

# Synthesis of Maximum Margin and Multiview Learning using Unlabeled Data

Sandor Szedmak<sup>1</sup> and John Shawe-Taylor<sup>1</sup> \*

<sup>1</sup> - Electronics and Computer Science, ISIS Group  
University of Southampton, SO17 1BJ, United Kingdom

**Abstract.** In this presentation we show the semi-supervised learning with two input sources can be transformed into a maximum margin problem to be similar to a binary SVM. Our formulation exploits the unlabeled data to reduce the complexity of the class of the learning functions. In order to measure how the complexity is decreased we use the Rademacher Complexity Theory. The corresponding optimization problem is convex and it is efficiently solvable for large-scale applications as well.

## 1 Introduction

Semi-supervised learning belongs to the main directions of the recent machine learning research. The exploitation of the unlabeled data is an attractive approach either to extend the capability of the known methods or to derive novel learning devices. In this presentation we give a synthesis of some earlier approaches and show an optimization framework with good statistical generalization capability. Our idea consists of:

**Multiview learning**, when two or more sources of the inputs are given with the same output, see in papers of Blum et al. [4] and Dasgupta et al. [5],

**Maximum margin learning**, where the unlabeled cases within the margin considered as errors, see in Bennett et al. [3],

**Reduction of the learning class complexity**, where the “closeness” of the learning functions is assumed on the unlabeled data, see in the conference paper about skeleton learning of Lugosi et al. [6] and Balcan et al. [1],

**Combination of the multiview and maximum margin learning**, which introduces new constraints to the optimization, see in Meng et al. [7].

Blending these ideas we arrive at an optimization framework, which can be solved efficiently at the level of computational complexity of a binary SVM problem. In the first part the optimization problem is presented and then the generalization performance is analyzed. We apply the Rademacher complexity to give an upper bound on the expected error. At the end, experimental results will illustrate the usefulness of the approach. The underlying theory is summarized in Appendix A and B. Further details can be found in [2].

---

\*This work is supported by PASCAL Network of Excellence (IST-2002-506778) European Community IST Programmes. The authors want to thank Gabor Lugosi for very fruitful discussions.

## 2 Optimization Framework

In the semi-supervised, multiview learning, with two views, we are given a compound sample  $S_L$  of pairs of outputs and inputs  $\{(y_i, (\mathbf{x}_i^1, \mathbf{x}_i^2)) : y_i \in \{-1, 1\}, \mathbf{x}_i^k \in \mathcal{X}_k, i = 1, \dots, m_L, k = 1, 2\}$  independently and identically generated by an unknown multivariate distribution  $\mathcal{P}(Y, X)$ , and a compound sample  $S_U = \{(\mathbf{u}_i^1, \mathbf{u}_i^2) : \mathbf{u}_i^k \in \mathcal{X}_k, i = 1, \dots, m_U, k = 1, 2\}$  of unlabeled cases independently and identically chosen from the marginal distribution  $\mathcal{P}(X)$  of  $\mathcal{P}(Y, X)$ . Furthermore, there is given a set of embedding of the inputs into Hilbert spaces by the functions  $\phi_k : \mathcal{X}_k \rightarrow \mathcal{H}_{\phi_k}, k = 1, 2$ . The image vectors of these mappings are called feature vectors in the sequel.

The objective is to find linear functions  $f_k(\mathbf{x}_k) = \mathbf{w}_k^T \phi_k(\mathbf{x}_k) + b_k$  which can predict the potential label value for any labeled and unlabeled cases. The decision function is then defined as  $\frac{1}{2} \sum_k \text{sign}(f_k)$ . In order to exploit the unlabeled data in the optimization problem we choose particular solutions of the SVM subproblems where the predictors give similar solutions on the unlabeled data.

The matrix  $\mathbf{Y}$  is a diagonal matrix of the labels  $\{y_i\}$ , and the matrices  $\mathbf{X}_k$  and  $\mathbf{U}_k$  comprise the labeled and unlabeled inputs in their rows for  $k = 1, 2$ . Applying the embedding functions  $\phi_k$  on  $\mathbf{X}_k$  and  $\mathbf{U}_k$  gives matrices with the feature vectors in their rows, otherwise all vectors are column vectors.

The optimization problem formulating our learning framework is given by

$$\begin{array}{ll}
 \min & \frac{1}{2} \sum_k \mathbf{w}_k^T \mathbf{w}_k + \mathbf{1}^T \sum_k C_k \boldsymbol{\xi}_k + C_\eta \mathbf{1}^T (\boldsymbol{\eta}^+ + \boldsymbol{\eta}^-) \\
 \text{w.r.t.} & (\mathbf{w}_k, b_k, \boldsymbol{\xi}_k), k = 1, 2, (\boldsymbol{\eta}^+, \boldsymbol{\eta}^-), \\
 \text{Synthesis:} & \text{s.t.} \quad \sum_k (-1)^{k-1} (\phi_k(\mathbf{U}_k) \mathbf{w}_k + b_k) = \boldsymbol{\eta}^+ - \boldsymbol{\eta}^-, \\
 \text{Subproblems:} & \mathbf{Y}(\phi_k(\mathbf{X}_k) \mathbf{w}_k + b_k) \geq \mathbf{1} - \boldsymbol{\xi}_k, \\
 & \boldsymbol{\xi}_k \geq \mathbf{0}, (\boldsymbol{\eta}^+, \boldsymbol{\eta}^-) \geq \mathbf{0}, k = 1, 2.
 \end{array} \tag{1}$$

We use the acronym SVM.2K for the problem (1) in the sequel.

Introducing Lagrangian multipliers  $\boldsymbol{\alpha}_k$  for any  $k$  to the SVM subproblems and  $\boldsymbol{\beta}$  to the synthesis constraints then we can express the normal vector of the separating hyperplanes for any  $k$  by

$$\mathbf{w}_k = [\phi_k(\mathbf{X}_k)^T, \phi_k(\mathbf{U}_k)^T] \mathbf{g}_k, \text{ where } \mathbf{g}_k^T = (\boldsymbol{\alpha}_k^T \mathbf{Y}, -(-1)^k \boldsymbol{\beta}^T),$$

then unfolding the Karush-Kuhn-Tucker conditions we can derive the corresponding dual problem. The kernel matrices in the dual have the structure

$$\mathbf{K}_k = \begin{bmatrix} \phi_k(\mathbf{X}_k) \phi_k(\mathbf{X}_k)^T & \phi_k(\mathbf{X}_k) \phi_k(\mathbf{U}_k)^T \\ \phi_k(\mathbf{U}_k) \phi_k(\mathbf{X}_k)^T & \phi_k(\mathbf{U}_k) \phi_k(\mathbf{U}_k)^T \end{bmatrix} = \begin{bmatrix} \mathbf{K}_k^{LL} & \mathbf{K}_k^{LU} \\ \mathbf{K}_k^{UL} & \mathbf{K}_k^{UU} \end{bmatrix} = \begin{bmatrix} \mathbf{K}_k^L \\ \mathbf{K}_k^U \end{bmatrix}.$$

Thus, the dual reads as

$$\begin{array}{ll}
 \min_{\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \boldsymbol{\beta}} & \frac{1}{2} \sum_k \mathbf{g}_k^T \mathbf{K}_k \mathbf{g}_k - \sum_k \mathbf{1}^T \boldsymbol{\alpha}_k \\
 \text{s.t.} & \mathbf{1}^T \mathbf{g}_k = 0, \text{ where } \mathbf{g}_k^T = (\boldsymbol{\alpha}_k^T \mathbf{Y} - (-1)^k \boldsymbol{\beta}^T), \\
 & 0 \leq \boldsymbol{\alpha}_k \leq C_k, -C_\eta \leq \boldsymbol{\beta} \leq C_\eta, k = 1, 2.
 \end{array}$$

### 3 Theoretical Foundation

To illuminate the theoretical background of the presented learning method we use the theory of the Rademacher Complexity, see in [2]. We assume the notation taken from Section 2.

#### 3.1 Analyzing the Semisupervised Multiview Learning

For SVM\_2K, the two feature sets are  $((\phi_1(\mathbf{X}_1), \phi_1(\mathbf{U}_1))$  and  $((\phi_2(\mathbf{X}_2), \phi_2(\mathbf{U}_2))$ .

An application of Theorem 3, see in the Appendix A, shows that

$$f_D(\mathbf{u}, f_1, f_2) \stackrel{\text{def}}{=} \mathbb{E}_{\mathbf{u}}[|\sum_k (-1)^{k-1} (\mathbf{w}_k^T \phi(\mathbf{u}_k) + b_k)|] \\ \leq \frac{1}{m_U} \mathbf{1}^T (\boldsymbol{\eta}^+ + \boldsymbol{\eta}^-) + \frac{2C}{m_U} \sqrt{\sum_k \text{tr}(\mathbf{K}_k)} + 3\sqrt{\frac{2\ln(2/\delta)}{2m_U}} =: D.$$

with probability at least  $1 - \delta$ . We have assumed that  $\|\mathbf{w}_k\|^2 + b_k^2 \leq C^2$  for some prefixed  $C$  and any  $k = 1, 2$ . Hence, the class of functions we are considering when applying SVM\_2K to this problem can be restricted to

$$\mathcal{F}_{C,D} = \left\{ f \mid f : (\mathbf{x}_k) \rightarrow \frac{1}{2} \sum_{k=1} (\mathbf{K}_k(x_k) \mathbf{g}_k + b_k), \right. \\ \left. \mathbf{g}_k^T \mathbf{K}_k \mathbf{g}_k + b_k^2 \leq C^2, \quad k = 1, 2, \quad f_D(\mathbf{u}, f_1, f_2) \leq D \right\},$$

where  $\mathbf{K}_k(\mathbf{x}_k) = [\phi(\mathbf{x}_k)^T \phi(\mathbf{X}_k)^T, \phi(\mathbf{x}_k)^T \phi(\mathbf{U}_k)^T]$ .

The class  $\mathcal{F}_{C,D}$  is clearly closed under negation.

Applying the usual Rademacher techniques for margin bounds on generalization we obtain the following result.

**Theorem 1.** *Fix  $\delta \in (0, 1)$  and let  $\mathcal{F}_{C,D}$  be the class of functions described above. Let  $(\mathbf{X}_k)$  labeled and  $(\mathbf{U}_k)$  unlabeled samples be drawn independently according to a probability distribution  $\mathcal{P}(X)$  for  $k = 1, 2$ . Then with probability at least  $1 - \delta$  over random draws of labeled samples of size  $m_L$ , every  $f \in \mathcal{F}_{C,D}$  satisfies*

$$P_{(x,y) \sim \mathcal{D}}(\text{sign}(f(\mathbf{x})) \neq y) \leq \frac{1}{2m_L} \mathbf{1}^T \sum_k \boldsymbol{\xi}_k + \hat{R}_\ell(\mathcal{F}_{C,D}) + 3\sqrt{\frac{\ln(2/\delta)}{2m_L}}.$$

It therefore remains to compute the empirical Rademacher complexity of  $\mathcal{F}_{C,D}$ , which is the critical discriminator between the bounds for the individual SVMs and that of the SVM\_2K. The details are presented in the Appendix B.

### 4 Experiments

In the experiment we used the image dataset<sup>2</sup> being commonly used for generic object recognition. In the cross-validation 5% percent of the cases are randomly

<sup>2</sup>Available at <http://www.robots.ox.ac.uk/~vgg/data/>

chosen as labeled cases and all others were unlabeled. This random selection of the labeled cases was ten times repeated.

Table 1 shows the mean and the standard deviation of the accuracies and the Rademacher complexities computed for each image class by using SVM working on the concatenation of two feature sets and the same ones in case of the SVM\_2K. It shows that the smaller Rademacher complexity can increase the mean and decrease the standard deviation of the accuracies in classification.

		Image classes		
		Airplanes	Faces	Motorbikes
SVM on two features	mean(std)	91.4(2.8)	96.8(1.3)	94.1(1.3)
	Rad. comp.	588.2	339.1	574.3
SVM_2K	mean(std)	<b>91.9(2.1)</b>	<b>97.4(1.3)</b>	<b>94.3(1.1)</b>
	Rad. comp.	236.4	34.4	197.4

Table 1: Accuracies(%), standard deviation and estimation of Rademacher Complexities of the SVM acting on two feature sets and the SVM\_2K in three image classes. 5% of the cases were labeled.

## A Appendix: Short Introduction of Rademacher Complexity Theory

We begin with the definitions required for Rademacher complexity, see for example Bartlett and Mendelson [2] (see also [8] for an introductory exposition).

**Definition 2.** For a sample  $S = \{x_1, \dots, x_\ell\}$  generated by a distribution  $D$  on a set  $X$  and a real-valued function class  $\mathcal{F}$  with a domain  $X$ , the empirical Rademacher complexity of  $\mathcal{F}$  is the random variable

$$\hat{R}_\ell(\mathcal{F}) = \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \left| \frac{2}{\ell} \sum_{i=1}^{\ell} \sigma_i f(x_i) \right| \middle| x_1, \dots, x_\ell \right]$$

where  $\sigma = \{\sigma_1, \dots, \sigma_\ell\}$  are independent uniform  $\{\pm 1\}$ -valued Rademacher random variables. The Rademacher complexity of  $\mathcal{F}$  is

$$R_\ell(\mathcal{F}) = \mathbb{E}_S [\hat{R}_\ell(\mathcal{F})] = \mathbb{E}_{S, \sigma} \left[ \sup_{f \in \mathcal{F}} \left| \frac{2}{\ell} \sum_{i=1}^{\ell} \sigma_i f(x_i) \right| \right]$$

We use  $\mathbb{E}_D$  to denote expectation with respect to a distribution  $D$  and  $\mathbb{E}_S$  when the distribution is the uniform (empirical) distribution on a sample  $S$ .

**Theorem 3.** Fix  $\delta \in (0, 1)$  and let  $\mathcal{F}$  be a class of functions mapping from  $S$  to  $[0, 1]$ . Let  $(x_i)_{i=1}^{\ell}$  be drawn independently according to a probability distribution

$\mathcal{D}$ . Then with probability at least  $1 - \delta$  over random draws of samples of size  $\ell$ , every  $f \in \mathcal{F}$  satisfies

$$\mathbb{E}_{\mathcal{D}} [f(x)] \leq \mathbb{E}_S [f(x)] + R_{\ell}(\mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2\ell}} \leq \mathbb{E}_S [f(x)] + \hat{R}_{\ell}(\mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2\ell}}$$

Given a training set  $S$  the class of functions that we will primarily be considering are linear functions with bounded norm

$$\left\{ x \rightarrow \sum_{i=1}^{\ell} \alpha_i \kappa(x_i, x) : \alpha' K \alpha \leq B^2 \right\} \subseteq \{x \rightarrow \langle w, \phi(x) \rangle : \|w\| \leq B\} = \mathcal{F}_B,$$

where  $\phi$  is the feature mapping corresponding to the kernel  $\kappa$  and  $K$  is the corresponding kernel matrix for the sample  $S$ . The following result bounds the Rademacher complexity of linear function classes.

**Theorem 4.** [2] *If  $\kappa : X \times X \rightarrow R$  is a kernel, and  $S = \{x_1, \dots, x_{\ell}\}$  is a i.i.d. sample from  $X$ , then the empirical Rademacher complexity of the class  $\mathcal{F}_B$  satisfies  $\hat{R}_{\ell}(\mathcal{F}) \leq \frac{2B}{\ell} \sqrt{\text{tr}(K)}$ .*

## B Appendix: Empirical Rademacher Complexity of $\mathcal{F}_{C,D}$

We define an auxiliary function of the weight vectors  $\bar{\mathbf{w}}_k = (\mathbf{w}_k, b_k)$ ,  $k = 1, 2$ ,

$$D(\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2) \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}} \left[ \left| \sum_k (-1)^{k-1} (\phi_k(\mathbf{U}_k) \mathbf{w}_k + b_k) \right| \right],$$

and the Rademacher complexity of the class  $\mathcal{F}_{C,D}$

$$\begin{aligned} \hat{R}_{\ell}(\mathcal{F}_{C,D}) &= \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}_{C,D}} \left| \frac{2}{m_L} \sum_{i=1}^{m_L} \sigma_i f(\mathbf{x}_i) \right| \right] \\ &= \mathbb{E}_{\sigma} \left[ \sup_{\substack{\|\bar{\mathbf{w}}_k\| \leq C, k=1,2 \\ D(\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2) \leq D}} \left| \frac{2}{m_L} \sum_{k=1}^2 \sigma^T(\mathbf{K}_k(\mathbf{x}_k) \mathbf{g}_k + b_k) \right| \right]. \end{aligned}$$

Based on the reversed version of the basic Rademacher complexity theorem where the roles of the empirical and true expectations are swapped:

**Theorem 5.** *Fix  $\delta \in (0, 1)$  and let  $\mathcal{F}$  be a class of functions mapping from  $S$  to  $[0, 1]$ . Let  $(x_i)_{i=1}^{\ell}$  be drawn independently according to a probability distribution  $\mathcal{D}$ . Then with probability at least  $1 - \delta$  over random draws of samples of size  $\ell$ , every  $f \in \mathcal{F}$  satisfies*

$$\mathbb{E}_S [f(x)] \leq \mathbb{E}_{\mathcal{D}} [f(x)] + R_{\ell}(\mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2\ell}} \leq \mathbb{E}_{\mathcal{D}} [f(x)] + \hat{R}_{\ell}(\mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2\ell}}.$$

The proof tracks that of Theorem 3 but is omitted through lack of space.

For weight vectors  $\{\bar{\mathbf{w}}_k\}$  satisfying  $D(\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2) \leq D$ , an application of Theorem 5 shows that with probability at least  $1 - \delta$  we have

$$\begin{aligned} \hat{D}(\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2) &\stackrel{def}{=} \mathbb{E}_S \left[ \left| \sum_k (-1)^{k-1} (\phi(\mathbf{U}_k) \mathbf{w}_k + b_k) \right| \right] \\ &\leq D + \frac{2C}{m_U} \sqrt{\sum_k \text{tr}(\mathbf{K}_k)} + 3 \sqrt{\frac{2 \ln(2/\delta)}{2m_U}} \\ &\leq \frac{1}{m_U} \mathbf{1}^T (\eta^+ + \eta^-) + \frac{2C}{m_U} \sqrt{\sum_k \text{tr}(\mathbf{K}_k)} + 3 \sqrt{\frac{\ln(2/\delta)}{2m_U}} =: \hat{D}. \end{aligned}$$

The above result shows that the Rademacher complexity of  $\mathcal{F}_{C,D}$  with probability greater than  $1 - \delta$  satisfies

$$\hat{R}_\ell(\mathcal{F}_{C,D}) \leq \mathbb{E}_\sigma \left[ \sup_{\substack{\|\bar{\mathbf{w}}_k\| \leq C, k=1,2, \\ \hat{D}(\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2) \leq \hat{D}}} \left| \frac{2}{m_L} \boldsymbol{\sigma}^T \sum_k [\phi_k(\mathbf{X}_k) \mathbf{w}_k + \mathbf{1} b_k] \right| \right],$$

where  $\boldsymbol{\sigma} \in \{-1, +1\}^{m_L}$ . Note that the expression in square brackets is concentrated under the uniform distribution of Rademacher variables. Hence, we can estimate the complexity for randomly chosen instantiation  $\hat{\sigma}$  of the the Rademacher variables  $\sigma$ . We now must find the value of  $\{\bar{\mathbf{w}}_k\}$  that maximizes

$$\begin{aligned} \max \quad & \frac{1}{m_L} \left| \left[ \sum_k \boldsymbol{\sigma}^T \phi_k(\mathbf{X}_k) \mathbf{w}_k + \sum_k \boldsymbol{\sigma}^T \mathbf{1} b_k \right] \right| = \frac{1}{m_L} \left| \sum_k \boldsymbol{\sigma}^T (\mathbf{K}_k^L \mathbf{g}_k + b_k) \right| \\ \text{s.t.} \quad & \mathbf{g}_k^T \mathbf{K}_k \mathbf{g}_k + b_k^2 \leq C^2, \quad k = 1, 2, \\ & \frac{1}{m_U} \mathbf{1}^T \left| \left( \sum_k (-1)^{k-1} (\mathbf{K}_k^U \mathbf{g}_k + b_k) \right) \right| \leq \hat{D}. \end{aligned}$$

The expected value of the objective function computed on randomly chosen  $\hat{\sigma}$ 's is the estimate of the Rademacher complexity.

## References

- [1] M.F. Balcan and A. Blum. A pac-style model for learning from labeled and unlabeled data. In *COLT*, pages 111–126, 2005.
- [2] P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- [3] K. Bennett and A. Demirez. Semi-supervised support vector machines. In M.S. Kearns, S.A. Solla, and D.A. Cohn, editors, *Advances in Neural Information Processing Systems*, volume 12, pages 368–374. MIT Press, Cambridge, MA, 1998.
- [4] 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.
- [5] S. Dasgupta, M.L. Littman, and D. McAllester. Pac generalization bounds for co-training. In *Advances in Neural Information Processing Systems (NIPS)*. 2001.
- [6] G. Lugosi and M. Pinter. A data-dependent skeleton estimate for learning. In *Proceedings of the 9th Annual Conference on Computational Learning Theory, New York*,, pages 51–58. Association for Computing Machinery, 1996.
- [7] H. Meng, J. Shawe-Taylor, S. Szedmak, and J.R.D. Farquhar. Support vector machine to synthesise kernels. In *Sheffield Machine Learning Workshop Proceedings, Lecture Notes in Computer Science*. Springer, 2005.
- [8] J. Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.