  $n_c = \#\{\mathbf{x}_i \in X|y_i = c\}$  is the number of samples with the class label equal to  $c$ , and

$$\begin{aligned} H(f|y = c) &= - \int p(f|y = c) \log p(f|y = c) df \\ &\approx - \sum_{i:y_i=c} p(f_i|y = c) \log p(f_i|y = c) \end{aligned} \quad (13)$$

is the conditional entropy of  $f$  given the class label  $y = c$  [49, 50, 51]. We also use the KDE to estimate the conditional probability density function

$$p(f|y = c) = \frac{1}{n_c} \sum_{i:y_i=c} K(f - f_i, \sigma) \quad (14)$$

Substituting it to (12), we have the estimated  $H(f|y)$ ,

$$H(f|y) \approx - \sum_{c \in \{+1, -1\}} \frac{n_c}{n} \left( \sum_{i:y_i=c} p(f_i|y = c) \log p(f_i|y = c) \right) \quad (15)$$

With the estimated entropy  $H(f)$  and the conditional entropy  $H(f|y)$ , the mutual information between the variable  $f$  and  $y$  could be rewritten as the function of parameter  $\mathbf{w}$  by substituting  $f_i = \mathbf{w}^\top \mathbf{x}_i$ ,

$$\begin{aligned} \tilde{I}(f, y; \mathbf{w}) &= H(f) - H(f|y) \\ &= - \sum_{i=1}^n p(f_i) \log p(f_i) \\ &\quad + \sum_{c \in \{+1, -1\}} \frac{n_c}{n} \left( \sum_{i:y_i=c} p(f_i|y = c) \log p(f_i|y = c) \right) \\ &= - \sum_{i=1}^n p(\mathbf{w}^\top \mathbf{x}_i) \log p(\mathbf{w}^\top \mathbf{x}_i) \\ &\quad + \sum_{c \in \{+1, -1\}} \frac{n_c}{n} \left( \sum_{i:y_i=c} p(\mathbf{w}^\top \mathbf{x}_i|y = c) \log p(\mathbf{w}^\top \mathbf{x}_i|y = c) \right) \end{aligned} \quad (16)$$

To learn the classifier parameter  $\mathbf{w}$ , we maximize the mutual information with regard to  $\mathbf{w}$ ,

$$\max_{\mathbf{w}} \tilde{I}(f, y; \mathbf{w}) \quad (17)$$

**Remark:** It should be noted that similar to our method, the algorithm proposed in [38] maximizes KL-divergence between the class PDF,  $p(f|y = c)$ ,

and the total PDF,  $p(f)$ , therefore, [38] has relation to method in Kullback-Leibler boosting. However, different from our method, it uses KL-divergence as a criterion to select the most discriminating features, whereas our method uses mutual information as a criterion to learn the classifier parameter.

#### 2.1.4. Overall Optimization Problem

By combining the optimization problems proposed in (7), (8) and (17), the optimization problem for the proposed classifier parameter learning method is obtained as

$$\min_{\mathbf{w}} \left\{ O(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n L(f_i, y_i; \mathbf{w}) + \alpha \frac{1}{2} \|\mathbf{w}\|_2^2 - \beta \tilde{I}(f, y; \mathbf{w}) \right\}, \quad (18)$$

where  $\alpha$  and  $\beta$  are tradeoff parameters. In the objective function, there are three terms. The first one is optimized so that the prediction error is minimized, the second term is used to control the complexity of the classifier, and the last term is introduced so that the mutual information between the classification response and the true class label can be maximized.

#### 2.2. Optimization

Direct optimization to (18) is difficult. Instead of seeking a closed-form solution, we try to optimize it using gradient descent method in an iterative algorithm [52]. In each iteration, we employ the gradient descent method to update  $\mathbf{w}$ . According to the optimization theory, if  $Q(\mathbf{w})$  is defined and differentiable in a neighborhood of a point  $\mathbf{w}^{old}$ , then  $Q(\mathbf{w})$  decreases faster if  $\mathbf{w}$  goes from  $\mathbf{w}^{old}$  in the direction of the negative gradient of  $Q(\mathbf{w})$  at  $\mathbf{w}^{old}$ ,  $-\nabla Q(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^{old}}$ . Thus the new  $\mathbf{w}$  is obtained by

$$\mathbf{w}^{new} \leftarrow \mathbf{w}^{old} - \eta \nabla Q(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^{old}}, \quad (19)$$

where  $\eta$  is the descent step.

The key step is to compute the gradient of  $Q(\mathbf{w})$ , which is calculated as

$$\nabla Q(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \nabla L(f_i, y_i; \mathbf{w}) + \alpha \mathbf{w} - \beta \nabla \tilde{I}(f, y; \mathbf{w}), \quad (20)$$

where  $\nabla L(f_i, y_i; \mathbf{w})$  and  $\nabla \tilde{I}(f, y; \mathbf{w})$  are the gradient of  $L(f_i, y_i; \mathbf{w})$  and  $\tilde{I}(f, y; \mathbf{w})$  respectively. They are given analytically as follows.

##### 2.2.1. Computation of $\nabla L(f_i, y_i; \mathbf{w})$

We give the analytical gradients of different definitions of  $L(f_i, y_i; \mathbf{w})$  as follows:

**Hinge Loss** is not a smooth function, but we can first update  $\tau_i$  using previous  $\mathbf{w}$  as in (3), and then fix it when we derivate  $\nabla L(f_i, y_i; \mathbf{w})$ ,

$$\nabla L(f_i, y_i; \mathbf{w}) = -\tau_i \times (y_i \mathbf{x}_i). \quad (21)$$

**Squared Loss** is a smooth function, and its gradient is

$$\nabla L(f_i, y_i; \mathbf{w}) = -(1 - y_i \mathbf{w}^\top \mathbf{x}_i) \times (y_i \mathbf{x}_i). \quad (22)$$

**Logistic Loss** is also smooth with its gradient as

$$\nabla L(f_i, y_i; \mathbf{w}) = -\frac{\exp(-y_i \mathbf{w}^\top \mathbf{x}_i)}{1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)} \times (y_i \mathbf{x}_i). \quad (23)$$

**Exponential loss** is also smooth, and its gradient can be obtained as

$$\nabla L(f_i, y_i; \mathbf{w}) = -\exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \times (y_i \mathbf{x}_i). \quad (24)$$

### 2.2.2. Computation of $\nabla \tilde{I}(f, y; \mathbf{w})$

The gradient of  $\tilde{I}(f, y; \mathbf{w})$  is computed as

$$\begin{aligned} \nabla \tilde{I}(f, y; \mathbf{w}) &= -\sum_{i=1}^n (\nabla p(\mathbf{w}^\top \mathbf{x}_i) \log p(\mathbf{w}^\top \mathbf{x}_i) + \nabla p(\mathbf{w}^\top \mathbf{x}_i)) \\ &\quad + \sum_{c \in \{+1, -1\}} \frac{n_c}{n} \left[ \sum_{i: y_i = c} (\nabla p(\mathbf{w}^\top \mathbf{x}_i | y = c) \log p(\mathbf{w}^\top \mathbf{x}_i | y = c) + \nabla p(\mathbf{w}^\top \mathbf{x}_i | y = c)) \right] \\ &= -\sum_{i=1}^n (\log p(\mathbf{w}^\top \mathbf{x}_i) + 1) \nabla p(\mathbf{w}^\top \mathbf{x}_i) \\ &\quad + \sum_{c \in \{+1, -1\}} \frac{n_c}{n} \left( \sum_{i: y_i = c} (\log p(\mathbf{w}^\top \mathbf{x}_i | y = c) + 1) \nabla p(\mathbf{w}^\top \mathbf{x}_i | y = c) \right), \end{aligned} \quad (25)$$

where the gradients of  $p(\mathbf{w}^\top \mathbf{x}_i)$  and  $p(\mathbf{w}^\top \mathbf{x}_i | y = c)$  are computed as

$$\begin{aligned} \nabla p(\mathbf{w}^\top \mathbf{x}_i) &= \frac{1}{n\sigma^2} \sum_{j=1}^n \exp\left(-\frac{(\mathbf{w}^\top \mathbf{x}_i - \mathbf{w}^\top \mathbf{x}_j)^2}{2\sigma^2}\right) (\mathbf{w}^\top \mathbf{x}_i - \mathbf{w}^\top \mathbf{x}_j)(\mathbf{x}_j - \mathbf{x}_i), \\ \nabla p(\mathbf{w}^\top \mathbf{x}_i | y = c) &= \frac{1}{n_c \sigma^2} \sum_{j: y_j = c} \exp\left(-\frac{(\mathbf{w}^\top \mathbf{x}_i - \mathbf{w}^\top \mathbf{x}_j)^2}{2\sigma^2}\right) (\mathbf{w}^\top \mathbf{x}_i - \mathbf{w}^\top \mathbf{x}_j)(\mathbf{x}_j - \mathbf{x}_i). \end{aligned} \quad (26)$$

## 3. Experiments

In this section, we evaluate the proposed classification method on two real world pattern classification problems.

### 3.1. Experiment I: Zinc-binding Site Prediction

Zinc is an important element for many biological processes of an organism, and it is closely related to many different diseases. Moreover, it is also critical for proteins to play their functional roles [33, 53, 54]. Thus functional annotation of zinc-binding proteins is necessary to biological process control and disease treatment. To this end, predicting zinc-binding sites of proteins shows its importance in bioinformatics problems. In the first experiment, we evaluate the proposed classification method on the problem of predicting zinc-binding sites.

#### 3.1.1. Data Set and Protocol

For the purpose of experiment, we collected a set of amino acids of four types, which are CYS, HIS, GLU and ASP (CHED). These four types are the most common zinc-binding site types, which take up roughly 96% of the known zinc-binding sites. In the collected data set, there are 1,937 zinc-binding CHEDs and 11,049 non-zinc-binding CHEDs, resulting a data set of 13,986 data samples. Given a candidate CHED, the problem of zinc-binding site prediction is to predict if it is a zinc-binding site or a non-zinc-binding site. In this experiment, we treated a zinc-binding CHED as a positive sample, and a non-zinc-binding CHED as a negative sample. To extract features from a CHED, we computed the position specific substitution matrices (PSSM) [55], the relative weight of gapless real matches to pseudocounts (RW-GRMTP) [56], Shannon entropy [33], and composition of  $k$ -spaced amino acid pairs (CKSAAP) [57], and concatenated them to form a feature vector for each data sample. Please note that the value of each feature was scaled to the range between -1 and 1, so that the performance does not depend on the selection of scaling.

To conduct the experiment, we used the 10-fold cross validation protocol [58]. The entire data set was split into ten non-overlapping folds, and each fold was used as a test set in turn, while the remaining nine folds were combined and used as a training set. The proposed algorithm was performed to the training set to learn a classifier from the feature vectors of the training samples, and then the learned classifier was used to predict the class labels of the test samples. Please note that the tradeoff parameters of the proposed algorithm was tuned within the training set. The averaged value of the hyper-parameters  $\alpha$  and  $\beta$  are 5.8 and 44.8. The parameter  $\sigma$  was computed as  $\sigma = \zeta \times dist$ , where  $dist$  was the median value of distances between pairs of training samples, and the averaged value of  $\zeta$  was 0.451. The classification performance was measured by comparing the predicted labels against the true labels. The receiver operating characteristic (ROC) and recall-precision curves were used as performance metrics. The ROC curve was obtained by plotting true positive rates (TPRs) against false positive rates (FPRs), while recall-precision curve was obtained by plotting precision against recall values. TPR, FPR, recall, and precision are defined as

$$\begin{aligned}
TPR &= \frac{TP}{TP + FN}, FPR = \frac{FP}{FP + TN}, \\
recall &= \frac{TP}{TP + FN}, precision = \frac{TP}{TP + FP},
\end{aligned}
\tag{27}$$

where  $TP$ ,  $FP$ ,  $FN$  and  $TN$  represent the number of true positives, false positives, false negatives and true negatives, respectively. Moreover, area under ROC curve (AUC) [59] was used as a single performance measure. A good classifier should achieve a ROC curve close to the top left corner of the figure, a recall-precision curve close to the top right corner, and also a high AUC value.

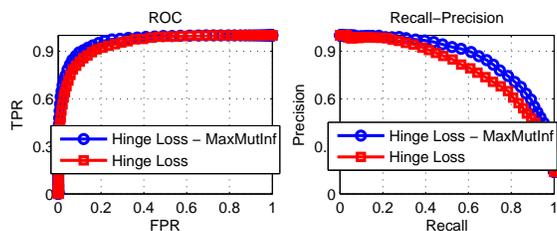
### 3.1.2. Results

In this experiment, we compared the proposed mutual information regularized classifier against the original loss functions based classifier without mutual information regularization, so that the improvement achieved by maximum mutual information regularization could be verified. The four different loss functions listed in Section 2 were considered, and the corresponding classifiers were evaluated here. The ROC and recall-precision curves of four loss functions based classification methods are given in Fig. 1. The proposed maximum mutual information regularized method is denoted as ‘‘MaxMutInf’’ after a loss function title in the figure. It turns out that maximum mutual information regularization improves all the four loss functions based classification methods significantly. Although various loss functions achieved different performances, all of them could be boosted by reducing the uncertainty about true class labels, which could be measured by the mutual information between class labels and classification responses. Therefore, the results show that maximizing mutual information is highly effective in reducing uncertainty of true class labels, and hence it can significantly improve the quality of classification.

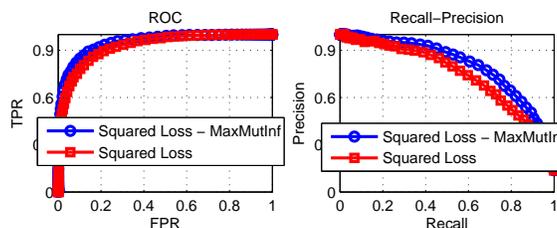
Moreover, we also plotted AUCs of different methods in Fig. 2. Again, we observe that maximum mutual information regularization improves different loss functions based classifiers. We also can see that among these four loss functions, hinge loss achieves the highest AUC values, while squared loss achieves the lowest. The AUC value of classifiers regularized by both hinge loss and mutual information is 0.9635, while that of squared loss and mutual information is even lower than 0.95. The performances of logistic and exponential loss functions are similar, and they are between the performances of hinge loss and squared loss.

Since the mutual information is used as a new regularization technique, we are also interested in how the proposed regularization alone works. We therefore compared the following three cases.

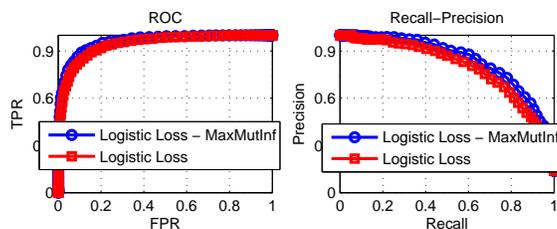
1. **Conventional case** which only uses the classification loss regularization. This case is corresponding to setting  $\beta = 0$  in (18). In this case, we only used the hinge loss since it has been shown that this loss function obtains better accuracy than other loss functions.
2. **Mutual information regularization case** which is corresponding to the problem in (18) when the first term is ignored.



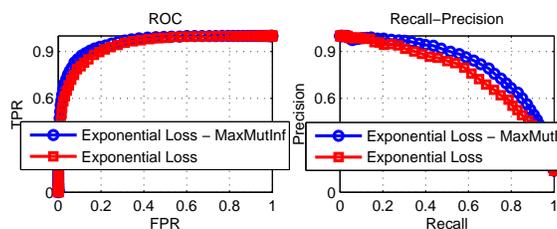
(a) Hinge Loss



(b) Squared Loss



(c) Logistic Loss



(d) Exponential Loss

Figure 1: ROC and recall curves of experiments on zinc-binding site prediction. “Loss” stands for combination of classification loss and  $\ell_2$ -regularization, and “Loss - MaxMutInf” stands for combination of classification loss,  $\ell_2$  and maximum mutual information-regularization.

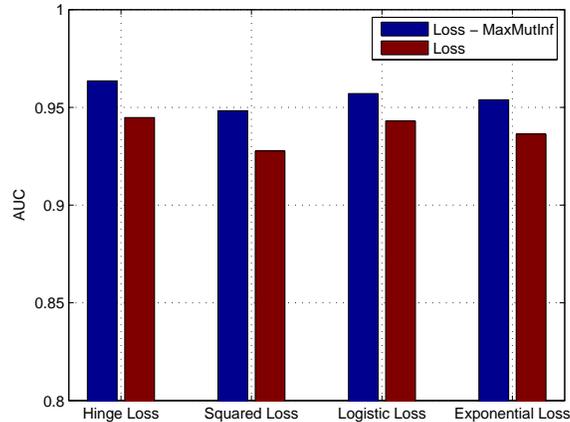


Figure 2: AUC values of experiments on zinc-binding site prediction. “Loss” stands for combination of classification loss and  $\ell_2$ -regularization, and “Loss - MaxMutInf” stands for combination of classification loss,  $\ell_2$  and maximum mutual information-regularization.

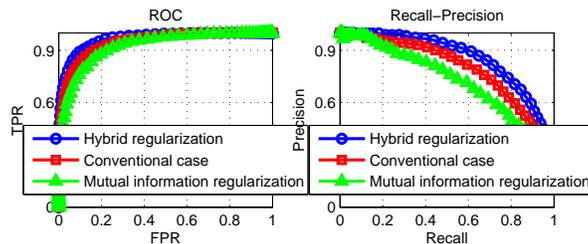


Figure 3: Comparison results of three regularization cases on zinc-binding site prediction.

### 3. Hybrid regularization case which is the proposed framework which combines the classification loss minimization and mutual information regularization.

The comparison results are given in Fig. 3. It can be seen that the conventional case which only uses the hinge loss function achieved better results than the method with only mutual information regularization, and the hybrid regularization achieved the best results. This means mutual information regularization cannot obtain good performance by itself and should be used with traditional loss functions.

#### 3.2. Experiment II: HEP-2 Cell Image Classification

Antinuclear Autoantibodies (ANA) test is a technology used to determine whether a human immune system is creating antibodies to fight against infections. ANA is usually done by a specific fluorescence pattern of HEP-2 cell images [60]. Recently, there is a great need for computer based HEP-3 cell

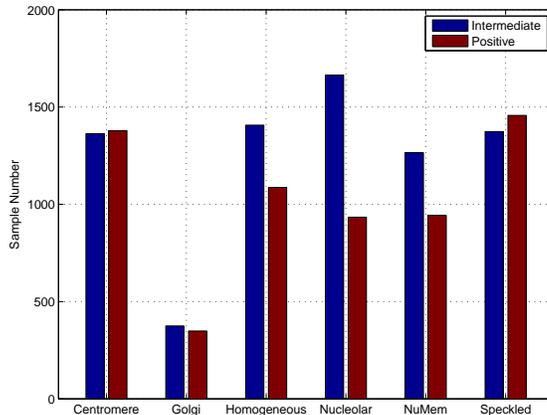


Figure 4: Number of samples of different classes of the HEP-2 cell image data set.

image classification, since manual classification is time-consuming and not accurate enough. In the second experiment, we will evaluate the performance of the proposed classifier on the problem of classifying HEP-2 cell images.

### 3.2.1. Data Set and Protocol

In this experiment, we used the database of HEP-2 cell images of the ICIP 2014 competition of cell classification by fluorescent image analysis [61]. In this data set, there are 13,596 cell images, and they belong to six cell classes, which are namely Centromere, Golgi, Homogeneous, Nucleolar, NuclearMembrane, and Speckled. Each cell image is segmented by a mask image showing the boundary of the cell. Moreover, the entire data set is composed of two groups of different density types, which are Intermediate and Positive. Overall, the intermediate group outnumbers the positive group, with an exception that, for the cases of Centromere and Speckled, the latter marginally outnumbers the former. The number of images in different classes of two groups are given in Fig. 4. To present each image for the classification problem, we extracted shape and texture features and concatenated them to form a visual feature vector [60].

Experiments were conducted in two groups respectively. We also adopted the 10-fold cross validation for the experiment. To handle the problem of multiple class problem, we used the one-against-all strategy. Each class was treated as a positive class in turn, while all remaining five classes were combined to form a negative class. A classifier was learned for each class to discriminate it from other classes. A test sample was assigned to a class with the largest classification response. The classification accuracy was used as a classification performance metric.

### 3.2.2. Results

The boxplots of accuracies of the 10-fold cross validation on the two groups of HEP-2 cell image data set are given in Fig. 5. From this figure, it could be observed that for both two groups of data sets, the proposed regularization method can improve the classification performances significantly, despite of the variety of loss functions. It can also be seen that the performances on the second group (Positive) is inferior to that of the first group (Intermediate). This indicates that it is more difficult to classify cell images when their contrast is low. However, the improvement achieved by mutual information regularization is consistent over these two groups.

## 4. Conclusions

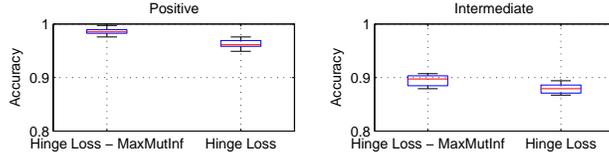
Can knowing the classification response of a data sample reduce uncertainty about its true class label? In this paper, we proposed this question and tried to answer it by learning an optimal classifier to reduce such uncertainty. Inspired by the fact that the reduced uncertainty can be measured by the mutual information between classification responses and true class labels, we proposed a new classifier learning algorithm, by maximizing the mutual information when learning the classifier. Particularly, our algorithm adds a maximum mutual information regularization term. We investigated the classification performances when maximum mutual information was used to regularize the classifier learning based on four different loss functions. The experimental results show that the proposed regularization can improve the classification performances of all these four loss function based classifiers. In the future, we will study how to apply the proposed algorithm on large scale dataset based on some distributed big data platforms [62, 63, 64, 65] and use it to signal and power integrity applications [66, 67, 68, 69, 70, 71, 72].

## Acknowledgements

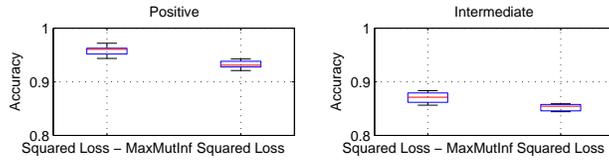
This work was supported by grants from King Abdullah University of Science and Technology (KAUST), Saudi Arabia.

## References

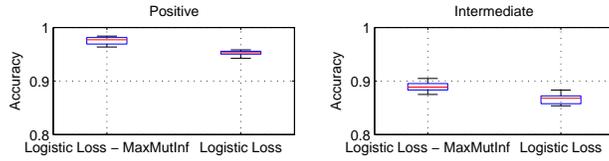
- [1] Q. Cai, H. He, H. Man, Imbalanced evolving self-organizing learning, *Neurocomputing* 133 (0) (2014) 258 – 270.
- [2] T. Ojala, M. Pietikäinen, T. Mäenpää, Multiresolution gray-scale and rotation invariant texture classification with local binary patterns, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24 (7) (2002) 971–987.
- [3] Q. Sun, F. Hu, Q. Hao, Mobile target scenario recognition via low-cost pyroelectric sensing system: Toward a context-enhanced accurate identification, *IEEE transactions on systems, man, and cybernetics. Systems* 44 (3) (2014) 375–384.



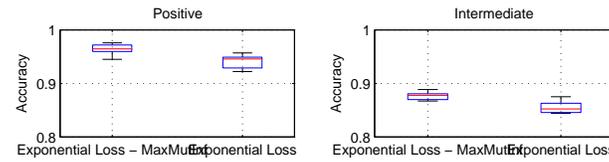
(a) Hinge Loss



(b) Squared Loss



(c) Logistic Loss



(d) Exponential Loss

Figure 5: Experiment results on HEP-2 cell image data set. Please note that “Loss” stands for combination of classification loss and  $\ell_2$ -regularization, and “Loss - MaxMutInf” stands for combination of classification loss,  $\ell_2$  and maximum mutual information-regularization.

- [4] L. Li, J. Yang, Y. Xu, Z. Qin, H. Zhang, Document clustering based on max-correntropy non-negative matrix factorization, 2014.
- [5] L. Li, J. Yang, K. Zhao, Y. Xu, H. Zhang, Z. Fan, Graph regularized non-negative matrix factorization by maximizing correntropy, arXiv preprint arXiv:1405.2246.
- [6] B. Alipanahi, X. Gao, E. Karakoc, L. Donaldson, M. Li, Picky: a novel svd-based nmr spectra peak picking method, *Bioinformatics* 25 (12) (2009) i268–i275.
- [7] J. J.-Y. Wang, X. Wang, X. Gao, Non-negative matrix factorization by maximizing correntropy for cancer clustering, *BMC Bioinformatics* 14 (1) (2013) 107.
- [8] J. J. Wang, H. Bensmail, X. Gao, Multiple graph regularized protein domain ranking, *BMC bioinformatics* 13 (1) (2012) 307.
- [9] J. Wang, X. Gao, Q. Wang, Y. Li, Prodis-contshc: learning protein dissimilarity measures and hierarchical context coherently for protein-protein comparison in protein database retrieval, *BMC bioinformatics* 13 (Suppl 7) (2012) S2.
- [10] Z. Liu, A. Abbas, B.-Y. Jing, X. Gao, Wavpeak: picking nmr peaks through wavelet-based smoothing and volume-based filtering, *Bioinformatics* 28 (7) (2012) 914–920.
- [11] P. Wang, Intelligent pattern recognition and applications to biometrics in an interactive environment, in: *GRAPP 2009 - Proceedings of the 4th International Conference on Computer Graphics Theory and Applications*, 2009, pp. IS21–IS22.
- [12] K. Roy, P. Bhattacharya, C. Y. Suen, Towards nonideal iris recognition based on level set method, genetic algorithms and adaptive asymmetrical svms, *Engineering Applications of Artificial Intelligence* 24 (3) (2011) 458–475.
- [13] F. Tafazzoli, R. Safabakhsh, Model-based human gait recognition using leg and arm movements, *Engineering applications of artificial intelligence* 23 (8) (2010) 1237–1246.
- [14] J. Yang, B. Payne, M. Hitz, Z. Fei, L. Li, T. Wei, Location aided energy balancing strategy in green cellular networks, arXiv preprint arXiv:1406.5258.
- [15] L. Xu, Z. Zhan, S. Xu, K. Ye, An evasion and Counter-Evasion study in malicious websites detection, in: *2014 IEEE Conference on Communications and Network Security (CNS) (IEEE CNS 2014)*, San Francisco, USA, 2014.

- [16] L. Xu, Z. Zhan, S. Xu, K. Ye, Cross-layer detection of malicious websites, in: Proceedings of the third ACM conference on Data and application security and privacy, ACM, 2013, pp. 141–152.
- [17] Q. Cai, Y. Yin, H. Man, Dspm: Dynamic structure preserving map for action recognition, in: Multimedia and Expo (ICME), 2013 IEEE International Conference on, 2013, pp. 1–6. doi:10.1109/ICME.2013.6607606.
- [18] J. J.-Y. Wang, Y. Sun, X. Gao, Sparse structure regularized ranking, Multimedia Tools and Applications (2014) 1–20.
- [19] J. J.-Y. Wang, H. Bensmail, X. Gao, Joint learning and weighting of visual vocabulary for bag-of-feature based tissue classification, Pattern Recognition 46 (12) (2013) 3249–3255.
- [20] Y. Zhou, L. Li, T. Zhao, H. Zhang, Region-based high-level semantics extraction with cedd, in: Network Infrastructure and Digital Content, 2010 2nd IEEE International Conference on, IEEE, 2010, pp. 404–408.
- [21] P. Jonathon Phillips, H. Moon, S. Rizvi, P. Rauss, The feret evaluation methodology for face-recognition algorithms, IEEE Transactions on Pattern Analysis and Machine Intelligence 22 (10) (2000) 1090–1104.
- [22] W. Zhao, R. Chellappa, P. Phillips, A. Rosenfeld, Face recognition: A literature survey, ACM Computing Surveys 35 (4) (2003) 399–458.
- [23] Q. Sun, P. Wu, Y. Wu, M. Guo, J. Lu, Unsupervised multi-level non-negative matrix factorization model: Binary data case, Journal of Information Security 3 (2012) 245.
- [24] Y. Zhou, Y. Liu, H. Li, W. Teng, Z. Li, Fault feature extraction for gear crack based on bispectral entropy, Zhongguo Jixie Gongcheng/China Mechanical Engineering 24 (2) (2013) 190–194.
- [25] J. J.-Y. Wang, X. Gao, Beyond cross-domain learning: Multiple domain nonnegative matrix factorization, Engineering Applications of Artificial Intelligence 28 (0) (2013) 181 – 189.
- [26] M. Al-Shedivat, J. J.-Y. Wang, M. Alzahrani, J. Z. Huang, X. Gao, Supervised transfer sparse coding, in: Twenty-Eighth AAAI Conference on Artificial Intelligence, 2014, pp. 1665 – 1672.
- [27] J. J.-Y. Wang, H. Bensmail, X. Gao, Feature selection and multi-kernel learning for sparse representation on a manifold, Neural Networks 51 (0) (2014) 9 – 16.
- [28] J. J.-Y. Wang, H. Bensmail, N. Yao, X. Gao, Discriminative sparse coding on multi-manifolds, Knowledge-Based Systems 54 (2013) 199–206.

- [29] J.-Y. Wang, I. Almasri, X. Gao, Adaptive graph regularized nonnegative matrix factorization via feature selection, in: Pattern Recognition (ICPR), 2012 21st International Conference on, IEEE, 2012, pp. 963–966.
- [30] J. J.-Y. Wang, H. Bensmail, X. Gao, Multiple graph regularized nonnegative matrix factorization, Pattern Recognition 46 (10) (2013) 2840 – 2847.
- [31] Y. Zhou, L. Li, H. Zhang, Adaptive learning of region-based pls model for total scene annotation, arXiv preprint arXiv:1311.5590.
- [32] T. Subbulakshmi, A. Afroze, Multiple learning based classifiers using layered approach and feature selection for attack detection, in: 2013 IEEE International Conference on Emerging Trends in Computing, Communication and Nanotechnology, ICE-CCN 2013, 2013, pp. 308–314.
- [33] Z. Chen, Y. Wang, Y.-F. Zhai, J. Song, Z. Zhang, Zincexplorer: An accurate hybrid method to improve the prediction of zinc-binding sites from protein sequences, Molecular BioSystems 9 (9) (2013) 2213–2222.
- [34] E. Liu, E. K. P. Chong, L. L. Scharf, Greedy adaptive linear compression in signal-plus-noise models, IEEE Transactions on Information Theory 60 (4) (2014) 2269–2280.
- [35] R. Battiti, Using mutual information for selecting features in supervised neural net learning, IEEE Transactions on Neural Networks 5 (4) (1994) 537–550.
- [36] L. Sun, J. Xu, Feature selection using mutual information based uncertainty measures for tumor classification, Bio-Medical Materials and Engineering 24 (1) (2014) 763–770.
- [37] P. J. Moreno, P. P. Ho, N. Vasconcelos, A kullback-leibler divergence based kernel for svm classification in multimedia applications, in: Advances in Neural Information Processing Systems 16, MIT Press, 2004, pp. 1385–1392.
- [38] C. Liu, H.-Y. Shum, Kullback-leibler boosting, in: Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on, Vol. 1, IEEE, 2003, pp. I–587.
- [39] Y. Wu, Y. Liu, Robust truncated hinge loss support vector machines, Journal of the American Statistical Association 102 (479) (2007) 974–983.
- [40] O. Yildiz, E. Alpaydin, Statistical tests using hinge/-sensitive loss, in: Computer and Information Sciences III - 27th International Symposium on Computer and Information Sciences, ISCIS 2012, 2013, pp. 153–160.
- [41] S. Bach, B. Huang, B. London, L. Getoor, Hinge-loss markov random fields: Convex inference for structured prediction, in: Uncertainty in Artificial Intelligence - Proceedings of the 29th Conference, UAI 2013, 2013, pp. 32–41.

- [42] X. Wang, Y. Jiang, M. Huang, H. Zhang, Robust variable selection with exponential squared loss, *Journal of the American Statistical Association* 108 (502) (2013) 632–643.
- [43] J. Luo, Regression learning in decision guidance systems: Models, languages, and algorithms, Ph.D. thesis, George Mason University (2012).
- [44] J. Luo, A. Brodsky, Y. Li, An em-based ensemble learning algorithm on piecewise surface regression problem, *International Journal of Applied Mathematics and Statistics* 28 (4) (2012) 59–74.
- [45] C. Park, J.-Y. Koo, P. Kim, J. Lee, Stepwise feature selection using generalized logistic loss, *Computational Statistics and Data Analysis* 52 (7) (2008) 3709–3718.
- [46] J. Bekenstein, Black holes and entropy, *Physical Review D* 7 (8) (1973) 2333–2346.
- [47] A. Elgammal, R. Duraiswami, D. Harwood, L. Davis, Background and foreground modeling using nonparametric kernel density estimation for visual surveillance, *Proceedings of the IEEE* 90 (7) (2002) 1151–1162.
- [48] S. Zhong, D. Chen, Q. Xu, T. Chen, Optimizing the gaussian kernel function with the formulated kernel target alignment criterion for two-class pattern classification, *Pattern Recognition* 46 (7) (2013) 2045–2054.
- [49] A. Carvalho, P. Adão, P. Mateus, Efficient approximation of the conditional relative entropy with applications to discriminative learning of bayesian network classifiers, *Entropy* 15 (7) (2013) 2716–2735.
- [50] A. Porta, G. Baselli, D. Liberati, N. Montano, C. Cogliati, T. Gnecchi-Ruscione, A. Malliani, S. Cerutti, Measuring regularity by means of a corrected conditional entropy in sympathetic outflow, *Biological Cybernetics* 78 (1) (1998) 71–78.
- [51] G.-Y. Wang, H. Yu, D.-C. Yang, Decision table reduction based on conditional information entropy, *Jisuanji Xuebao/Chinese Journal of Computers* 25 (7) (2002) 759–766.
- [52] A. P. Dempster, N. M. Laird, D. B. Rubin, Maximum likelihood from incomplete data via the em algorithm, *Journal of the Royal Statistical Society. Series B (Methodological)* (1977) 1–38.
- [53] A. Kumar, D. Hati, T. Thaker, L. Miah, P. Cunningham, C. Domene, T. Bui, A. Drake, L. McDermott, Strong and weak zinc binding sites in human zinc-2-glycoprotein, *FEBS Letters* 587 (24) (2013) 3949–3954.
- [54] Z. Liu, Y. Wang, C. Zhou, Y. Xue, W. Zhao, H. Liu, Computationally characterizing and comprehensive analysis of zinc-binding sites in proteins, *Biochimica et Biophysica Acta - Proteins and Proteomics* 1844 (1 PART B) (2014) 171–180.

- [55] L. Kelley, R. MacCallum, M. Sternberg, Enhanced genome annotation using structural profiles in the program 3d-pssm, *Journal of Molecular Biology* 299 (2) (2000) 499–520.
- [56] S. Menchetti, A. Passerini, P. Frasconi, C. Andreini, A. Rosato, Improving prediction of zinc binding sites by modeling the linkage between residues close in sequence, in: *Research in Computational Molecular Biology*, Springer, 2006, pp. 309–320.
- [57] W. Zhang, X. Xu, M. Yin, N. Luo, J. Zhang, J. Wang, Prediction of methylation sites using the composition of k-spaced amino acid pairs, *Protein and Peptide Letters* 20 (8) (2013) 911–917.
- [58] P. Burman, A comparative study of ordinary cross-validation, v-fold cross-validation and the repeated learning-testing methods, *Biometrika* 76 (3) (1989) 503–514.
- [59] G. Zou, L. Yue, Using confidence intervals to compare several correlated areas under the receiver operating characteristic curves, *Statistics in Medicine* 32 (29) (2013) 5077–5090.
- [60] P. Agrawal, M. Vatsa, R. Singh, Hep-2 cell image classification: A comparative analysis, in: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, Vol. 8184 LNCS, 2013, pp. 195 – 202.
- [61] M. V. A. W. Peter Hobson, Gennaro Percannella, Competition on cells classification by fluorescent image analysis, <http://nerone.diiiie.unisa.it/contest-icip-2013/index.shtml> (2013).
- [62] Y. Su, Y. Wang, G. Agrawal, R. Kettimuthu, SDQuery DSI: integrating data management support with a wide area data transfer protocol, in: *SC*, 2013, p. 47.
- [63] Y. Wang, Y. Su, G. Agrawal, Supporting a Light-Weight Data Management Layer Over HDF5, in: *Cluster, Cloud and Grid Computing (CCGrid)*, 2013 13th IEEE/ACM International Symposium on, IEEE, 2013, pp. 335–342.
- [64] Y. Wang, W. Jiang, G. Agrawal, SciMATE: A novel mapreduce-like framework for multiple scientific data formats, in: *Cluster, Cloud and Grid Computing (CCGrid)*, 2012 12th IEEE/ACM International Symposium on, IEEE, 2012, pp. 443–450.
- [65] Q. Sun, F. Hu, H. Qi, Context awareness emergence for distributed binary pyroelectric sensors, in: *Multisensor Fusion and Integration for Intelligent Systems (MFI)*, 2010 IEEE Conference on, IEEE, 2010, pp. 162–167.

- [66] H. Liu, F. Shi, Y. Wang, N. Wong, Frequency-domain transient analysis of multitime partial differential equation systems, in: VLSI and System-on-Chip (VLSI-SoC), 2011 IEEE/IFIP 19th International Conference on, IEEE, 2011, pp. 160–163.
- [67] Y. Wang, Z. Zhang, C.-K. Koh, G. Shi, G. K. Pang, N. Wong, Passivity enforcement for descriptor systems via matrix pencil perturbation, *Computer-Aided Design of Integrated Circuits and Systems*, IEEE Transactions on 31 (4) (2012) 532–545.
- [68] C.-U. Lei, Y. Wang, Q. Chen, N. Wong, On vector fitting methods in signal/power integrity applications, in: Proceedings of the International MultiConference of Engineers and Computer Scientists 2010, IMECS 2010, Newswood Limited., 2010, pp. 1407–1412.
- [69] Y. Wang, C.-U. Lei, G. K. Pang, N. Wong, Mfti: matrix-format tangential interpolation for modeling multi-port systems, in: Proceedings of the 47th Design Automation Conference, ACM, 2010, pp. 683–686.
- [70] Y. Wang, Z. Zhang, C.-K. Koh, G. K. Pang, N. Wong, Peds: Passivity enforcement for descriptor systems via hamiltonian-symplectic matrix pencil perturbation, in: Proceedings of the International Conference on Computer-Aided Design, IEEE Press, 2010, pp. 800–807.
- [71] E. G. Lim, Z. Wang, C.-U. Lei, Y. Wang, K. Man, Ultra wideband antennas—past and present, *International Journal of Computer Science* 37 (3) (2010) 304–314.
- [72] C.-U. Lei, Y. Wang, Q. Chen, N. Wong, A decade of vector fitting development: Applications on signal/power integrity, *IAENG Transactions on Engineering Technologies* (2010) 435–449.