

# Manifold Regularization: A Geometric Framework for Learning from Examples

Mikhail Belkin\*, Partha Niyogi<sup>†</sup>, Vikas Sindhwani<sup>‡</sup>  
The University of Chicago  
Hyde Park, Chicago, IL 60637  
{misha, niyogi, vikass}@cs.uchicago.edu

September 9, 2004

## Abstract

We propose a family of learning algorithms based on a new form of regularization that allows us to exploit the geometry of the marginal distribution. We focus on a semi-supervised framework that incorporates labeled and unlabeled data in a general-purpose learner. Some transductive graph learning algorithms and standard methods including Support Vector Machines and Regularized Least Squares can be obtained as special cases. We utilize properties of Reproducing Kernel Hilbert spaces to prove new Representer theorems that provide theoretical basis for the algorithms. As a result (in contrast to purely graph based approaches) we obtain a natural out-of-sample extension to novel examples and so are able to handle both transductive and truly semi-supervised settings. We present experimental evidence suggesting that our semi-supervised algorithms are able to use unlabeled data effectively. Finally we have a brief discussion of unsupervised and fully supervised learning within our general framework.

## 1 Introduction

In this paper, we introduce a new framework for data-dependent regularization that exploits the geometry of the probability distribution. While this framework allows us to approach the full range of learning problems from unsupervised to supervised, we focus on the problem of semi-supervised learning. The problem of learning

---

\*Department of Computer Science, also TTI Chicago

<sup>†</sup>Departments of Computer Science and Statistics

<sup>‡</sup>Department of Computer Science, also TTI Chicago

from labeled and unlabeled data (*semi-supervised* and *transductive* learning) has attracted considerable attention in recent years. Some recently proposed methods include Transductive SVM [35, 22], Cotraining [13], and a variety of graph based methods [12, 14, 32, 37, 38, 24, 23, 4]. We also note two regularization based techniques [16, 7]. The latter reference is closest in spirit to the intuitions of our paper.

The idea of regularization has a rich mathematical history going back to [34], where it is used for solving ill-posed inverse problems. Regularization is a key idea in the theory of splines (e.g., [36]) and is widely used in machine learning (e.g., [20]). Many machine learning algorithms, including Support Vector Machines, can be interpreted as instances of regularization.

Our framework exploits the geometry of the probability distribution that generates the data and incorporates it as an additional regularization term. Hence, there are two regularization terms — one controlling the complexity of the classifier in the *ambient space* and the other controlling the complexity as measured by the *geometry* of the distribution. We consider in some detail the special case where this probability distribution is supported on a submanifold of the ambient space.

The points below highlight several aspects of the current paper:

1. Our general framework brings together three distinct concepts that have received some independent recent attention in machine learning:
  - i. The first of these is the technology of *spectral graph theory* (e.g., see [15]) that has been applied to a wide range of clustering and classification tasks over the last two decades. Such methods typically reduce to certain eigenvalue problems.
  - ii. The second is the geometric point of view embodied in a class of algorithms that can be termed as *manifold learning* (see webpage [39] for a fairly long list of references). These methods attempt to use the geometry of the probability distribution by assuming that its support has the geometric structure of a Riemannian manifold.
  - iii. The third important conceptual framework is the set of ideas surrounding regularization in Reproducing Kernel Hilbert Spaces. This leads to the class of *kernel based algorithms* for classification and regression (e.g., see [31], [36], [20]).

We show how to bring these ideas together in a coherent and natural way to incorporate geometric structure in a kernel based regularization framework. As far as we know, these ideas have not been unified in a similar fashion before.

2. Within this general framework, we propose two specific families of algo-

rithms: the Laplacian Regularized Least Squares (hereafter LapRLS) and the Laplacian Support Vector Machines (hereafter LapSVM). These are natural extensions of RLS and SVM respectively. In addition, several recently proposed transductive methods (e.g., [38, 4]) are also seen to be special cases of this general approach.

3. We elaborate on the RKHS foundations of our algorithms and show how geometric knowledge of the probability distribution may be incorporated in such a setting. In particular, a new Representer theorem provides a functional form of the solution when the distribution is known and an empirical version which involves an expansion over labeled and unlabeled points when the distribution is unknown. These Representer theorems provide the basis for our algorithms.
4. Our framework with an ambiently defined RKHS and the associated Representer theorems result in a natural out-of-sample extension from the data set (labeled and unlabeled) to novel examples. This is in contrast to the variety of purely graph based approaches that have been considered in the last few years. Such graph based approaches work in a transductive setting and do not naturally extend to the semi-supervised case where novel test examples need to be classified (predicted). Also see [8, 11] for some recent related work on out-of-sample extensions.

## 1.1 The Significance of Semi-Supervised Learning

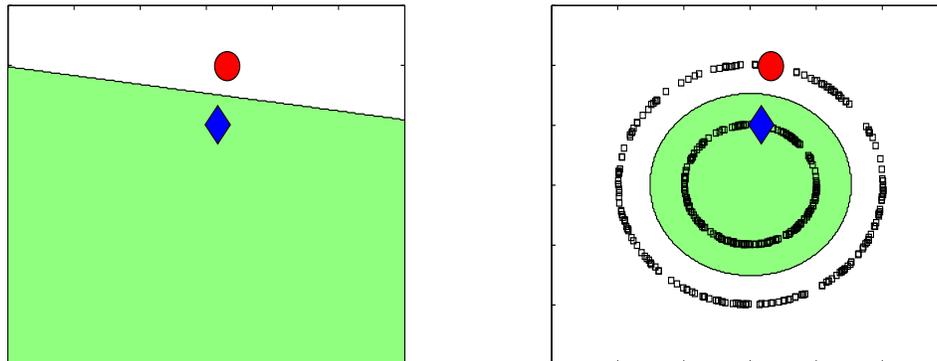
From an engineering standpoint, it is clear that collecting labeled data is generally more involved than collecting unlabeled data. As a result, an approach to pattern recognition that is able to make better use of unlabeled data to improve recognition performance is of potentially great practical significance.

However, the significance of semi-supervised learning extends beyond purely utilitarian considerations. Arguably, most natural (human or animal) learning occurs in the semi-supervised regime. We live in a world where we are constantly exposed to a stream of natural stimuli. These stimuli comprise the unlabeled data that we have easy access to. For example, in phonological acquisition contexts, a child is exposed to many acoustic utterances. These utterances do not come with identifiable phonological markers. Corrective feedback is the main source of directly labeled examples. In many cases, a small amount of feedback is sufficient to allow the child to master the acoustic-to-phonetic mapping of any language.

The ability of humans to do unsupervised learning (e.g. learning clusters and categories of objects) suggests that unlabeled data can be usefully processed to learn natural invariances, to form categories, and to develop classifiers. In most

pattern recognition tasks, humans have access only to a small number of labeled examples. Therefore the success of human learning in this “small sample” regime is plausibly due to effective utilization of the large amounts of unlabeled data to extract information that is useful for generalization.

Figure 1: Unlabeled Data and Prior Beliefs



Consequently, if we are to make progress in understanding how natural learning comes about, we need to think about the basis of semi-supervised learning. Figure 1 illustrates how unlabeled examples may force us to restructure our hypotheses during learning. Imagine a situation where one is given two labeled examples — one positive and one negative — as shown in the left panel. If one is to induce a classifier on the basis of this, a natural choice would seem to be the linear separator as shown. Indeed, a variety of theoretical formalisms (Bayesian paradigms, Regularization, Minimum Description Length or Structural Risk Minimization principles, and the like) have been constructed to rationalize such a choice. In most of these formalisms, one structures the set of one’s hypothesis functions by a prior notion of simplicity and one may then justify why the linear separator is the simplest structure consistent with the data.

Now consider the situation where one is given additional unlabeled examples as shown in the right panel. We argue that it is self-evident that in the light of this new unlabeled set, one must re-evaluate one’s prior notion of simplicity. The particular geometric structure of the marginal distribution suggests that the most natural classifier is now the circular one indicated in the right panel. Thus the geometry of the marginal distribution must be incorporated in our regularization principle to impose structure on the space of functions in nonparametric classifi-

cation or regression. This is the intuition we formalize in the rest of the paper. The success of our approach depends on whether we can extract structure from the marginal distribution and on the extent to which such structure may reveal the underlying truth.

## 1.2 Outline of the Paper

The paper is structured as follows: in Sec. 2, we develop the basic framework for semi-supervised learning where we ultimately formulate an objective function that can utilize both labeled and unlabeled data. The framework is developed in an RKHS setting and we state two kinds of Representer theorems describing the functional form of the solutions. In Sec. 3, we elaborate on the theoretical underpinnings of this framework and prove the Representer theorems of Sec. 2. While the Representer theorem for the finite sample case can be proved using standard orthogonality arguments, the Representer theorem for the known marginal distribution requires more subtle considerations. In Sec. 4, we derive the different algorithms for semi-supervised learning that arise out of our framework. Connections to related algorithms are stated. In Sec. 5, we describe experiments that evaluate the algorithms against each other and demonstrate the usefulness of unlabeled data. In Sec. 6, we consider the cases of fully supervised and unsupervised learning. In Sec. 7 we conclude.

# 2 The Semi-Supervised Learning Framework

## 2.1 Background

Recall the standard framework of learning from examples. There is a probability distribution  $P$  on  $X \times \mathbb{R}$  according to which examples are generated for function learning. Labeled examples are  $(x, y)$  pairs generated according to  $P$ . Unlabeled examples are simply  $x \in X$  drawn according to the marginal distribution  $\mathcal{P}_X$  of  $P$ .

One might hope that knowledge of the marginal  $\mathcal{P}_X$  can be exploited for better function learning (e.g. in classification or regression tasks). Of course, if there is no identifiable relation between  $\mathcal{P}_X$  and the conditional  $\mathcal{P}(y|x)$ , the knowledge of  $\mathcal{P}_X$  is unlikely to be of much use.

Therefore, we will make a specific assumption about the connection between the marginal and the conditional distributions. We will assume that if two points  $x_1, x_2 \in X$  are *close* in the *intrinsic* geometry of  $\mathcal{P}_X$ , then the conditional distributions  $\mathcal{P}(y|x_1)$  and  $\mathcal{P}(y|x_2)$  are similar. In other words, the conditional probability

distribution  $\mathcal{P}(y|x)$  varies smoothly along the geodesics in the intrinsic geometry of  $\mathcal{P}_X$ .

We utilize these geometric intuitions to extend an established framework for function learning. A number of popular algorithms such as SVM, Ridge regression, splines, Radial Basis Functions may be broadly interpreted as regularization algorithms with different empirical cost functions and complexity measures in an appropriately chosen Reproducing Kernel Hilbert Space (RKHS).

For a Mercer kernel  $K : X \times X \rightarrow \mathbb{R}$ , there is an associated RKHS  $\mathcal{H}_K$  of functions  $X \rightarrow \mathbb{R}$  with the corresponding norm  $\|\cdot\|_K$ . Given a set of labeled examples  $(x_i, y_i)$ ,  $i = 1, \dots, l$  the standard framework estimates an unknown function by minimizing

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f) + \gamma \|f\|_K^2 \quad (1)$$

where  $V$  is some loss function, such as squared loss  $(y_i - f(x_i))^2$  for RLS or the soft margin loss function for SVM. Penalizing the RKHS norm imposes smoothness conditions on possible solutions. The classical Representer Theorem states that the solution to this minimization problem exists in  $\mathcal{H}_K$  and can be written as

$$f^*(x) = \sum_{i=1}^l \alpha_i K(x_i, x) \quad (2)$$

Therefore, the problem is reduced to optimizing over the finite dimensional space of coefficients  $\alpha_i$ , which is the algorithmic basis for SVM, Regularized Least Squares and other regression and classification schemes.

We first consider the case when the marginal distribution is already known.

## 2.2 Marginal $\mathcal{P}_X$ is known

Our goal is to extend this framework by incorporating additional information about the geometric structure of the marginal  $\mathcal{P}_X$ . We would like to ensure that the solution is smooth with respect to both the ambient space and the marginal distribution  $\mathcal{P}_X$ . To achieve that, we introduce an additional regularizer:

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f) + \gamma_A \|f\|_K^2 + \gamma_I \|f\|_I^2 \quad (3)$$

where  $\|f\|_I^2$  is an appropriate penalty term that should reflect the intrinsic structure of  $\mathcal{P}_X$ . Here  $\gamma_A$  controls the complexity of the function in the *ambient* space while  $\gamma_I$  controls the complexity of the function in the *intrinsic* geometry of  $\mathcal{P}_X$ . It turns

out that one can derive an explicit functional form for the solution  $f^*$  as shown in the following theorem.

**Theorem 2.1.** *Assume that the penalty term  $\|f\|_I$  is sufficiently smooth with respect to the RKHS norm  $\|f\|_K$ . Then the solution  $f^*$  to the optimization problem in Eqn. 3 above exists and admits the following representation*

$$f^*(x) = \sum_{i=1}^l \alpha_i K(x_i, x) + \int_{\mathcal{M}} \alpha(y) K(x, y) d\mathcal{P}_X(y) \quad (4)$$

where  $\mathcal{M} = \text{supp}\{\mathcal{P}_X\}$  is the support of the marginal  $\mathcal{P}_X$ .

We postpone the proof and the exact formulation of smoothness conditions on the norm  $\|\cdot\|_I$  until the next section.

The Representer Theorem above allows us to express the solution  $f^*$  directly in terms of the labeled data, the (ambient) kernel  $K$ , and the marginal  $\mathcal{P}_X$ . If  $\mathcal{P}_X$  is unknown, we see that the solution may be expressed in terms of an empirical estimate of  $\mathcal{P}_X$ . Depending on the nature of this estimate, different approximations to the solution may be developed. In the next section, we consider a particular approximation scheme that leads to a simple algorithmic framework for learning from labeled and unlabeled data.

### 2.3 Marginal $\mathcal{P}_X$ Unknown

In most applications the marginal  $\mathcal{P}_X$  is not known. Therefore we must attempt to get empirical estimates of  $\mathcal{P}_X$  and  $\|\cdot\|_I$ . Note that in order to get such empirical estimates it is sufficient to have *unlabeled* examples.

A case of particular recent interest (e.g., see [27, 33, 5, 19] for a discussion on dimensionality reduction) is when the support of  $\mathcal{P}_X$  is a compact submanifold  $\mathcal{M} \subset X = \mathbb{R}^n$ . In that case, one natural choice for  $\|f\|_I$  is  $\int_{\mathcal{M}} \langle \nabla_{\mathcal{M}} f, \nabla_{\mathcal{M}} f \rangle$ .

The optimization problem becomes

$$f^* = \underset{f \in \mathcal{H}_K}{\text{argmin}} \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f) + \gamma_A \|f\|_K^2 + \gamma_I \int_{\mathcal{M}} \langle \nabla_{\mathcal{M}} f, \nabla_{\mathcal{M}} f \rangle$$

The term  $\int_{\mathcal{M}} \langle \nabla_{\mathcal{M}} f, \nabla_{\mathcal{M}} f \rangle$  may be approximated on the basis of labeled and unlabeled data using the graph Laplacian ([4], also see [7]). Thus, given a set of  $l$  labeled examples  $\{(x_i, y_i)\}_{i=1}^l$  and a set of  $u$  unlabeled examples  $\{x_j\}_{j=l+1}^{l+u}$ , we consider the following optimization problem :

$$f^* = \underset{f \in \mathcal{H}_K}{\text{argmin}} \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f) + \gamma_A \|f\|_K^2 + \frac{\gamma_I}{(u+l)^2} \sum_{i,j=1}^{l+u} (f(x_i) - f(x_j))^2 W_{ij}$$

$$= \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f) + \gamma_A \|f\|_K^2 + \frac{\gamma_I}{(u+l)^2} \mathbf{f}^T L \mathbf{f} \quad (5)$$

where  $W_{ij}$  are edge weights in the data adjacency graph,  $\mathbf{f} = [f(x_1), \dots, f(x_{l+u})]^T$ , and  $L$  is the graph Laplacian given by  $L = D - W$ . Here, the diagonal matrix  $D$  is given by  $D_{ii} = \sum_{j=1}^{l+u} W_{ij}$ . The normalizing coefficient  $\frac{1}{(u+l)^2}$  is the natural scale factor for the empirical estimate of the Laplace operator. We note that on a sparse adjacency graph it may be replaced by  $\sum_{i,j=1}^{l+u} W_{ij}$ .

The following version of the Representer Theorem shows that the minimizer has an expansion in terms of both labeled and unlabeled examples and is a key to our algorithms.

**Theorem 2.2.** *The minimizer of optimization problem 5 admits an expansion*

$$f^*(x) = \sum_{i=1}^{l+u} \alpha_i K(x_i, x) \quad (6)$$

*in terms of the labeled and unlabeled examples.*

The proof is a variation of the standard orthogonality argument and is presented in Section 3.4.

**Remark 1:** Several natural choices of  $\|\cdot\|_I$  exist. Some examples are:

1. Iterated Laplacians  $\mathcal{L}^k$ . Differential operators  $\mathcal{L}^k$  and their linear combinations provide a natural family of smoothness penalties.
2. Heat semigroup  $e^{-\mathcal{L}t}$  is a family of smoothing operators corresponding to the process of diffusion (Brownian motion) on the manifold. One can take  $\|f\|_I^2 = \int_{\mathcal{M}} f e^{\mathcal{L}t}(f)$ . We note that for small values of  $t$  the corresponding Green's function (the heat kernel of  $\mathcal{M}$ ) can be approximated by a sharp Gaussian in the ambient space.
3. Squared norm of the Hessian (cf. [19]). While the Hessian  $\mathbf{H}(f)$  (the matrix of second derivatives of  $f$ ) generally depends on the coordinate system, it can be shown that the Frobenius norm (the sum of squared eigenvalues) of  $\mathbf{H}$  is the same in any geodesic coordinate system and hence is invariantly defined for a Riemannian manifold  $\mathcal{M}$ . Using the Frobenius norm of  $\mathbf{H}$  as a regularizer presents an intriguing generalization of thin-plate splines. We also note that  $\mathcal{L}(f) = \operatorname{tr}(\mathbf{H}(f))$ .

**Remark 2:** Note that  $K$  restricted to  $\mathcal{M}$  (denoted by  $K_{\mathcal{M}}$ ) is also a kernel defined on  $\mathcal{M}$  with an associated RKHS  $\mathcal{H}_{\mathcal{M}}$  of functions  $\mathcal{M} \rightarrow \mathbb{R}$ . While this might suggest  $\|f\|_I = \|f_{\mathcal{M}}\|_{K_{\mathcal{M}}}$  ( $f_{\mathcal{M}}$  is  $f$  restricted to  $\mathcal{M}$ ) as a reasonable choice for  $\|f\|_I$ ,

it turns out, that for the minimizer  $f^*$  of the corresponding optimization problem we get  $\|f^*\|_I = \|f^*\|_K$ , yielding the same solution as standard regularization, although with a different  $\gamma$ . This observation follows from the restriction properties of RKHS discussed in the next section and is formally stated as Proposition 3.4. Therefore it is impossible to have an out-of-sample extension without two *different* measures of smoothness. On the other hand, a different ambient kernel restricted to  $\mathcal{M}$  can potentially serve as the intrinsic regularization term. For example, a sharp Gaussian kernel can be used as an approximation to the heat kernel on  $\mathcal{M}$ .

### 3 Theoretical Underpinnings and Results

In this section we briefly review the theory of Reproducing Kernel Hilbert spaces and their connection to integral operators. We proceed to establish the Representer theorems from the previous section.

#### 3.1 General Theory of RKHS

We start by recalling some basic properties of Reproducing Kernel Hilbert Spaces (see the original work [2] and also [17] for a nice discussion in the context of learning theory) and their connections to integral operators. We say that a Hilbert space  $\mathcal{H}$  of functions  $X \rightarrow \mathbb{R}$  has the *reproducing property*, if  $\forall x \in X$  the evaluation functional  $f \rightarrow f(x)$  is continuous. For the purposes of this discussion we will assume that  $X$  is compact. By Riesz Representation theorem it follows that for a given  $x \in X$ , there is a function  $h_x \in \mathcal{H}$ , s.t.

$$\forall f \in \mathcal{H} \quad \langle h_x, f \rangle_{\mathcal{H}} = f(x)$$

We can therefore define the corresponding *kernel function*

$$K(x, y) = \langle h_x, h_y \rangle_{\mathcal{H}}$$

It follows that  $h_x(y) = \langle h_x, h_y \rangle_{\mathcal{H}} = K(x, y)$  and thus  $\langle K(x, \cdot), f \rangle = f(x)$ . It is clear that  $K(x, \cdot) \in \mathcal{H}$ .

It can be easily shown that  $K(x, y)$  is a positive semi-definite kernel, i.e., that given  $n$  points  $x_1, \dots, x_n$ , the corresponding matrix  $K$  with  $K_{ij} = K(x_i, x_j)$  is positive semi-definite. Moreover, if the space  $\mathcal{H}$  is sufficiently rich, that is if for any  $x_1, \dots, x_n$  there is a function  $f$ , s.t.  $f(x_1) = 1, f(x_i) = 0, i > 1$ , then  $K$  is strictly positive definite. Conversely any function  $K(x, y)$  with such property gives rise to an RKHS. For simplicity we will assume that all our RKHS are rich (the corresponding kernels are sometimes called *universal*).

We proceed to endow  $X$  with a measure  $\mu$  (supported on all of  $X$ ). We denote the corresponding Hilbert norm by  $\langle \cdot, \cdot \rangle_\mu$ .

We can now consider the integral operator  $L_K$  corresponding the kernel  $K$ :

$$(L_K f)(x) = \int_X f(y)K(x, y) d\mu$$

It is well-known that this operator is compact and self-adjoint with respect to  $\mathcal{L}_\mu^2$  and, by the Spectral Theorem, its eigenfunctions  $e_1(x), e_2(x), \dots$  form an orthogonal basis of  $\mathcal{L}_\mu^2$  and the corresponding eigenvalues  $\lambda_1, \lambda_2, \dots$  are discrete with finite multiplicity,  $\lim_{i \rightarrow \infty} \lambda_i = 0$ .

We see that

$$\langle K(x, \cdot), e_i(\cdot) \rangle_\mu = \lambda_i e_i(x)$$

and therefore  $K(x, y) = \sum_i \lambda_i e_i(x) e_i(y)$ . Writing a function  $f$  in that basis, we have  $f = \sum a_i e_i(x)$  and  $\langle K(x, \cdot), f(\cdot) \rangle_\mu = \sum_i \lambda_i a_i e_i(x)$ . Assuming that  $e_i \in \mathcal{H}$ , which can be ensured by extending  $\mathcal{H}$ , if necessary, we see that

$$e_j(x) = \langle K(x, \cdot), e_j(\cdot) \rangle_{\mathcal{H}} = \sum_i \lambda_i e_i(x) \langle e_i, e_j \rangle_{\mathcal{H}}$$

Therefore  $\langle e_i, e_j \rangle_{\mathcal{H}} = 0$ , if  $i \neq j$ , and  $\langle e_i, e_i \rangle_{\mathcal{H}} = \frac{1}{\lambda_i}$ . On the other hand  $\langle e_i, e_j \rangle_\mu = 0$ , if  $i \neq j$ , and  $\langle e_i, e_i \rangle_\mu = 1$ . This observation establishes a simple relationship between the Hilbert norms in  $\mathcal{H}$  and  $\mathcal{L}_\mu^2$ . We also see that  $f = \sum a_i e_i(x) \in \mathcal{H}$  if and only if  $\sum \frac{a_i^2}{\lambda_i} < \infty$ .

Consider now the operator  $L_K^{1/2}$ . It can be defined as the only positive definite self-adjoint operator, s.t.  $L_K = L_K^{1/2} \circ L_K^{1/2}$ . Assuming that the series  $\tilde{K}(x, y) = \sum_i \sqrt{\lambda_i} e_i(x) e_i(y)$  converges, we can write

$$(L_K^{1/2} f)(x) = \int_X f(y) \tilde{K}(x, y) d\mu$$

It is easy to check that  $L_K^{1/2}$  is an isomorphism between  $\mathcal{H}$  and  $\mathcal{L}_\mu^2$ , that is

$$\forall f, g \in \mathcal{H} \quad \langle f, g \rangle_\mu = \langle L_K^{1/2} f, L_K^{1/2} g \rangle_{\mathcal{H}}$$

Therefore  $\mathcal{H}$  is the image  $L_K^{1/2}$  acting on  $\mathcal{L}_\mu^2$ .

**Lemma 3.1.** *A function  $f(x) = \sum_i a_i e_i(x)$  can be represented as  $f = L_K g$  for some  $g$  if and only if*

$$\sum_{i=1}^{\infty} \frac{a_i^2}{\lambda_i} < \infty \tag{7}$$

*Proof.* Suppose  $f = L_K g$ . Write  $g(x) = \sum_i b_i e_i(x)$ . We know that  $g \in L^2_\mu$  if and only if  $\sum_i b_i^2 < \infty$ . Since  $L_K(\sum_i b_i e_i) = \sum_i b_i \lambda_i e_i = \sum_i a_i e_i$ , we obtain  $a_i = b_i \lambda_i$ . Therefore  $\sum_{i=1}^\infty \frac{a_i^2}{\lambda_i^2} < \infty$ .

Conversely, if the condition in the inequality 7 is satisfied,  $f = L_K g$ , where  $g = \sum \frac{a_i}{\lambda_i} e_i$ .  $\square$

### 3.2 Proof of Theorems

Now let us recall the Eqn. 3:

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f) + \gamma_A \|f\|_K^2 + \gamma_I \|f\|_I^2 \quad (8)$$

We have an RKHS  $\mathcal{H}_K$  and the probability distribution  $\mu$  which is supported on  $\mathcal{M} \subset X$ . We define a linear space  $\mathcal{S}$  to be the closure with respect to the RKHS norm of  $\mathcal{H}_K$  of the linear span of kernels centered at points of  $\mathcal{M}$ :

$$\mathcal{S} = \overline{\operatorname{span}\{K(x, \cdot) \mid x \in \mathcal{M}\}}$$

Before we start we need to introduce some notation. By a subscript  $\mathcal{M}$  we will denote the restriction to  $\mathcal{M}$ . For example, by  $\mathcal{S}_\mathcal{M}$  we denote the restriction of functions from  $\mathcal{S}$  to  $\mathcal{M}$ . It can be shown ([2], p. 350) that the space  $(\mathcal{H}_K)_\mathcal{M}$  of functions from  $\mathcal{H}_K$  restricted to  $\mathcal{M}$  is an RKHS with the kernel  $K_\mathcal{M}$ , in other words  $(\mathcal{H}_K)_\mathcal{M} = \mathcal{H}_{K_\mathcal{M}}$ .

**Lemma 3.2.** *The following properties of  $\mathcal{S}$  hold:*

1.  $\mathcal{S}$  with the inner product induced by  $\mathcal{H}_K$  is a Hilbert space.
2.  $\mathcal{S}_\mathcal{M} = (\mathcal{H}_K)_\mathcal{M}$ .
3. The orthogonal complement  $\mathcal{S}^\perp$  to  $\mathcal{S}$  in  $\mathcal{H}_K$  consists of all functions vanishing on  $\mathcal{M}$ .

*Proof.* 1. From the definition of  $\mathcal{S}$  it is clear by that  $\mathcal{S}$  is a complete subspace of  $\mathcal{H}_K$ .

2. We give a convergence argument similar to one found in [2]. Since  $(\mathcal{H}_K)_\mathcal{M} = \mathcal{H}_{K_\mathcal{M}}$  any function  $f_\mathcal{M}$  in it can be written as  $f_\mathcal{M} = \lim_{n \rightarrow \infty} f_{\mathcal{M},n}$ , where  $f_{\mathcal{M},n} = \sum_i \alpha_{in} K_\mathcal{M}(x_{in}, \cdot)$  is the sum of some kernel functions.

Consider the corresponding sum  $f_n = \sum_i \alpha_{in} K(x_{in}, \cdot)$ . From the definition of the norm we see that  $\|f_n - f_k\|_K = \|f_{\mathcal{M},n} - f_{\mathcal{M},k}\|_{K_\mathcal{M}}$  and therefore  $f_n$  is a Cauchy sequence. Thus  $f = \lim_{n \rightarrow \infty} f_n$  exists and its restriction to  $\mathcal{M}$  must equal

$f_{\mathcal{M}}$ . This shows that  $(\mathcal{H}_K)_{\mathcal{M}} \subset \mathcal{S}_{\mathcal{M}}$ . The other direction follows by the same argument.

3. Let  $g \in \mathcal{S}^{\perp}$ . By the reproducing property for any  $x \in \mathcal{M}$ ,  $g(x) = \langle K(x, \cdot), g(\cdot) \rangle_K = 0$  and therefore any function in  $\mathcal{S}^{\perp}$  vanishes on  $\mathcal{M}$ . On the other hand, if  $g$  vanishes on  $\mathcal{M}$  it is perpendicular to each  $K(x, \cdot)$ ,  $x \in \mathcal{M}$  and is therefore perpendicular to the closure of their span  $\mathcal{S}$ .  $\square$

**Lemma 3.3.** *Assume that the intrinsic norm is such that for any  $f, g \in \mathcal{H}_K$ ,  $(f - g)|_{\mathcal{M}} \equiv 0$  implies that  $\|f\|_I = \|g\|_I$ . Then if the solution  $f^*$ , of the optimization problem in Eqn. 8 exists and belongs to  $\mathcal{S}$ .*

*Proof.* Any  $f \in \mathcal{H}_K$  can be written as  $f = f_{\mathcal{S}} + f_{\mathcal{S}^{\perp}}$ , where  $f_{\mathcal{S}}$  is the projection of  $f$  to  $\mathcal{S}$  and  $f_{\mathcal{S}^{\perp}}$  is its orthogonal complement.

For any  $x \in \mathcal{M}$  we have  $K(x, \cdot) \in \mathcal{S}$ . By the previous Lemma  $f_{\mathcal{S}^{\perp}}$  vanishes on  $\mathcal{M}$ . We have  $f(x_i) = f_{\mathcal{S}}(x_i) \forall_i$  and by assumption  $\|f_{\mathcal{S}}\|_I = \|f\|_I$ .

On the other hand,  $\|f\|_K^2 = \|f_{\mathcal{S}}\|_K^2 + \|f_{\mathcal{S}^{\perp}}\|_K^2$  and therefore  $\|f\|_K \geq \|f_{\mathcal{S}}\|_K$ . It follows that the minimizer  $f^*$  is in  $\mathcal{S}$ .  $\square$

As a direct corollary of these consideration, we obtain the following

**Proposition 3.4.** *If  $\|f\|_I = \|f\|_{K_{\mathcal{M}}}$  then the minimizer of Eqn. 8 is identical to that of the usual regularization problem (Eqn. 1) although with a different regularization parameter.*

We can now restrict our attention to the study of  $\mathcal{S}$ . While it is clear that the right-hand side of Eqn. 4 lies in  $\mathcal{S}$ , not every element in  $\mathcal{S}$  can be written like that. For example,  $K(x, \cdot)$ , where  $x$  is not one of the data points  $x_i$  is generally not of that form.

We will now assume that for  $f \in \mathcal{S}$

$$\|f\|_I^2 = \langle f, Df \rangle_{\mathcal{L}_{\mu}^2}$$

where  $D$  is an appropriate smoothness operator, such as an inverse integral operator or a differential operator, e.g.,  $Df = \text{grad}_{\mathcal{M}} f$ . The Representer theorem, however, holds for even more general  $D$ :

**Theorem 3.5.** *Let  $D$  be a bounded operator  $D : \mathcal{S} \rightarrow \mathcal{L}_{\mathcal{P}_X}^2$ . Then the solution  $f^*$  of the optimization problem in Eqn. 3 exists and can be written*

$$f^*(x) = \sum_{i=1}^l \alpha_i K(x_i, x) + \int_{\mathcal{M}} \alpha(y) K(x, y) d\mathcal{P}_X(y) \quad (9)$$

*Proof.* For simplicity we will assume that the loss function  $V$  is differentiable. This condition can ultimately be eliminated by approximating a non-differentiable function appropriately and passing to the limit.

Put

$$H(f) = \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f(y_i)) + \gamma_A \|f\|_K^2 + \gamma_I \|f\|_I^2 \quad (10)$$

We first show that the solution to Eqn. 3  $f^*$  exists and by Lemma 3.3 belongs to  $\mathcal{S}$ . It follows easily from Cor. 3.7 and standard results about compact embeddings of Sobolev spaces (e.g., [1]) that a ball  $\mathcal{B}_r \subset \mathcal{H}_K$ ,  $\mathcal{B}_r = \{f \in \mathcal{S}, s.t. \|f\|_K \leq r\}$  is compact in  $\mathcal{L}_X^\infty$ . Therefore for any such ball the minimizer in that ball  $f_r^*$  must exist and belong to  $\mathcal{B}_r$ . On the other hand, by substituting the zero function

$$H(f_r^*) \leq H(0) = \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, 0)$$

If the loss is actually zero, then zero function is a solution, otherwise

$$\gamma_A \|f_r^*\|_K^2 < \sum_{i=1}^l V(x_i, y_i, 0)$$

and hence  $f_r^* \in \mathcal{B}_r$ , where

$$r = \sqrt{\frac{\sum_{i=1}^l V(x_i, y_i, 0)}{\gamma_A}}$$

Therefore we cannot decrease  $H(f^*)$  by increasing  $r$  beyond a certain point, which shows that  $f^* = f_r^*$  with  $r$  as above, which completes the proof of existence. If  $V$  is convex, such solution will also be unique.

We proceed to derive the Eqn. 9. As before, let  $e_1, e_2, \dots$  be the basis associated to the integral operator  $(L_K f)(x) = \int_{\mathcal{M}} f(y) K(x, y) d\mathcal{P}_X(y)$ . Write  $f^* = \sum_i a_i e_i(x)$ . By substituting  $f^*$  into  $H(f)$  we obtain:

$$H(f^*) = \frac{1}{l} \sum_{j=1}^l V(x_j, y_j, \sum_i a_i e_i(x_i)) + \gamma_A \|f^*\|_K^2 + \gamma_I \|f^*\|_I^2$$

Assume that  $V$  is differentiable with respect to each  $a_k$ . We have  $\|\sum_i a_i e_i(x)\|_K^2 = \sum_i \frac{a_i^2}{\lambda_i}$ . Differentiating with respect to the coefficients  $a_i$  yields the following set of equations:

$$\frac{\partial H(f^*)}{\partial a_k} = \frac{1}{l} \sum_{j=1}^l e_k(x_j) \partial_3 V(x_j, y_j, \sum_i a_i e_i) + 2\gamma_A \frac{a_k}{\lambda_k} + \gamma_I \langle Df, e_k \rangle + \gamma_I \langle f, De_k \rangle = 0$$

where  $\partial_3 V$  denotes the derivative with respect to the third argument of  $V$ .

$\langle Df, e_k \rangle + \langle f, De_k \rangle = \langle (D + D^*)f, e_k \rangle$  and hence

$$a_k = -\frac{1}{\gamma_A l} \sum_{j=1}^l e_k(x_j) \partial_3 V(x_j, y_j, f^*) - \frac{\gamma_I}{2\gamma_A} \lambda_k \langle Df^* + D^*f^*, e_k \rangle$$

Since  $f^*(x) = \sum_k a_k e_k(x)$  and recalling that  $K(x, y) = \sum_i \lambda_i e_i(x) e_i(y)$

$$\begin{aligned} f^*(x) &= -\frac{1}{\gamma_A l} \sum_k \sum_{j=1}^l \lambda_k e_k(x) e_k(x_j) \partial_3 V(x_j, y_j, f^*) - \frac{\gamma_I}{2\gamma_A} \sum_k \lambda_k \langle Df^* + D^*f^*, e_k \rangle e_k \\ &= -\frac{1}{\gamma_A l} \sum_{j=1}^l K(x, x_j) \partial_3 V(x_j, y_j, f^*) - \frac{\gamma_I}{2\gamma_A} \sum_k \lambda_k \langle Df^* + D^*f^*, e_k \rangle e_k \end{aligned}$$

We see that the first summand is a sum of the kernel functions centered at data points. It remains to show that the second summand has an integral representation, i.e. can be written as  $\int_{\mathcal{M}} \alpha(y) K(x, y) d\mathcal{P}_X(y)$ , which is equivalent to being in the image of  $L_K$ . To verify this we apply Lemma 3.1. We need that

$$\sum_k \frac{\lambda_k^2 \langle Df^* + D^*f^*, e_k \rangle^2}{\lambda_k^2} = \sum_k \langle Df^* + D^*f^*, e_k \rangle^2 < \infty \quad (11)$$

Since  $D$ , its adjoint operator  $D^*$  and hence their sum are bounded the inequality 11 is satisfied for any function in  $\mathcal{S}$ .  $\square$

### 3.3 Manifold Setting<sup>1</sup>

We now show that for the case when  $\mathcal{M}$  is a manifold and  $D$  is a differential operator, such as the Laplace-Beltrami operator  $\mathcal{L}$ , the boundedness condition of Theorem 3.5 is satisfied. While we consider the case when the manifold has no boundary, the same argument goes through for manifold with boundary, with, for example, Dirichet's boundary conditions (vanishing at the boundary). Thus the setting of Theorem 3.5 is very general, applying, among other things, to arbitrary differential operators on compact domains in Euclidean space.

Let  $\mathcal{M}$  be a  $\mathcal{C}^\infty$  manifold without boundary with an infinitely differentiable embedding in some ambient space  $X$ ,  $D$  a differential operator with  $\mathcal{C}^\infty$  coefficients and let  $\mu$ , be the measure corresponding to some  $\mathcal{C}^\infty$  nowhere vanishing volume

---

<sup>1</sup>We thank Peter Constantin and Todd Dupont for help with this section.

form on  $\mathcal{M}$ . We assume that the kernel  $K(x, y)$  is also infinitely differentiable.<sup>2</sup> As before for an operator  $A$ ,  $A^*$  denotes the adjoint operator.

*Proof.* First note that it is enough to show that  $D$  is bounded on  $\mathcal{H}_{K_{\mathcal{M}}}$ , since  $D$  only depends on the restriction  $f_{\mathcal{M}}$ . As before, let  $L_{K_{\mathcal{M}}}(f)(x) = \int_{\mathcal{M}} f(y)K_{\mathcal{M}}(x, y) d\mu$  is the integral operator associated to  $K_{\mathcal{M}}$ . Note that  $D^*$  is also a differential operator of the same degree as  $D$ . The integral operator  $L_{K_{\mathcal{M}}}$  is bounded (compact) from  $L_{\mu}^2$  to any Sobolev space  $H^{sob}$ . Therefore the operator  $L_{K_{\mathcal{M}}}D$  is also bounded. We therefore see that  $DL_{K_{\mathcal{M}}}D^*$  is bounded  $L_{\mu}^2 \rightarrow L_{\mu}^2$ . Therefore there is a constant  $C$ , s.t.  $\langle DL_{K_{\mathcal{M}}}D^*f, f \rangle_{L_{\mu}^2} \leq C\|f\|_{L_{\mu}^2}$ .

The square root  $T = L_{K_{\mathcal{M}}}^{1/2}$  of the self-adjoint positive definite operator  $L_{K_{\mathcal{M}}}$  is a self-adjoint positive definite operator as well. Thus  $(DT)^* = TD^*$ . By definition of the operator norm, for any  $\epsilon > 0$  there exists  $f \in L_{\mu}^2$ ,  $\|f\|_{L_{\mu}^2} \leq 1 + \epsilon$ , such that

$$\begin{aligned} \|DT\|_{L_{\mu}^2}^2 &= \|TD^*\|_{L_{\mu}^2}^2 \leq \langle TD^*f, TD^*f \rangle_{L_{\mu}^2} = \\ &= \langle DLD^*f, f \rangle_{L_{\mu}^2} \leq \|DLD^*\|_{L_{\mu}^2} \|f\|_{L_{\mu}^2}^2 \leq C(1 + \epsilon)^2 \end{aligned}$$

Therefore the operator  $DT : L_{\mu}^2 \rightarrow L_{\mu}^2$  is bounded (and also  $\|DT\|_{L_{\mu}^2} \leq C$ , since  $\epsilon$  is arbitrary).

Now recall that  $T$  provides an isometry between  $L_{\mu}^2$  and  $\mathcal{H}_{K_{\mathcal{M}}}$ . That means that for any  $g \in \mathcal{H}_{K_{\mathcal{M}}}$  there is  $f \in L_{\mu}^2$ , such that  $Tf = g$  and  $\|f\|_{L_{\mu}^2} = \|g\|_{\mathcal{H}_{K_{\mathcal{M}}}}$ . Thus  $\|Dg\|_{L_{\mu}^2} = \|DTf\|_{L_{\mu}^2} \leq C\|g\|_{\mathcal{H}_{K_{\mathcal{M}}}}$ , which shows that  $T : \mathcal{H}_{K_{\mathcal{M}}} \rightarrow L_{\mu}^2$  is bounded and concludes the proof.  $\square$

Since  $\mathcal{S}$  is a subspace of  $\mathcal{H}_K$  the main result follows immediately:

**Corollary 3.6.**  *$D$  is a bounded operator  $\mathcal{S} \rightarrow L_{\mu}^2$  and the conditions of Theorem 3.5 hold.*

Before finishing the theoretical discussion we obtain a useful

**Corollary 3.7.** *The operator  $T = L_K^{1/2}$  on  $L_{\mu}^2$  is a bounded (and in fact compact) operator  $L_{\mu}^2 \rightarrow H^{sob}$ , where  $H^{sob}$  is an arbitrary Sobolev space.*

*Proof.* Follows from the fact that  $DT$  is bounded operator  $L_{\mu}^2 \rightarrow L_{\mu}^2$  for an arbitrary differential operator  $D$  and standard results on compact embeddings of Sobolev spaces (see, e.g. [1]).  $\square$

---

<sup>2</sup>While we have assumed that all objects are infinitely differentiable, it is not hard to specify the precise differentiability conditions. Roughly speaking, a degree  $k$  differential operator  $D$  is bounded as an operator  $\mathcal{H}_K \rightarrow L_{\mu}^2$ , if the kernel  $K(x, y)$  has  $2k$  derivatives.

### 3.4 The Representer Theorem for the Empirical Case

In the case when  $\mathcal{M}$  is unknown and sampled via labeled and unlabeled examples, the Laplace-Beltrami operator on  $\mathcal{M}$  may be approximated by the Laplacian of the data adjacency graph (see [3, 7] for some discussion). A regularizer based on the graph Laplacian leads to the optimization problem posed in Eqn. 5. We now provide a proof of Theorem 2.2 which states that the solution to this problem admits a representation in terms of an expansion over labeled and unlabeled points. The proof is based on a simple orthogonality argument (e.g., [31]).

*Proof. (Theorem 2.2)* Any function  $f \in \mathcal{H}_K$  can be uniquely decomposed into a component  $f_{\parallel}$  in the linear subspace spanned by the kernel functions  $\{K(x_i, \cdot)\}_{i=1}^{l+u}$ , and a component  $f_{\perp}$  orthogonal to it. Thus,

$$f = f_{\parallel} + f_{\perp} = \sum_{i=1}^{l+u} \alpha_i K(x_i, \cdot) + f_{\perp}$$

By the reproducing property, as the following arguments show, the evaluation of  $f$  on any data point  $x_j$ ,  $1 \leq j \leq l+u$  is independent of the orthogonal component  $f_{\perp}$ :

$$f(x_j) = \langle f, K(x_j, \cdot) \rangle = \left\langle \sum_{i=1}^{l+u} \alpha_i K(x_i, \cdot), K(x_j, \cdot) \right\rangle + \langle f_{\perp}, K(x_j, \cdot) \rangle$$

Since the second term vanishes, and  $\langle K(x_i, \cdot), K(x_j, \cdot) \rangle = K(x_i, x_j)$ , it follows that  $f(x_j) = \sum_{i=1}^{l+u} \alpha_i K(x_i, x_j)$ . Thus, the empirical terms involving the loss function and the intrinsic norm in the optimization problem in Eqn. 5 depend only on the value of the coefficients  $\{\alpha_i\}_{i=1}^{l+u}$  and the gram matrix of the kernel function.

Indeed, since the orthogonal component only increases the norm of  $f$  in  $\mathcal{H}_K$ :  $\|f\|_K^2 = \|\sum_{i=1}^{l+u} \alpha_i K(x_i, \cdot)\|_K^2 + \|f_{\perp}\|_H^2 \geq \|\sum_{i=1}^{l+u} \alpha_i K(x_i, \cdot)\|_K^2$ , it follows that the minimizer of problem 5 must have  $f_{\perp} = 0$ , and therefore admits a representation  $f^*(\cdot) = \sum_{i=1}^{l+u} \alpha_i K(x_i, \cdot)$ .  $\square$

The simple form of the minimizer, given by this theorem, allows us to translate our extrinsic and intrinsic regularization framework into optimization problems over the finite dimensional space of coefficients  $\{\alpha_i\}_{i=1}^{l+u}$ , and invoke the machinery of kernel based algorithms. In the next section, we derive these algorithms, and explore their connections to other related work.

## 4 Algorithms

We now discuss standard regularization algorithms (RLS and SVM) and present their extensions (LapRLS and LapSVM respectively). These are obtained by solving the optimization problems posed in Eqn. (5) for different choices of cost function  $V$  and regularization parameters  $\gamma_A, \gamma_I$ . To fix notation, we assume we have  $l$  labeled examples  $\{(x_i, y_i)\}_{i=1}^l$  and  $u$  unlabeled examples  $\{x_j\}_{j=l+1}^{j=l+u}$ . We use  $K$  interchangeably to denote the kernel function or the Gram matrix.

### 4.1 Regularized Least Squares

The Regularized Least Squares algorithm is a fully supervised method where we solve :

$$\min_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l (y_i - f(x_i))^2 + \gamma \|f\|_K^2 \quad (12)$$

The classical Representer Theorem can be used to show that the solution is of the following form:

$$f^*(x) = \sum_{i=1}^l \alpha_i^* K(x, x_i) \quad (13)$$

Substituting this form in the problem above, we arrive at following convex differentiable objective function of the  $l$ -dimensional variable  $\alpha = [\alpha_1 \dots \alpha_l]^T$ :

$$\alpha^* = \operatorname{argmin} \frac{1}{l} (Y - K\alpha)^T (Y - K\alpha) + \gamma \alpha^T K \alpha \quad (14)$$

where  $K$  is the  $l \times l$  gram matrix  $K_{ij} = K(x_i, x_j)$  and  $Y$  is the label vector  $Y = [y_1 \dots y_l]^T$ .

The derivative of the objective function vanishes at the minimizer :

$$\frac{1}{l} (Y - K\alpha^*)^T (-K) + \gamma K \alpha^* = 0$$

which leads to the following solution.

$$\alpha^* = (K + \gamma l I)^{-1} Y \quad (15)$$

### 4.2 Laplacian Regularized Least Squares (LapRLS)

The Laplacian Regularized Least Squares algorithm solves problem (5) with the squared loss function:

$$\min_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l (y_i - f(x_i))^2 + \gamma_A \|f\|_K^2 + \frac{\gamma_I}{(u+l)^2} \mathbf{f}^T L \mathbf{f}$$

As before, the Representer Theorem can be used to show that the solution is an expansion of kernel functions over both the labeled and the unlabeled data :

$$f^*(x) = \sum_{i=1}^{l+u} \alpha_i^* K(x, x_i) \quad (16)$$

Substituting this form in the problem above, as before, we arrive at a convex differentiable objective function of the  $l+u$ -dimensional variable  $\alpha = [\alpha_1 \dots \alpha_{l+u}]^T$ :

$$\alpha^* = \operatorname{argmin}_{\alpha \in \mathbb{R}^{l+u}} \frac{1}{l} (Y - JK\alpha)^T (Y - JK\alpha) + \gamma_A \alpha^T K \alpha + \frac{\gamma_I}{(u+l)^2} \alpha^T K L K \alpha$$

where  $K$  is the  $(l+u) \times (l+u)$  Gram matrix over labeled and unlabeled points;  $Y$  is an  $(l+u)$  dimensional label vector given by :  $Y = [y_1, \dots, y_l, 0, \dots, 0]$  and  $J$  is an  $(l+u) \times (l+u)$  diagonal matrix given by  $J = \operatorname{diag}(1, \dots, 1, 0, \dots, 0)$  with the first  $l$  diagonal entries as 1 and the rest 0.

The derivative of the objective function vanishes at the minimizer :

$$\frac{1}{l} (Y - JK\alpha)^T (-JK) + (\gamma_A K + \frac{\gamma_I l}{(u+l)^2} K L K) \alpha = 0$$

which leads to the following solution.

$$\alpha^* = (JK + \gamma_A l I + \frac{\gamma_I l}{(u+l)^2} L K)^{-1} Y \quad (17)$$

Note that when  $\gamma_I = 0$ , Eqn. (17) gives zero coefficients over unlabeled data, and the coefficients over the labeled data are exactly those for standard RLS.

### 4.3 Support Vector Machine Classification

Here we outline the SVM approach to binary classification problems. For SVMs, the following problem is solved :

$$\min_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l (1 - y_i f(x_i))_+ + \gamma \|f\|_K^2$$

where the hinge loss is defined as:  $(1 - yf(x))_+ = \max(0, 1 - yf(x))$  and the labels  $y_i \in \{-1, +1\}$ .

Again, the solution is given by:

$$f^*(x) = \sum_{i=1}^l \alpha_i^* K(x, x_i) \quad (18)$$

Following SVM expositions, the above problem can be equivalently written as:

$$\begin{aligned} \min_{f \in \mathcal{H}_K, \xi_i \in \mathbb{R}} & \frac{1}{l} \sum_{i=1}^l \xi_i + \gamma \|f\|_K^2 \\ \text{subject to : } & y_i f(x_i) \geq 1 - \xi_i \quad i = 1, \dots, l \\ & \xi_i \geq 0 \quad i = 1, \dots, l \end{aligned} \quad (19)$$

Using the Lagrange multipliers technique, and benefiting from strong duality, the above problem has a simpler quadratic dual program in the Lagrange multipliers  $\beta = [\beta_1, \dots, \beta_l]^T \in \mathbb{R}^l$ :

$$\begin{aligned} \beta^* &= \max_{\beta \in \mathbb{R}^l} \sum_{i=1}^l \beta_i - \frac{1}{2} \beta^T Q \beta \\ \text{subject to : } & \sum_{i=1}^l y_i \beta_i = 0 \\ & 0 \leq \beta_i \leq \frac{1}{l} \quad i = 1, \dots, l \end{aligned} \quad (20)$$

where the equality constraint arises due to an unregularized bias term that is often added to the sum in Eqn (18), and the following notation is used :

$$\begin{aligned} Y &= \text{diag}(y_1, y_2, \dots, y_l) \\ Q &= Y \left( \frac{K}{2\gamma} \right) Y \\ \alpha^* &= \frac{Y \beta^*}{2\gamma} \end{aligned} \quad (21)$$

Here again,  $K$  is the gram matrix over labeled points. SVM practitioners may be familiar with a slightly different parameterization involving the  $C$  parameter :  $C = \frac{1}{2\gamma l}$  is the weight on the hinge loss term (instead of using a weight  $\gamma$  on the norm term in the optimization problem). The  $C$  parameter appears as the upper bound (instead of  $\frac{1}{l}$ ) on the values of  $\beta$  in the quadratic program. For additional details on the derivation and alternative formulations of SVMs, see [31], [26].

#### 4.4 Laplacian Support Vector Machines

By including the intrinsic smoothness penalty term, we can extend SVMs by solving the following problem:

$$\min_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l (1 - y_i f(x_i))_+ + \gamma_A \|f\|_K^2 + \frac{\gamma_I}{(u+l)^2} \mathbf{f}^T L \mathbf{f}$$

By the representer theorem, as before, the solution to the problem above is given by:

$$f^*(x) = \sum_{i=1}^{l+u} \alpha_i^* K(x, x_i)$$

Often in SVM formulations, an unregularized bias term  $b$  is added to the above form. Again, the primal problem can be easily seen to be the following:

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^{l+u}, \xi \in \mathbb{R}^l} \quad & \frac{1}{l} \sum_{i=1}^l \xi_i + \gamma_A \alpha^T K \alpha + \frac{\gamma_I}{(u+l)^2} \alpha^T K L K \alpha \\ \text{subject to: } \quad & y_i \left( \sum_{j=1}^{l+u} \alpha_j K(x_i, x_j) + b \right) \geq 1 - \xi_i, \quad i = 1, \dots, l \\ & \xi_i \geq 0 \quad i = 1, \dots, l \end{aligned}$$

Introducing the Lagrangian, with  $\beta_i, \zeta_i$  as Lagrange multipliers:

$$\begin{aligned} L(\alpha, \xi, b, \beta, \zeta) = \quad & \frac{1}{l} \sum_{i=1}^l \xi_i + \frac{1}{2} \alpha^T \left( 2\gamma_A K + 2 \frac{\gamma_I}{(l+u)^2} K L K \right) \alpha \\ & - \sum_{i=1}^l \beta_i \left( y_i \left( \sum_{j=1}^{l+u} \alpha_j K(x_i, x_j) + b \right) - 1 + \xi_i \right) - \sum_{i=1}^l \zeta_i \xi_i \end{aligned}$$

Passing to the dual requires the following steps:

$$\begin{aligned} \frac{\partial L}{\partial b} = 0 \quad & \implies \sum_{i=1}^l \beta_i y_i = 0 \\ \frac{\partial L}{\partial \xi_i} = 0 \quad & \implies \frac{1}{l} - \beta_i - \zeta_i = 0 \\ & \implies 0 \leq \beta_i \leq \frac{1}{l} \quad (\xi_i, \zeta_i \text{ are non-negative}) \end{aligned}$$

Using above identities, we formulate a reduced Lagrangian:

$$\begin{aligned}
L^R(\alpha, \beta) &= \frac{1}{2}\alpha^T(2\gamma_A K + 2\frac{\gamma I}{(u+l)^2}K L K)\alpha - \sum_{i=1}^l \beta_i(y_i \sum_{j=1}^{l+u} \alpha_j K(x_i, x_j) - 1) \\
&= \frac{1}{2}\alpha^T(2\gamma_A K + 2\frac{\gamma I}{(u+l)^2}K L K)\alpha - \alpha^T K J^T Y \beta + \sum_{i=1}^l \beta_i \quad (22)
\end{aligned}$$

where  $J = [I \ 0]$  is an  $l \times (l+u)$  matrix with  $I$  as the  $l \times l$  identity matrix (assuming the first  $l$  points are labeled) and  $Y = \text{diag}(y_1, y_2, \dots, y_l)$ .

Taking derivative of the reduced Lagrangian wrt  $\alpha$ :

$$\frac{\partial L^R}{\partial \alpha} = (2\gamma_A K + 2\frac{\gamma I}{(u+l)^2}K L K)\alpha - K J^T Y \beta$$

This implies:

$$\alpha = (2\gamma_A I + 2\frac{\gamma I}{(u+l)^2}L K)^{-1} J^T Y \beta^* \quad (23)$$

Note that the relationship between  $\alpha$  and  $\beta$  is no longer as simple as the SVM algorithm. In particular, the  $(l+u)$  expansion coefficients are obtained by solving a linear system involving the  $l$  dual variables that will appear in the SVM dual problem.

Substituting back in the reduced Lagrangian we get:

$$\beta^* = \max_{\beta \in \mathbb{R}^l} \sum_{i=1}^l \beta_i - \frac{1}{2}\beta^T Q \beta \quad (24)$$

$$\begin{aligned}
\text{subject to :} \quad & \sum_{i=1}^l \beta_i y_i = 0 \\
& 0 \leq \beta_i \leq \frac{1}{l} \quad i = 1, \dots, l \quad (25)
\end{aligned}$$

where

$$Q = Y J K (2\gamma_A I + 2\frac{\gamma I}{(l+u)^2}L K)^{-1} J^T Y$$

Laplacian SVMs can be implemented by using a standard SVM solver with the quadratic form induced by the above matrix, and using the solution to obtain the expansion coefficients by solving the linear system in Eqn (23).

Note that when  $\gamma_I = 0$ , the SVM QP and Eqns (24,23), give zero expansion coefficients over the unlabeled data. The expansion coefficients over the labeled data and the Q matrix are as in standard SVM, in this case.

The Manifold Regularization algorithms are summarized in the Table 1.

Table 1: A Summary of the algorithms

<i>Manifold Regularization algorithms</i>	
<b>Input:</b>	$l$ labeled examples $\{(x_i, y_i)\}_{i=1}^l$ , $u$ unlabeled examples $\{x_j\}_{j=l+1}^{l+u}$
<b>Output:</b>	Estimated function $f : \mathbb{R}^n \rightarrow \mathbb{R}$
<b>Step 1</b>	► Construct data adjacency graph with $(l + u)$ nodes using, e.g, $k$ nearest neighbors or a graph kernel. Choose edge weights $W_{ij}$ , e.g. binary weights or heat kernel weights $W_{ij} = e^{-\ x_i - x_j\ ^2/4t}$ .
<b>Step 2</b>	► Choose a kernel function $K(x, y)$ . Compute the Gram matrix $K_{ij} = K(x_i, x_j)$ .
<b>Step 3</b>	► Compute graph Laplacian matrix : $L = D - W$ where $D$ is a diagonal matrix given by $D_{ii} = \sum_{j=1}^{l+u} W_{ij}$ .
<b>Step 4</b>	► Choose $\gamma_A$ and $\gamma_I$ .
<b>Step 5</b>	► Compute $\alpha^*$ using Eqn. (17) for squared loss (Laplacian RLS) or using Eqn.s (24,23) together with the SVM QP solver for soft margin loss (Laplacian SVM).
<b>Step 6</b>	► Output function $f^*(x) = \sum_{i=1}^{l+u} \alpha_i^* K(x_i, x)$ .

## 4.5 Related Work and Connections to Other Algorithms

In this section we survey various approaches to semi-supervised and transductive learning and highlight connections of Manifold Regularization to other algorithms.

**Transductive SVM (TSVM)** [35], [22]: TSVMs are based on the following optimization principle :

$$f^* = \underset{f \in \mathcal{H}_K, y_{l+1}, \dots, y_{l+u}}{\operatorname{argmin}} \quad C \sum_{i=0}^l (1 - y_i f(x_i))_+ + C^* \sum_{i=l+1}^{l+u} (1 - y_i f(x_i))_+ + \|f\|_K^2$$

which proposes a joint optimization of the SVM objective function over binary-valued labels on the unlabeled data and functions in the RKHS. Here,  $C, C^*$  are parameters that control the relative hinge-loss over labeled and unlabeled sets. The

joint optimization is implemented in [22] by first using an inductive SVM to label the unlabeled data and then iteratively solving SVM quadratic programs, at each step switching labels to improve the objective function. However this procedure is susceptible to local minima and requires an unknown, possibly very large number of label switches before converging. Note that even though TSVM were inspired by transductive inference, they do provide an out-of-sample extension.

**Semi-Supervised SVMs (S<sup>3</sup>VM)** [10], [21] : S<sup>3</sup>VM incorporate unlabeled data by including the minimum hinge-loss for the two choices of labels for each unlabeled example. This is formulated as a mixed-integer program for linear SVMs in [10] and is found to be intractable for large amounts of unlabeled data. [21] reformulate this approach as a concave minimization problem which is solved by a successive linear approximation algorithm. The presentation of these algorithms is restricted to the linear case.

**Measure-Based Regularization** [7]: The conceptual framework of this work is closest to our approach. The authors consider a gradient based regularizer that penalizes variations of the function more in high density regions and less in low density regions leading to the following optimization principle:

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=0}^l V(f(x_i), y_i) + \gamma \int_X \langle \nabla f(x), \nabla f(x) \rangle p(x) dx$$

where  $p$  is the density of the marginal distribution  $\mathcal{P}_X$ . The authors observe that it is not straightforward to find a kernel for arbitrary densities  $p$ , whose associated RKHS norm is  $\int \langle \nabla f(x), \nabla f(x) \rangle p(x) dx$ . Thus, in the absence of a representer theorem, the authors propose to perform minimization of the regularized loss on a fixed set of basis functions chosen apriori, i.e,  $\mathcal{F} = \{ \sum_{i=1}^q \alpha_i \phi_i \}$ . For the hinge loss, this paper derives an SVM quadratic program in the coefficients  $\{ \alpha_i \}_{i=1}^q$  whose  $Q$  matrix is calculated by computing  $q^2$  integrals over gradients of the basis functions. However the algorithm does not demonstrate performance improvements in real world experiments. It is also worth noting that while [7] use the gradient  $\nabla f(x)$  in the ambient space, we use the gradient over a submanifold  $\operatorname{grad}_{\mathcal{M}} f$  for penalizing the function. In a situation where the data truly lies on or near a submanifold  $\mathcal{M}$ , the difference between these two penalizers can be significant since smoothness in the normal direction to the data manifold is irrelevant to classification or regression.

**Graph Based Approaches** See e.g., [12, 14, 32, 37, 38, 24, 23, 4]: A variety of graph based methods have been proposed for transductive inference. However, these methods do not provide an out-of-sample extension. In [38], nearest neighbor labeling for test examples is proposed once unlabeled examples have been labeled by transductive learning. In [14], test points are approximately represented as a

linear combination of training and unlabeled points in the feature space induced by the kernel. Manifold regularization provides natural out-of-sample extensions to several graph based approaches. These connections are summarized in Table 2. We also note the very recent work [8] on out-of-sample extensions for semi-supervised learning. For Graph Regularization and Label Propagation see [29, 6, 38].

**Cotraining** [13] The Co-training algorithm was developed to integrate abundance of unlabeled data with availability of multiple sources of information in domains like web-page classification. Weak learners are trained on labeled examples and their predictions on subsets of unlabeled examples are used to mutually expand the training set. Note that this setting may not be applicable in several cases of practical interest where one does not have access to multiple information sources.

**Bayesian Techniques** See e.g., [25, 28, 16]. An early application of semi-supervised learning to Text classification appeared in [25] where a combination of EM algorithm and Naive-Bayes classification is proposed to incorporate unlabeled data. [28] provides a detailed overview of Bayesian frameworks for semi-supervised learning. The recent work in [16] formulates a new information-theoretic principle to develop a regularizer for conditional log-likelihood.

Table 2: Connections of Manifold Regularization to other algorithms

Parameters	Corresponding algorithms (square loss or hinge loss)
$\gamma_A \geq 0 \quad \gamma_I \geq 0$	Manifold Regularization
$\gamma_A \geq 0 \quad \gamma_I = 0$	Standard Regularization (RLS or SVM)
$\gamma_A \rightarrow 0 \quad \gamma_I > 0$	Out-of-sample extension for Graph Regularization (RLS or SVM)
$\gamma_A \rightarrow 0 \quad \gamma_I \rightarrow 0$ $\gamma_I \gg \gamma_A$	Out-of-sample extension for Label Propagation (RLS or SVM)
$\gamma_A \rightarrow 0 \quad \gamma_I = 0$	Hard margin SVM or Interpolated RLS

## 5 Experiments

We performed experiments on a synthetic dataset and three real world classification problems arising in visual and speech recognition, and text categorization. Comparisons are made with inductive methods (SVM, RLS). Based on a survey of related approaches, as summarized in Section 4.5, we chose to also compare Laplacian SVM with Transductive SVM. Other approaches lack out-of-sample extension, use different base-classifiers or paradigms, or are implementationally not preferable. All software and datasets used for these experiments will be made

available at:

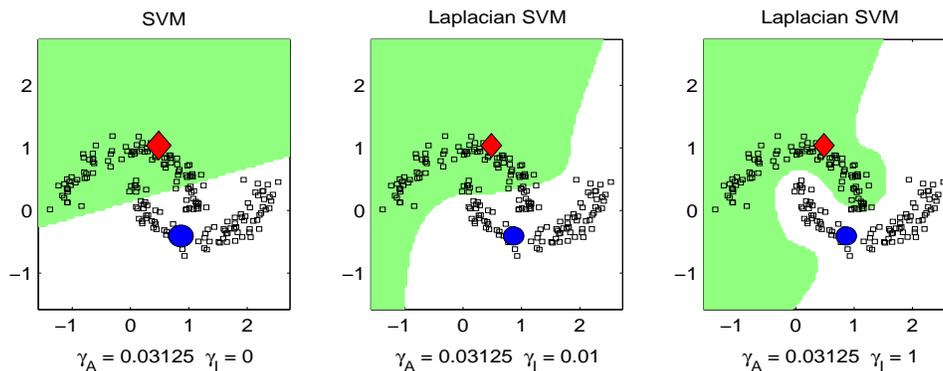
[http://manifold.cs.uchicago.edu/manifold\\_regularization/manifold.html](http://manifold.cs.uchicago.edu/manifold_regularization/manifold.html).

## 5.1 Synthetic Data : Two Moons Dataset

The two moons dataset is shown in Figure 2. The dataset contains 200 examples with only 1 labeled example for each class. Also shown are the decision surfaces of Laplacian SVM for increasing values of the intrinsic regularization parameter  $\gamma_I$ . When  $\gamma_I = 0$ , Laplacian SVM disregards unlabeled data and returns the SVM decision boundary which is fixed by the location of the two labeled points. As  $\gamma_I$  is increased, the intrinsic regularizer incorporates unlabeled data and causes the decision surface to appropriately adjust according to the geometry of the two classes.

In Figure 3, the best decision surfaces across a wide range of parameter settings are also shown for SVM, Transductive SVM and Laplacian SVM. Figure 3 demonstrates how TSVM fails to find the optimal solution. The Laplacian SVM decision boundary seems to be intuitively most satisfying.

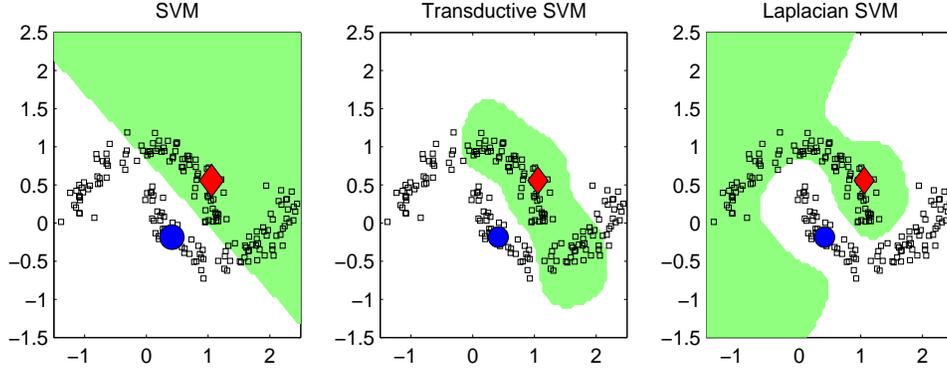
Figure 2: Laplacian SVM with RBF Kernels for various values of  $\gamma_I$ . Labeled points are shown in color, other points are unlabeled.



## 5.2 Handwritten Digit Recognition

In this set of experiments we applied Laplacian SVM and Laplacian RLSC algorithms to 45 binary classification problems that arise in pairwise classification of handwritten digits. The first 400 images for each digit in the USPS training set (preprocessed using PCA to 100 dimensions) were taken to form the training

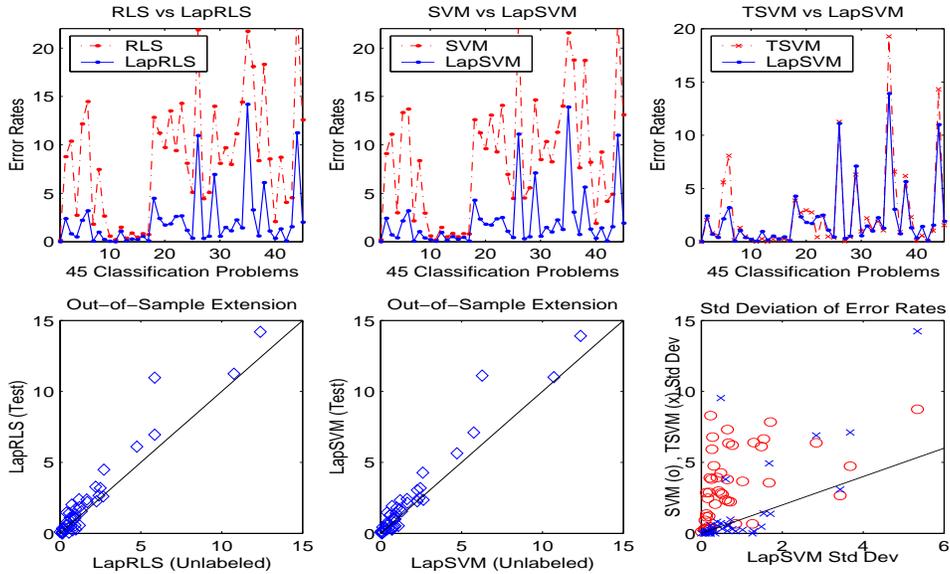
Figure 3: Two Moons Dataset: Best decision surfaces using RBF kernels for SVM, TSVM and Laplacian SVM. Labeled points are shown in color, other points are unlabeled.



set. The remaining images formed the test set. 2 images for each class were randomly labeled ( $l=2$ ) and the rest were left unlabeled ( $u=398$ ). Following [30], we chose to train classifiers with polynomial kernels of degree 3, and set the weight on the regularization term for inductive methods as  $\gamma l = 0.05 (C = 10)$ . For manifold regularization, we chose to split the same weight in the ratio 1 : 9 so that  $\gamma_A l = 0.005$ ,  $\frac{\gamma l}{(u+l)^2} = 0.045$ . The observations reported in this section hold consistently across a wide choice of parameters.

In Figure 4, we compare the error rates of manifold regularization algorithms, inductive classifiers and TSVM, at the precision-recall breakeven points in the ROC curves for the 45 binary classification problems. These results are averaged over 10 random choices of labeled examples. The following comments can be made: (a) Manifold regularization results in significant improvements over inductive classification, for both RLS and SVM, and either compares well or significantly outperforms TSVM across the 45 classification problems. Note that TSVM solves multiple quadratic programs in the size of the labeled and unlabeled sets whereas LapSVM solves a single QP in the size of the labeled set, followed by a linear system. This resulted in substantially faster training times for LapSVM in this experiment. (b) Scatter plots of performance on test and unlabeled data sets confirm that the out-of-sample extension is good for both LapRLS and LapSVM. (c) As shown in the scatter plot in Figure 4 on standard deviation of error rates, we found Laplacian algorithms to be significantly more stable than the inductive methods and TSVM, with respect to choice of the labeled data

Figure 4: USPS Experiment - Error Rates at Precision-Recall Breakeven points for 45 binary classification problems

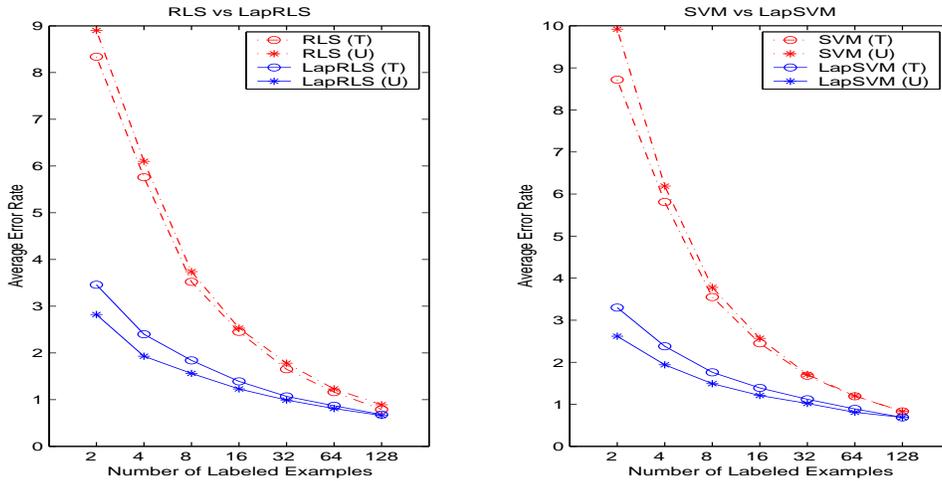


In Figure 5, we demonstrate the benefit of unlabeled data as a function of the number of labeled examples.

### 5.3 Spoken Letter Recognition

This experiment was performed on the Isolet database of letters of the English alphabet spoken in isolation (available from the UCI machine learning repository). The data set contains utterances of 150 subjects who spoke the name of each letter of the English alphabet twice. The speakers are grouped into 5 sets of 30 speakers each, referred to as isolet1 through isolet5. For the purposes of this experiment, we chose to train on the first 30 speakers (isolet1) forming a training set of 1560 examples, and test on isolet5 containing 1559 examples (1 utterance is missing in the database due to poor recording). We considered the task of classifying the first 13 letters of the English alphabet from the last 13. The experimental set-up is meant to simulate a real-world situation: we considered 30 binary classification problems corresponding to 30 splits of the training data where all 52 utterances of one speaker were labeled and all the rest were left unlabeled. The test set is composed of entirely new speakers, forming the separate group isolet5.

Figure 5: USPS Experiment - Mean Error Rate at Precision-Recall Breakeven points as a function of number of labeled points (T: Test Set, U: Unlabeled Set)



We chose to train with RBF kernels of width  $\sigma = 10$  (this was the best value among several settings with respect to 5-fold cross-validation error rates for the fully supervised problem using standard SVM). For SVM and RLSC we set  $\gamma l = 0.05$  ( $C = 10$ ) (this was the best value among several settings with respect to mean error rates over the 30 splits). For Laplacian RLS and Laplacian SVM we set  $\gamma_{Al} = \frac{\gamma l}{(u+l)^2} = 0.005$ . In Figure 6, we compare these algorithms. The following comments can be made: (a) LapSVM and LapRLS make significant performance improvements over inductive methods and TSVM, for predictions on unlabeled speakers that come from the same group as the labeled speaker, over all choices of the labeled speaker. (b) On Isolet5 which comprises of a separate group of speakers, performance improvements are smaller but consistent over the choice of the labeled speaker. This can be expected since there appears to be a systematic bias that affects all algorithms, in favor of same-group speakers. To test this hypothesis, we performed another experiment in which the training and test utterances are both drawn from Isolet1. Here, the second utterance of each letter for each of the 30 speakers in Isolet1 was taken away to form the test set containing 780 examples. The training set consisted of the first utterances for each letter. As before, we considered 30 binary classification problems arising when all utterances of one speaker are labeled and other training speakers are left unlabeled. The scatter plots in Figure 7 confirm our hypothesis, and show high correlation between in-sample and out-of-sample performance of our algorithms in this experiment. It is encourag-

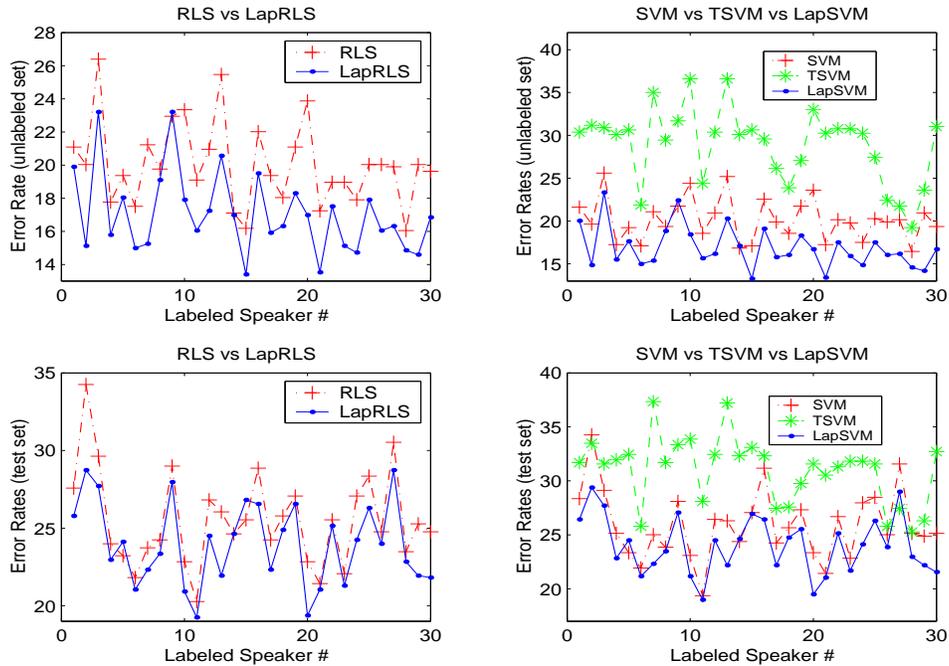


Figure 6: Isolet Experiment - Error Rates at precision-recall breakeven points of 30 binary classification problems

ing to note performance improvements with unlabeled data in Experiment 1 where the test data comes from a slightly different distribution. This robustness is often desirable in real-world applications.

#### 5.4 Text Categorization

We performed Text Categorization experiments on the WebKB dataset which consists of 1051 web pages collected from Computer Science department web-sites of various universities. The task is to classify these webpages into two categories: *course* or *non-course*. We considered learning classifiers using only textual content of the webpages, ignoring link information. A bag-of-words vector space representation for documents is built using the the top 3000 words (skipping HTML headers) having highest mutual information with the class variable, followed by TFIDF mapping. Feature vectors are normalized to unit length. 9 documents were found to contain none of these words and were removed from the dataset.

For the first experiment, we ran LapRLS and LapSVM in a transductive set-

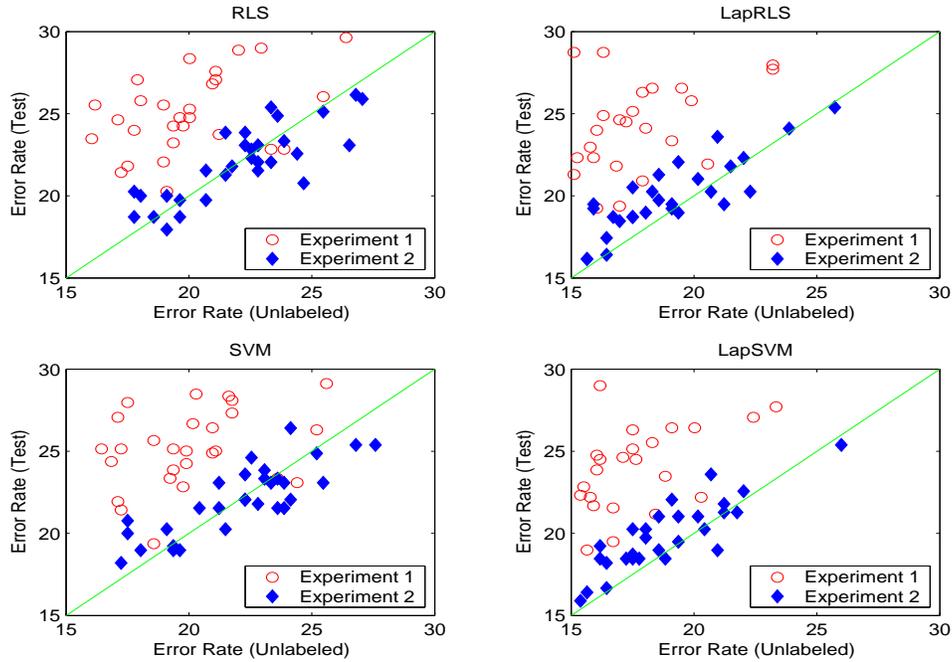


Figure 7: Isolet Experiment - Error Rates at precision-recall breakeven points on Test set Versus Unlabeled Set. In Experiment 1, the training data comes from Isolet 1 and the test data comes from Isolet5; in Experiment 2, both training and test sets come from Isolet1.

ting, with 12 randomly labeled examples (3 course and 9 non-course) and the rest unlabeled. In Table 1, we report the precision and error rates at the precision-recall breakeven point averaged over 100 realizations of the data, and include results reported in [23] for Spectral Graph Transduction, and the Cotraining algorithm [12] for comparison. We used 15 nearest neighbor graphs, weighted by cosine distances and used iterated laplacians of degree 3. For inductive methods,  $\gamma_A l$  was set to 0.01 for RLS and 1.00 for SVM. For LapRLS and LapSVM,  $\gamma_A$  was set as in inductive methods, with  $\frac{\gamma l}{(l+u)^2} = 100\gamma_A l$ . These parameters were chosen based on a simple grid search for best performance over the first 5 realizations of the data. Linear Kernels and cosine distances were used since these have found wide-spread applications in text classification problems [18].

Since the exact datasets on which these algorithms were run, somewhat differ in preprocessing, preparation and experimental protocol, these results are only meant to suggest that Manifold Regularization algorithms perform similar to state-

Table 3: Precision and Error Rates at the Precision-Recall Breakeven Points of supervised and transductive algorithms.

Method	PRBEP	Error
k-NN [23]	73.2	13.3
SGT [23]	86.2	6.2
Naive-Bayes [13]	—	12.9
Cotraining [13]	—	6.20
SVM	76.39 (5.6)	10.41 (2.5)
TSVM <sup>3</sup>	88.15 (1.0)	5.22 (0.5)
LapSVM	87.73 (2.3)	5.41 (1.0)
RLS	73.49 (6.2)	11.68 (2.7)
LapRLS	86.37 (3.1)	5.99 (1.4)

of-the-art methods for transductive inference in text classification problems. The following comments can be made: (a) Transductive categorization with LapSVM and LapRLS leads to significant improvements over inductive categorization with SVM and RLS. (b) [23] reports 91.4% precision-recall breakeven point, and 4.6% error rate for TSVM. Results for TSVM reported in the table were obtained when we ran the TSVM implementation using SVM-Light software on this particular dataset. The average training time for TSVM was found to be more than 10 times slower than for LapSVM (c) The Co-training results were obtained on unseen test datasets utilizing additional hyperlink information, which was excluded in our experiments. This additional information is known to improve performance, as demonstrated in [23] and [13].

In the next experiment, we randomly split the webkb data into a test set of 263 examples and a training set of 779 examples. We noted the performance of inductive and semi-supervised classifiers on unlabeled and test sets as a function of the number of labeled examples in the training set. The performance measure is the precision-recall breakeven point (PRBEP), averaged over 100 random data splits. Results are presented in the top panel of Figure 8. The benefit of unlabeled data can be seen by comparing the performance curves of inductive and semi-supervised classifiers.

We also performed experiments with different sizes of the training set, keeping a randomly chosen test set of 263 examples. The bottom panel in Figure 8 presents the quality of transduction and semi-supervised learning with Laplacian SVM (Laplacian RLS performed similarly) as a function of the number of labeled

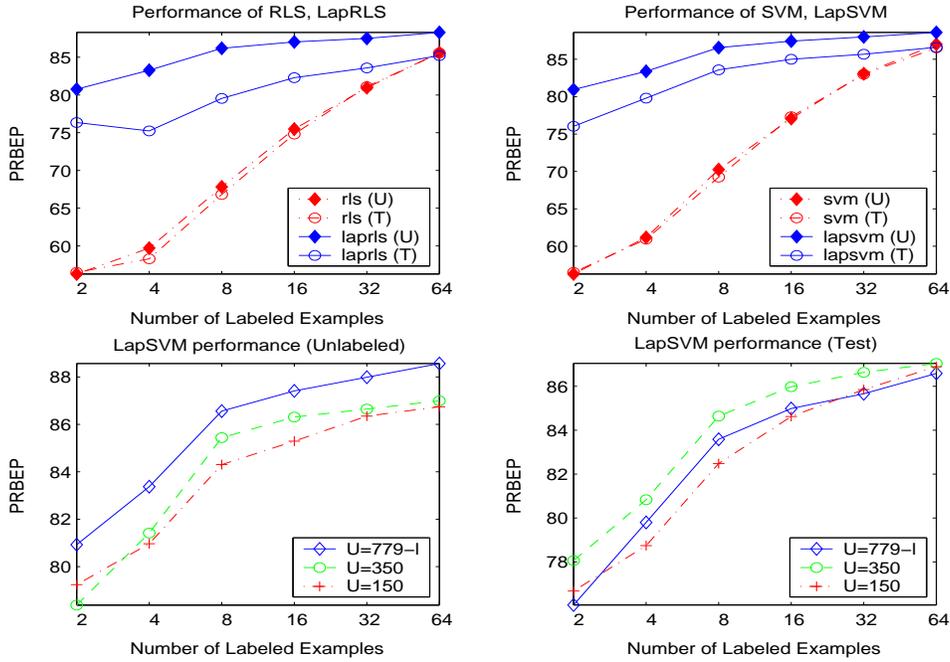


Figure 8: WebKb Text Classification Experiment : The top panel presents performance in terms of precision-recall breakeven points (PRBEP) of RLS,SVM,Laplacian RLS and Laplacian SVM as a function of number of labeled examples, on Test (marked as T) set and Unlabeled set (marked as U and of size 779-number of labeled examples). The bottom panel presents performance curves of Laplacian SVM for different number of unlabeled points.

examples for different amounts of unlabeled data. We find that transduction improves with increasing unlabeled data. We expect this to be true for test set performance as well, but do not observe this consistently since we use a fixed set of parameters that become suboptimal as unlabeled data is increased. The optimal choice of the regularization parameters depends on the amount of labeled and unlabeled data, and should be adjusted by the model selection protocol accordingly.

## 6 Unsupervised and Fully Supervised Cases

While the previous discussion concentrated on the semi-supervised case, our framework covers both unsupervised and fully supervised cases as well. We briefly dis-

cuss each in turn.

## 6.1 Unsupervised Learning: Clustering and Data Representation

In the unsupervised case one is given a collection of unlabeled data points  $x_1, \dots, x_u$ . Our basic algorithmic framework embodied in the optimization problem in Eqn. 3 has three terms: (i) fit to labeled data, (ii) extrinsic regularization and (iii) intrinsic regularization. Since no labeled data is available, the first term does not arise anymore. Therefore we are left with the following optimization problem:

$$\min_{f \in \mathcal{H}_K} \gamma_A \|f\|_K^2 + \gamma_I \|f\|_I^2 \quad (26)$$

Of course, only the ratio  $\frac{\gamma_A}{\gamma_I}$  matters. As before  $\|f\|_I^2$  can be approximated using the unlabeled data. Choosing  $\|f\|_I^2 = \int_{\mathcal{M}} \langle \text{grad}_{\mathcal{M}} f, \text{grad}_{\mathcal{M}} f \rangle$  and approximating it by the empirical Laplacian, we are left with the following optimization problem:

$$f^* = \underset{\substack{\sum_i f(x_i)=0; \sum_i f(x_i)^2=1 \\ f \in \mathcal{H}_K}}{\text{argmin}} \gamma \|f\|_K^2 + \sum_{i \sim j} (f(x_i) - f(x_j))^2 \quad (27)$$

Note that to avoid degenerate solutions we need to impose some additional conditions (cf. [5]). It turns out that a version of Representer theorem still holds showing that the solution to Eqn. 27 admits a representation of the form

$$f^* = \sum_{i=1}^u \alpha_i K(x_i, \cdot)$$

By substituting back in Eqn. 27, we come up with the following optimization problem:

$$\alpha = \underset{\substack{\mathbf{1}^T K \alpha = 0 \\ \alpha^T K^2 \alpha = 1}}{\text{argmin}} \gamma \|f\|_K^2 + \sum_{i \sim j} (f(x_i) - f(x_j))^2$$

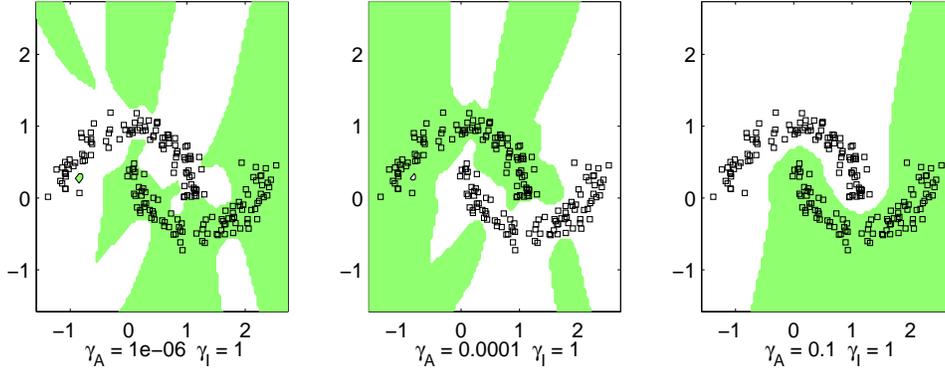
where  $\mathbf{1}$  is the vector of all ones and  $\alpha = (\alpha_1, \dots, \alpha_u)$  and  $K$  is the corresponding Gram matrix.

Letting  $P$  be the projection onto the subspace of  $\mathbb{R}^u$  orthogonal to  $K\mathbf{1}$ , one obtains the solution for the constrained quadratic problem, which is given by the generalized eigenvalue problem

$$P(\gamma K + K L K) P \mathbf{v} = \lambda P K^2 P \mathbf{v} \quad (28)$$

The final solution is given by  $\alpha = P \mathbf{v}$ , where  $\mathbf{v}$  is the eigenvector corresponding to the smallest eigenvalue.

Figure 9: Two Moons Dataset: Regularized Clustering



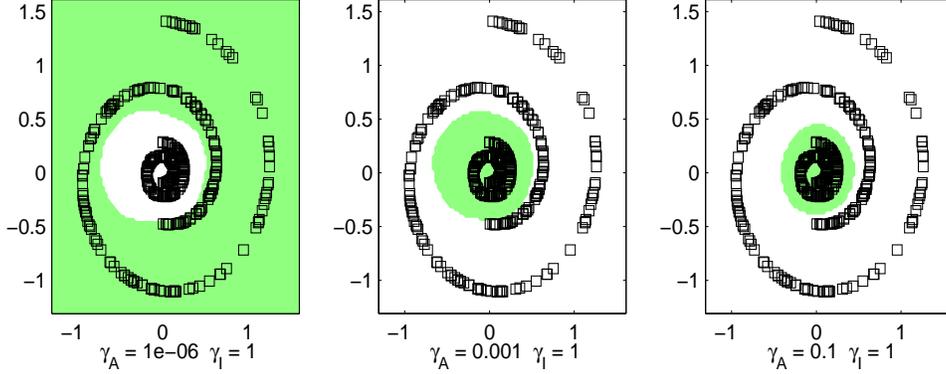
**Remark 1:** The framework for clustering sketched above provides a way of doing regularized spectral clustering, where  $\gamma$  controls the smoothness of the resulting function in the ambient space. We also obtain a natural out-of-sample extension for clustering points not in the original data set. Figures 9,10 show results of this method on two two-dimensional clustering problems. Unlike recent work [9, 11] on out-of-sample extensions, our method is based on a Representer theorem for RKHS.

**Remark 2:** By taking multiple eigenvectors of the system in Eqn. 28 we obtain a natural regularized out-of-sample extension of Laplacian eigenmaps. This leads to new method for dimensionality reduction and data representation. Further study of this approach is a direction of future research.

## 6.2 Fully Supervised Learning

The fully supervised case represents the other end of the spectrum of learning. Since standard supervised algorithms (SVM and RLS) are special cases of manifold regularization, our framework is also able to deal with a labeled dataset containing no unlabeled examples. Additionally, manifold regularization can augment supervised learning with intrinsic regularization, possibly in a class-dependent

Figure 10: Two Spirals Dataset: Regularized Clustering



manner, which suggests the following algorithm :

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f) + \gamma_A \|f\|_K^2 + \frac{\gamma_I^+}{(u+l)^2} \mathbf{f}_+^T L_+ \mathbf{f}_+ + \frac{\gamma_I^-}{(u+l)^2} \mathbf{f}_-^T L_- \mathbf{f}_- \quad (29)$$

Here we introduce two intrinsic regularization parameters  $\gamma_I^+$ ,  $\gamma_I^-$  and regularize separately for the two classes :  $\mathbf{f}_+$ ,  $\mathbf{f}_-$  are the vectors of evaluations of the function  $f$ , and  $L_+$ ,  $L_-$  are the graph Laplacians, on positive and negative examples respectively. The solution to the above problem for RLS and SVM can be obtained by replacing  $\gamma_I L$  by the block-diagonal matrix  $\begin{pmatrix} \gamma_I^+ L_+ & 0 \\ 0 & \gamma_I^- L_- \end{pmatrix}$  in the manifold regularization formulae given in Section 4.

Detailed experimental study of this approach to supervised learning is left for future work.

## 7 Conclusions and Further Directions

We have provided a novel framework for data-dependent geometric regularization. It is based on a new Representer theorem that provides a basis for several algorithms for unsupervised, semi-supervised and fully supervised learning. This framework brings together ideas from the theory of regularization in Reproducing Kernel Hilbert spaces, manifold learning and spectral methods.

There are several directions of future research:

- 1. Convergence and generalization error:** The crucial issue of dependence of generalization error on the number of labeled and unlabeled examples is still very poorly understood. Some very preliminary steps in that direction have been taken in [6].
- 2. Model selection:** Model selection involves choosing appropriate values for the extrinsic and intrinsic regularization parameters. We do not as yet have a good understanding of how to choose these parameters. More systematic procedures need to be developed.
- 3. Efficient algorithms:** It is worth noting that naive implementations of our optimization algorithms give rise to cubic time complexities, which might be impractical for large problems. Efficient algorithms for exact or approximate solutions need to be devised.
- 4. Additional structure:** In this paper we have shown how to incorporate the geometric structure of the marginal distribution into the regularization framework. We believe that this framework will extend to other structures that may constrain the learning task and bring about effective learnability. One important example of such structure is invariance under certain classes of natural transformations, such as invariance under lighting conditions in vision.

## Acknowledgments

We are grateful to Marc Coram, Steve Smale and Peter Bickel for intellectual support and to NSF funding for financial support. We would like to acknowledge the Toyota Technological Institute for its support for this work.

## References

- [1] R. A. Adams, *Sobolev Spaces*, Academic Press, New York, 1975.
- [2] N. Aronszajn, *Theory of Reproducing Kernels*, Transactions of the American mathematical Society, Vol. 68, Issue 3, p337-404, 1950.
- [3] M. Belkin, *Problems of Learning on Manifolds*, The University of Chicago, Ph.D. Dissertation, 2003.
- [4] M. Belkin, P. Niyogi, *Using Manifold Structure for Partially Labeled Classification*, NIPS 2002.
- [5] M. Belkin, P. Niyogi. (2003). *Laplacian Eigenmaps for Dimensionality Reduction and Data Representation*, Neural Computation, Vol. 15, No. 6, 1373-1396.

- [6] M. Belkin, I. Matveeva, P. Niyogi, *Regression and Regularization on Large Graphs*, COLT 2004.
- [7] O. Bousquet, O. Chapelle, M. Hein, *Measure Based Regularization*, NIPS 2003.
- [8] Y. Bengio, O. Delalleau and N. Le Roux, *Efficient Non-Parametric Function Induction in Semi-Supervised Learning*, Technical Report 1247, DIRO, University of Montreal, 2004.
- [9] Y. Bengio, J-F. Paiement, and P. Vincent, *Out-of-Sample Extensions for LLE, Isomap, MDS, Eigenmaps, and Spectral Clustering*, NIPS 2003.
- [10] K. Bennett and A. Demirez, *Semi-Supervised Support Vector Machines*, Advances in Neural Information Processing Systems, 12, M. S. Kearns, S. A. Solla, D. A. Cohn, editors, MIT Press, Cambridge, MA, 1998, pp 368-374
- [11] M. Brand, *Nonlinear dimensionality reduction by kernel eigenmaps*. Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence, pp. 547-552, Acapulco, Mexico, 9-15 August 2003.
- [12] A. Blum, S. Chawla, *Learning from Labeled and Unlabeled Data using Graph Min-cuts*, ICML 2001.
- [13] A. Blum, T. Mitchell, *Combining Labeled and Unlabeled Data with Co-Training*, Proceedings of the 11th Annual Conference on Computational Learning Theory, pages 92–100, 1998
- [14] Chapelle, O., J. Weston and B. Schoelkopf, *Cluster Kernels for Semi-Supervised Learning*, NIPS 2002.
- [15] F. R. K. Chung. (1997). *Spectral Graph Theory*. Regional Conference Series in Mathematics, number 92.
- [16] A. Corduneanu, T. Jaakkola, *On Information Regularization*, UAI 2003.
- [17] F. Cucker, S. Smale, *On the Mathematical Foundations of Learning*, Bull. Amer. Math. Soc. 39 (2002), 1-49.
- [18] S. T. Dumais, J. Platt, D. Heckerman and M. Sahami (1998). *Inductive learning algorithms and representations for text categorization*, In Proceedings of ACM-CIKM98, Nov. 1998, pp. 148-155
- [19] D. L. Donoho, C. E. Grimes, *Hessian Eigenmaps: new locally linear embedding techniques for high-dimensional data*, Proceedings of the National Academy of Arts and Sciences vol. 100 pp. 5591-5596.
- [20] T. Evgeniou, M. Pontil and T. Poggio, *Regularization Networks and Support Vector Machines*, Advances in Computational Mathematics, Vol. 13, 1-50, 2000.
- [21] G. Fung and O. L. Mangasarian, *Semi-Supervised Support Vector Machines for Unlabeled Data Classification*, Optimization Methods and Software 15, 2001, 29-44
- [22] T. Joachims, *Transductive Inference for Text Classification using Support Vector Machines*, ICML 1999.

- [23] T. Joachims, *Transductive Learning via Spectral Graph Partitioning*, Proceedings of the International Conference on Machine Learning (ICML), 2003
- [24] C.C. Kemp, T.L. Griffiths, S. Stromsten, J.B. Tenenbaum, *Semi-supervised Learning with Trees*, NIPS 2003.
- [25] K. Nigam, A. McCallum, S. Thrun and T. Mitchell. *Text Classification from Labeled and Unlabeled Documents using EM*. Machine Learning, 39(2/3). pp. 103-134. 2000.
- [26] R. Rifkin, *Everything Old Is New Again: A Fresh Look at Historical Approaches in Machine Learning*, PhD Thesis, MIT, 2002.
- [27] Sam T. Roweis, Lawrence K. Saul. (2000). *Nonlinear Dimensionality Reduction by Locally Linear Embedding*, Science, vol 290.
- [28] M. Seeger *Learning with Labeled and Unlabeled Data*, Technical Report. Edinburgh University (2000)
- [29] A. Smola and R. Kondor, *Kernels and Regularization on Graphs*, COLT/KW 2003.
- [30] B. Schoelkopf, C.J.C. Burges, V. Vapnik, *Extracting Support Data for a Given Task*, KDD95.
- [31] B. Schoelkopf, A. Smola, *Learning with Kernels*, 644, MIT Press, Cambridge, MA (2002).
- [32] Martin Szummer, Tommi Jaakkola, *Partially labeled classification with Markov random walks*, NIPS 2001.
- [33] J.B.Tenenbaum, V. de Silva, J. C. Langford. (2000). *A Global Geometric Framework for Nonlinear Dimensionality Reduction*, Science, Vol 290.
- [34] A.N. Tikhonov, *Regularization of Incorrectly Posed Problems*, Soviet Math. Doklady 4, 1963 (English Translation).
- [35] V. Vapnik, *Statistical Learning Theory*, Wiley-Interscience, 1998.
- [36] G. Wahba. (1990). *Spline Models for Observational Data*, Society for Industrial and Applied Mathematics.
- [37] D. Zhou, O. Bousquet, T.N. Lal, J. Weston and B. Schoelkopf, *Learning with Local and Global Consistency*, NIPS 2003.
- [38] X. Zhu, J. Lafferty and Z. Ghahramani, *Semi-supervised learning using Gaussian fields and harmonic functions*, ICML 2003.
- [39] <http://www.cse.msu.edu/~lawhiu/manifold/>