

---

# Squared-loss Mutual Information Regularization: A Novel Information-theoretic Approach to Semi-supervised Learning

---

Gang Niu  
Wittawat Jitkrittum

Department of Computer Science, Tokyo Institute of Technology, Tokyo, 152-8552, Japan

GANG@SG.CS.TITECH.AC.JP  
WITTAWATJ@GMAIL.COM

Bo Dai

College of Computing, Georgia Institute of Technology, Atlanta, GA 30332, USA

BDAI6@GATECH.EDU

Hiroataka Hachiya  
Masashi Sugiyama

Department of Computer Science, Tokyo Institute of Technology, Tokyo, 152-8552, Japan

HACCHAN@GMAIL.COM  
SUGI@CS.TITECH.AC.JP

## Abstract

We propose *squared-loss mutual information regularization* (SMIR) for multi-class probabilistic classification, following the *information maximization principle*. SMIR is convex under mild conditions and thus improves the nonconvexity of mutual information regularization. It offers all of the following four abilities to semi-supervised algorithms: Analytical solution, out-of-sample/multi-class classification, and probabilistic output. Furthermore, novel generalization error bounds are derived. Experiments show SMIR compares favorably with state-of-the-art methods.

## 1. Introduction

Semi-supervised learning, which utilizes both labeled and unlabeled data for training, has attracted much attention over the last decade. Many semi-supervised assumptions have been made to extract information from unlabeled data. Among them, the *manifold assumption* (Belkin et al., 2006) is of vital importance. Its origin is the *low-density separation principle*.

However, this low-density separation principle is not the only way to go. A useful alternative is the *information maximization principle* (IMP). IMP comes from information maximization clustering (Agakov & Barber, 2006; Gomes et al., 2010; Sugiyama et al., 2011),

where a probabilistic classifier is trained in an unsupervised manner, so that a given information measure between data and cluster assignments is maximized. These clustering methods have shown IMP is reasonable and powerful.

Following IMP, we propose an information-theoretic approach to semi-supervised learning. Specifically, the *squared-loss mutual information* (SMI) (Suzuki et al., 2009) is designated as the information measure to be maximized. Then, we introduce an SMI approximator with no logarithm inside (Sugiyama et al., 2011), and propose the model of *SMI regularization* (SMIR). Unlike maximizing the mutual information, SMIR is strictly convex under mild conditions and the unique globally optimal solution is accessible. Albeit we can employ any convex loss in principle, SMIR can get rid of logarithm in the involved optimization and guarantees the analytic expression of the globally optimal solution if we use the *squared difference of two probabilities* (Sugiyama, 2010). SMIR aims at *multi-class probabilistic classifiers* that possess the innate ability of multi-class classification with the probabilistic output, and no reduction from the multi-class case to the binary case (cf. Allwein et al., 2000) is needed. These classifiers can also naturally handle unseen data and need no explicit out-of-sample extension. To the best of our knowledge, SMIR is the only framework up to the present which leads to semi-supervised algorithms equipped with all these properties.

Furthermore, we establish two *data-dependent generalization error bounds* for a reduced SMIR algorithm based on the theory of *Rademacher averages* (Bartlett & Mendelson, 2002). Our error bounds can consider

not only labeled data but also unlabeled data. Thus, they can reflect the properties of the particular mechanism generating the data. Thanks to the analytical solution, our bounds also have closed-form expression even though they depend on the data in terms of the Rademacher complexity. Notice that previous bounds (Belkin et al., 2004; Cortes et al., 2008) just focus on the regression error, and none of semi-supervised algorithms hitherto have similar theoretical results.

The rest of this paper is organized as follows. First of all, we present preliminaries, and propose the model and algorithm of SMIR in Section 2. In Section 3, we derive the generalization error bounds. The comparisons to related works are in Section 4, and then the experiments are in Section 5.

## 2. Squared-loss Mutual Information Regularization (SMIR)

In this section, we propose the SMIR approach.

### 2.1. Preliminaries

Let  $\mathcal{X} \subseteq \mathbb{R}^d$  and  $\mathcal{Y} = \{1, \dots, c\}$  where  $d$  and  $c$  are natural numbers,  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$  have an underlying  $p(\mathbf{x}, y)$  and  $p(\mathbf{x}) > 0$  over  $\mathcal{X}$ . Given i.i.d.  $\{(\mathbf{x}_i, y_i)\}_{i=1}^l$  and  $\{\mathbf{x}_i\}_{i=l+1}^n$  where  $n = l + u$  and  $l \ll u$ , we aim at estimating  $p(y | \mathbf{x})$ . Then, we can classify any  $\mathbf{x} \in \mathcal{X}$  to  $\hat{y} = \arg \max_{y \in \mathcal{Y}} p(y | \mathbf{x})$ .

As an information measure, *squared-loss mutual information* (SMI) (Suzuki et al., 2009) between random variables  $X$  and  $Y$  is defined by

$$\text{SMI} := \frac{1}{2} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(\mathbf{x}) p(y) \left( \frac{p(\mathbf{x}, y)}{p(\mathbf{x}) p(y)} - 1 \right)^2 d\mathbf{x}.$$

SMI is the *Pearson divergence* (Pearson, 1900) from  $p(\mathbf{x}, y)$  to  $p(\mathbf{x}) p(y)$ , while the *mutual information* (Shannon, 1948) is the *Kullback-Leibler divergence* (Kullback & Leibler, 1951) from  $p(\mathbf{x}, y)$  to  $p(\mathbf{x}) p(y)$ . They both belong to  $f$ -divergence (Ali & Silvey, 1966; Csiszár, 1967), and thus share similar properties. For instance, both of them are nonnegative, and take zero if and only if  $X$  and  $Y$  are independent.

In Sugiyama et al. (2011), a computationally-efficient unsupervised SMI approximator was proposed. By assuming a uniform class-prior probability  $p(y) = 1/c$ , SMI becomes

$$\text{SMI} = \frac{c}{2} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} (p(y | \mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} - \frac{1}{2}. \quad (1)$$

Then,  $p(y | \mathbf{x})$  is approximated by a kernel model:

$$q(y | \mathbf{x}; \boldsymbol{\alpha}) := \sum_{i=1}^n \alpha_{y,i} k(\mathbf{x}, \mathbf{x}_i), \quad (2)$$

where  $\boldsymbol{\alpha} = \{\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_c\}$  and  $\boldsymbol{\alpha}_y = (\alpha_{y,1}, \dots, \alpha_{y,n})^\top$  are model parameters, and  $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  is a kernel. After approximating the expectation w.r.t.  $p(\mathbf{x})$  in Eq. (1) by the empirical average, an SMI approximator is derived as

$$\widehat{\text{SMI}} = \frac{c}{2n} \sum_{y \in \mathcal{Y}} \boldsymbol{\alpha}_y^\top \mathbf{K}^2 \boldsymbol{\alpha}_y - \frac{1}{2},$$

where  $\mathbf{K} \in \mathbb{R}^{n \times n}$  is the kernel matrix.

### 2.2. Basic model

Instead of Eq. (2), we introduce an alternative kernel model for SMIR (the reason will be explained in Remark 1). Let the empirical kernel map (Schölkopf & Smola, 2001) be

$$\Phi_n : \mathcal{X} \mapsto \mathbb{R}^n, \mathbf{x} \mapsto (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n))^\top,$$

the degree of  $\mathbf{x}_i$  be  $d_i = \sum_{j=1}^n k(\mathbf{x}_i, \mathbf{x}_j)$ , and the degree matrix be  $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$ . We approximate the class-posterior probability  $p(y | \mathbf{x})$  by<sup>1</sup>

$$q(y | \mathbf{x}; \boldsymbol{\alpha}) := \langle \mathbf{K}^{-1/2} \Phi_n(\mathbf{x}), \mathbf{D}^{-1/2} \boldsymbol{\alpha}_y \rangle, \quad (3)$$

where  $\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{j=1}^n a_j b_j$  is the inner product. Plugging (3) into Eq. (1) gives us an alternative SMI approximator:

$$\widehat{\text{SMI}} = \frac{c}{2n} \text{tr} \left( \mathbf{A}^\top \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2} \mathbf{A} \right) - \frac{1}{2}, \quad (4)$$

where  $\mathbf{A} = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_c) \in \mathbb{R}^{n \times c}$  is the matrix representation of model parameters.

Subsequently, we employ Eq. (4) to regularize a loss function  $\Delta(p, q)$  that is convex w.r.t.  $q$ . More specifically, we have three objectives: (i) Minimize  $\Delta(p, q)$ ; (ii) Maximize  $\widehat{\text{SMI}}$ ; (iii) Regularize  $\boldsymbol{\alpha}$ . Therefore, we formulate the optimization problem of SMIR as

$$\min_{\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_c \in \mathbb{R}^n} \Delta(p, q) - \gamma \widehat{\text{SMI}} + \lambda \sum_{y \in \mathcal{Y}} \frac{1}{2} \|\boldsymbol{\alpha}_y\|_2^2, \quad (5)$$

where  $\gamma, \lambda > 0$  are regularization parameters.

A remarkable characteristic of optimization (5) is its convexity, as long as the kernel function  $k$  is nonnegative and  $\lambda > \gamma c/n$ :

**Theorem 1.** *Assume that  $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}_+$  and  $\lambda > \gamma c/n$ . Then optimization (5) is strictly convex, and there exists a unique globally optimal solution.*<sup>2</sup>

<sup>1</sup>Assume that  $\mathbf{K}$  is full-rank, and then  $\mathbf{K}^{-1/2}$  is well-defined. The Gaussian kernel matrix is full-rank as long as  $\forall i \neq j, \mathbf{x}_i \neq \mathbf{x}_j$ .

<sup>2</sup>In the rest of this paper, we will assume that  $k$  is nonnegative and  $\lambda > \gamma c/n$ . See Appendix A for the proof.

*Remark 1.* We introduced Eq. (3) due to the following reasons: (i) In principle, any kernel model linear w.r.t.  $\alpha_y$  may be used to approximate  $p(y | \mathbf{x})$ , and maximizing  $\widehat{\text{SMI}}$  alone must be non-convex. However, optimization (5) becomes convex if  $\lambda$  is large enough. Hence, only  $\lambda$  above a certain threshold is acceptable: The threshold of (3) is  $\gamma c/n$ . The threshold of (2) is  $\|K\|_2^2 \cdot \gamma c/n$  where  $\|K\|_2$  is the spectral norm of  $K$ . It depends upon all the training data thoroughly and is usually much larger than  $\gamma c/n$ . (ii) We found that (3) experimentally outperformed (2).

### 2.3. Proposed algorithm

Due to limited space, we give a brief derivation here.

We choose the squared difference of probabilities  $p$  and  $q$  as the loss function (Sugiyama, 2010):

$$\Delta_2(p, q) := \frac{1}{2} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} (p(y | \mathbf{x}) - q(y | \mathbf{x}; \alpha))^2 p(\mathbf{x}) d\mathbf{x}.$$

It enables the analytical solution and facilitates our future theoretical analysis. Its empirical version is

$$\widehat{\Delta}_2 = \text{Const.} - \frac{1}{l} \sum_{i=1}^l q(y_i | \mathbf{x}_i) + \frac{1}{2l} \sum_{i=1}^l \sum_{y=1}^c (q(y | \mathbf{x}_i))^2. \quad (6)$$

Let  $\mathbf{Y} \in \mathbb{R}^{l \times c}$  be the class indicator matrix for  $l$  labeled data and  $\mathbf{B} = (\mathbf{I}_l; \mathbf{0}_{u \times l}) \in \mathbb{R}^{n \times l}$ . Subsequently, Eq. (6) can be expressed by

$$\begin{aligned} \widehat{\Delta}_2 = \text{Const.} & - \frac{1}{l} \text{tr}(\mathbf{Y}^\top \mathbf{B}^\top \mathbf{K}^{1/2} \mathbf{D}^{-1/2} \mathbf{A}) \\ & + \frac{1}{2l} \text{tr}(\mathbf{A}^\top \mathbf{D}^{-1/2} \mathbf{K}^{1/2} \mathbf{B} \mathbf{B}^\top \mathbf{K}^{1/2} \mathbf{D}^{-1/2} \mathbf{A}). \end{aligned} \quad (7)$$

Substituting Eq. (7) into optimization (5), we will get the following objective function:

$$\begin{aligned} \mathcal{F}(\mathbf{A}) = & -\frac{1}{l} \text{tr}(\mathbf{Y}^\top \mathbf{B}^\top \mathbf{K}^{1/2} \mathbf{D}^{-1/2} \mathbf{A}) \\ & + \frac{1}{2l} \text{tr}(\mathbf{A}^\top \mathbf{D}^{-1/2} \mathbf{K}^{1/2} \mathbf{B} \mathbf{B}^\top \mathbf{K}^{1/2} \mathbf{D}^{-1/2} \mathbf{A}) \\ & - \frac{\gamma c}{2n} \text{tr}(\mathbf{A}^\top \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2} \mathbf{A}) + \frac{\lambda}{2} \text{tr}(\mathbf{A}^\top \mathbf{A}). \end{aligned}$$

At last, by equating  $\nabla \mathcal{F}$  to the zero matrix, we obtain the analytical solution to unconstrained optimization problem (5):

$$\begin{aligned} \mathbf{A}_{\mathcal{F}}^* = & n \left( n \mathbf{D}^{-1/2} \mathbf{K}^{1/2} \mathbf{B} \mathbf{B}^\top \mathbf{K}^{1/2} \mathbf{D}^{-1/2} + \lambda n l \mathbf{I}_n \right. \\ & \left. - \gamma l c \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2} \right)^{-1} \mathbf{D}^{-1/2} \mathbf{K}^{1/2} \mathbf{B} \mathbf{Y}. \end{aligned}$$

We recommend to post-process the model parameters as

$$\beta_y = n \pi_y \cdot \frac{\mathbf{K}^{-1/2} \mathbf{D}^{-1/2} \alpha_y^*}{\mathbf{1}_n^\top \mathbf{K}^{1/2} \mathbf{D}^{-1/2} \alpha_y^*},$$

where  $\beta_y$  is a normalized version of  $\alpha_y^*$ , and  $\pi_y$  is an estimate of  $p(y)$  based on labeled data. In addition, probability estimates should be nonnegative and thus our final solution can be expressed as follows (cf. Yamada et al., 2011):

$$\hat{p}(y | \mathbf{x}) = \frac{\max(0, \langle \Phi_n(\mathbf{x}), \beta_y \rangle)}{\sum_{y'=1}^c \max(0, \langle \Phi_n(\mathbf{x}), \beta_{y'} \rangle)}.$$

Although  $q(y | \mathbf{x}; \alpha^*)$  might be negative or unnormalized, Kanamori et al. (2012) implies that minimizing  $\Delta_2$  could achieve the optimal non-parametric convergence rate from  $q$  to  $p$ , and when we have enough data  $q$  is automatically a probability (i.e., non-negative and normalized).

### 3. Generalization Error Bounds

To elucidate the generalization capability, we reduce SMIR to binary classification. Now, a class label  $y$  is  $\pm 1$ , a single vector  $\alpha \in \mathbb{R}^n$  is enough to construct a discriminative model, and we classify any  $x \in \mathcal{X}$  to

$$\hat{y} = \text{sign}(\langle \mathbf{K}^{-1/2} \Phi_n(\mathbf{x}), \mathbf{D}^{-1/2} \alpha \rangle).$$

Let us encode the information of class labels into  $\mathbf{y} = (y_1, \dots, y_l)^\top \in \mathbb{R}^l$ . The solution is then

$$\begin{aligned} \alpha_{\mathcal{F}}^* = & n \left( n \mathbf{D}^{-1/2} \mathbf{K}^{1/2} \mathbf{B} \mathbf{B}^\top \mathbf{K}^{1/2} \mathbf{D}^{-1/2} + \lambda n l \mathbf{I}_n \right. \\ & \left. - \gamma l c \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2} \right)^{-1} \mathbf{D}^{-1/2} \mathbf{K}^{1/2} \mathbf{B} \mathbf{y}, \end{aligned} \quad (8)$$

and for convenience, we define the decision function

$$f(\mathbf{x}) = \langle \mathbf{K}^{-1/2} \Phi_n(\mathbf{x}), \mathbf{D}^{-1/2} \alpha_{\mathcal{F}}^* \rangle. \quad (9)$$

Let  $\mathbb{E}$  and  $\hat{\mathbb{E}}$  stand for the true and empirical expectations,  $\ell(z) = (1 - \text{sign}(z))/2$  be the *indicator loss*, and  $\ell_\eta(z) = \min(1, \max(0, 1 - z/\eta))$  be the *surrogate loss*. We bound  $\mathbb{E}\ell(yf)$  using the theory of *Rademacher averages* (Bartlett & Mendelson, 2002). If all labels are available for evaluation, we can evaluate  $\hat{\mathbb{E}}\ell_\eta(yf)$  over all training data and bound  $\mathbb{E}\ell(yf)$  more tightly. We state the theoretical result in Theorem 2 and prove it in Appendix B.

**Theorem 2.** *Assume that*

$$\exists B_k > 0, \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}, k(\mathbf{x}, \mathbf{x}') \leq B_k^2.$$

Let  $\alpha_{\mathcal{F}}^*$  and  $f(\mathbf{x})$  be the optimal solution and the decision function defined in Eqs. (8) and (9) respectively, and

$$B_{\mathcal{F}} = \|\mathbf{D}^{-1/2} \alpha_{\mathcal{F}}^*\|_2, B'_{\mathcal{F}} = \|\mathbf{K}^{-1/2} \mathbf{D}^{-1/2} \alpha_{\mathcal{F}}^*\|_1.$$

For any  $\eta > 0$  and  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , we have

$$\begin{aligned} \mathbb{E}\ell(yf(\mathbf{x})) &\leq \frac{1}{l} \sum_{i=1}^l \ell_{\eta}(y_i f(\mathbf{x}_i)) + \frac{2B_k B_{\mathcal{F}}}{\eta\sqrt{l}} \\ &+ \min\left(3, 1 + \frac{4B_k^2 B'_{\mathcal{F}}}{\eta}\right) \sqrt{\frac{\ln(2/\delta)}{2l}}. \end{aligned} \quad (10)$$

If the ground truth class labels  $y_{l+1}, \dots, y_n$  are also available for evaluation, with probability at least  $1 - \delta$ , we have

$$\begin{aligned} \mathbb{E}\ell(yf(\mathbf{x})) &\leq \frac{1}{n} \sum_{i=1}^n \ell_{\eta}(y_i f(\mathbf{x}_i)) + \frac{2B_k B_{\mathcal{F}}}{\eta\sqrt{n}} \\ &+ \min\left(3, 1 + \frac{4B_k^2 B'_{\mathcal{F}}}{\eta}\right) \sqrt{\frac{\ln(2/\delta)}{2n}}. \end{aligned} \quad (11)$$

Theorem 2 gives the tightest upper bounds (i.e., the coefficients of  $1/\sqrt{l}$  and  $1/\sqrt{n}$  are smallest under each given scenario) based on the inductive Rademacher complexity. The bound in Eq. (10) is asymptotically  $O(1/\sqrt{l})$ , if we only know the first  $l$  labels. In such cases, we may benefit from unlabeled data by a lower empirical error. It becomes  $O(1/\sqrt{n})$  in Eq. (11) if we can access the other  $u$  labels, even though they are not used for training. Due to the smaller deviation of the empirical error and the empirical Rademacher complexity when they are estimated over all training data, we can improve the order from  $O(1/\sqrt{l})$  to  $O(1/\sqrt{n})$ . Nevertheless, there is no free lunch: In (11), the empirical error is evaluated over all training data, and it may be significantly higher than that evaluated over labeled data. Basically, (10) or (11) which right-hand side is smaller reflects whether the information maximization principle befits the data set or not.

## 4. Related Works

Information-theoretic semi-supervised approaches directly constrain  $p(y | \mathbf{x})$  by unlabeled data or some  $p(\mathbf{x})$  given as the prior knowledge. *Information regularization* (IR; Szummer & Jaakkola, 2002) is the pioneer for this purpose. Compared with later information maximization methods, IR minimizes the mutual information (MI) based on a key observation: Within a small region  $Q \subset \mathcal{X}$ ,  $MI_Q$  is low/high if the label information is pure/chaotic. Subsequently, IR estimates a cover  $\mathcal{C}$  of  $\mathcal{X}$  from  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , and minimizes the maximal  $MI_Q$  for  $Q \in \mathcal{C}$ , subject to class constraints provided by labeled data. The advantage of IR is its flexibility and convexity, while the drawback is that it is unclear how to estimate  $\mathcal{C}$  properly. Each region should be small enough to preserve the locality of the

label information in a single region; each pair of regions should be connected to ensure the dependence of  $p(y | \mathbf{x})$  over all regions, and this implies a great number of tiny regions.

By employing the Shannon entropy of  $p(y | \mathbf{x})$  as a measure of class overlap, *entropy regularization* (ER; Grandvalet & Bengio, 2004) minimizes the entropy from a viewpoint of *maximum a posteriori* estimation. More specifically, ER regularizes the maximum log-likelihood estimation of a logistic regression or kernel logistic regression model by an entropy term:

$$\begin{aligned} \max_{\alpha} &\sum_{i=1}^l \ln q(y_i | \mathbf{x}_i; \alpha) \\ &+ \gamma \sum_{i=l+1}^n \sum_{y \in \mathcal{Y}} q(y | \mathbf{x}_i; \alpha) \ln q(y | \mathbf{x}_i; \alpha). \end{aligned}$$

ER favors low-density separations, since the low/high entropy means that the class overlap is mild/intensive. ER and IR seem opposite at a first glance, because MI equals the difference of the entropies of class prior and posterior. However, IR minimizes MI *locally* and ER minimizes the entropy *globally*, so both of them highly penalize the variations of the class-posterior probability in high-density regions. A recent framework called *regularized information maximization* (RIM; Gomes et al., 2010) follows ER and further maximizes the entropy of the class-prior probability to encourage balanced classes. ER and RIM do not model  $p(\mathbf{x})$  explicitly which is a major improvement, but the disadvantage is the non-convexity of their optimizations.

*Expectation regularization* (XR; Mann & McCallum, 2007) goes one step further such that it does not use  $p(\mathbf{x})$  at all. Therefore, XR does not favor low-density separations and can handle highly overlapped classes. XR encourages the predictions on unlabeled data to match a designer-provided expectation by minimizing the KL-divergence between the expectations predicted by the model and provided as the prior knowledge. If there is no prior knowledge, XR will match the class prior of unlabeled data with that of labeled data:

$$\begin{aligned} \max_{\alpha} &\sum_{i=1}^l \ln q(y_i | \mathbf{x}_i; \alpha) - \lambda \sum_{y \in \mathcal{Y}} \frac{1}{2} \|\alpha_y\|_2^2 \\ &+ \gamma \sum_{y \in \mathcal{Y}} \pi_y \ln \left( \sum_{i=l+1}^n q(y | \mathbf{x}_i; \alpha) \right), \end{aligned}$$

where  $\pi_y$  is an estimate of  $p(y)$  through labeled data, and  $q(y | \mathbf{x}; \alpha)$  is a logistic or kernel logistic regression model. Unlike IR and ER, XR does not prefer low-density separations. As a result, XR cannot deal with low-dimensional data with nonlinear structures (such as the famous *two-moons* or *two-circles*), if there are not enough labeled data.

On the other hand, there are lots of geometric methods for semi-supervised learning. Please see Table 1

Table 1. Summary of existing semi-supervised learning methods.

	AS	OC	MC	PO
Geometric				
Transductive SVM (Joachims, 1999)	×	○	△	×
Semi-supervised SVM (Bennett & Demiriz, 1998)	×	○	△	×
Laplacian SVM (Belkin et al., 2006)	×	○	△	×
Laplacian Regularized Least Squares (Belkin et al., 2006)	○	○	△	×
Markov Random Walks (Szummer & Jaakkola, 2001)	×	×	○	○
Local and Global Consistency (Zhou et al., 2003)	○	△	○	×
Spectral Graph Transducer (Joachims, 2003)	○	×	×	×
Harmonic Energy Minimization (Zhu et al., 2003)	○	×	×	○
Sparse Eigenfunction Bases (Sinha & Belkin, 2009)	×	○	×	×
Information-theoretic				
Information Regularization (Szummer & Jaakkola, 2002)	×	○	○	○
Entropy Regularization (Grandvalet & Bengio, 2004)	×	○	○	○
Expectation Regularization (Mann & McCallum, 2007)	×	○	○	○
Regularized Information Maximization (Gomes et al., 2010)	×	○	○	○
Squared-loss Mutual Information Regularization	○	○	○	○

AS: analytical solution    OC: out-of-sample classification    MC: multi-class classification    PO: probabilistic output  
 ○: Yes    ×: No    △: Extension has been proposed

as a list of representative methods. Note that all geometric methods in Table 1 are in the style of either large margins or similarity graphs. According to Table 1, we could know that many methods based on similarity graphs (Szummer & Jaakkola, 2001; Zhou et al., 2003; Joachims, 2003; Zhu et al., 2003) are transductive, while the information-theoretic methods are all inductive; only two geometric methods (Szummer & Jaakkola, 2001; Zhou et al., 2003) could deal with multi-class data directly, while it is an inherent property of all information-theoretic methods. However, none of previous information-theoretic methods have analytical solutions, due to the logarithms in the entropy, MI or KL-divergence. Thanks to SMI, the proposed SMIR involves a strictly convex optimization problem with no logarithm inside and consequently it has the analytic expression of the unique globally optimal solution.

The similarity between ER and SMIR is intriguing. RIM followed ER historically. Nonetheless, if we start from MI maximization with the uniform  $p(y)$ , we will get ER as

$$\max_{\alpha} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} q(y | \mathbf{x}; \alpha) \ln q(y | \mathbf{x}; \alpha) p(\mathbf{x}) d\mathbf{x}.$$

Recall that SMI maximization under the assumption of the uniform  $p(y)$  is expressed by

$$\max_{\alpha} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} q(y | \mathbf{x}; \alpha) q(y | \mathbf{x}; \alpha) p(\mathbf{x}) d\mathbf{x}.$$

As a consequence, they have the similar preference as the logarithm is strictly monotonically increasing. The

vital difference is the convexity and the analytical solution: SMIR is convex and the globally optimal solution can be obtained analytically, whereas ER is non-convex so any locally optimal solution has to be found numerically.<sup>3</sup>

## 5. Experiments

In this section, we numerically evaluate SMIR. The specification of benchmark data sets is summarized in Table 2. Besides the four well-tried benchmarks in the first block (i.e., USPS, MNIST, 20Newsgroups and Isolet), there are eight benchmarks from a book entitled *Semi-Supervised Learning* (Chapelle et al., 2006)<sup>4</sup> in the second block, and eight benchmarks from the *UCI machine learning repository*<sup>5</sup> in the third block except that Senseval-2 is from a workshop for *word sense disambiguation*<sup>6</sup>. Detailed explanation of benchmarks is omitted due to lack of space. Our experiments consist of three parts:

Firstly, we compare SMIR with entropy regularization (ER; Grandvalet & Bengio, 2004) and expectation regularization (XR; Mann & McCallum, 2007). The probabilistic models are the logistic regression

$$q(y | \mathbf{x}; \alpha) \propto \exp\langle \mathbf{x}, \alpha_y \rangle, \alpha_y \in \mathbb{R}^d,$$

and the kernel logistic regression (Ker)

$$q(y | \mathbf{x}; \alpha) \propto \exp\langle \Phi_n(\mathbf{x}), \alpha_y \rangle, \alpha_y \in \mathbb{R}^n,$$

<sup>3</sup>SMIR may also be solved numerically in consideration of the computational efficiency for large  $n$  in practice.

<sup>4</sup><http://olivier.chapelle.cc/ssl-book/benchmarks.html>.

<sup>5</sup><http://archive.ics.uci.edu/ml/>.

<sup>6</sup><http://www.senseval.org/>.



Table 2. Specification of benchmark data sets.

	# Classes	# Dimensions	# Data	Balance of classes (in %)
USPS	10	256	11000	10 per class
MNIST	10	784	70000	11.3 / 10.0 / 10.2 / 9.8 / 9.0 / 9.8 / 10.4 / 9.8 / 9.9 / 9.9
20Newsgroups	7	53975	11269	4.3 / 25.8 / 5.2 / 21.1 / 21.0 / 5.3 / 17.3
Isolet	26	617	7797	3.85 per class
g241c	2	241	1500	50.0 / 50.0
g241n	2	241	1500	50.1 / 49.9
Digit1	2	241	1500	51.1 / 48.9
USPS	2	241	1500	80.0 / 20.0
COIL	6	241	1500	16.7 per class
COIL2	2	241	1500	50.0 / 50.0
BCI	2	117	400	50.0 / 50.0
Text	2	11960	1500	50.0 / 50.0
Diabetes	2	8	768	65.1 / 34.9
Wine	3	13	178	33.1 / 39.9 / 27.0
Vowel	11	13	990	9.1 per class
Image	2	18	1155	42.9 / 57.1
Vehicle	4	18	846	25.1 / 25.7 / 25.8 / 23.5
German	2	20	1000	70.0 / 30.0
Satimage	6	36	6435	23.8 / 10.9 / 21.1 / 9.7 / 11.0 / 23.4
Senseval-2	3	50	534	33.3 per class

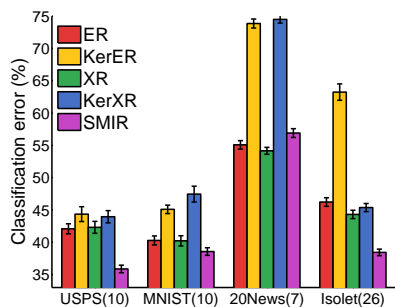
where  $\langle \cdot, \cdot \rangle$  is the inner product,  $\Phi_n$  is the empirical kernel map for the Gaussian kernel. SMIR also applies the Gaussian kernel, so there are three kernel methods which allow nonlinear decision boundaries in  $\mathbb{R}^d$ . The two-fold cross-validation is performed to select the hyperparameters. The kernel width is the median of all pairwise distances times the best value among  $\{1/15, 1/10, 1/5, 1/2, 1\}$ . A Gaussian prior of parameters, which is same as the third term of optimization (5), is included for XR and KerXR (Mann & McCallum, 2007). No extra prior is added to ER or KerER, since ER itself is a prior from a viewpoint of maximum a posteriori estimation (Grandvalet & Bengio, 2004). Therefore, ER/KerER has one regularization parameter whereas XR/KerXR and SMIR have two. The candidate list of regularization parameters is  $10^{\{-7, -3, -1, 1, 3\}}$ , except that  $\lambda$  is chosen from  $\gamma c/n + 10^{\{-10, -8, -6, -4, -2\}}$  for SMIR to ensure the convexity. The *minFunc*<sup>7</sup> package for unconstrained optimization using line-search methods (the quasi-Newton limited-memory BFGS updates, by default) is utilized to solve ER/KerER and XR/KerXR. Since minimizing the entropy is non-convex, we initialize ER/KerER with the globally optimal solution of its supervised part.

We evaluated them on USPS, MNIST, 20Newsgroups and Isolet. Pearson’s correlation (Hall, 2000) was used to select 1000 most informative features for 20Newsgroups. For each data set, we prepared a multi-class task, namely, the tasks using 10 classes of USPS and

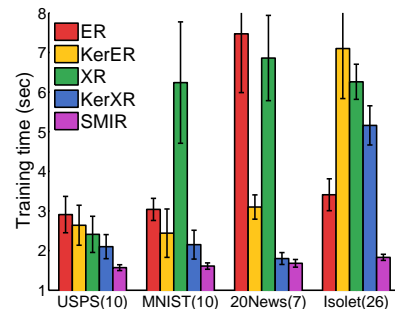
MNIST, 7 classes of 20Newsgroups, and 26 classes of Isolet. In addition, extensive experiments of simple classification tasks were conducted, including 45 binary tasks of USPS, 45 binary tasks of MNIST and 21 binary tasks of 20Newsgroups. Isolet may lead to too many binary tasks and these tasks are often too easy, and thus we combined 26 letters into 13 groups (e.g., ‘a’ with ‘b’, ‘c’ with ‘d’ etc.) and treated each group as a single class resulting in 78 simple classification tasks. For each task, we repeatedly ran all methods on 100 random samplings, where the sample size was fixed to 500. Each random sampling was partitioned into a training set and a test set with 80% and 20% data, and 10% class labels of training data were revealed to construct labeled data.

Figure 1 reports the experimental results of the multi-class tasks, Figure 2 reports the experimental results of the simple tasks, and Table 3 summarizes the experimental results. We can see from Figure 1 that SMIR outperformed others on the multi-class tasks of USPS, MNIST and Isolet. Likewise Figure 1 indicates that SMIR was the most computationally-efficient algorithm on all four multi-class tasks. According to Figure 2, SMIR was the best on the simple tasks of USPS, 20Newsgroups and Isolet, but was slightly inferior to plain ER on MNIST. Note that there were 12 highly imbalanced tasks among 21 simple tasks of 20Newsgroups, which implies that the uniform class-prior assumption will not affect the performance of SMIR essentially, if the tasks are not so complicated. The experiments of Isolet further imply that SMIR is

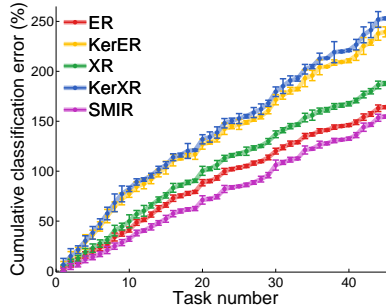
<sup>7</sup><http://www.di.ens.fr/~mschmidt/Software/minFunc/>.



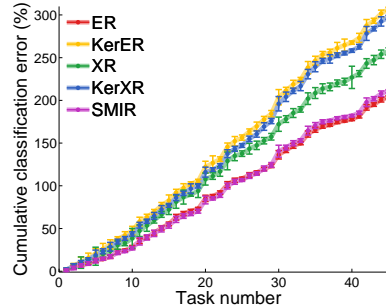
(a) Classification error



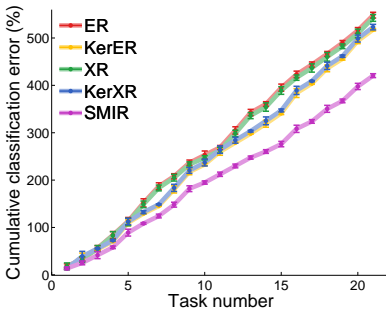
(b) Training time



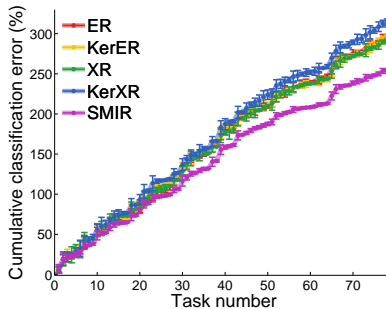
(a) USPS (45 tasks)



(b) MNIST (45 tasks)



(c) 20News (21 tasks)



(d) Isolet (78 tasks)

Figure 1. Experimental results of the multi-class classification tasks. Means with standard errors are shown by bar charts.

Figure 2. Experimental results of the simple classification tasks. The cumulative classification error at the  $k$ -th task is the sum of classification errors from the first to  $k$ -th tasks. Non-cumulative standard deviations are shown along the curves.

Table 3. Summary of all experimental results on USPS, MNIST, 20Newsgroups and Isolet. For each method, we measure how frequently it is the best or a comparable method based on the unpaired  $t$ -test at the significance level 5%, and the training time is averaged over all samplings of all tasks. The most accurate method and the most computationally-efficient method are highlighted in boldface.

	ER	KerER	XR	KerXR	SMIR
USPS, best or comparable (%)	45.65	15.22	21.74	17.39	<b>73.91</b>
MNIST, best or comparable (%)	<b>86.95</b>	0.00	19.57	2.17	80.43
20News, best or comparable (%)	36.36	18.18	36.36	18.18	<b>63.64</b>
Isolet, best or comparable (%)	60.76	62.03	68.35	48.10	<b>81.01</b>
USPS, training time (sec)	1.545	1.906	<b>1.149</b>	1.770	1.608
MNIST, training time (sec)	2.367	1.676	2.060	<b>1.536</b>	1.575
20News, training time (sec)	3.987	2.023	4.144	1.917	<b>1.654</b>
Isolet, training time (sec)	2.377	1.842	2.194	1.728	<b>1.723</b>

fairly good at multi-modal data, since all classes there had two clusters. Compared with KerER and KerXR, the plain ER and XR were better on USPS, MNIST and Isolet, but worse on 20Newsgroups. Nonetheless, ER/XR always outperformed KerER/KerXR in Table 3. Even though other algorithms often converged quite quickly on the simple tasks, SMIR was still a computationally-efficient algorithm after taking these simple tasks into account.

Secondly, we compare SMIR with two well-known geometric methods: Laplacian regularized least squares (LapRLS; Belkin et al., 2006) with a multi-class exten-

Table 4. Comparisons of LapRLS, LGC and SMIR, by means with standard errors of the classification error (in %) on the multi-class tasks. The best method and comparable ones based on the 5% unpaired  $t$ -test are highlighted in boldface.

	LapRLS	LGC	SMIR
USPS	39.64 $\pm$ 0.55	<b>36.53 <math>\pm</math> 0.53</b>	<b>35.87 <math>\pm</math> 0.59</b>
MNIST	42.34 $\pm$ 0.67	42.70 $\pm$ 0.60	<b>38.56 <math>\pm</math> 0.59</b>
20News	64.85 $\pm$ 0.61	73.03 $\pm$ 0.24	<b>56.90 <math>\pm</math> 0.68</b>
Isolet	39.98 $\pm$ 0.56	40.62 $\pm$ 0.47	<b>38.43 <math>\pm</math> 0.51</b>

Table 5. Means with standard errors of the classification error (in %) on benchmarks from Chapelle et al. (2006). The best method and comparable ones based on the unpaired  $t$ -test at the significance level 5% are highlighted in boldface.

	ER	KerER	XR	KerXR	LapRLS	LGC	SMIR
g241c	30.14 ± 0.55	<b>24.86 ± 0.66</b>	31.66 ± 0.81	<b>24.42 ± 0.69</b>	34.12 ± 0.69	36.53 ± 0.74	31.69 ± 0.66
g241n	<b>33.07 ± 0.58</b>	35.65 ± 0.98	<b>33.90 ± 0.83</b>	36.67 ± 0.99	35.07 ± 0.65	38.15 ± 0.72	<b>33.76 ± 0.65</b>
Digit1	12.12 ± 0.41	<b>9.31 ± 0.32</b>	12.47 ± 0.49	<b>9.68 ± 0.62</b>	11.44 ± 0.43	11.87 ± 0.46	10.23 ± 0.40
USPS	26.60 ± 0.59	17.58 ± 0.33	27.07 ± 0.90	18.02 ± 0.72	12.45 ± 0.34	<b>10.27 ± 0.33</b>	12.23 ± 0.40
COIL	46.16 ± 0.78	38.58 ± 0.98	50.55 ± 1.20	39.81 ± 1.06	37.03 ± 0.81	<b>32.95 ± 0.88</b>	<b>33.62 ± 0.82</b>
COIL2	28.83 ± 0.72	25.81 ± 0.75	31.54 ± 1.02	27.73 ± 0.98	26.52 ± 0.65	<b>23.39 ± 0.71</b>	<b>24.12 ± 0.69</b>
BCI	<b>40.58 ± 0.67</b>	47.76 ± 0.45	43.21 ± 0.70	48.35 ± 0.46	43.46 ± 0.63	48.70 ± 0.44	47.27 ± 0.57
Text	<b>34.92 ± 0.56</b>	44.36 ± 0.58	<b>35.38 ± 0.54</b>	43.79 ± 0.65	44.50 ± 0.54	49.53 ± 0.18	38.80 ± 0.64

Table 6. Means with standard errors of the classification error (in %) on seven UCI benchmarks and Senseval-2. The best method and comparable ones based on the unpaired  $t$ -test at the significance level 5% are highlighted in boldface.

	ER	KerER	XR	KerXR	LapRLS	LGC	SMIR
Diabetes	<b>27.26 ± 0.41</b>	29.70 ± 0.50	28.41 ± 0.53	30.16 ± 0.72	32.01 ± 0.62	32.32 ± 0.42	29.87 ± 0.57
Wine	8.09 ± 0.44	<b>4.21 ± 0.44</b>	10.56 ± 1.21	6.56 ± 0.95	8.21 ± 0.45	7.71 ± 0.44	6.91 ± 0.54
Vowel	70.65 ± 0.78	63.03 ± 0.78	69.70 ± 0.77	<b>61.32 ± 0.68</b>	63.90 ± 0.65	64.13 ± 0.66	<b>62.77 ± 0.65</b>
Image	27.32 ± 0.64	22.38 ± 0.67	26.91 ± 0.75	23.07 ± 0.90	<b>18.80 ± 0.66</b>	<b>19.45 ± 0.65</b>	<b>19.82 ± 0.67</b>
Vehicle	39.43 ± 0.90	45.61 ± 0.78	48.44 ± 1.10	46.86 ± 0.91	<b>38.22 ± 0.79</b>	43.01 ± 0.54	<b>37.48 ± 0.74</b>
German	32.30 ± 0.55	<b>29.31 ± 0.31</b>	32.76 ± 0.65	<b>29.45 ± 0.35</b>	30.96 ± 0.42	30.94 ± 0.33	30.62 ± 0.43
Satimage	31.01 ± 0.73	22.59 ± 0.58	34.79 ± 0.68	25.12 ± 1.43	20.15 ± 0.40	<b>18.75 ± 0.34</b>	<b>18.96 ± 0.39</b>
Senseval-2	<b>32.72 ± 0.62</b>	35.56 ± 0.73	37.14 ± 1.10	36.37 ± 0.83	34.66 ± 0.71	37.77 ± 0.67	<b>33.11 ± 0.74</b>

sion, as well as learning with local and global consistency (LGC; Zhou et al., 2003) with an out-of-sample extension. They represent the state-of-the-art manifold regularization and similarity graph transduction respectively. Similarly to SMIR, their optimizations are convex and can be solved analytically. LapRLS is extended using the one-vs-rest trick, and LGC is extended via the Nadaraya-Watson estimator (Delalleau et al., 2005). The experimental setup and the candidates of hyperparameters for LapRLS and LGC are same as SMIR, except that the regularization parameter  $\alpha$  of LGC is chosen from  $\{0.2, 0.4, 0.6, 0.8, 0.99\}$ . SMIR was always best or tie in Table 4, and thus it is fairly competitive with those pure geometric methods on these benchmarks.

Finally, we take all seven methods and compare their performance on the sixteen benchmarks listed in Table 2. The experimental results are reported in Tables 5 and 6, where the experimental setup and the candidates of hyperparameters are same as previous experiments. To be clear, there are two benchmarks, BCI and Wine, whose sample size is less than 500. As a result, each of their random samplings included the whole set, and the randomness or the difference of the classification error was actually from how the training, test and cross-validation data were split and also how labeled data were selected. We can see from Table 5 that ER, LGC and SMIR were best or comparable on three benchmarks, and KerER, XR and KerXR were best or comparable on two benchmarks. Moreover, in

Table 6, SMIR won or tied five times, while all other methods except XR won or tied twice. Therefore, it is reasonable and practical to maximize SMI following the information maximization principle, and SMIR is a promising information-theoretic approach to semi-supervised learning.

## 6. Conclusions

In this paper, we proposed squared-loss mutual information regularization (SMIR). Compared with other information-theoretic regularization, SMIR is convex with no logarithm in the involved optimization problem, and thus enables the analytic expression of the globally optimal solution. We established novel data-dependent generalization error bounds that even incorporate the information of unlabeled data. We then evaluated SMIR on twenty benchmark data sets, and the results demonstrated that SMIR compared favorably with entropy regularization, expectation regularization, manifold regularization, and similarity graph transduction.

## Acknowledgments

GN was supported by the MEXT scholarship 103250, WJ was supported by the Okazaki Kaheita International Scholarship Foundation, HH was supported by the FIRST program, and MS was supported by the MEXT KAKENHI 25700022.



## References

- Agakov, F. and Barber, D. Kernelized infomax clustering. In *NIPS*, 2006.
- Ali, S. M. and Silvey, S. D. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society, Series B*, 28(1):131–142, 1966.
- Allwein, E., Schapire, R., and Singer, Y. Reducing multi-class to binary: a unifying approach for margin classifiers. *Journal of Machine Learning Research*, 1:113–141, 2000.
- Bartlett, P. and Mendelson, S. Rademacher and Gaussian complexities: risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- Belkin, M., Matveeva, I., and Niyogi, P. Regularization and semi-supervised learning on large graphs. In *ALT*, 2004.
- Belkin, M., Niyogi, P., and Sindhvani, V. Manifold regularization: a geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research*, 7:2399–2434, 2006.
- Bennett, K. and Demiriz, A. Semi-supervised support vector machines. In *NIPS*, 1998.
- Chapelle, O., Schölkopf, B., and Zien, A. (eds.). *Semi-Supervised Learning*. MIT Press, 2006.
- Cortes, C., Mohri, M., Pechyony, D., and Rastogi, A. Stability of transductive regression algorithms. In *ICML*, 2008.
- Csiszár, I. Information-type measures of difference of probability distributions and indirect observation. *Studia Scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.
- Delalleau, O., Bengio, Y., and Le Roux, N. Efficient non-parametric function induction in semi-supervised learning. In *AISTATS*, 2005.
- El-Yaniv, R. and Pechyony, D. Transductive Rademacher complexity and its applications. *Journal of Artificial Intelligence Research*, 35:193–234, 2009.
- Gomes, R., Krause, A., and Perona, P. Discriminative clustering by regularized information maximization. In *NIPS*, 2010.
- Grandvalet, Y. and Bengio, Y. Semi-supervised learning by entropy minimization. In *NIPS*, 2004.
- Hall, M. Correlation-based feature selection for discrete and numeric class machine learning. In *ICML*, 2000.
- Joachims, T. Transductive inference for text classification using support vector machines. In *ICML*, 1999.
- Joachims, T. Transductive learning via spectral graph partitioning. In *ICML*, 2003.
- Kanamori, T., Suzuki, T., and Sugiyama, M. Statistical analysis of kernel-based least-squares density-ratio estimation. *Machine Learning*, 86(3):335–367, 2012.
- Koltchinskii, V. Rademacher penalties and structural risk minimization. *IEEE Transactions on Information Theory*, 47(5):1902–1914, 2001.
- Kullback, S. and Leibler, R. A. On information and sufficiency. *Annals of Mathematical Statistics*, 22:79–86, 1951.
- Mann, G. and McCallum, A. Simple, robust, scalable semi-supervised learning via expectation regularization. In *ICML*, 2007.
- McDiarmid, C. On the method of bounded differences. In Siemons, J. (ed.), *Surveys in Combinatorics*, pp. 148–188. Cambridge University Press, 1989.
- Meir, R. and Zhang, T. Generalization error bounds for Bayesian mixture algorithms. *Journal of Machine Learning Research*, 4:839–860, 2003.
- Pearson, K. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. *Philosophical Magazine*, 50:157–175, 1900.
- Schölkopf, B. and Smola, A. *Learning with Kernels*. MIT Press, 2001.
- Shannon, C. E. A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423 & 623–656, 1948.
- Sinha, K. and Belkin, M. Semi-supervised learning using sparse eigenfunction bases. In *NIPS*, 2009.
- Sugiyama, M. Superfast-trainable multi-class probabilistic classifier by least-squares posterior fitting. *IEICE Transactions on Information and Systems*, E93-D(10):2690–2701, 2010.
- Sugiyama, M., Yamada, M., Kimura, M., and Hachiya, H. On information-maximization clustering: Tuning parameter selection and analytic solution. In *ICML*, 2011.
- Suzuki, T., Sugiyama, M., Kanamori, T., and Sese, J. Mutual information estimation reveals global associations between stimuli and biological processes. *BMC Bioinformatics*, 10(1):S52, 2009.
- Szummer, M. and Jaakkola, T. Partially labeled classification with Markov random walks. In *NIPS*, 2001.
- Szummer, M. and Jaakkola, T. Information regularization with partially labeled data. In *NIPS*, 2002.
- Yamada, M., Sugiyama, M., Wichern, G., and Simm, J. Improving the accuracy of least-squares probabilistic classifiers. *IEICE Transactions on Information and Systems*, E94-D(6):1337–1340, 2011.
- Zhou, D., Bousquet, O., Navin Lal, T., Weston, J., and Schölkopf, B. Learning with local and global consistency. In *NIPS*, 2003.
- Zhu, X., Ghahramani, Z., and Lafferty, J. Semi-supervised learning using Gaussian fields and harmonic functions. In *ICML*, 2003.