

# PAC-Bayes Analysis of Multi-view Learning

**Shiliang Sun**

SHILIANGSUN@GMAIL.COM

*Shanghai Key Laboratory of Multidimensional Information Processing  
Department of Computer Science and Technology  
East China Normal University  
500 Dongchuan Road, Shanghai 200241, China*

**John Shawe-Taylor**

J.SHAWE-TAYLOR@UCL.AC.UK

*Department of Computer Science  
University College London  
Gower Street, London WC1E 6BT, United Kingdom*

**Liang Mao**

LMAO14@OUTLOOK.COM

*Shanghai Key Laboratory of Multidimensional Information Processing  
Department of Computer Science and Technology  
East China Normal University  
500 Dongchuan Road, Shanghai 200241, China*

**Editor:**

## Abstract

This paper presents eight PAC-Bayes bounds to analyze the generalization performance of multi-view classifiers. These bounds adopt data dependent Gaussian priors which emphasize classifiers with high view agreements. The center of the prior for the first two bounds is the origin, while the center of the prior for the third and fourth bounds is given by a data dependent vector. An important technique to obtain these bounds is two derived logarithmic determinant inequalities whose difference lies in whether the dimensionality of data is involved. The centers of the fifth and sixth bounds are calculated on a separate subset of the training set. The last two bounds use unlabeled data to represent view agreements and are thus applicable to semi-supervised multi-view learning. We evaluate all the presented multi-view PAC-Bayes bounds on benchmark data and compare them with previous single-view PAC-Bayes bounds. The usefulness and performance of the multi-view bounds are discussed.

**Keywords:** PAC-Bayes bound, statistical learning theory, support vector machine, multi-view learning

## 1. Introduction

Multi-view learning is a promising research direction with prevalent applicability (Sun, 2013). For instance, in multimedia content understanding, multimedia segments can be described by both their video and audio signals, and the video and audio signals are regarded as the two views. Learning from data relies on collecting data that contain a sufficient signal and encoding our prior knowledge in increasingly sophisticated regularization schemes that enable the signal to be extracted. With certain co-regularization schemes, multi-view learning performs well on various learning tasks.

Statistical learning theory (SLT) provides a general framework to analyze the generalization performance of machine learning algorithms. The theoretical outcomes can be used to motivate algorithm design, select models or give insights on the effects and behaviors of some interesting quantities. For example, the well-known large margin principle in support vector machines (SVMs) is well supported by various SLT bounds (Vapnik, 1998; Bartlett and Mendelson, 2002; Sun and Shawe-Taylor, 2010). Different from early bounds that often rely on the complexity measures of the considered function classes, the recent PAC-Bayes bounds (McAllester, 1999; Seeger, 2002; Langford, 2005) give the tightest predictions of the generalization performance, for which the prior and posterior distributions of learners are involved on top of the PAC (Probably Approximately Correct) learning setting (Catoni, 2007; Germain et al., 2009). Beyond the common supervised learning, PAC-Bayes analysis has also been applied to other tasks, e.g., density estimation (Seldin and Tishby, 2010; Higgs and Shawe-Taylor, 2010) and reinforcement learning (Seldin et al., 2012).

Although the field of multi-view learning has enjoyed a great success with algorithms and applications and is provided with some theoretical results, PAC-Bayes analysis of multi-view learning is still absent. In this paper, we attempt to fill the gap between the developments in theory and practice by proposing new PAC-Bayes bounds for multi-view learning.

An earlier attempt to analyze the generalization of two-view learning was made using Rademacher complexity (Farquhar et al., 2006; Rosenberg and Bartlett, 2007). The bound relied on estimating the empirical Rademacher complexity of the class of pairs of functions from the two views that are matched in expectation under the data generating distribution. Hence, this approach also implicitly relied on the data generating distribution to define the function class (and hence prior). The current paper makes the definition of the prior in terms of the data generating distribution explicit through the PAC-Bayes framework and provides several bounds. However, the main advantage is that it defines a framework that makes explicit the definition of the prior in terms of the data generating distribution, setting a template for other related approaches to encoding complex prior knowledge that relies on the data generating distribution.

Kakade and Foster (2007) characterized the expected regret of a semi-supervised multi-view regression algorithm. The results given by Sridharan and Kakade (2008) take an information theoretic approach that involves a number of assumptions that may be difficult to check in practice. With these assumptions theoretical results including PAC-style analysis to bound expected losses were given, which involve some Bayes optimal predictor and but cannot provide computable classification error bounds since the data generating distribution is usually unknown. These results therefore represent a related but distinct approach.

We adopt a PAC-Bayes analysis where we encode our assumptions through priors defined in terms of the data generating distribution. Such priors have been studied by Catoni (2007) under the name of localized priors and more recently by Lever et al. (2013) as data distribution dependent priors. Both papers considered schemes for placing a prior over classifiers defined through their true generalization errors. In contrast, the prior that we consider is mainly used to encode the assumption about the relationship between the two views in the data generating distribution. Such data distribution dependent priors cannot be subjected to traditional Bayesian analysis since we do not have an explicit form for the

prior, making inference impossible. Hence, this paper illustrates one of the advantages that arise from the PAC-Bayes framework.

The PAC-Bayes theorem bounds the true error of the distribution of classifiers in terms of a term from the sample complexity and the KL divergence between the posterior and the prior distributions of classifiers. The key technical innovations of the paper enable the bounding of the KL divergence term in terms of empirical quantities despite involving priors that cannot be computed. This approach was adopted in Parrado-Hernández et al. (2012) for some simple priors such as the Gaussian centered at  $\mathbb{E}[y\phi(\mathbf{x})]$ . The current paper treats a significantly more sophisticated case where the priors encode our expectation that good weight vectors can be found that give similar outputs from both views.

Specifically, we first provide four PAC-Bayes bounds using priors that reflect how well the two views agree on average over all examples. The first two bounds use a Gaussian prior centered at the origin, while the third and fourth ones adopt a different prior whose center is not the origin. However, the formulations of the priors involve mathematical expectations with respect to the unknown data distributions. We manage to bound the expectation related terms with their empirical estimations on a finite sample of data. Then, we further provide two PAC-Bayes bounds using a part of the training data to determine priors, and two PAC-Bayes bounds for semi-supervised multi-view learning where unlabeled data are involved in the definition of the priors.

When a natural feature split does not exist, multi-view learning could still obtain performance improvements with manufactured splits, provided that each of the views contains not only enough information for the learning task itself, but some knowledge that other views do not have. It is therefore important that people should split features into views satisfying the assumptions. However, data split is still an open question and beyond the scope of this paper.

The rest of this paper is organized as follows. After briefly reviewing the PAC-Bayes bound for SVMs in Section 2, we give and derive four multi-view PAC-Bayes bounds involving only empirical quantities in Section 3 and Section 4. Then we give two bounds whose centers are calculated on a separate subset of the training data in Section 5. After that, we present two semi-supervised multi-view PAC-Bayes bounds in Section 6. The optimization formulations of the related single-view and multi-view SVMs as well as semi-supervised multi-view SVMs are given in Section 7. After evaluating the usefulness and performance of the bounds in Section 8, we give concluding remarks in Section 9.

## 2. PAC-Bayes Bound and Specialization to SVMs

Consider a binary classification problem. Let  $\mathcal{D}$  be the distribution of feature  $\mathbf{x}$  lying in an input space  $\mathcal{X}$  and the corresponding output label  $y$  where  $y \in \{-1, 1\}$ . Suppose  $Q$  is a posterior distribution over the parameters of the classifier  $c$ . Define the true error and empirical error of a classifier as

$$e_{\mathcal{D}} = Pr_{(\mathbf{x}, y) \sim \mathcal{D}}(c(\mathbf{x}) \neq y),$$

$$\hat{e}_S = Pr_{(\mathbf{x}, y) \sim S}(c(\mathbf{x}) \neq y) = \frac{1}{m} \sum_{i=1}^m I(c(\mathbf{x}_i) \neq y_i),$$

where  $S$  is a sample including  $m$  examples, and  $I(\cdot)$  is the indicator function. With the distribution  $Q$ , we can then define the average true error  $E_{Q,\mathcal{D}} = \mathbb{E}_{c \sim Q} e_{\mathcal{D}}$ , and the average empirical error  $\hat{E}_{Q,S} = \mathbb{E}_{c \sim Q} \hat{e}_S$ . The following lemma provides the PAC-Bayes bound on  $E_{Q,\mathcal{D}}$  in the current context of binary classification.

**Theorem 1 (PAC-Bayes Bound (Langford, 2005))** *For any data distribution  $\mathcal{D}$ , for any prior  $P(c)$  over the classifier  $c$ , for any  $\delta \in (0, 1]$ :*

$$Pr_{S \sim \mathcal{D}^m} \left( \forall Q(c) : KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{KL(Q \| P) + \ln(\frac{m+1}{\delta})}{m} \right) \geq 1 - \delta, \quad (1)$$

where  $KL(Q \| P) = \mathbb{E}_{c \sim Q} \ln \frac{Q(c)}{P(c)}$  is the KL divergence between  $Q$  and  $P$ , and  $KL_+(q \| p) = q \ln \frac{q}{p} + (1 - q) \ln \frac{1-q}{1-p}$  for  $p > q$  and 0 otherwise.

Suppose from the  $m$  training examples we learn an SVM classifier represented by  $c_{\mathbf{u}}(\mathbf{x}) = \text{sign}(\mathbf{u}^\top \phi(\mathbf{x}))$ , where  $\phi(\mathbf{x})$  is a projection of the original feature to a certain feature space induced by some kernel function. Define the prior and the posterior of the classifier to be Gaussian with  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\mathbf{u} \sim \mathcal{N}(\mu \mathbf{w}, \mathbf{I})$ , respectively. Note that here  $\|\mathbf{w}\| = 1$ , and thus the distance between the center of the posterior and the origin is  $\mu$ . With this specialization, we give the PAC-Bayes bound for SVMs (Langford, 2005; Parrado-Hernández et al., 2012) below.

**Theorem 2** *For any data distribution  $\mathcal{D}$ , for any  $\delta \in (0, 1]$ , we have*

$$Pr_{S \sim \mathcal{D}^m} \left( \forall \mathbf{w}, \mu : KL_+(\hat{E}_{Q,S}(\mathbf{w}, \mu) \| E_{Q,\mathcal{D}}(\mathbf{w}, \mu)) \leq \frac{\frac{\mu^2}{2} + \ln(\frac{m+1}{\delta})}{m} \right) \geq 1 - \delta, \quad (2)$$

where  $\|\mathbf{w}\| = 1$ .

All that remains is calculating the empirical stochastic error rate  $\hat{E}_{Q,S}$ . It can be shown that for a posterior  $Q = \mathcal{N}(\mu \mathbf{w}, \mathbf{I})$  with  $\|\mathbf{w}\| = 1$ , we have

$$\hat{E}_{Q,S} = \mathbb{E}_S \left[ \tilde{F}(\mu \gamma(\mathbf{x}, y)) \right], \quad (3)$$

where  $\mathbb{E}_S$  is the average over the  $m$  training examples,  $\gamma(\mathbf{x}, y)$  is the normalized margin of the example

$$\gamma(\mathbf{x}, y) = y \mathbf{w}^\top \phi(\mathbf{x}) / \|\phi(\mathbf{x})\|, \quad (4)$$

and  $\tilde{F}(x)$  is the Gaussian cumulative distribution

$$\tilde{F}(x) = \int_x^\infty \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx. \quad (5)$$

The generalization error of the original SVM classifier  $c_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \phi(\mathbf{x}))$  can be bounded by at most twice the average true error  $E_{Q,\mathcal{D}}(\mathbf{w}, \mu)$  of the corresponding stochastic classifier (Langford and Shawe-Taylor, 2002). That is, for any  $\mu$  we have

$$Pr_{(\mathbf{x}, y) \sim \mathcal{D}} \left( \text{sign}(\mathbf{w}^\top \phi(\mathbf{x})) \neq y \right) \leq 2E_{Q,\mathcal{D}}(\mathbf{w}, \mu). \quad (6)$$

### 3. Multi-view PAC-Bayes Bounds

We propose a new data dependent prior for PAC-Bayes analysis of multi-view learning. In particular, we take the distribution on the concatenation of the two weight vectors  $\mathbf{u}_1$  and  $\mathbf{u}_2$  as their individual product:  $\tilde{P}([\mathbf{u}_1^\top, \mathbf{u}_2^\top]^\top) = P_1(\mathbf{u}_1)P_2(\mathbf{u}_2)$  but then weight it in some manner associated with how well the two weights agree averagely on all examples. That is, the prior is

$$P([\mathbf{u}_1^\top, \mathbf{u}_2^\top]^\top) \propto P_1(\mathbf{u}_1)P_2(\mathbf{u}_2)V(\mathbf{u}_1, \mathbf{u}_2),$$

where  $P_1(\mathbf{u}_1)$  and  $P_2(\mathbf{u}_2)$  are Gaussian with zero mean and identity covariance, and

$$V(\mathbf{u}_1, \mathbf{u}_2) = \exp \left\{ -\frac{1}{2\sigma^2} \mathbb{E}_{(\mathbf{x}_1, \mathbf{x}_2)} (\mathbf{x}_1^\top \mathbf{u}_1 - \mathbf{x}_2^\top \mathbf{u}_2)^2 \right\}.$$

To specialize the PAC-Bayes bound for multi-view learning, we consider classifiers of the form

$$c(\mathbf{x}) = \text{sign}(\mathbf{u}^\top \phi(\mathbf{x})), \quad (7)$$

where  $\mathbf{u} = [\mathbf{u}_1^\top, \mathbf{u}_2^\top]^\top$  is the concatenated weight vector from two views, and  $\phi(\mathbf{x})$  can be the concatenated  $\mathbf{x} = [\mathbf{x}_1^\top, \mathbf{x}_2^\top]^\top$  itself or a concatenation of maps of  $\mathbf{x}$  to kernel-induced feature spaces. Note that  $\mathbf{x}_1$  and  $\mathbf{x}_2$  indicate features of one example from the two views, respectively. For simplicity, here we use the original features to derive our results, though kernel maps can be implicitly employed as well. Our dimensionality independent bounds work even when the dimension of the kernelized feature space goes to infinity.

According to our setting, the classifier prior is fixed to be

$$P(\mathbf{u}) \propto \mathcal{N}(\mathbf{0}, \mathbf{I}) \times V(\mathbf{u}_1, \mathbf{u}_2), \quad (8)$$

Function  $V(\mathbf{u}_1, \mathbf{u}_2)$  makes the prior place large probability mass on parameters with which the classifiers from two views agree well on all examples averagely. The posterior is chosen to be of the form

$$Q(\mathbf{u}) = \mathcal{N}(\mu \mathbf{w}, \mathbf{I}), \quad (9)$$

where  $\|\mathbf{w}\| = 1$ .

Define  $\tilde{\mathbf{x}} = [\mathbf{x}_1^\top, -\mathbf{x}_2^\top]^\top$ . We have

$$\begin{aligned} P(\mathbf{u}) &\propto \mathcal{N}(\mathbf{0}, \mathbf{I}) \times V(\mathbf{u}_1, \mathbf{u}_2) \\ &\propto \exp \left\{ -\frac{1}{2} \mathbf{u}^\top \mathbf{u} \right\} \times \exp \left\{ -\frac{1}{2\sigma^2} \mathbb{E}_{(\mathbf{x}_1, \mathbf{x}_2)} (\mathbf{x}_1^\top \mathbf{u}_1 - \mathbf{x}_2^\top \mathbf{u}_2)^2 \right\} \\ &= \exp \left\{ -\frac{1}{2} \mathbf{u}^\top \mathbf{u} \right\} \times \exp \left\{ -\frac{1}{2\sigma^2} \mathbb{E}_{\tilde{\mathbf{x}}} (\mathbf{u}^\top \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \mathbf{u}) \right\} \\ &= \exp \left\{ -\frac{1}{2} \mathbf{u}^\top \mathbf{u} \right\} \times \exp \left\{ -\frac{1}{2\sigma^2} \mathbf{u}^\top \mathbb{E}(\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top) \mathbf{u} \right\} \\ &= \exp \left\{ -\frac{1}{2} \mathbf{u}^\top \left( \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top)}{\sigma^2} \right) \mathbf{u} \right\}. \end{aligned}$$

That is,  $P(\mathbf{u}) = \mathcal{N}(\mathbf{0}, \Sigma)$  with  $\Sigma = \left( \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top)}{\sigma^2} \right)^{-1}$ .

Suppose  $\dim(\mathbf{u}) = d$ . Given the above prior and posterior, we have the following theorem to characterize their divergence.

**Theorem 3**

$$KL(Q(\mathbf{u})\|P(\mathbf{u})) = \frac{1}{2} \left( -\ln \left( \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \right) + \frac{1}{\sigma^2} \mathbb{E}[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2(\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \right). \quad (10)$$

**Proof** It is easy to show that the KL divergence between two Gaussians (Rasmussen and Williams, 2006) in an  $N$ -dimensional space is

$$KL(\mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0)\|\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)) = \frac{1}{2} \left( \ln \left( \frac{|\Sigma_1|}{|\Sigma_0|} \right) + \text{tr}(\Sigma_1^{-1}\Sigma_0) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \Sigma_1^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - d \right). \quad (11)$$

The KL divergence between the posterior and prior is thus

$$\begin{aligned} KL(Q(\mathbf{u})\|P(\mathbf{u})) &= \frac{1}{2} \left( -\ln \left( \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \right) + \text{tr} \left( \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right) + \mu^2 \mathbf{w}^\top \left( \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right) \mathbf{w} - d \right) \\ &= \frac{1}{2} \left( -\ln \left( \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \right) + \text{tr} \left( \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right) + \mu^2 \mathbf{w}^\top \left( \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right) \mathbf{w} + \mu^2 \right) \\ &= \frac{1}{2} \left( -\ln \left( \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \right) + \frac{1}{\sigma^2} \mathbb{E}[\text{tr}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)] + \frac{\mu^2}{\sigma^2} \mathbb{E}[(\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \right) \\ &= \frac{1}{2} \left( -\ln \left( \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \right) + \frac{1}{\sigma^2} \mathbb{E}[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}}] + \frac{\mu^2}{\sigma^2} \mathbb{E}[(\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \right) \\ &= \frac{1}{2} \left( -\ln \left( \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \right) + \frac{1}{\sigma^2} \mathbb{E}[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2(\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \right), \end{aligned}$$

which completes the proof. ■

The problem with this expression is that it contains expectations over the input distribution that we are unable to compute. This is because we have defined the prior distribution in terms of the input distribution via the  $V$  function. Such priors are referred to as localized by Catoni (2007). While his work considered specific examples of such priors that satisfy certain optimality conditions, the definition we consider here is encoding natural prior assumptions about the link between the input distribution and the classification function, namely that it will have a simple representation in both views. This is an example of luckiness (Shawe-Taylor et al., 1998), where generalization is estimated making assumptions that if proven true lead to tighter bounds, as for example in the case of a large margin classifier.

We now develop methods that estimate the relevant quantities in (10) from empirical data, so that there will be additional empirical estimations involved in the final bounds besides the usual empirical error.

We proceed to provide and prove two inequalities on the involved logarithmic determinant function, which are very important for the subsequent multi-view PAC-Bayes bounds.

**Theorem 4**

$$-\ln \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \leq -d \ln \mathbb{E} \left[ \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right|^{1/d} \right], \quad (12)$$

$$-\ln \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \leq -\mathbb{E} \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right|. \quad (13)$$

**Proof** According to the Minkowski determinant theorem, for  $n \times n$  positive semi-definite matrices  $A$  and  $B$ , the following inequality holds

$$|A + B|^{1/n} \geq |A|^{1/n} + |B|^{1/n}, \quad (14)$$

which implies that the function  $A \mapsto |A|^{1/n}$  is concave on the set of  $n \times n$  positive semi-definite matrices. Therefore, with Jensen's inequality we have

$$\begin{aligned} -\ln \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| &= -d \ln \left| \mathbb{E} \left( \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right) \right|^{1/d} \\ &\leq -d \ln \mathbb{E} \left[ \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right|^{1/d} \right]. \end{aligned}$$

Since the natural logarithm is concave, we further have

$$-d \ln \mathbb{E} \left[ \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right|^{1/d} \right] \leq -d \mathbb{E} \left[ \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right|^{1/d} \right] = -\mathbb{E} \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right|,$$

and thereby

$$-\ln \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \leq -\mathbb{E} \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right|. \quad (15)$$

■

Denote  $R = \sup_{\tilde{\mathbf{x}}} \|\tilde{\mathbf{x}}\|$ . From inequality (12), we can finally prove the following theorem, as detailed in Appendix A.

**Theorem 5 (Multi-view PAC-Bayes bound 1)** *Consider a classifier prior given in (8) and a classifier posterior given in (9). For any data distribution  $\mathcal{D}$ , for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$ , the following inequality holds*

$$\begin{aligned} \forall \mathbf{w}, \mu : KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) &\leq \\ \frac{-\frac{d}{2} \ln \left[ f_m - (\sqrt[d]{(R/\sigma)^2 + 1} - 1) \sqrt{\frac{1}{2m} \ln \frac{3}{\delta}} \right]_+ + \frac{H_m}{2\sigma^2} + \frac{(1+\mu^2)R^2}{2\sigma^2} \sqrt{\frac{1}{2m} \ln \frac{3}{\delta}} + \frac{\mu^2}{2} + \ln \left( \frac{m+1}{\delta/3} \right)}{m}, \end{aligned}$$

where

$$\begin{aligned} f_m &= \frac{1}{m} \sum_{i=1}^m \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right|^{1/d}, \\ H_m &= \frac{1}{m} \sum_{i=1}^m [\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2], \end{aligned}$$

and  $\|\mathbf{w}\| = 1$ .

From the bound formulation, we see that if  $(\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2$  is small, that is, if the two view outputs tend to agree, the bound will be tight.

Note that, although the formulation of  $f_m$  involves the outer product of feature vectors, it can actually be represented by the inner product, which is obvious through the following determinant equality

$$\left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right| = \frac{\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i}{\sigma^2} + 1, \quad (16)$$

where we have used the fact that matrix  $\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top$  has rank 1 and has only one nonzero eigenvalue.

We can use inequality (13) instead of (12) to derive a  $d$ -independent bound (see Theorem 6 below), which is independent of the dimensionality of the feature representation space.

**Theorem 6 (Multi-view PAC-Bayes bound 2)** *Consider a classifier prior given in (8) and a classifier posterior given in (9). For any data distribution  $\mathcal{D}$ , for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$ , the following inequality holds*

$$\forall \mathbf{w}, \mu : KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{\tilde{f}/2 + \frac{1}{2} \left( \frac{(1+\mu^2)R^2}{\sigma^2} + \ln(1 + \frac{R^2}{\sigma^2}) \right) \sqrt{\frac{1}{2m} \ln \frac{2}{\delta}} + \frac{\mu^2}{2} + \ln \left( \frac{m+1}{\delta/2} \right)}{m},$$

where

$$\tilde{f} = \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{\sigma^2} [\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2] - \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right| \right), \quad (17)$$

and  $\|\mathbf{w}\| = 1$ .

The proof of this theorem is given in Appendix B.

Since this bound is independent with  $d$  and the term  $\left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right|$  involving the outer product can be represented by the inner product through (16), this bound can be employed when the dimension of the kernelized feature space goes to infinity.

#### 4. Another Two Multi-view PAC-Bayes Bounds

We further propose a new prior whose center is not located at the origin, inspired by Parrado-Hernández et al. (2012). The new classifier prior is

$$P(\mathbf{u}) \propto \mathcal{N}(\eta \mathbf{w}_p, \mathbf{I}) \times V(\mathbf{u}_1, \mathbf{u}_2), \quad (18)$$

and the posterior is still

$$Q(\mathbf{u}) = \mathcal{N}(\mu \mathbf{w}, \mathbf{I}), \quad (19)$$

where  $\eta > 0$ ,  $\|\mathbf{w}\| = 1$  and  $\mathbf{w}_p = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[y \mathbf{x}]$  (or  $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[y \phi(\mathbf{x})]$  in a predefined kernel space) with  $\mathbf{x} = [\mathbf{x}_1^\top, \mathbf{x}_2^\top]^\top$ .

We have

$$\begin{aligned} P(\mathbf{u}) &\propto \mathcal{N}(\eta \mathbf{w}_p, \mathbf{I}) \times V(\mathbf{u}_1, \mathbf{u}_2) \\ &\propto \exp \left\{ -\frac{1}{2} (\mathbf{u} - \eta \mathbf{w}_p)^\top (\mathbf{u} - \eta \mathbf{w}_p) \right\} \times \exp \left\{ -\frac{1}{2\sigma^2} \mathbf{u}^\top \mathbb{E}(\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top) \mathbf{u} \right\}. \end{aligned}$$

That is,  $P(\mathbf{u}) = \mathcal{N}(\mathbf{u}_p, \Sigma)$  with  $\Sigma = \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right)^{-1}$  and  $\mathbf{u}_p = \eta\Sigma\mathbf{w}_p$ .

With  $d$  being the dimensionality of  $\mathbf{u}$ , the KL divergence between the posterior and prior is

$$\begin{aligned} & KL(Q(\mathbf{u})\|P(\mathbf{u})) \\ &= \frac{1}{2} \left( -\ln\left|\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right| + \text{tr}\left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right) + (\mathbf{u}_p - \mu\mathbf{w})^\top \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right) (\mathbf{u}_p - \mu\mathbf{w}) - d \right) \\ &= \frac{1}{2} \left( -\ln\left|\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right| + \frac{1}{\sigma^2}\mathbb{E}[\tilde{\mathbf{x}}^\top\tilde{\mathbf{x}}] + (\mathbf{u}_p - \mu\mathbf{w})^\top \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right) (\mathbf{u}_p - \mu\mathbf{w}) \right). \end{aligned} \quad (20)$$

We have

$$\begin{aligned} & (\mathbf{u}_p - \mu\mathbf{w})^\top \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right) (\mathbf{u}_p - \mu\mathbf{w}) \\ &= \eta^2\mathbf{w}_p^\top \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right)^{-1} \mathbf{w}_p - 2\eta\mu\mathbf{w}_p^\top \mathbf{w} + \mu^2\mathbf{w}^\top \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right) \mathbf{w} \\ &= \eta^2\mathbf{w}_p^\top \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right)^{-1} \mathbf{w}_p - 2\eta\mu\mathbf{w}_p^\top \mathbf{w} + \frac{\mu^2}{\sigma^2}\mathbb{E}[(\mathbf{w}^\top\tilde{\mathbf{x}})^2] + \mu^2 \\ &= \eta^2\mathbf{w}_p^\top \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right)^{-1} \mathbf{w}_p - 2\eta\mu\mathbb{E}[y(\mathbf{w}^\top\mathbf{x})] + \frac{\mu^2}{\sigma^2}\mathbb{E}[(\mathbf{w}^\top\tilde{\mathbf{x}})^2] + \mu^2 \\ &\leq \eta^2\mathbf{w}_p^\top \mathbf{w}_p - 2\eta\mu\mathbb{E}[y(\mathbf{w}^\top\mathbf{x})] + \frac{\mu^2}{\sigma^2}\mathbb{E}[(\mathbf{w}^\top\tilde{\mathbf{x}})^2] + \mu^2, \end{aligned} \quad (21)$$

where for the last inequality we have used the fact that matrix  $\mathbf{I} - \left(\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right)^{-1}$  is symmetric and positive semi-definite.

Define  $\hat{\mathbf{w}}_p = \mathbb{E}_{(\mathbf{x}, y) \sim S}[y\mathbf{x}] = \frac{1}{m} \sum_{i=1}^m [y_i\mathbf{x}_i]$ . We have

$$\begin{aligned} \eta^2\mathbf{w}_p^\top \mathbf{w}_p &= \|\eta\mathbf{w}_p - \mu\mathbf{w} + \mu\mathbf{w}\|^2 \\ &= \|\eta\mathbf{w}_p - \mu\mathbf{w}\|^2 + \mu^2 + 2(\eta\mathbf{w}_p - \mu\mathbf{w})^\top \mu\mathbf{w} \\ &\leq \|\eta\mathbf{w}_p - \mu\mathbf{w}\|^2 + \mu^2 + 2\mu\|\eta\mathbf{w}_p - \mu\mathbf{w}\| \\ &= (\|\eta\mathbf{w}_p - \mu\mathbf{w}\| + \mu)^2. \end{aligned} \quad (22)$$

Moreover, we have

$$\|\eta\mathbf{w}_p - \mu\mathbf{w}\| = \|\eta\mathbf{w}_p - \eta\hat{\mathbf{w}}_p + \eta\hat{\mathbf{w}}_p - \mu\mathbf{w}\| \leq \|\eta\mathbf{w}_p - \eta\hat{\mathbf{w}}_p\| + \|\eta\hat{\mathbf{w}}_p - \mu\mathbf{w}\|. \quad (23)$$

From (20), (21), (22) and (23), it follows that

$$\begin{aligned} KL(Q(\mathbf{u})\|P(\mathbf{u})) &\leq -\frac{1}{2} \ln\left|\mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right| + \frac{1}{2} (\|\eta\mathbf{w}_p - \eta\hat{\mathbf{w}}_p\| + \|\eta\hat{\mathbf{w}}_p - \mu\mathbf{w}\| + \mu)^2 + \\ &\quad \frac{1}{2\sigma^2}\mathbb{E} \left[ \tilde{\mathbf{x}}^\top\tilde{\mathbf{x}} - 2\eta\mu\sigma^2 y(\mathbf{w}^\top\mathbf{x}) + \mu^2(\mathbf{w}^\top\tilde{\mathbf{x}})^2 \right] + \frac{\mu^2}{2}. \end{aligned} \quad (24)$$

By using inequalities (12) and (13), we get the following two theorems, whose proofs are detailed in Appendix C and Appendix D, respectively.

**Theorem 7 (Multi-view PAC-Bayes bound 3)** Consider a classifier prior given in (18) and a classifier posterior given in (19). For any data distribution  $\mathcal{D}$ , for any  $\mathbf{w}$ , positive  $\mu$ , and positive  $\eta$ , for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$  the following multi-view PAC-Bayes bound holds

$$KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{-\frac{d}{2} \ln \left[ f_m - (\sqrt[4]{(R/\sigma)^2 + 1} - 1) \sqrt{\frac{1}{2m} \ln \frac{4}{\delta}} \right]_+}{m} + \frac{\frac{1}{2} \left( \frac{\eta R}{\sqrt{m}} \left( 2 + \sqrt{2 \ln \frac{4}{\delta}} \right) + \|\eta \hat{\mathbf{w}}_p - \mu \mathbf{w}\| + \mu \right)^2 + \frac{\hat{H}_m}{2\sigma^2} + \frac{R^2 + \mu^2 R^2 + 4\eta\mu\sigma^2 R}{2\sigma^2} \sqrt{\frac{1}{2m} \ln \frac{4}{\delta}} + \frac{\mu^2}{2} + \ln \left( \frac{m+1}{\delta/4} \right)}{m},$$

where

$$f_m = \frac{1}{m} \sum_{i=1}^m \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right|^{1/d},$$

$$\hat{H}_m = \frac{1}{m} \sum_{i=1}^m [\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i - 2\eta\mu\sigma^2 y_i (\mathbf{w}^\top \mathbf{x}_i) + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2],$$

and  $\|\mathbf{w}\| = 1$ .

Besides the term  $(\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2$  that appears in the previous bounds, we can see that if  $\|\eta \hat{\mathbf{w}}_p - \mu \mathbf{w}\|$  is small, that is, the centers of the prior and posterior tend to overlap, the bound will be tight.

**Theorem 8 (Multi-view PAC-Bayes bound 4)** Consider a classifier prior given in (18) and a classifier posterior given in (19). For any data distribution  $\mathcal{D}$ , for any  $\mathbf{w}$ , positive  $\mu$ , and positive  $\eta$ , for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$  the following multi-view PAC-Bayes bound holds

$$KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{\frac{1}{2} \left( \frac{\eta R}{\sqrt{m}} \left( 2 + \sqrt{2 \ln \frac{3}{\delta}} \right) + \|\eta \hat{\mathbf{w}}_p - \mu \mathbf{w}\| + \mu \right)^2}{m} + \frac{\frac{\tilde{H}_m}{2} + \frac{R^2 + 4\eta\mu\sigma^2 R + \mu^2 R^2 + \sigma^2 \ln(1 + \frac{R^2}{\sigma^2})}{2\sigma^2} \sqrt{\frac{1}{2m} \ln \frac{3}{\delta}} + \frac{\mu^2}{2} + \ln \left( \frac{m+1}{\delta/3} \right)}{m},$$

where

$$\tilde{H}_m = \frac{1}{m} \sum_{i=1}^m \left[ \frac{\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i - 2\eta\mu\sigma^2 y_i (\mathbf{w}^\top \mathbf{x}_i) + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2}{\sigma^2} - \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right| \right], \quad (25)$$

and  $\|\mathbf{w}\| = 1$ .

## 5. Separate Training Data Dependent Multi-view PAC-Bayes Bounds

We attempt to improve our bounds by using a separate set of training data to determine new priors, inspired by Ambroladze et al. (2007) and Parrado-Hernández et al. (2012). We consider a spherical Gaussian whose center is calculated on a subset  $T$  of training set

comprising  $r$  training patterns and labels. In the experiments this is taken as a random subset, but for simplicity of the presentation we will assume  $T$  comprises the last  $r$  examples  $\{\mathbf{x}_k, y_k\}_{k=m-r+1}^m$ .

The new prior is

$$P(\mathbf{u}) = \mathcal{N}(\eta \mathbf{w}_p, \mathbf{I}), \quad (26)$$

and the posterior is again

$$Q(\mathbf{u}) = \mathcal{N}(\mu \mathbf{w}, \mathbf{I}). \quad (27)$$

One reasonable choice of  $\mathbf{w}_p$  is

$$\mathbf{w}_p = \left( \mathbb{E}_{\tilde{\mathbf{x}}}[\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top] \right)^{-1} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[y\mathbf{x}], \quad (28)$$

which is the solution to the following optimization problem

$$\max_{\mathbf{w}} \frac{\mathbb{E}_{\mathbf{x}_1, y}[y\mathbf{w}_1^\top \mathbf{x}_1] + \mathbb{E}_{\mathbf{x}_2, y}[y\mathbf{w}_2^\top \mathbf{x}_2]}{\mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2}[(\mathbf{w}_1^\top \mathbf{x}_1 - \mathbf{w}_2^\top \mathbf{x}_2)^2]}, \quad (29)$$

where  $\mathbf{w} = [\mathbf{w}_1^\top, \mathbf{w}_2^\top]^\top$ . We use the subset  $T$  to approximate  $\mathbf{w}_p$ , that is, let

$$\begin{aligned} \mathbf{w}_p &= \left( \mathbb{E}_{\tilde{\mathbf{x}} \sim T}[\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top] \right)^{-1} \mathbb{E}_{(\mathbf{x}, y) \sim T}[y\mathbf{x}] \\ &= \left( \frac{1}{m-r} \sum_{k=r}^{m-r+1} [\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top] \right)^{-1} \frac{1}{m-r} \sum_{k=r}^{m-r+1} [y_k \mathbf{x}_k]. \end{aligned} \quad (30)$$

The KL divergence between the posterior and prior is

$$KL(Q(\mathbf{u}) \| P(\mathbf{u})) = KL(\mathcal{N}(\mu \mathbf{w}, \mathbf{I}) \| \mathcal{N}(\eta \mathbf{w}_p, \mathbf{I})) = \|\eta \mathbf{w}_p - \mu \mathbf{w}\|^2. \quad (31)$$

Since we separate  $r$  examples to calculate the prior, the actual size of training set that we apply bound to is  $m - r$ . We have the following bound.

**Theorem 9 (Multi-view PAC-Bayes bound 5)** *Consider a classifier prior given in (26) and a classifier posterior given in (27), with  $\mathbf{w}_p$  given in (30). For any data distribution  $\mathcal{D}$ , for any  $\mathbf{w}$ , positive  $\mu$ , and positive  $\eta$ , for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$  the following multi-view PAC-Bayes bound holds*

$$KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{\frac{1}{2} \|\eta \mathbf{w}_p - \mu \mathbf{w}\|^2 + \ln \frac{m-r+1}{\delta}}{m-r} \quad (32)$$

and  $\|\mathbf{w}\| = 1$ .

Another choice of  $\mathbf{w}_p$  is to learn a multi-view SVM classifier with the subset  $T$ , leading to the following bound.

**Theorem 10 (Multi-view PAC-Bayes bound 6)** *Consider a classifier prior given in (26) and a classifier posterior given in (27). Classifier  $\mathbf{w}_p$  has been learned from a subset  $T$  of  $r$  examples a priori separated from a training set  $S$  of  $m$  samples. For any data*

distribution  $\mathcal{D}$ , for any  $\mathbf{w}$ , positive  $\mu$ , and positive  $\eta$ , for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$  the following multi-view PAC-Bayes bound holds

$$KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{\frac{1}{2} \|\eta \mathbf{w}_p - \mu \mathbf{w}\|^2 + \ln \frac{m-r+1}{\delta}}{m-r} \quad (33)$$

and  $\|\mathbf{w}\| = 1$ .

Although the above two bounds look similar, they are essentially different in that the priors are determined differently. We will see in the experimental results that they also perform differently when applied in our experiments.

## 6. Semi-supervised Multi-view PAC-Bayes Bounds

Now we consider PAC-Bayes analysis for semi-supervised multi-view learning, where besides the  $m$  labeled examples we are further provided with  $u$  unlabeled examples  $U = \{\tilde{\mathbf{x}}_j\}_{j=m+1}^{m+u}$ . We replace  $V(\mathbf{u}_1, \mathbf{u}_2)$  with  $\hat{V}(\mathbf{u}_1, \mathbf{u}_2)$ , which has the form

$$\hat{V}(\mathbf{u}_1, \mathbf{u}_2) = \exp \left\{ -\frac{1}{2\sigma^2} \mathbf{u}^\top \mathbb{E}_U(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top) \mathbf{u} \right\}, \quad (34)$$

where  $\mathbb{E}_U$  means the empirical average over the unlabeled set  $U$ .

### 6.1 Noninformative Prior Center

Under a similar setting with Section 3, that is,  $P(\mathbf{u}) \propto \mathcal{N}(\mathbf{0}, \mathbf{I}) \times \hat{V}(\mathbf{u}_1, \mathbf{u}_2)$ , we have  $P(\mathbf{u}) = \mathcal{N}(\mathbf{0}, \Sigma)$  with  $\Sigma = \left( \mathbf{I} + \frac{\mathbb{E}_U(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right)^{-1}$ . Therefore, according to Theorem 3, we have

$$KL(Q(\mathbf{u}) \| P(\mathbf{u})) = \frac{1}{2} \left( -\ln \left( \left| \mathbf{I} + \frac{\mathbb{E}_U(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \right) + \frac{1}{\sigma^2} \mathbb{E}_U[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \right). \quad (35)$$

Substituting (35) into Theorem 1, we reach the following semi-supervised multi-view PAC-Bayes bound.

**Theorem 11 (Semi-supervised multi-view PAC-Bayes bound 1)** *Consider a classifier prior given in (8) with  $\hat{V}$  defined in (34), a classifier posterior given in (9) and an unlabeled set  $U = \{\tilde{\mathbf{x}}_j\}_{j=m+1}^{m+u}$ . For any data distribution  $\mathcal{D}$ , for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$ , the following inequality holds*

$$\forall \mathbf{w}, \mu : \frac{KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{\frac{1}{2} \left( -\ln \left( \left| \mathbf{I} + \frac{\mathbb{E}_U(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \right) + \frac{1}{\sigma^2} \mathbb{E}_U[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \right) + \ln \left( \frac{m+1}{\delta} \right)}{m},$$

where  $\|\mathbf{w}\| = 1$ .

## 6.2 Informative Prior Center

Similar to Section 4, we take the classifier prior to be

$$P(\mathbf{u}) \propto \mathcal{N}(\eta\mathbf{w}_p, \mathbf{I}) \times \hat{V}(\mathbf{u}_1, \mathbf{u}_2), \quad (36)$$

where  $\hat{V}(\mathbf{u}_1, \mathbf{u}_2)$  is given by (34),  $\eta > 0$  and  $\mathbf{w}_p = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[y\mathbf{x}]$  with  $\mathbf{x} = [\mathbf{x}_1^\top, \mathbf{x}_2^\top]^\top$ . We have  $P(\mathbf{u}) = \mathcal{N}(\mathbf{u}_p, \Sigma)$  with  $\Sigma = \left(\mathbf{I} + \frac{\mathbb{E}_U(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right)^{-1}$  and  $\mathbf{u}_p = \eta\Sigma\mathbf{w}_p$ .

By similar reasoning, we get

$$\begin{aligned} KL(Q(\mathbf{u})\|P(\mathbf{u})) &\leq -\frac{1}{2} \ln\left|\mathbf{I} + \frac{\mathbb{E}_U(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right| + \frac{1}{2}(\|\eta\mathbf{w}_p - \eta\hat{\mathbf{w}}_p\| + \|\eta\hat{\mathbf{w}}_p - \mu\mathbf{w}\| + \mu)^2 + \\ &\quad \frac{1}{2\sigma^2} \mathbb{E}_U \left[ \tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2(\mathbf{w}^\top \tilde{\mathbf{x}})^2 \right] - \eta\mu \mathbb{E} \left[ y(\mathbf{w}^\top \mathbf{x}) \right] + \frac{\mu^2}{2}, \end{aligned} \quad (37)$$

which is analogous to (24).

Then, we can give the following semi-supervised multi-view PAC-Bayes bound, whose proof is provided in Appendix E.

**Theorem 12 (Semi-supervised multi-view PAC-Bayes bound 2)** *Consider a classifier prior given in (36) with  $\hat{V}$  defined in (34), a classifier posterior given in (19) and an unlabeled set  $U = \{\tilde{\mathbf{x}}_j\}_{j=m+1}^{m+u}$ . For any data distribution  $\mathcal{D}$ , for any  $\mathbf{w}$ , positive  $\mu$ , and positive  $\eta$ , for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$ , the following inequality holds*

$$\begin{aligned} KL_+(\hat{E}_{Q,S}\|E_{Q,\mathcal{D}}) &\leq \frac{\frac{1}{2} \left( \frac{\eta R}{\sqrt{m}} \left( 2 + \sqrt{2 \ln \frac{3}{\delta}} \right) + \|\eta\hat{\mathbf{w}}_p - \mu\mathbf{w}\| + \mu \right)^2}{m} + \\ &\quad \frac{\frac{1}{2} \left( -\ln\left|\mathbf{I} + \frac{\mathbb{E}_U(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2}\right| + \frac{1}{\sigma^2} \mathbb{E}_U[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2(\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \right) + \bar{S}_m + \eta\mu R \sqrt{\frac{2}{m} \ln \frac{3}{\delta}} + \ln\left(\frac{m+1}{\delta/3}\right)}{m}, \end{aligned}$$

where

$$\bar{S}_m = \frac{1}{m} \sum_{i=1}^m [-\eta\mu y_i(\mathbf{w}^\top \mathbf{x}_i)],$$

and  $\|\mathbf{w}\| = 1$ .

## 7. Learning Algorithms

Below we provide the optimization formulations for the single-view and multi-view SVMs as well as semi-supervised multi-view SVMs that are adopted to train classifiers and calculate PAC-Bayes bounds. Note that the augmented vector representation is used by appending a scalar 1 at the end of the feature representations, in order to formulate the classifier in a simple form without the explicit bias term.

## 7.1 SVMs

The optimization problem (Cristianini and Shawe-Taylor, 2000; Shawe-Taylor and Sun, 2011) is formulated as

$$\begin{aligned} \min_{\mathbf{w}, \boldsymbol{\xi}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^\top \mathbf{x}_i) \geq 1 - \xi_i, \quad i = 1, \dots, n, \\ & \xi_i \geq 0, \quad i = 1, \dots, n, \end{aligned} \quad (38)$$

where scalar  $C$  controls the balance between the margin and empirical loss. This problem is a differentiable convex problem with affine constraints. The constraint qualification is satisfied by the refined Slater's condition.

The Lagrangian of problem (38) is

$$\begin{aligned} L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\gamma}) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \lambda_i \left[ y_i(\mathbf{w}^\top \mathbf{x}_i) - 1 + \xi_i \right] \\ & - \sum_{i=1}^n \gamma_i \xi_i, \quad \lambda_i \geq 0, \quad \gamma_i \geq 0, \end{aligned} \quad (39)$$

where  $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_n]^\top$  and  $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_n]^\top$  are the associated Lagrange multipliers. From the optimality conditions, we obtain

$$\partial_{\mathbf{w}} L(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*, \boldsymbol{\lambda}^*, \boldsymbol{\gamma}^*) = \mathbf{w}^* - \sum_{i=1}^n \lambda_i^* y_i \mathbf{x}_i = 0, \quad (40)$$

$$\partial_{\xi_i} L(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*, \boldsymbol{\lambda}^*, \boldsymbol{\gamma}^*) = C - \lambda_i^* - \gamma_i^* = 0, \quad i = 1, \dots, n. \quad (41)$$

The dual optimization problem is derived as

$$\begin{aligned} \min_{\boldsymbol{\lambda}} \quad & \frac{1}{2} \boldsymbol{\lambda}^\top D \boldsymbol{\lambda} - \boldsymbol{\lambda}^\top \mathbf{1} \\ \text{s.t.} \quad & \boldsymbol{\lambda} \succeq 0, \\ & \boldsymbol{\lambda} \preceq C \mathbf{1}, \end{aligned} \quad (42)$$

where  $D$  is a symmetric  $n \times n$  matrix with entries  $D_{ij} = y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$ . Once the solution  $\boldsymbol{\lambda}^*$  is given, the SVM decision function is given by

$$c^*(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^n y_i \lambda_i^* \mathbf{x}^\top \mathbf{x}_i \right). \quad (43)$$

Using the kernel trick, the optimization problem for SVMs is still (42). However, now  $D_{ij} = y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j)$  with the kernel function  $\kappa(\cdot, \cdot)$ , and the solution for the SVM classifier is formulated as

$$c^*(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^n y_i \lambda_i^* \kappa(\mathbf{x}_i, \mathbf{x}) \right). \quad (44)$$

## 7.2 MvSVMs

Denote the classifier weights from two views by  $\mathbf{w}_1$  and  $\mathbf{w}_2$  which are not assumed to be unit vectors at the moment. Inspired by semi-supervised multi-view SVMs (Sindhwani et al., 2005; Sindhwani and Rosenberg, 2008; Sun and Shawe-Taylor, 2010), the objective function of the multi-view SVMs (MvSVMs) can be given by

$$\begin{aligned} \min_{\mathbf{w}_1, \mathbf{w}_2, \xi_1, \xi_2} \quad & \frac{1}{2}(\|\mathbf{w}_1\|^2 + \|\mathbf{w}_2\|^2) + C_1 \sum_{i=1}^n (\xi_1^i + \xi_2^i) + C_2 \sum_{i=1}^n (\mathbf{w}_1^\top \mathbf{x}_1^i - \mathbf{w}_2^\top \mathbf{x}_2^i)^2 \\ \text{s.t.} \quad & y_i \mathbf{w}_1^\top \mathbf{x}_1^i \geq 1 - \xi_1^i, \quad i = 1, \dots, n, \\ & y_i \mathbf{w}_2^\top \mathbf{x}_2^i \geq 1 - \xi_2^i, \quad i = 1, \dots, n, \\ & \xi_1^i, \xi_2^i \geq 0, \quad i = 1, \dots, n. \end{aligned} \quad (45)$$

If kernel functions are used, the solution of the above optimization problem can be given by  $\mathbf{w}_1 = \sum_{i=1}^n \alpha_1^i k_1(\mathbf{x}_1^i, \cdot)$ , and  $\mathbf{w}_2 = \sum_{i=1}^n \alpha_2^i k_2(\mathbf{x}_2^i, \cdot)$ . Since a function defined on view  $j$  only depends on the  $j$ th feature set, the solution is given by

$$\mathbf{w}_1 = \sum_{i=1}^n \alpha_1^i k_1(\mathbf{x}_i, \cdot), \quad \mathbf{w}_2 = \sum_{i=1}^n \alpha_2^i k_2(\mathbf{x}_i, \cdot). \quad (46)$$

It can be shown that

$$\begin{aligned} \|\mathbf{w}_1\|^2 &= \boldsymbol{\alpha}_1^\top K_1 \boldsymbol{\alpha}_1, \quad \|\mathbf{w}_2\|^2 = \boldsymbol{\alpha}_2^\top K_2 \boldsymbol{\alpha}_2, \\ \sum_{i=1}^n (\mathbf{w}_1^\top \mathbf{x}_i - \mathbf{w}_2^\top \mathbf{x}_i)^2 &= (K_1 \boldsymbol{\alpha}_1 - K_2 \boldsymbol{\alpha}_2)^\top (K_1 \boldsymbol{\alpha}_1 - K_2 \boldsymbol{\alpha}_2), \end{aligned}$$

where  $K_1$  and  $K_2$  are kernel matrices from two views.

The optimization problem (45) can be reformulated as the following

$$\begin{aligned} \min_{\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \xi_1, \xi_2} \quad & F_0 = \frac{1}{2}(\boldsymbol{\alpha}_1^\top K_1 \boldsymbol{\alpha}_1 + \boldsymbol{\alpha}_2^\top K_2 \boldsymbol{\alpha}_2) + C_2 (K_1 \boldsymbol{\alpha}_1 - K_2 \boldsymbol{\alpha}_2)^\top (K_1 \boldsymbol{\alpha}_1 - K_2 \boldsymbol{\alpha}_2) + \\ & C_1 \sum_{i=1}^n (\xi_1^i + \xi_2^i) \\ \text{s.t.} \quad & y_i \left( \sum_{j=1}^n \alpha_1^j k_1(\mathbf{x}_j, \mathbf{x}_i) \right) \geq 1 - \xi_1^i, \quad i = 1, \dots, n, \\ & y_i \left( \sum_{j=1}^n \alpha_2^j k_2(\mathbf{x}_j, \mathbf{x}_i) \right) \geq 1 - \xi_2^i, \quad i = 1, \dots, n, \\ & \xi_1^i, \xi_2^i \geq 0, \quad i = 1, \dots, n. \end{aligned} \quad (47)$$

The derivation of the dual optimization formulation is detailed in Appendix F. Table 1 summarizes the MvSVM algorithm.

Table 1: The MvSVM Algorithm

**Input:**

A training set with  $n$  examples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  (each example has two views).  
 Kernel function  $k_1(\cdot, \cdot)$  and  $k_2(\cdot, \cdot)$  for two views, respectively.  
 Regularization coefficients  $C_1, C_2$ .

**Algorithm:**

- 1 Calculate Gram matrices  $K_1$  and  $K_2$  from two views.
- 2 Calculate  $A, B, D$  according to (90).
- 3 Solve the quadratic optimization problem (92) to get  $\lambda_1, \lambda_2$ .
- 4 Calculate  $\alpha_1$  and  $\alpha_2$  using (86) and (87).

**Output:** Classifier parameters  $\alpha_1$  and  $\alpha_2$  used by (46).

### 7.3 Semi-supervised MvSVMs (SMvSVMs)

Next we give the optimization formulation for semi-supervised MvSVMs (SMvSVMs) (Sindhwani et al., 2005; Sindhwani and Rosenberg, 2008; Sun and Shawe-Taylor, 2010), where besides the  $n$  labeled examples we further have  $u$  unlabeled examples.

Denote the classifier weights from two views by  $\mathbf{w}_1$  and  $\mathbf{w}_2$  which are not assumed to be unit vectors. The objective function of SMvSVMs is

$$\begin{aligned} \min_{\mathbf{w}_1, \mathbf{w}_2, \xi_1, \xi_2} \quad & \frac{1}{2}(\|\mathbf{w}_1\|^2 + \|\mathbf{w}_2\|^2) + C_1 \sum_{i=1}^n (\xi_1^i + \xi_2^i) + C_2 \sum_{i=1}^{n+u} (\mathbf{w}_1^\top \mathbf{x}_1^i - \mathbf{w}_2^\top \mathbf{x}_2^i)^2 \\ \text{s.t.} \quad & y_i \mathbf{w}_1^\top \mathbf{x}_1^i \geq 1 - \xi_1^i, \quad i = 1, \dots, n, \\ & y_i \mathbf{w}_2^\top \mathbf{x}_2^i \geq 1 - \xi_2^i, \quad i = 1, \dots, n, \\ & \xi_1^i, \xi_2^i \geq 0, \quad i = 1, \dots, n. \end{aligned} \quad (48)$$

If kernel functions are used, the solution can be expressed by  $\mathbf{w}_1 = \sum_{i=1}^{n+u} \alpha_1^i k_1(\mathbf{x}_1^i, \cdot)$ , and  $\mathbf{w}_2 = \sum_{i=1}^{n+u} \alpha_2^i k_2(\mathbf{x}_2^i, \cdot)$ . Since a function defined on view  $j$  only depends on the  $j$ th feature set, the solution is given by

$$\mathbf{w}_1 = \sum_{i=1}^{n+u} \alpha_1^i k_1(\mathbf{x}_i, \cdot), \quad \mathbf{w}_2 = \sum_{i=1}^{n+u} \alpha_2^i k_2(\mathbf{x}_i, \cdot). \quad (49)$$

It is straightforward to show that

$$\begin{aligned} \|\mathbf{w}_1\|^2 &= \alpha_1^\top K_1 \alpha_1, \quad \|\mathbf{w}_2\|^2 = \alpha_2^\top K_2 \alpha_2, \\ \sum_{i=1}^{n+u} (\mathbf{w}_1^\top \mathbf{x}_i - \mathbf{w}_2^\top \mathbf{x}_i)^2 &= (K_1 \alpha_1 - K_2 \alpha_2)^\top (K_1 \alpha_1 - K_2 \alpha_2), \end{aligned}$$

where  $(n+u) \times (n+u)$  matrices  $K_1$  and  $K_2$  are kernel matrices from two views.

The optimization problem (48) can be reformulated as

$$\min_{\alpha_1, \alpha_2, \xi_1, \xi_2} \tilde{F}_0 = \frac{1}{2}(\alpha_1^\top K_1 \alpha_1 + \alpha_2^\top K_2 \alpha_2) + C_2 (K_1 \alpha_1 - K_2 \alpha_2)^\top (K_1 \alpha_1 - K_2 \alpha_2) +$$

$$\begin{aligned}
& C_1 \sum_{i=1}^n (\xi_1^i + \xi_2^i) \\
\text{s.t. } & y_i \left( \sum_{j=1}^{n+u} \alpha_1^j k_1(\mathbf{x}_j, \mathbf{x}_i) \right) \geq 1 - \xi_1^i, \quad i = 1, \dots, n, \\
& y_i \left( \sum_{j=1}^{n+u} \alpha_2^j k_2(\mathbf{x}_j, \mathbf{x}_i) \right) \geq 1 - \xi_2^i, \quad i = 1, \dots, n, \\
& \xi_1^i, \xi_2^i \geq 0, \quad i = 1, \dots, n.
\end{aligned} \tag{50}$$

The derivation of the dual optimization formulation is detailed in Appendix G. Table 2 summarizes the SMvSVM algorithm.

Table 2: The SMvSVM Algorithm

---

**Input:**

A training set with  $n$  examples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  (each example has two views) and  $u$  unlabeled examples.

Kernel function  $k_1(\cdot, \cdot)$  and  $k_2(\cdot, \cdot)$  for two views, respectively.

Regularization coefficients  $C_1, C_2$ .

**Algorithm:**

- 1 Calculate Gram matrices  $K_1$  and  $K_2$  from two views.
- 2 Calculate  $A, B, D$  according to (104).
- 3 Solve the quadratic optimization problem (106) to get  $\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2$ .
- 4 Calculate  $\boldsymbol{\alpha}_1$  and  $\boldsymbol{\alpha}_2$  using (100) and (101).

**Output:** Classifier parameters  $\boldsymbol{\alpha}_1$  and  $\boldsymbol{\alpha}_2$  used by (49).

---

## 8. Experiments

The new bounds are evaluated on one synthetic and three real-world multi-view data sets where the learning task is binary classification. Below we first introduce the used data and the experimental settings. Then we report the test errors of the involved variants of the SVM algorithms, and evaluate the usefulness and relative performance of the new PAC-Bayes bounds.

### 8.1 Data Sets

The four multi-view data sets are introduced as follows.

#### SYNTHETIC

The synthetic data include 2000 examples half of which belong to the positive class. The dimensionality for each of the two views is 50. We first generate two random direction vectors one for each view, and then for each view sample 2000 points to make the inner products between the direction and the feature vector of half of the points be positive and the inner products for the other half of the points be negative. For the same point, the

corresponding inner products calculated from the two views are made identical. Finally, we add Gaussian white noise to the generated data to form the synthetic data set.

#### HANDWRITTEN

The handwritten digit data set is taken from the UCI machine learning repository (Bache and Lichman, 2013), which includes features of ten handwritten digits (0 ~ 9) extracted from a collection of Dutch utility maps. It consists of 2000 examples (200 examples per class) with the first view being the 76 Fourier coefficients, and the second view being the 64 Karhunen-Loève coefficients of each image. Binary classification between digits (1, 2, 3) and (4, 5, 6) is used for experiments.

#### ADS

The ads data are used for classifying web images into ads and non-ads (Kushmerick, 1999). This data set consists of 3279 examples with 459 of them being ads. 1554 binary attributes (weights of text terms related to an image using Boolean model) are used for classification, whose values can be 0 and 1. These attributes are divided into two views: one view describes the image itself (terms in the image’s caption, URL and alt text) and the other view contains features from other information (terms in the page and destination URLs). The two views have 587 and 967 features, respectively.

#### COURSE

The course data set consists of 1051 two-view web pages collected from computer science department web sites at four universities: Cornell University, University of Washington, University of Wisconsin, and University of Texas. There are 230 course pages and 821 non-course pages. The two views are words occurring in a web page and words appearing in the links pointing to that page (Blum and Mitchell, 1998; Sun and Shawe-Taylor, 2010). The document vectors are normalized to *tf-idf* (term frequency-inverse document frequency) features and then principal component analysis is used to perform dimensionality reduction. The dimensions of the two views are 500 and 87, respectively.

## 8.2 Experimental Settings

Our experiments include algorithm test error evaluation and PAC-Bayes bound evaluation for single-view learning, multi-view learning, supervised learning and semi-supervised learning. For single-view learning, SVMs are trained separately on each of the two views and the third view (concatenating the previous two views to form a long view), providing three supervised classifiers which are called SVM-1, SVM-2 and SVM-3, respectively. Evaluating the performance of the third view is interesting to compare single-view and multi-view learning methods, since single-view learning on the third view can exploit the same data as the usual multi-view learning algorithms. The MvSVMs and SMvSVMs are supervised multi-view learning and semi-supervised multi-view learning algorithms, respectively. The linear kernel is used for all the algorithms.

For each data set, four experimental settings are used. All the settings use 20% of all the examples as the unlabeled examples. For the remaining examples, the four settings use

20%, 40%, 60% and 80% of them as the labeled training set, respectively, and the rest forms the test set. Supervised algorithms will not use the unlabeled training data. For multi-view PAC-Bayes bound 5 and 6, we use 20% of the labeled training set to calculate the prior, and evaluate the bounds on the remaining 80% of training set. Each setting involves 10 random partitions of the above subsets. The reported performance is the average test error and standard deviation over these random partitions.

Model parameters, i.e.,  $C$  in SVMs, and  $C_1, C_2$  in MvSVMs and SMvSVMs, are selected by three-fold cross-validation on each labeled training set, where  $C_1, C_2$  are selected from  $\{10^{-6}, 10^{-4}, 10^{-2}, 1, 10, 100\}$  and  $C$  is selected from  $\{10^{-8}, 5 \times 10^{-8}, 10^{-7}, 5 \times 10^{-7}, 10^{-6}, 5 \times 10^{-6}, 10^{-5}, 5 \times 10^{-5}, 10^{-4}, 5 \times 10^{-4}, 10^{-3}, 5 \times 10^{-3}, 10^{-2}, 5 \times 10^{-2}, 10^{-1}, 5 \times 10^{-1}, 1, 5, 10, 20, 25, 30, 40, 50, 55, 60, 70, 80, 85, 90, 100, 300, 500, 700, 900, 1000\}$ . All the PAC-Bayes bounds are evaluated with a confidence of  $\delta = 0.05$ . We normalize  $\mathbf{w}$  in the posterior when we calculate the bounds. For multi-view PAC-Bayes bounds,  $\sigma$  is fixed to 100,  $\eta$  is set to 1, and  $R$  is equal to 1 which is clear from the augmented feature representation and data normalization preprocessing (all the training examples after feature augmentation are divided by a common value to make the maximum feature vector length be one).

We evaluate the following eleven PAC-Bayes bounds where the last eight bounds are presented in this paper.

- PB-1: The PAC-Bayes bound given by Theorem 2 and the SVM algorithm on the first view.
- PB-2: The PAC-Bayes bound given by Theorem 2 and the SVM algorithm on the second view.
- PB-3: The PAC-Bayes bound given by Theorem 2 and the SVM algorithm on the third view.
- MvPB-1: Multi-view PAC-Bayes bound 1 with the MvSVM algorithm.
- MvPB-2: Multi-view PAC-Bayes bound 2 with the MvSVM algorithm.
- MvPB-3: Multi-view PAC-Bayes bound 3 with the MvSVM algorithm.
- MvPB-4: Multi-view PAC-Bayes bound 4 with the MvSVM algorithm.
- MvPB-5: Multi-view PAC-Bayes bound 5 with the MvSVM algorithm.
- MvPB-6: Multi-view PAC-Bayes bound 6 with the MvSVM algorithm.
- SMvPB-1: Semi-supervised multi-view PAC-Bayes bound 1 with the SMvSVM algorithm.
- SMvPB-2: Semi-supervised multi-view PAC-Bayes bound 2 with the SMvSVM algorithm.

### 8.3 Test Errors

The prediction performances of SVMs, MvSVMs and SMvSVMs for the four experimental settings are reported in Table 3, Table 4, Table 5 and Table 6, respectively. For each data set, the best performance is indicated with boldface numbers. From all of these results, we see that MvSVMs and SMvSVMs have the best overall performance and sometimes single-view SVMs can have the best performances. SMvSVMs often perform better than MvSVMs since additional unlabeled examples are used, especially when the labeled training data set is small. Moreover, as expected, with more labeled training data the prediction performance of the algorithms will usually increase.

### 8.4 PAC-Bayes Bounds

Table 7, Table 8, Table 9 and Table 10 show the values of various PAC-Bayes bounds under different settings, where for each data set the best bound is indicated in bold and the best multi-view bound is indicated with underline.

From all the bound results, we find that the best single-view bound is usually tighter than the best multi-view bound, except on the synthetic data set. One possible explanation for this is that, the synthetic data set is ideal and in accordance with the assumptions for multi-view learning encoded in the prior, while the real world data sets are not. This also indicates that there is much space and possibility for further developments of multi-view PAC-Bayes analysis. In addition, with more labeled training data the corresponding bound will usually become tighter. Last but not least, among the eight presented multi-view PAC-Bayes bounds on real world data sets, the tightest one is often the first semi-supervised multi-view bound which exploits unlabeled data to calculate the function  $\hat{V}(\mathbf{u}_1, \mathbf{u}_2)$  and needs no further relaxation. The results also show that the second multi-view PAC-Bayes bound (dimensionality-independent bound with the prior distribution centered at the origin) is sometimes very good.

## 9. Conclusion

The paper lays the foundation of a theoretical and practical framework for defining priors that encode non-trivial interactions between data distributions and classifiers and translating them into sophisticated regularization schemes and associated generalization bounds. Specifically, we have presented eight new multi-view PAC-Bayes bounds, which integrate

Test Error	Synthetic	Handwritten	Ads	Course
SVM-1	17.20 ± 1.39	5.66 ± 0.94	5.84 ± 0.56	19.15 ± 1.54
SVM-2	19.98 ± 0.76	3.98 ± 0.68	5.25 ± 0.79	<b>10.15</b> ± 1.60
SVM-3	16.55 ± 2.04	<b>1.65</b> ± 0.53	4.62 ± 0.80	10.33 ± 1.34
MvSVM	10.54 ± 0.73	2.17 ± 0.64	<b>4.55</b> ± 0.66	10.55 ± 1.47
SMvSVM	<b>10.30</b> ± 0.79	2.04 ± 0.69	4.70 ± 0.70	10.28 ± 1.63

Table 3: Average error rates (%) and standard deviations for different learning algorithms under the 20% training setting.

Test Error	Synthetic	Handwritten	Ads	Course
SVM-1	14.49 ± 0.98	5.57 ± 0.41	5.04 ± 0.83	14.23 ± 1.27
SVM-2	16.88 ± 1.06	3.75 ± 0.99	4.14 ± 0.40	7.64 ± 0.80
SVM-3	10.31 ± 0.82	<b>1.51</b> ± 0.39	3.61 ± 0.54	7.68 ± 0.97
MvSVM	7.72 ± 0.78	1.98 ± 0.61	3.56 ± 0.54	7.00 ± 0.93
SMvSVM	<b>7.48</b> ± 0.66	2.03 ± 0.61	<b>3.44</b> ± 0.54	<b>6.81</b> ± 0.98

Table 4: Average error rates (%) and standard deviations for different learning algorithms under the 40% training setting.

Test Error	Synthetic	Handwritten	Ads	Course
SVM-1	14.23 ± 1.24	5.16 ± 0.61	4.32 ± 0.50	11.28 ± 1.30
SVM-2	16.11 ± 0.94	3.46 ± 0.94	3.90 ± 0.58	6.53 ± 1.44
SVM-3	9.08 ± 1.07	1.77 ± 0.85	3.43 ± 0.51	6.62 ± 1.33
MvSVM	<b>7.30</b> ± 0.85	<b>1.67</b> ± 0.63	3.45 ± 0.32	<b>5.82</b> ± 1.73
SMvSVM	7.31 ± 0.80	1.82 ± 0.70	<b>3.36</b> ± 0.38	5.93 ± 1.63

Table 5: Average error rates (%) and standard deviations for different learning algorithms under the 60% training setting.

Test Error	Synthetic	Handwritten	Ads	Course
SVM-1	13.06 ± 2.00	5.42 ± 1.51	4.47 ± 0.60	9.70 ± 1.64
SVM-2	16.03 ± 1.73	3.54 ± 1.33	3.59 ± 0.66	5.62 ± 1.68
SVM-3	8.06 ± 1.11	1.93 ± 0.66	<b>2.96</b> ± 0.51	5.56 ± 1.72
MvSVM	<b>6.28</b> ± 1.20	<b>1.82</b> ± 0.75	3.19 ± 0.63	4.20 ± 1.51
SMvSVM	<b>6.28</b> ± 1.19	1.93 ± 0.77	3.15 ± 0.75	<b>3.96</b> ± 1.59

Table 6: Average error rates (%) and standard deviations for different learning algorithms under the 80% training setting.

PAC-Bayes Bound	Synthetic	Handwritten	Ads	Course
PB-1	60.58 ± 0.12	54.61 ± 1.59	40.49 ± 2.09	<b>58.93</b> ± 8.90
PB-2	60.72 ± 0.09	<b>45.17</b> ± 3.74	<b>40.44</b> ± 2.12	61.64 ± 1.49
PB-3	<b>60.49</b> ± 0.12	47.62 ± 3.42	43.75 ± 3.15	59.67 ± 2.32
MvPB-1	61.27 ± 0.07	51.63 ± 2.89	40.87 ± 2.77	63.54 ± 0.45
MvPB-2	61.04 ± 0.07	51.45 ± 2.89	40.80 ± 2.77	63.26 ± 0.47
MvPB-3	62.35 ± 0.01	63.44 ± 0.62	56.38 ± 1.49	66.37 ± 0.06
MvPB-4	62.17 ± 0.01	63.23 ± 0.61	56.29 ± 1.48	66.14 ± 0.06
MvPB-5	61.84 ± 0.09	52.52 ± 3.01	43.21 ± 2.94	64.36 ± 0.43
MvPB-6	63.74 ± 0.08	58.65 ± 7.09	54.94 ± 4.68	67.75 ± 0.25
SMvPB-1	<u>60.60</u> ± 0.06	<u>49.84</u> ± 2.87	<u>40.65</u> ± 3.25	<u>62.77</u> ± 0.49
SMvPB-2	62.17 ± 0.01	62.94 ± 0.62	56.28 ± 1.30	66.14 ± 0.06

Table 7: Average PAC-Bayes bounds (%) and standard deviations for different learning algorithms under the 20% training setting.

PAC-Bayes Bound	Synthetic	Handwritten	Ads	Course
PB-1	57.20 ± 0.05	45.26 ± 1.48	33.11 ± 3.89	59.68 ± 0.52
PB-2	57.40 ± 0.11	<b>35.45</b> ± 3.22	<b>28.85</b> ± 3.26	<b>55.26</b> ± 1.97
PB-3	57.15 ± 0.07	35.48 ± 2.26	32.74 ± 4.29	56.12 ± 0.78
MvPB-1	57.69 ± 0.09	40.85 ± 3.23	33.36 ± 2.17	59.17 ± 0.51
MvPB-2	57.54 ± 0.08	<u>40.76</u> ± 3.22	<u>33.32</u> ± 2.17	58.99 ± 0.50
MvPB-3	58.97 ± 0.02	57.26 ± 1.17	51.68 ± 1.38	61.91 ± 0.07
MvPB-4	58.85 ± 0.02	57.15 ± 1.16	51.62 ± 1.37	61.77 ± 0.10
MvPB-5	57.44 ± 0.13	42.56 ± 3.36	35.86 ± 2.23	59.91 ± 0.48
MvPB-6	<b>52.67</b> ± 2.36	42.57 ± 5.93	47.34 ± 3.05	62.86 ± 0.09
SMvPB-1	57.27 ± 0.06	<u>40.76</u> ± 3.26	34.26 ± 3.00	<u>58.69</u> ± 0.44
SMvPB-2	58.85 ± 0.01	57.22 ± 1.18	52.16 ± 1.50	61.77 ± 0.09

Table 8: Average PAC-Bayes bounds (%) and standard deviations for different learning algorithms under the 40% training setting.

PAC-Bayes Bound	Synthetic	Handwritten	Ads	Course
PB-1	55.45 ± 0.08	42.07 ± 2.35	29.65 ± 1.93	57.52 ± 0.22
PB-2	55.71 ± 0.08	30.70 ± 2.05	<b>28.59</b> ± 3.71	<b>53.71</b> ± 2.27
PB-3	55.39 ± 0.16	<b>30.50</b> ± 3.31	30.49 ± 4.35	53.78 ± 1.01
MvPB-1	55.89 ± 0.08	34.16 ± 1.88	31.72 ± 4.13	56.90 ± 0.46
MvPB-2	55.78 ± 0.07	34.09 ± 1.88	<u>31.69</u> ± 4.13	56.75 ± 0.45
MvPB-3	57.38 ± 0.01	52.82 ± 1.08	49.77 ± 2.49	59.82 ± 0.07
MvPB-4	57.29 ± 0.01	52.73 ± 1.07	49.74 ± 2.48	59.69 ± 0.07
MvPB-5	55.60 ± 0.08	36.17 ± 1.88	34.11 ± 4.26	57.56 ± 0.42
MvPB-6	<b>39.20</b> ± 5.03	<u>31.76</u> ± 4.17	47.56 ± 3.81	60.67 ± 0.05
SMvPB-1	55.58 ± 0.06	33.93 ± 2.00	32.33 ± 3.37	<u>56.53</u> ± 0.43
SMvPB-2	57.28 ± 0.01	52.76 ± 1.15	50.51 ± 1.64	59.69 ± 0.07

Table 9: Average PAC-Bayes bounds (%) and standard deviations for different learning algorithms under the 60% training setting.

PAC-Bayes Bound	Synthetic	Handwritten	Ads	Course
PB-1	54.64 ± 0.77	37.52 ± 1.42	<b>28.97</b> ± 1.51	56.21 ± 0.18
PB-2	54.59 ± 0.04	28.47 ± 2.07	30.28 ± 1.83	<b>51.28</b> ± 2.97
PB-3	54.21 ± 0.08	<b>26.50</b> ± 2.15	29.74 ± 3.42	52.00 ± 0.85
MvPB-1	54.65 ± 0.05	30.25 ± 0.86	29.69 ± 0.84	55.77 ± 1.09
MvPB-2	54.63 ± 0.05	30.19 ± 0.86	<u>29.67</u> ± 0.84	55.38 ± 0.50
MvPB-3	56.41 ± 0.00	49.51 ± 0.52	48.12 ± 0.94	58.55 ± 0.07
MvPB-4	56.32 ± 0.01	49.43 ± 0.54	48.09 ± 0.92	58.44 ± 0.07
MvPB-5	54.36 ± 0.05	32.39 ± 0.88	31.44 ± 0.98	56.22 ± 0.41
MvPB-6	<b>26.89</b> ± 2.05	31.52 ± 3.33	46.31 ± 1.50	59.23 ± 0.18
SMvPB-1	54.41 ± 0.03	<u>30.15</u> ± 0.79	30.55 ± 2.28	<u>55.24</u> ± 0.43
SMvPB-2	56.32 ± 0.01	49.43 ± 0.46	48.77 ± 1.38	58.44 ± 0.06

Table 10: Average PAC-Bayes bounds (%) and standard deviations for different learning algorithms under the 80% training setting.

the view agreement as a key measure to modulate the prior distributions of classifiers. As extensions of PAC-Bayes analysis to the multi-view learning scenario, the proposed theoretical results are promising to fill the gap between the developments in theory and practice of multi-view learning, and are also possible to serve as the underpinnings to explain the effectiveness of multi-view learning. We have validated the theoretical superiority of multi-view learning in the ideal case of synthetic data, though this is not so evident for real world data which may not well meet our assumptions on the priors for multi-view learning.

The usefulness of the proposed bounds has been shown. Although often the current bounds are not the tightest, they indeed open the possibility of applying PAC-Bayes analysis to multi-view learning. We think the set of bounds could be further tightened in the future by adopting other techniques. It is also possible to study algorithms whose co-regularization term pushes towards the minimization of the multi-view PAC-Bayes bounds. In addition, we may use the work in this paper to motivate PAC-Bayes analysis for other learning tasks such as multi-task learning and domain adaptation, since these tasks are closely related to the current multi-view learning.

## Acknowledgments

This work is supported by the National Natural Science Foundation of China under Project 61370175, the Scientific Research Foundation for the Returned Overseas Chinese Scholars, State Education Ministry, and Shanghai Knowledge Service Platform Project (No. ZF1213).

## Appendix A. Proof of Theorem 5

Define

$$f(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m) = \frac{1}{m} \sum_{i=1}^m \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right|^{1/d}. \quad (51)$$

Since the rank of matrix  $\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top / \sigma^2$  is 1 with the nonzero eigenvalue being  $\|\tilde{\mathbf{x}}_i\|^2 / \sigma^2$  and the determinant of a positive semi-definite matrix is equal to the product of its eigenvalues, it follows that

$$\begin{aligned} & \sup_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m, \tilde{\mathbf{x}}_i} |f(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m) - f(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_{i+1}, \dots, \tilde{\mathbf{x}}_m)| \\ &= \frac{1}{m} \left| \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right|^{1/d} - \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right|^{1/d} \right| \\ &\leq \frac{1}{m} (\sqrt[d]{(R/\sigma)^2 + 1} - 1). \end{aligned}$$

By McDiarmid's inequality (Shawe-Taylor and Cristianini, 2004), we have for all  $\epsilon > 0$ ,

$$P \left\{ \mathbb{E} \left[ \left| \mathbf{I} + \frac{\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top}{\sigma^2} \right|^{1/d} \right] \geq f(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m) - \epsilon \right\} \geq 1 - \exp \left( \frac{-2m\epsilon^2}{(\sqrt[d]{(R/\sigma)^2 + 1} - 1)^2} \right). \quad (52)$$

Setting the right hand side equal to  $1 - \frac{\delta}{3}$ , we have with probability at least  $1 - \frac{\delta}{3}$ ,

$$\mathbb{E} \left[ \left| \mathbf{I} + \frac{\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top}{\sigma^2} \right|^{1/d} \right] \geq f(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m) - (\sqrt[d]{(R/\sigma)^2 + 1} - 1) \sqrt{\frac{1}{2m} \ln \frac{3}{\delta}}, \quad (53)$$

and

$$-\ln \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \leq -d \ln \left[ f(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m) - (\sqrt[d]{(R/\sigma)^2 + 1} - 1) \sqrt{\frac{1}{2m} \ln \frac{3}{\delta}} \right]_+, \quad (54)$$

where to reach (54) we have used (12) and defined  $[\cdot]_+ = \max(\cdot, 0)$ .

Denote  $H_m = \frac{1}{m} \sum_{i=1}^m [\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2]$ . It is clear that

$$\mathbb{E}[H_m] = \mathbb{E} \left\{ \frac{1}{m} \sum_{i=1}^m [\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2] \right\} = \mathbb{E}[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2]. \quad (55)$$

Recall  $R = \sup_{\tilde{\mathbf{x}}} \|\tilde{\mathbf{x}}\|$ . By McDiarmid's inequality, we have for all  $\epsilon > 0$ ,

$$P \{ \mathbb{E}[H_m] \leq H_m + \epsilon \} \geq 1 - \exp \left( \frac{-2m\epsilon^2}{(1 + \mu^2)^2 R^4} \right). \quad (56)$$

Setting the right hand side equal to  $1 - \frac{\delta}{3}$ , we have with probability at least  $1 - \frac{\delta}{3}$ ,

$$\mathbb{E}[H_m] \leq H_m + (1 + \mu^2) R^2 \sqrt{\frac{1}{2m} \ln \frac{3}{\delta}}. \quad (57)$$

In addition, from Lemma 1, we have

$$Pr_{S \sim \mathcal{D}^m} \left( \forall Q(c) : KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{KL(Q \| P) + \ln \left( \frac{m+1}{\delta/3} \right)}{m} \right) \geq 1 - \delta/3. \quad (58)$$

According to the union bound ( $Pr(A \text{ or } B \text{ or } C) \leq Pr(A) + Pr(B) + Pr(C)$ ), the probability that at least one of the inequalities in (54), (57) and (58) fails is no larger than

$\delta/3 + \delta/3 + \delta/3 = \delta$ . Hence, the probability that all of the three inequalities hold is no less than  $1 - \delta$ . That is, with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^m$ , the following inequality holds

$$\forall \mathbf{w}, \mu : \frac{KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq -\frac{d}{2} \ln \left[ f_m - \left( \sqrt{(R/\sigma)^2 + 1} - 1 \right) \sqrt{\frac{1}{2m} \ln \frac{3}{\delta}} \right]_+ + \frac{H_m}{2\sigma^2} + \frac{(1+\mu^2)R^2}{2\sigma^2} \sqrt{\frac{1}{2m} \ln \frac{3}{\delta}} + \frac{\mu^2}{2} + \ln \left( \frac{m+1}{\delta/3} \right)}{m},$$

where  $f_m$  is a shorthand for  $f(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m)$ , and  $\|\mathbf{w}\| = 1$ .

## Appendix B. Proof of Theorem 6

Now the KL divergence between the posterior and prior becomes

$$\begin{aligned} KL(Q(\mathbf{u}) \| P(\mathbf{u})) &= \frac{1}{2} \left( -\ln \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top) \right| \right) + \frac{1}{\sigma^2} \mathbb{E}[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \\ &\leq \frac{1}{2} \left( -\mathbb{E} \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right| \right) + \frac{1}{\sigma^2} \mathbb{E}[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2] + \mu^2 \\ &= \frac{1}{2} \left( \mathbb{E} \left( \frac{1}{\sigma^2} [\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2] - \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right| \right) + \mu^2 \right). \end{aligned}$$

Define

$$\tilde{f}(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m) = \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{\sigma^2} [\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2] - \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right| \right). \quad (59)$$

Recall  $R = \sup_{\tilde{\mathbf{x}}} \|\tilde{\mathbf{x}}\|$ . Since the rank of matrix  $\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top / \sigma^2$  is 1 with the nonzero eigenvalue being  $\|\tilde{\mathbf{x}}_i\|^2 / \sigma^2$  and the determinant of a positive semi-definite matrix is equal to the product of its eigenvalues, it follows that

$$\begin{aligned} &\sup_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m, \tilde{\mathbf{x}}_i} |\tilde{f}(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m) - \tilde{f}(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_{i+1}, \dots, \tilde{\mathbf{x}}_m)| \\ &\leq \frac{1}{m} \left( \frac{(1+\mu^2)R^2}{\sigma^2} + \ln \left( 1 + \frac{R^2}{\sigma^2} \right) \right). \end{aligned}$$

By McDiarmid's inequality, we have for all  $\epsilon > 0$ ,

$$P \left\{ \mathbb{E} \left( \frac{1}{\sigma^2} [\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2] - \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right| \right) \leq \tilde{f} + \epsilon \right\} \geq 1 - \exp \left( \frac{-2m\epsilon^2}{\Delta^2} \right), \quad (60)$$

where  $\tilde{f}$  is short for  $\tilde{f}(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m)$ , and  $\Delta = \frac{(1+\mu^2)R^2}{\sigma^2} + \ln \left( 1 + \frac{R^2}{\sigma^2} \right)$ . Setting the right hand size of (60) equal to  $1 - \frac{\delta}{2}$ , we have with probability at least  $1 - \frac{\delta}{2}$ ,

$$\mathbb{E} \left( \frac{1}{\sigma^2} [\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2] - \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2} \right| \right) \leq \tilde{f} + \Delta \sqrt{\frac{1}{2m} \ln \frac{2}{\delta}}. \quad (61)$$

Meanwhile, from Lemma 1, we have

$$Pr_{S \sim \mathcal{D}^m} \left( \forall Q(c) : KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{KL(Q \| P) + \ln \left( \frac{m+1}{\delta/2} \right)}{m} \right) \geq 1 - \delta/2. \quad (62)$$

According to the union bound, we can complete the proof for the dimensionality-independent PAC-Bayes bound.

### Appendix C. Proof of Theorem 7

It is clear that from  $R = \sup_{\tilde{\mathbf{x}}} \|\tilde{\mathbf{x}}\|$ , we have  $\sup_{\mathbf{x}} \|\mathbf{x}\| = R$  and  $\sup_{(\mathbf{x}, y)} \|y\mathbf{x}\| = R$ .

From (54), it follows that with probability at least  $1 - \frac{\delta}{4}$ ,

$$-\ln \left| \mathbf{I} + \frac{\mathbb{E}(\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top)}{\sigma^2} \right| \leq -d \ln \left[ f_m - \left( \sqrt[4]{(R/\sigma)^2 + 1} - 1 \right) \sqrt{\frac{1}{2m} \ln \frac{4}{\delta}} \right]_+. \quad (63)$$

With reference to a bounding result on estimating the center of mass (Shawe-Taylor and Cristianini, 2004), it follows that with probability at least  $1 - \delta/4$  the following inequality holds

$$\|\mathbf{w}_p - \hat{\mathbf{w}}_p\| \leq \frac{R}{\sqrt{m}} \left( 2 + \sqrt{2 \ln \frac{4}{\delta}} \right). \quad (64)$$

Denote  $\hat{H}_m = \frac{1}{m} \sum_{i=1}^m [\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i - 2\eta\mu\sigma^2 y_i (\mathbf{w}^\top \mathbf{x}_i) + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2]$ . It is clear that

$$\mathbb{E}[\hat{H}_m] = \mathbb{E}[\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} - 2\eta\mu\sigma^2 y (\mathbf{w}^\top \mathbf{x}) + \mu^2 (\mathbf{w}^\top \tilde{\mathbf{x}})^2]. \quad (65)$$

By McDiarmid's inequality, we have for all  $\epsilon > 0$ ,

$$P \left\{ \mathbb{E}[\hat{H}_m] \leq \hat{H}_m + \epsilon \right\} \geq 1 - \exp \left( \frac{-2m\epsilon^2}{(R^2 + 4\eta\mu\sigma^2 R + \mu^2 R^2)^2} \right). \quad (66)$$

Setting the right hand side equal to  $1 - \frac{\delta}{4}$ , we have with probability at least  $1 - \frac{\delta}{4}$ ,

$$\mathbb{E}[\hat{H}_m] \leq \hat{H}_m + (R^2 + \mu^2 R^2 + 4\eta\mu\sigma^2 R) \sqrt{\frac{1}{2m} \ln \frac{4}{\delta}}. \quad (67)$$

In addition, according to Lemma 1, we have

$$Pr_{S \sim \mathcal{D}^m} \left( \forall Q(c) : KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{KL(Q \| P) + \ln \left( \frac{m+1}{\delta/4} \right)}{m} \right) \geq 1 - \delta/4. \quad (68)$$

Therefore, from the union bound, we get the result.

## Appendix D. Proof of Theorem 8

Applying (13) to (24), we obtain

$$\begin{aligned}
KL(Q(\mathbf{u})\|P(\mathbf{u})) &\leq -\frac{1}{2}\mathbb{E}\ln\left|\mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2}\right| + \frac{1}{2}(\|\eta\mathbf{w}_p - \eta\hat{\mathbf{w}}_p\| + \|\eta\hat{\mathbf{w}}_p - \mu\mathbf{w}\| + \mu)^2 + \\
&\quad \frac{1}{2\sigma^2}\mathbb{E}\left[\tilde{\mathbf{x}}^\top\tilde{\mathbf{x}} - 2\eta\mu\sigma^2y(\mathbf{w}^\top\mathbf{x}) + \mu^2(\mathbf{w}^\top\tilde{\mathbf{x}})^2\right] + \frac{\mu^2}{2} \\
&= \frac{1}{2}(\|\eta\mathbf{w}_p - \eta\hat{\mathbf{w}}_p\| + \|\eta\hat{\mathbf{w}}_p - \mu\mathbf{w}\| + \mu)^2 + \\
&\quad \frac{1}{2}\mathbb{E}\left[\frac{\tilde{\mathbf{x}}^\top\tilde{\mathbf{x}} - 2\eta\mu\sigma^2y(\mathbf{w}^\top\mathbf{x}) + \mu^2(\mathbf{w}^\top\tilde{\mathbf{x}})^2}{\sigma^2} - \ln\left|\mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2}\right|\right] + \frac{\mu^2}{2}.
\end{aligned}$$

Following Shawe-Taylor and Cristianini (2004), we have with probability at least  $1 - \delta/3$

$$\|\mathbf{w}_p - \hat{\mathbf{w}}_p\| \leq \frac{R}{\sqrt{m}} \left(2 + \sqrt{2\ln\frac{3}{\delta}}\right). \quad (69)$$

Denote  $\tilde{H}_m = \frac{1}{m} \sum_{i=1}^m \left[ \frac{\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i - 2\eta\mu\sigma^2 y_i(\mathbf{w}^\top \mathbf{x}_i) + \mu^2(\mathbf{w}^\top \tilde{\mathbf{x}}_i)^2}{\sigma^2} - \ln \left| \mathbf{I} + \frac{\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top}{\sigma^2} \right| \right]$ . It is clear that

$$\mathbb{E}[\tilde{H}_m] = \mathbb{E}\left[\frac{\tilde{\mathbf{x}}^\top\tilde{\mathbf{x}} - 2\eta\mu\sigma^2y(\mathbf{w}^\top\mathbf{x}) + \mu^2(\mathbf{w}^\top\tilde{\mathbf{x}})^2}{\sigma^2} - \ln\left|\mathbf{I} + \frac{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top}{\sigma^2}\right|\right]. \quad (70)$$

By McDiarmid's inequality, we have for all  $\epsilon > 0$ ,

$$P\left\{\mathbb{E}[\tilde{H}_m] \leq \tilde{H}_m + \epsilon\right\} \geq 1 - \exp\left(\frac{-2m\epsilon^2}{\left(\frac{R^2 + 4\eta\mu\sigma^2 R + \mu^2 R^2}{\sigma^2} + \ln\left(1 + \frac{R^2}{\sigma^2}\right)\right)^2}\right). \quad (71)$$

Setting the right hand side equal to  $1 - \frac{\delta}{3}$ , we have with probability at least  $1 - \frac{\delta}{3}$ ,

$$\mathbb{E}[\tilde{H}_m] \leq \tilde{H}_m + \left(\frac{R^2 + 4\eta\mu\sigma^2 R + \mu^2 R^2}{\sigma^2} + \ln\left(1 + \frac{R^2}{\sigma^2}\right)\right) \sqrt{\frac{1}{2m} \ln\frac{3}{\delta}}. \quad (72)$$

In addition, from Lemma 1, we have

$$Pr_{S \sim \mathcal{D}^m} \left( \forall Q(c) : KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{KL(Q\|P) + \ln\left(\frac{m+1}{\delta/3}\right)}{m} \right) \geq 1 - \delta/3. \quad (73)$$

By applying the union bound, we complete the proof.

## Appendix E. Proof of Theorem 12

We already have  $\sup_{\mathbf{x}} \|\mathbf{x}\| = R$  and  $\sup_{(\mathbf{x}, y)} \|y\mathbf{x}\| = R$  from the definition  $R = \sup_{\tilde{\mathbf{x}}} \|\tilde{\mathbf{x}}\|$ .

Following Shawe-Taylor and Cristianini (2004), we have with probability at least  $1 - \delta/3$

$$\|\mathbf{w}_p - \hat{\mathbf{w}}_p\| \leq \frac{R}{\sqrt{m}} \left(2 + \sqrt{2\ln\frac{3}{\delta}}\right). \quad (74)$$

Denote  $\bar{S}_m = \frac{1}{m} \sum_{i=1}^m [-\eta\mu y_i(\mathbf{w}^\top \mathbf{x}_i)]$ . It is clear that

$$\mathbb{E}[\bar{S}_m] = -\eta\mu \mathbb{E} \left[ y(\mathbf{w}^\top \mathbf{x}) \right]. \quad (75)$$

By McDiarmid's inequality, we have for all  $\epsilon > 0$ ,

$$P \left\{ \mathbb{E}[\bar{S}_m] \leq \bar{S}_m + \epsilon \right\} \geq 1 - \exp \left( \frac{-2m\epsilon^2}{(2\eta\mu R)^2} \right). \quad (76)$$

Setting the right hand side equal to  $1 - \frac{\delta}{3}$ , we have with probability at least  $1 - \frac{\delta}{3}$ ,

$$\mathbb{E}[\bar{S}_m] \leq \bar{S}_m + \eta\mu R \sqrt{\frac{2}{m} \ln \frac{3}{\delta}}. \quad (77)$$

In addition, from Lemma 1, we have

$$Pr_{S \sim \mathcal{D}^m} \left( \forall Q(c) : KL_+(\hat{E}_{Q,S} \| E_{Q,\mathcal{D}}) \leq \frac{KL(Q \| P) + \ln \left( \frac{m+1}{\delta/3} \right)}{m} \right) \geq 1 - \delta/3. \quad (78)$$

After applying the union bound, the proof is completed.

## Appendix F. Dual Optimization Derivation for MvSVMs

To optimize (47), here we derive the Lagrange dual function.

Let  $\lambda_1^i, \lambda_2^i, \nu_1^i, \nu_2^i \geq 0$  be the Lagrange multipliers associated with the inequality constraints of problem (47). The Lagrangian  $L(\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2, \boldsymbol{\nu}_1, \boldsymbol{\nu}_2)$  can be written as

$$\begin{aligned} L = F_0 - \sum_{i=1}^n & \left[ \lambda_1^i \left( y_i \left( \sum_{j=1}^n \alpha_1^j k_1(x_j, x_i) \right) - 1 + \xi_1^i \right) + \right. \\ & \left. \lambda_2^i \left( y_i \left( \sum_{j=1}^n \alpha_2^j k_2(x_j, x_i) \right) - 1 + \xi_2^i \right) + \nu_1^i \xi_1^i + \nu_2^i \xi_2^i \right]. \end{aligned}$$

To obtain the Lagrangian dual function,  $L$  has to be minimized with respect to the primal variables  $\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \boldsymbol{\xi}_1, \boldsymbol{\xi}_2$ . To eliminate these variables, we compute the corresponding partial derivatives and set them to 0, obtaining the following conditions

$$(K_1 + 2C_2 K_1 K_1) \boldsymbol{\alpha}_1 - 2C_2 K_1 K_2 \boldsymbol{\alpha}_2 = \Lambda_1, \quad (79)$$

$$(K_2 + 2C_2 K_2 K_2) \boldsymbol{\alpha}_2 - 2C_2 K_2 K_1 \boldsymbol{\alpha}_1 = \Lambda_2, \quad (80)$$

$$\lambda_1^i + \nu_1^i = C_1, \quad (81)$$

$$\lambda_2^i + \nu_2^i = C_1, \quad (82)$$

where we have defined

$$\begin{aligned} \Lambda_1 & \triangleq \sum_{i=1}^n \lambda_1^i y_i K_1(:, i), \\ \Lambda_2 & \triangleq \sum_{i=1}^n \lambda_2^i y_i K_2(:, i), \end{aligned}$$

with  $K_1(:, i)$  and  $K_2(:, i)$  being the  $i$ th columns of the corresponding Gram matrices.

Substituting (79)~(82) into  $L$  results in the following expression of the Lagrangian dual function  $g(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2, \boldsymbol{\nu}_1, \boldsymbol{\nu}_2)$

$$\begin{aligned}
g &= \frac{1}{2}(\boldsymbol{\alpha}_1^\top K_1 \boldsymbol{\alpha}_1 + \boldsymbol{\alpha}_2^\top K_2 \boldsymbol{\alpha}_2) + C_2(\boldsymbol{\alpha}_1^\top K_1 K_1 \boldsymbol{\alpha}_1 - 2\boldsymbol{\alpha}_1^\top K_1 K_2 \boldsymbol{\alpha}_2 + \\
&\quad \boldsymbol{\alpha}_2^\top K_2 K_2 \boldsymbol{\alpha}_2) - \boldsymbol{\alpha}_1^\top \Lambda_1 - \boldsymbol{\alpha}_2^\top \Lambda_2 + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i) \\
&= \frac{1}{2}\boldsymbol{\alpha}_1^\top \Lambda_1 + \frac{1}{2}\boldsymbol{\alpha}_2^\top \Lambda_2 - \boldsymbol{\alpha}_1^\top \Lambda_1 - \boldsymbol{\alpha}_2^\top \Lambda_2 + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i) \\
&= -\frac{1}{2}\boldsymbol{\alpha}_1^\top \Lambda_1 - \frac{1}{2}\boldsymbol{\alpha}_2^\top \Lambda_2 + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i). \tag{83}
\end{aligned}$$

Define

$$\begin{aligned}
\tilde{K}_1 &= K_1 + 2C_2 K_1 K_1, & \bar{K}_1 &= 2C_2 K_1 K_2, \\
\tilde{K}_2 &= K_2 + 2C_2 K_2 K_2, & \bar{K}_2 &= 2C_2 K_2 K_1.
\end{aligned}$$

Then, (79) and (80) become

$$\tilde{K}_1 \boldsymbol{\alpha}_1 - \bar{K}_1 \boldsymbol{\alpha}_2 = \Lambda_1, \tag{84}$$

$$\tilde{K}_2 \boldsymbol{\alpha}_2 - \bar{K}_2 \boldsymbol{\alpha}_1 = \Lambda_2. \tag{85}$$

From (84) and (85), we have

$$\begin{aligned}
(\tilde{K}_1 - \bar{K}_1 \tilde{K}_2^{-1} \bar{K}_2) \boldsymbol{\alpha}_1 &= \bar{K}_1 \tilde{K}_2^{-1} \Lambda_2 + \Lambda_1 \\
(\tilde{K}_2 - \bar{K}_2 \tilde{K}_1^{-1} \bar{K}_1) \boldsymbol{\alpha}_2 &= \bar{K}_2 \tilde{K}_1^{-1} \Lambda_1 + \Lambda_2.
\end{aligned}$$

Define  $M_1 \triangleq \tilde{K}_1 - \bar{K}_1 \tilde{K}_2^{-1} \bar{K}_2$  and  $M_2 \triangleq \tilde{K}_2 - \bar{K}_2 \tilde{K}_1^{-1} \bar{K}_1$ . It follows that

$$\boldsymbol{\alpha}_1 = M_1^{-1} [\bar{K}_1 \tilde{K}_2^{-1} \Lambda_2 + \Lambda_1], \tag{86}$$

$$\boldsymbol{\alpha}_2 = M_2^{-1} [\bar{K}_2 \tilde{K}_1^{-1} \Lambda_1 + \Lambda_2]. \tag{87}$$

Now with  $\boldsymbol{\alpha}_1$  and  $\boldsymbol{\alpha}_2$  substituted into (83), the Lagrange dual function  $g(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2, \boldsymbol{\nu}_1, \boldsymbol{\nu}_2)$  is

$$\begin{aligned}
g &= \inf_{\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \boldsymbol{\xi}_1, \boldsymbol{\xi}_2} L = -\frac{1}{2}\boldsymbol{\alpha}_1^\top \Lambda_1 - \frac{1}{2}\boldsymbol{\alpha}_2^\top \Lambda_2 + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i) \\
&= -\frac{1}{2}\Lambda_1^\top M_1^{-1} [\bar{K}_1 \tilde{K}_2^{-1} \Lambda_2 + \Lambda_1] - \frac{1}{2}\Lambda_2^\top M_2^{-1} [\bar{K}_2 \tilde{K}_1^{-1} \Lambda_1 + \Lambda_2] + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i).
\end{aligned}$$

The Lagrange dual problem is given by

$$\begin{aligned}
&\max_{\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2} g \\
&\text{s.t.} \quad \begin{cases} 0 \leq \lambda_1^i \leq C_1, & i = 1, \dots, n \\ 0 \leq \lambda_2^i \leq C_1, & i = 1, \dots, n. \end{cases} \tag{88}
\end{aligned}$$

As Lagrange dual functions are concave, we can formulate the Lagrange dual problem as a convex optimization problem

$$\begin{aligned} & \min_{\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2} -g \\ & \text{s.t.} \quad \begin{cases} 0 \leq \lambda_1^i \leq C_1, & i = 1, \dots, n \\ 0 \leq \lambda_2^i \leq C_1, & i = 1, \dots, n. \end{cases} \end{aligned} \quad (89)$$

Define matrix  $Y \triangleq \text{diag}(y_1, \dots, y_n)$ . Then,  $\Lambda_1 = K_1 Y \boldsymbol{\lambda}_1$  and  $\Lambda_2 = K_2 Y \boldsymbol{\lambda}_2$  with  $\boldsymbol{\lambda}_1 = (\lambda_1^1, \dots, \lambda_1^n)^\top$ , and  $\boldsymbol{\lambda}_2 = (\lambda_2^1, \dots, \lambda_2^n)^\top$ . It is clear that  $\tilde{K}_1$  and  $\tilde{K}_2$  are symmetric matrices, and  $\tilde{K}_1 = \tilde{K}_2^\top$ . Therefore, it follows that matrices  $M_1$  and  $M_2$  are also symmetric.

We have

$$\begin{aligned} -g &= \frac{1}{2} \Lambda_1^\top M_1^{-1} [\tilde{K}_1 \tilde{K}_2^{-1} \Lambda_2 + \Lambda_1] + \frac{1}{2} \Lambda_2^\top M_2^{-1} [\tilde{K}_2 \tilde{K}_1^{-1} \Lambda_1 + \Lambda_2] - \sum_{i=1}^n (\lambda_1^i + \lambda_2^i) \\ &= \frac{1}{2} \left\{ \boldsymbol{\lambda}_1^\top [Y K_1 M_1^{-1} K_1 Y] \boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_1^\top [Y K_1 M_1^{-1} \tilde{K}_1 \tilde{K}_2^{-1} K_2 Y] \boldsymbol{\lambda}_2 + \right. \\ & \quad \left. \boldsymbol{\lambda}_2^\top [Y K_2 M_2^{-1} \tilde{K}_2 \tilde{K}_1^{-1} K_1 Y] \boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_2^\top [Y K_2 M_2^{-1} K_2 Y] \boldsymbol{\lambda}_2 \right\} - \mathbf{1}^\top (\boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_2) \\ &= \frac{1}{2} (\boldsymbol{\lambda}_1^\top \quad \boldsymbol{\lambda}_2^\top) \begin{pmatrix} A & B \\ B^\top & D \end{pmatrix} \begin{pmatrix} \boldsymbol{\lambda}_1 \\ \boldsymbol{\lambda}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\lambda}_1 \\ \boldsymbol{\lambda}_2 \end{pmatrix}^\top \mathbf{1}_{2n}, \end{aligned}$$

where

$$A \triangleq Y K_1 M_1^{-1} K_1 Y, \quad B \triangleq Y K_1 M_1^{-1} \tilde{K}_1 \tilde{K}_2^{-1} K_2 Y, \quad D \triangleq Y K_2 M_2^{-1} K_2 Y, \quad (90)$$

$\mathbf{1}_{2n} = (1, \dots, 1_{(2n)})^\top$ , and we have used the fact that

$$Y K_1 M_1^{-1} \tilde{K}_1 \tilde{K}_2^{-1} K_2 Y = [Y K_2 M_2^{-1} \tilde{K}_2 \tilde{K}_1^{-1} K_1 Y]^\top. \quad (91)$$

Because of the convexity of function  $-g$ , we affirm that matrix  $\begin{pmatrix} A & B \\ B^\top & D \end{pmatrix}$  is positive semidefinite.

Hence, the optimization problem in (89) can be rewritten as

$$\begin{aligned} & \min_{\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2} \frac{1}{2} (\boldsymbol{\lambda}_1^\top \quad \boldsymbol{\lambda}_2^\top) \begin{pmatrix} A & B \\ B^\top & D \end{pmatrix} \begin{pmatrix} \boldsymbol{\lambda}_1 \\ \boldsymbol{\lambda}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\lambda}_1 \\ \boldsymbol{\lambda}_2 \end{pmatrix}^\top \mathbf{1}_{2n} \\ & \text{s.t.} \quad \begin{cases} 0 \leq \boldsymbol{\lambda}_1 \leq C_1 \mathbf{1}, \\ 0 \leq \boldsymbol{\lambda}_2 \leq C_1 \mathbf{1}. \end{cases} \end{aligned} \quad (92)$$

After solving this problem, we can then obtain classifier parameters  $\boldsymbol{\alpha}_1$  and  $\boldsymbol{\alpha}_2$  using (86) and (87), which are finally used by (46).

## Appendix G. Dual Optimization Derivation for SMvSVMs

To optimize (50), we first derive the Lagrange dual function following the same line of optimization derivations for MvSVMs. Although here some of the derivations are similar to those for MvSVMs, for completeness we include them.

Let  $\lambda_1^i, \lambda_2^i, \nu_1^i, \nu_2^i \geq 0$  be the Lagrange multipliers associated with the inequality constraints of problem (50). The Lagrangian  $L(\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2, \boldsymbol{\nu}_1, \boldsymbol{\nu}_2)$  can be formulated as

$$L = \tilde{F}_0 - \sum_{i=1}^n \left[ \lambda_1^i \left( y_i \left( \sum_{j=1}^{n+u} \alpha_1^j k_1(x_j, x_i) \right) - 1 + \xi_1^i \right) + \lambda_2^i \left( y_i \left( \sum_{j=1}^{n+u} \alpha_2^j k_2(x_j, x_i) \right) - 1 + \xi_2^i \right) + \nu_1^i \xi_1^i + \nu_2^i \xi_2^i \right].$$

To obtain the Lagrangian dual function,  $L$  will be minimized with respect to the primal variables  $\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \boldsymbol{\xi}_1, \boldsymbol{\xi}_2$ . To eliminate these variables, setting the corresponding partial derivatives to 0 results in the following conditions

$$(K_1 + 2C_2 K_1 K_1) \boldsymbol{\alpha}_1 - 2C_2 K_1 K_2 \boldsymbol{\alpha}_2 = \Lambda_1, \quad (93)$$

$$(K_2 + 2C_2 K_2 K_2) \boldsymbol{\alpha}_2 - 2C_2 K_2 K_1 \boldsymbol{\alpha}_1 = \Lambda_2, \quad (94)$$

$$\lambda_1^i + \nu_1^i = C_1, \quad (95)$$

$$\lambda_2^i + \nu_2^i = C_1, \quad (96)$$

where we have defined

$$\Lambda_1 \triangleq \sum_{i=1}^n \lambda_1^i y_i K_1(:, i),$$

$$\Lambda_2 \triangleq \sum_{i=1}^n \lambda_2^i y_i K_2(:, i),$$

with  $K_1(:, i)$  and  $K_2(:, i)$  being the  $i$ th columns of the corresponding Gram matrices.

Substituting (93)~(96) into  $L$  results in the Lagrangian dual function  $g(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2, \boldsymbol{\nu}_1, \boldsymbol{\nu}_2)$

$$\begin{aligned} g &= \frac{1}{2} (\boldsymbol{\alpha}_1^\top K_1 \boldsymbol{\alpha}_1 + \boldsymbol{\alpha}_2^\top K_2 \boldsymbol{\alpha}_2) + C_2 (\boldsymbol{\alpha}_1^\top K_1 K_1 \boldsymbol{\alpha}_1 - 2 \boldsymbol{\alpha}_1^\top K_1 K_2 \boldsymbol{\alpha}_2 + \\ &\quad \boldsymbol{\alpha}_2^\top K_2 K_2 \boldsymbol{\alpha}_2) - \boldsymbol{\alpha}_1^\top \Lambda_1 - \boldsymbol{\alpha}_2^\top \Lambda_2 + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i) \\ &= \frac{1}{2} \boldsymbol{\alpha}_1^\top \Lambda_1 + \frac{1}{2} \boldsymbol{\alpha}_2^\top \Lambda_2 - \boldsymbol{\alpha}_1^\top \Lambda_1 - \boldsymbol{\alpha}_2^\top \Lambda_2 + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i) \\ &= -\frac{1}{2} \boldsymbol{\alpha}_1^\top \Lambda_1 - \frac{1}{2} \boldsymbol{\alpha}_2^\top \Lambda_2 + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i). \end{aligned} \quad (97)$$

Define

$$\begin{aligned} \tilde{K}_1 &= K_1 + 2C_2 K_1 K_1, & \bar{K}_1 &= 2C_2 K_1 K_2, \\ \tilde{K}_2 &= K_2 + 2C_2 K_2 K_2, & \bar{K}_2 &= 2C_2 K_2 K_1. \end{aligned}$$

Then, (93) and (94) become

$$\tilde{K}_1 \boldsymbol{\alpha}_1 - \bar{K}_1 \boldsymbol{\alpha}_2 = \Lambda_1, \quad (98)$$

$$\tilde{K}_2 \boldsymbol{\alpha}_2 - \bar{K}_2 \boldsymbol{\alpha}_1 = \Lambda_2. \quad (99)$$

From (98) and (99), we have

$$\begin{aligned}(\tilde{K}_1 - \bar{K}_1 \tilde{K}_2^{-1} \bar{K}_2) \boldsymbol{\alpha}_1 &= \bar{K}_1 \tilde{K}_2^{-1} \Lambda_2 + \Lambda_1 \\ (\tilde{K}_2 - \bar{K}_2 \tilde{K}_1^{-1} \bar{K}_1) \boldsymbol{\alpha}_2 &= \bar{K}_2 \tilde{K}_1^{-1} \Lambda_1 + \Lambda_2.\end{aligned}$$

Define  $M_1 \triangleq \tilde{K}_1 - \bar{K}_1 \tilde{K}_2^{-1} \bar{K}_2$  and  $M_2 \triangleq \tilde{K}_2 - \bar{K}_2 \tilde{K}_1^{-1} \bar{K}_1$ . It is clear that

$$\boldsymbol{\alpha}_1 = M_1^{-1} \left[ \bar{K}_1 \tilde{K}_2^{-1} \Lambda_2 + \Lambda_1 \right], \quad (100)$$

$$\boldsymbol{\alpha}_2 = M_2^{-1} \left[ \bar{K}_2 \tilde{K}_1^{-1} \Lambda_1 + \Lambda_2 \right]. \quad (101)$$

With  $\boldsymbol{\alpha}_1$  and  $\boldsymbol{\alpha}_2$  substituted into (97), the Lagrange dual function  $g(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2, \boldsymbol{\nu}_1, \boldsymbol{\nu}_2)$  is then

$$\begin{aligned}g &= \inf_{\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \boldsymbol{\xi}_1, \boldsymbol{\xi}_2} L = -\frac{1}{2} \boldsymbol{\alpha}_1^\top \Lambda_1 - \frac{1}{2} \boldsymbol{\alpha}_2^\top \Lambda_2 + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i) \\ &= -\frac{1}{2} \Lambda_1^\top M_1^{-1} \left[ \bar{K}_1 \tilde{K}_2^{-1} \Lambda_2 + \Lambda_1 \right] - \frac{1}{2} \Lambda_2^\top M_2^{-1} \left[ \bar{K}_2 \tilde{K}_1^{-1} \Lambda_1 + \Lambda_2 \right] + \sum_{i=1}^n (\lambda_1^i + \lambda_2^i).\end{aligned}$$

The Lagrange dual problem is given by

$$\begin{aligned}\max_{\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2} \quad & g \\ \text{s.t.} \quad & \begin{cases} 0 \leq \lambda_1^i \leq C_1, & i = 1, \dots, n \\ 0 \leq \lambda_2^i \leq C_1, & i = 1, \dots, n. \end{cases}\end{aligned} \quad (102)$$

As Lagrange dual functions are concave, below we formulate the Lagrange dual problem as a convex optimization problem

$$\begin{aligned}\min_{\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2} \quad & -g \\ \text{s.t.} \quad & \begin{cases} 0 \leq \lambda_1^i \leq C_1, & i = 1, \dots, n \\ 0 \leq \lambda_2^i \leq C_1, & i = 1, \dots, n. \end{cases}\end{aligned} \quad (103)$$

Define matrix  $Y \triangleq \text{diag}(y_1, \dots, y_n)$ . Then,  $\Lambda_1 = K_{n1} Y \boldsymbol{\lambda}_1$  and  $\Lambda_2 = K_{n2} Y \boldsymbol{\lambda}_2$  with  $K_{n1} = K_1(:, 1:n)$ ,  $K_{n2} = K_2(:, 1:n)$ ,  $\boldsymbol{\lambda}_1 = (\lambda_1^1, \dots, \lambda_1^n)^\top$ , and  $\boldsymbol{\lambda}_2 = (\lambda_2^1, \dots, \lambda_2^n)^\top$ . It is clear that  $\tilde{K}_1$  and  $\tilde{K}_2$  are symmetric matrices, and  $\bar{K}_1 = \tilde{K}_2^\top$ . Therefore, it follows that matrices  $M_1$  and  $M_2$  are also symmetric.

We have

$$\begin{aligned}-g &= \frac{1}{2} \Lambda_1^\top M_1^{-1} \left[ \bar{K}_1 \tilde{K}_2^{-1} \Lambda_2 + \Lambda_1 \right] + \frac{1}{2} \Lambda_2^\top M_2^{-1} \left[ \bar{K}_2 \tilde{K}_1^{-1} \Lambda_1 + \Lambda_2 \right] - \sum_{i=1}^n (\lambda_1^i + \lambda_2^i) \\ &= \frac{1}{2} \left\{ \boldsymbol{\lambda}_1^\top [Y K_{n1}^\top M_1^{-1} K_{n1} Y] \boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_1^\top [Y K_{n1}^\top M_1^{-1} \bar{K}_1 \tilde{K}_2^{-1} K_{n2} Y] \boldsymbol{\lambda}_2 + \right. \\ &\quad \left. \boldsymbol{\lambda}_2^\top [Y K_{n2}^\top M_2^{-1} \bar{K}_2 \tilde{K}_1^{-1} K_{n1} Y] \boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_2^\top [Y K_{n2}^\top M_2^{-1} K_{n2} Y] \boldsymbol{\lambda}_2 \right\} - \mathbf{1}^\top (\boldsymbol{\lambda}_1 + \boldsymbol{\lambda}_2) \\ &= \frac{1}{2} (\boldsymbol{\lambda}_1^\top \quad \boldsymbol{\lambda}_2^\top) \begin{pmatrix} A & B \\ B^\top & D \end{pmatrix} \begin{pmatrix} \boldsymbol{\lambda}_1 \\ \boldsymbol{\lambda}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\lambda}_1 \\ \boldsymbol{\lambda}_2 \end{pmatrix}^\top \mathbf{1}_{2n},\end{aligned}$$

where

$$A \triangleq YK_{n1}^\top M_1^{-1} K_{n1} Y, \quad B \triangleq YK_{n1}^\top M_1^{-1} \bar{K}_1 \tilde{K}_2^{-1} K_{n2} Y, \quad D \triangleq YK_{n2}^\top M_2^{-1} K_{n2} Y, \quad (104)$$

$\mathbf{1}_{2n} = (1, \dots, 1_{(2n)})^\top$ , and we have used the fact that

$$YK_{n1}^\top M_1^{-1} \bar{K}_1 \tilde{K}_2^{-1} K_{n2} Y = [YK_{n2}^\top M_2^{-1} \bar{K}_2 \tilde{K}_1^{-1} K_{n1} Y]^\top. \quad (105)$$

Because of the convexity of function  $-g$ , we affirm that matrix  $\begin{pmatrix} A & B \\ B^\top & D \end{pmatrix}$  is positive semidefinite.

Hence, the optimization problem in (103) can be rewritten as

$$\begin{aligned} \min_{\lambda_1, \lambda_2} & \frac{1}{2} (\lambda_1^\top \ \lambda_2^\top) \begin{pmatrix} A & B \\ B^\top & D \end{pmatrix} \begin{pmatrix} \lambda_1 \\ \lambda_2 \end{pmatrix} - \begin{pmatrix} \lambda_1 \\ \lambda_2 \end{pmatrix}^\top \mathbf{1}_{2n} \\ \text{s.t.} & \begin{cases} 0 \preceq \lambda_1 \preceq C_1 \mathbf{1}, \\ 0 \preceq \lambda_2 \preceq C_1 \mathbf{1}. \end{cases} \end{aligned} \quad (106)$$

After solving this problem, we can then obtain classifier parameters  $\alpha_1$  and  $\alpha_2$  using (100) and (101), which are finally used by (49).

## References

- A. Ambroladze, E. Parrado-hernández, and J. Shawe-taylor. Tighter pac-bayes bounds. *Advances in Neural Information Processing Systems*, 19:9–16, 2007.
- K. Bache and M. Lichman. *UCI Machine Learning Repository*. University of California, Irvine, School of Information and Computer Sciences, 2013. URL <http://archive.ics.uci.edu/ml>.
- P. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- A. Blum and T. Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the 11th Annual Conference on Computational Learning Theory*, pages 92–100, 1998.
- O. Catoni. *PAC-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning*. Institute of Mathematical Statistics, Beachwood, Ohio, USA, 2007.
- N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines*. Cambridge University Press, Cambridge, UK, 2000.
- J. Farquhar, D. Hardoon, H. Meng, J. Shawe-Taylor, and S. Szedmak. Two view learning: SVM-2K, theory and practice. *Advances in Neural Information Processing Systems*, 18: 355–362, 2006.
- P. Germain, A. Lacasse, F. Laviolette, and M. Marchand. PAC-Bayesian learning of linear classifiers. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 353–360, 2009.

- M. Higgs and J. Shawe-Taylor. A PAC-Bayes bound for tailored density estimation. *Lecture Notes in Computer Science*, 6331:148–162, 2010.
- S. Kakade and D. Foster. Multi-view regression via canonical correlation analysis. In *Proceedings of the 20th Annual Conference on Learning Theory*, pages 82–96, 2007.
- N. Kushmerick. Learning to remove internet advertisements. In *Proceedings of the 3rd International Conference on Autonomous Agents*, pages 175–181, 1999.
- J. Langford. Tutorial on practical prediction theory for classification. *Journal of Machine Learning Research*, 6(Mar):273–306, 2005.
- J. Langford and J. Shawe-Taylor. PAC-Bayes & margins. *Advances in Neural Information Processing Systems*, 15:423–430, 2002.
- G. Lever, F. Laviolette, and J. Shawe-Taylor. Tighter PAC-Bayes bounds through distribution-dependent priors. *Theoretical Computer Science*, 473(Feb):4–28, 2013.
- D. McAllester. PAC-Bayesian model averaging. In *Proceedings of 12th Annual Conference on Computational Learning Theory*, pages 164–170, 1999.
- E. Parrado-Hernández, A. Ambroladze, J. Shawe-Taylor, and S. Sun. PAC-Bayes bounds with data dependent priors. *Journal of Machine Learning Research*, 13(Dec):3507–3531, 2012.
- C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, 2006.
- D. Rosenberg and P. Bartlett. The rademacher complexity of co-regularized kernel classes. In *Proceedings of the 11th International Conference on Artificial Intelligence and Statistics*, pages 396–403, 2007.
- M. Seeger. PAC-Bayesian generalisation error bounds for Gaussian process classification. *Journal of Machine Learning Research*, 3(Oct):233–269, 2002.
- Y. Seldin and N. Tishby. PAC-Bayesian analysis of co-clustering and beyond. *Journal of Machine Learning Research*, 11(Dec):3595–3646, 2010.
- Y. Seldin, F. Laviolette, N. Cesa-Bianchi, J. Shawe-Taylor, and P. Auer. PAC-Bayesian inequalities for martingales. *IEEE Transactions on Information Theory*, 58(12):7086–7093, 2012.
- J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, Cambridge, UK, 2004.
- J. Shawe-Taylor and S. Sun. A review of optimization methodologies in support vector machines. *Neurocomputing*, 74(17):3609–3618, 2011.
- J. Shawe-Taylor, P. Bartlett, R. Williamson, and M. Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 44(5):1926–1940, 1998.

- V. Sindhwani and D. Rosenberg. An rkhs for multi-view learning and manifold co-regularization. In *Proceedings of the 25th Annual International Conference on Machine Learning*, pages 976–983, 2008.
- V. Sindhwani, P. Niyogi, and M. Belkin. A co-regularization approach to semi-supervised learning with multiple views. In *Proceedings of ICML Workshop on Learning with Multiple Views*, pages 74–79, 2005.
- K. Sridharan and S. Kakade. An information theoretic framework for multi-view learning. In *Proceedings of the 21st Annual Conference on Learning Theory*, pages 403–414, 2008.
- S. Sun. A survey of multi-view machine learning. *Neural Computing and Applications*, 23(7-8):2031–2038, 2013.
- S. Sun and J. Shawe-Taylor. Sparse semi-supervised learning using conjugate functions. *Journal of Machine Learning Research*, 11(Sep):2423–2455, 2010.
- V. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.