

Computer Sciences Department

Semi-Supervised Learning Literature Survey

Xiaojin Zhu

Technical Report #1530

September 2005

UNIVERSITY OF
WISCONSIN
MADISON

Semi-Supervised Learning Literature Survey

Xiaojin Zhu

Last modified on September 7, 2005

Contents

1	FAQ	3
2	Generative Mixture Models and EM	5
2.1	Identifiability	6
2.2	Model Correctness	6
2.3	EM Local Maxima	7
2.4	Cluster and Label	8
3	Self-Training	8
4	Co-Training	8
5	Avoiding Changes in Dense Regions	10
5.1	Transductive SVM	10
5.2	Gaussian Processes	11
5.3	Information Regularization	11
5.4	Entropy Minimization	12
6	Graph-Based Methods	12
6.1	Regularization by Graph	12
6.1.1	Mincut	13
6.1.2	Discrete Markov Random Fields: Boltzmann Machines	13
6.1.3	Gaussian Random Fields and Harmonic Functions	14
6.1.4	Local and Global Consistency	14
6.1.5	Tikhonov Regularization	14
6.1.6	Manifold Regularization	15
6.1.7	Graph Kernels from the Spectrum of Laplacian	15
6.1.8	Spectral Graph Transducer	16
6.1.9	Tree-Based Bayes	16
6.1.10	Some Other Methods	17
6.2	Graph Construction	17
6.3	Fast Computation	18
6.4	Induction	19
6.5	Consistency	20
6.6	Directed Graphs	20
6.7	Connection to Standard Graphical Models	21
7	Metric-Based Model Selection	21

8	Computational Learning Theory	22
9	Related Areas	23
9.1	Spectral Clustering	23
9.2	Learning with Positive and Unlabeled Data	24
9.3	Clustering with Side Information	24
9.4	Nonlinear Dimensionality Reduction	25
9.5	Learning a Distance Metric	25
9.6	Inferring Label Sampling Mechanisms	27

1 FAQ

Q: What's in this Document?

A: We review some of the literature on semi-supervised learning. There has been a whole spectrum of interesting ideas on how to learn from both labeled and unlabeled data. This document is a chapter excerpt from the author's doctoral thesis (Zhu, 2005), but with recent updates. Please check the latest version at

http://www.cs.wisc.edu/~jerryzhu/pub/ssl_survey.pdf

The review is by no means comprehensive as the field of semi-supervised learning is evolving rapidly. It is difficult for one person to summarize the field. I apologize in advance for any missed papers and inaccuracies in descriptions. Corrections and comments are highly welcome. Please send them to jerryzhu@cs.wisc.edu.

Q: What is semi-supervised learning?

A: It's a special form of classification. Traditional classifiers need labeled data (feature / label pairs) to train. Labeled instances however are often difficult, expensive, or time consuming to obtain, as they require the efforts of experienced human annotators. Meanwhile unlabeled data may be relatively easy to collect, but there has been few ways to use them. Semi-supervised learning addresses this problem by using large amount of unlabeled data, together with the labeled data, to build better classifiers. Because semi-supervised learning requires less human effort and gives higher accuracy, it is of great interest both in theory and in practice.

Q: Can we really learn anything from unlabeled data? It looks like magic.

A: Yes we can – under certain assumptions. It's not magic, but good matching of problem structure with model assumption.

Q: Does unlabeled data always help?

A: No, there's no free lunch. Bad matching of problem structure with model as-

sumption can lead to degradation in classifier performance. For example, quite a few semi-supervised learning methods assume that the decision boundary should avoid regions with high $p(x)$. These methods include transductive support vector machines (SVMs), information regularization, Gaussian processes with null category noise model, graph-based methods if the graph weights is determined by pairwise distance. Nonetheless if the data is generated from two heavily overlapping Gaussian, the decision boundary would go right through the densest region, and these methods would perform badly. On the other hand EM with generative mixture models, another semi-supervised learning method, would have easily solved the problem. Detecting bad match in advance however is hard and remains an open question.

Q: How many semi-supervised learning methods are there?

A: Many. Some often-used methods include: EM with generative mixture models, self-training, co-training, transductive support vector machines, and graph-based methods. See the following sections for more methods.

Q: Which method should I use / is the best?

A: There is no direct answer to this question. Because labeled data is scarce, semi-supervised learning methods make strong model assumptions. Ideally one should use a method whose assumptions fit the problem structure. This may be difficult in reality. Nonetheless we can try the following checklist: Do the classes produce well clustered data? If yes, EM with generative mixture models may be a good choice; Do the features naturally split into two sets? If yes, co-training may be appropriate; Is it true that two points with similar features tend to be in the same class? If yes, graph-based methods can be used; Already using SVM? Transductive SVM is a natural extension; Is the existing supervised classifier complicated and hard to modify? Self-training is a practical wrapper method.

Q: How do semi-supervised learning methods use unlabeled data?

A: Semi-supervised learning methods use unlabeled data to either modify or reprioritize hypotheses obtained from labeled data alone. Although not all methods are probabilistic, it is easier to look at methods that represent hypotheses by $p(y|x)$, and unlabeled data by $p(x)$. Generative models have common parameters for the joint distribution $p(x, y)$. It is easy to see that $p(x)$ influences $p(y|x)$. Mixture models with EM is in this category, and to some extent self-training. Many other methods are discriminative, including transductive SVM, Gaussian processes, information regularization, and graph-based methods. Original discriminative training cannot be used for semi-supervised learning, since $p(y|x)$ is estimated ignoring $p(x)$. To solve the problem, $p(x)$ dependent terms are often brought into the ob-

jective function, which amounts to assuming $p(y|x)$ and $p(x)$ share parameters.

Q: What is the difference between ‘transductive learning’ and ‘semi-supervised learning’?

A: Different authors use slightly different names. In this survey I will use the following convention:

- ‘Semi-supervised learning’ refers to the use of both labeled and unlabeled data for training. It contrasts supervised learning (data all labeled) or unsupervised learning (data all unlabeled). Other names are ‘learning from labeled and unlabeled data’ or ‘learning from partially labeled/classified data’. Notice semi-supervised learning can be either transductive or inductive.
- ‘Transductive learning’ will be used to contrast inductive learning. A learner is transductive if it only works on the labeled and unlabeled training data, and cannot handle unseen data. The early graph-based methods are often transductive. Inductive learners can naturally handle unseen data. Notice under this convention *transductive support vector machines* (TSVMs) are in fact inductive learners, because the resulting classifiers are defined over the whole space. The name TSVM originates from the intention to work only on the observed data (though people use them for induction anyway).
- In this survey semi-supervised learning refers to ‘semi-supervised classification’, where one has additional unlabeled data and the goal is classification. Its cousin ‘semi-supervised clustering’, where one has unlabeled data with some pairwise constraints and the goal is clustering, is only briefly discussed later in the survey.

We will follow the above convention in the survey.

Q: Where can I learn more?

A: An existing survey can be found in (Seeger, 2001).

2 Generative Mixture Models and EM

This is perhaps the oldest semi-supervised learning method. It assumes a generative model $p(x, y) = p(y)p(x|y)$ where $p(x|y)$ is an identifiable mixture distribution, for example Gaussian mixture models. With large amount of unlabeled data, the mixture components can be identified; then ideally we only need one labeled example per component to fully determine the mixture distribution. One can think of the mixture components as ‘soft clusters’.

Nigam et al. (2000) apply the EM algorithm on mixture of multinomial for the task of text classification. They showed the resulting classifiers perform better than those trained only from L . Baluja (1998) uses the same algorithm on a face orientation discrimination task. Fujino et al. (2005) extend generative mixture models by including a ‘bias correction’ term and discriminative training using the maximum entropy principle.

One has to pay attention to a few things:

2.1 Identifiability

The mixture model ideally should be identifiable. In general let $\{p_\theta\}$ be a family of distributions indexed by a parameter vector θ . θ is identifiable if $\theta_1 \neq \theta_2 \Rightarrow p_{\theta_1} \neq p_{\theta_2}$, up to a permutation of mixture components. If the model family is identifiable, in theory with infinite U one can learn θ up to a permutation of component indices.

Here is an example showing the problem with unidentifiable models. The model $p(x|y)$ is uniform for $y \in \{+1, -1\}$. Assuming with large amount of unlabeled data U we know $p(x)$ is uniform in $[0, 1]$. We also have 2 labeled data points $(0.1, +1), (0.9, -1)$. Can we determine the label for $x = 0.5$? No. With our assumptions we cannot distinguish the following two models:

$$p(y = 1) = 0.2, p(x|y = 1) = \text{unif}(0, 0.2), p(x|y = -1) = \text{unif}(0.2, 1) \quad (1)$$

$$p(y = 1) = 0.6, p(x|y = 1) = \text{unif}(0, 0.6), p(x|y = -1) = \text{unif}(0.6, 1) \quad (2)$$

which give opposite labels at $x = 0.5$, see Figure 1. It is known that a mixture of Gaussian is identifiable. Mixture of multivariate Bernoulli (McCallum & Nigam, 1998) is not identifiable. More discussions on identifiability and semi-supervised learning can be found in e.g. (Ratsaby & Venkatesh, 1995) and (Corduneanu & Jaakkola, 2001).

2.2 Model Correctness

If the mixture model assumption is correct, unlabeled data is guaranteed to improve accuracy (Castelli & Cover, 1995) (Castelli & Cover, 1996) (Ratsaby & Venkatesh, 1995). However if the model is wrong, unlabeled data may actually hurt accuracy. Figure 2 shows an example. This has been observed by multiple researchers. Cozman et al. (2003) give a formal derivation on how this might happen.

It is thus important to carefully construct the mixture model to reflect reality. For example in text categorization a topic may contain several sub-topics, and will be better modeled by multiple multinomial instead of a single one (Nigam et al., 2000). Some other examples are (Shahshahani & Landgrebe, 1994) (Miller & Uyar, 1997). Another solution is to down-weighting unlabeled data (Corduneanu &

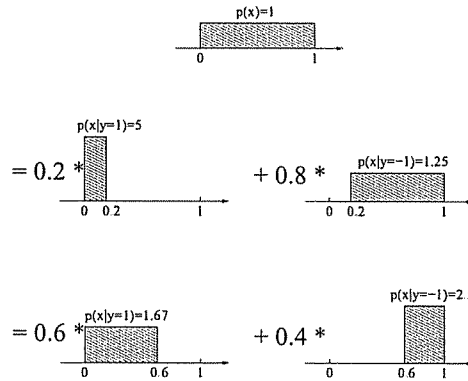


Figure 1: An example of unidentifiable models. Even if we know $p(x)$ (top) is a mixture of two uniform distributions, we cannot uniquely identify the two components. For instance, the mixtures on the second and third line give the same $p(x)$, but they classify $x = 0.5$ differently.

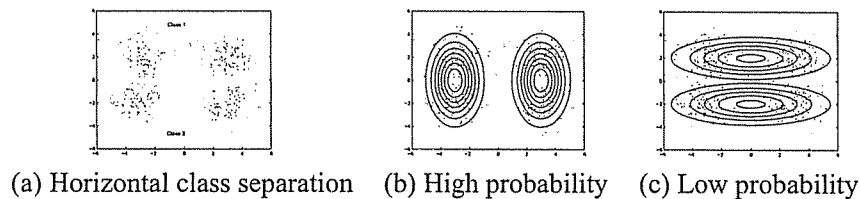


Figure 2: If the model is wrong, higher likelihood may lead to lower classification accuracy. For example, (a) is clearly not generated from two Gaussians. If we insist that each class is a single Gaussian, (b) will have higher probability than (c). But (b) has around 50% accuracy, while (c)'s is much better.

Jaakkola, 2001), which is also used by Nigam et al. (2000), and by Callison-Burch et al. (2004) who estimate word alignment for machine translation.

2.3 EM Local Maxima

Even if the mixture model assumption is correct, in practice mixture components are identified by the Expectation-Maximization (EM) algorithm (Dempster et al., 1977). EM is prone to local maxima. If a local maximum is far from the global maximum, unlabeled data may again hurt learning. Remedies include smart choice of starting point by active learning (Nigam, 2001).

2.4 Cluster and Label

We shall also mention that instead of using an probabilistic generative mixture model, some approaches employ various clustering algorithms to cluster the whole dataset, then label each cluster with labeled data, e.g. (Demiriz et al., 1999) (Dara et al., 2000). Although they can perform well if the particular clustering algorithms match the true data distribution, these approaches are hard to analyze due to their algorithmic nature.

3 Self-Training

Self-training is a commonly used technique for semi-supervised learning. In self-training a classifier is first trained with the small amount of labeled data. The classifier is then used to classify the unlabeled data. Typically the most confident unlabeled points, together with their predicted labels, are added to the training set. The classifier is re-trained and the procedure repeated. Note the classifier uses its own predictions to teach itself. The procedure is also called self-teaching or bootstrapping (not to be confused with the statistical procedure with the same name). The generative model and EM approach of section 2 can be viewed as a special case of ‘soft’ self-training. One can imagine that a classification mistake can reinforce itself. Some algorithms try to avoid this by ‘unlearn’ unlabeled points if the prediction confidence drops below a threshold.

Self-training has been applied to several natural language processing tasks. Yarowsky (1995) uses self-training for word sense disambiguation, e.g. deciding whether the word ‘plant’ means a living organism or a factory in a give context. Riloff et al. (2003) uses it to identify subjective nouns. Maeireizo et al. (2004) classify dialogues as ‘emotional’ or ‘non-emotional’ with a procedure involving two classifiers. Self-training has also been applied to parsing and machine translation. Rosenberg et al. (2005) apply self-training to object detection systems from images, and show the semi-supervised technique compares favorably with a state-of-the-art detector.

4 Co-Training

Co-training (Blum & Mitchell, 1998) (Mitchell, 1999) assumes that features can be split into two sets; Each sub-feature set is sufficient to train a good classifier; The two sets are conditionally independent given the class. Initially two separate classifiers are trained with the labeled data, on the two sub-feature sets respectively. Each classifier then classifies the unlabeled data, and ‘teaches’ the other classifier

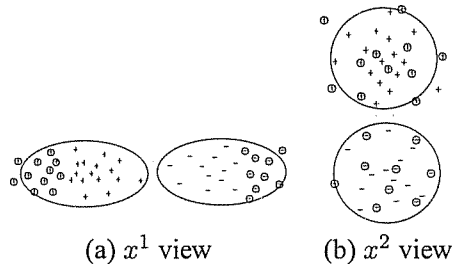


Figure 3: Co-Training: Conditional independent assumption on feature split. With this assumption the high confident data points in x^1 view, represented by circled labels, will be randomly scattered in x^2 view. This is advantageous if they are to be used to teach the classifier in x^2 view.

with the few unlabeled examples (and the predicted labels) they feel most confident. Each classifier is retrained with the additional training examples given by the other classifier, and the process repeats.

In co-training, unlabeled data helps by reducing the version space size. In other words, the two classifiers (or hypotheses) must agree on the much larger unlabeled data as well as the labeled data.

We need the assumption that sub-features are sufficiently good, so that we can trust the labels by each learner on U . We need the sub-features to be conditionally independent so that one classifier's high confident data points are *iid* samples for the other classifier. Figure 3 visualizes the assumption.

Nigam and Ghani (2000) perform extensive empirical experiments to compare co-training with generative mixture models and EM. Their result shows co-training performs well if the conditional independence assumption indeed holds. In addition, it is better to probabilistically label the entire U , instead of a few most confident data points. They name this paradigm co-EM. Finally, if there is no natural feature split, the authors create artificial split by randomly break the feature set into two subsets. They show co-training with artificial feature split still helps, though not as much as before. Jones (2005) used co-training, co-EM and other related methods for information extraction from text.

Co-training makes strong assumptions on the splitting of features. One might wonder if these conditions can be relaxed. Goldman and Zhou (2000) use two learners of different type but both takes the whole feature set, and essentially use one learner's high confidence data points, identified with a set of statistical tests, in U to teach the other learning and vice versa. Balcan et al. (2005b) relax the conditional independence assumption with a much weaker expansion condition, and justify the iterative co-training procedure.

More generally, we can define learning paradigms that utilize the agreement among different learners. Co-training can be viewed as a special case with two learners and a specific algorithm to enforce agreement. For instance, the work of Leskes (2005) is discussed in Section 8.

5 Avoiding Changes in Dense Regions

5.1 Transductive SVM

Discriminative methods work on $p(y|x)$ directly. This brings up the danger of leaving $p(x)$ outside of the parameter estimation loop, if $p(x)$ and $p(y|x)$ do not share parameters. Notice $p(x)$ is usually all we can get from unlabeled data. It is believed that if $p(x)$ and $p(y|x)$ do not share parameters, semi-supervised learning cannot help. This point is emphasized in (Seeger, 2001). Zhang and Oles (2000) give both theoretical and experimental evidence of the same point on transductive support vector machines (TSVMs). However empirically TSVMs seem to be beneficial.

TSVM is an extension of standard support vector machines with unlabeled data. In a standard SVM only the labeled data is used, and the goal is to find a maximum margin linear boundary in the Reproducing Kernel Hilbert Space. In a TSVM the unlabeled data is also used. The goal is to find a labeling of the unlabeled data, so that a linear boundary has the maximum margin on both the original labeled data and the (now labeled) unlabeled data. The decision boundary has the smallest generalization error bound on unlabeled data (Vapnik, 1998). Intuitively, unlabeled data guides the linear boundary away from dense regions. However

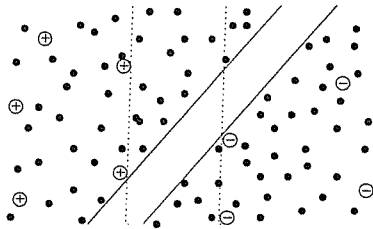


Figure 4: In TSVM, U helps to put the decision boundary in sparse regions. With labeled data only, the maximum margin boundary is plotted with dotted lines. With unlabeled data (black dots), the maximum margin boundary would be the one with solid lines.

finding the exact transductive SVM solution is NP-hard. Several approximation algorithms have been proposed and show positive results, see e.g. (Joachims, 1999)

(Bennett & Demiriz, 1999) (Demirez & Bennett, 2000) (Fung & Mangasarian, 1999) (Chapelle & Zien, 2005).

Xu and Schuurmans (2005) present a training method based on semi-definite programming (SDP, which applies to the completely unsupervised SVMs as well). In the simple binary classification case, the goal of finding a good labeling for unlabeled data is formulated as finding a positive semi-definite matrix M . M is meant to be the continuous relaxation of the label outer product matrix yy^T , and the SVM objective is expressed as semi-definite programming on M . There are effective (although still expensive) SDP solvers. Importantly, the authors propose multi-class version of the SDP, which results in multi-class SVM for semi-supervised learning.

The maximum entropy discrimination approach (Jaakkola et al., 1999) also maximizes the margin, and is able to take into account unlabeled data, with SVM as a special case.

The application of graph kernels (Zhu et al., 2005) to SVMs differs from TSVM. The graph kernels are special semi-supervised kernels applied to a standard SVM; TSVM is a special optimization criterion regardless of the kernel being used.

5.2 Gaussian Processes

Lawrence and Jordan (2005) proposed a Gaussian process approach, which can be viewed as the Gaussian process parallel of TSVM. The key difference to a standard Gaussian process is in the noise model. A ‘null category noise model’ maps the hidden continuous variable f to three instead of two labels, specifically to the never used label ‘0’ when f is around zero. On top of that, it is restricted that unlabeled data points cannot take the label 0. This pushes the posterior of f away from zero for the unlabeled points. It achieves the similar effect of TSVM where the margin avoids dense unlabeled data region. However nothing special is done on the process model. Therefore all the benefit of unlabeled data comes from the noise model. A very similar noise model is proposed in (Chu & Ghahramani, 2004) for ordinal regression.

This is different from the Gaussian processes in (Zhu et al., 2003b), where we have a semi-supervised Gram matrix, and semi-supervised learning originates from the process model, not the noise model.

5.3 Information Regularization

Szummer and Jaakkola (2002) propose the information regularization framework to control the label conditionals $p(y|x)$ by $p(x)$, where $p(x)$ may be estimated from unlabeled data. The idea is that labels shouldn’t change too much in regions where

$p(x)$ is high. The authors use the mutual information $I(x; y)$ between x and y as a measure of label complexity. $I(x; y)$ is small when the labels are homogeneous, and large when labels vary. This motivates the minimization of the product of $p(x)$ mass in a region with $I(x; y)$ (normalized by a variance term). The minimization is carried out on multiple overlapping regions covering the data space.

The theory is developed further in (Corduneanu & Jaakkola, 2003). Corduneanu and Jaakkola (2005) extend the work by formulating semi-supervised learning as a communication problem. Regularization is expressed as the rate of information, which again discourages complex conditionals $p(y|x)$ in regions with high $p(x)$. The problem becomes finding the unique $p(y|x)$ that minimizes a regularized loss on labeled data. The authors give a local propagation algorithm.

5.4 Entropy Minimization

The hyperparameter learning method in section 7.2 of (Zhu, 2005) uses entropy minimization. Grandvalet and Bengio (2005) used the label entropy on unlabeled data as a regularizer. By minimizing the entropy, the method assumes a prior which prefers minimal class overlap.

6 Graph-Based Methods

Graph-based semi-supervised methods define a graph where the nodes are labeled and unlabeled examples in the dataset, and edges (may be weighted) reflect the similarity of examples. These methods usually assume label smoothness over the graph. Graph methods are nonparametric, discriminative, and transductive in nature.

6.1 Regularization by Graph

Many graph-based methods can be viewed as estimating a function f on the graph. One wants f to satisfy two things at the same time: 1) it should be close to the given labels y_L on the labeled nodes, and 2) it should be smooth on the whole graph. This can be expressed in a regularization framework where the first term is a loss function, and the second term is a regularizer.

Several graph-based methods listed here are similar to each other. They differ in the particular choice of the loss function and the regularizer. We believe it is more important to construct a good graph than to choose among the methods. However graph construction, as we will see later, is not a well studied area.

6.1.1 Mincut

Blum and Chawla (2001) pose semi-supervised learning as a graph mincut (also known as *st-cut*) problem. In the binary case, positive labels act as sources and negative labels act as sinks. The objective is to find a minimum set of edges whose removal blocks all flow from the sources to the sinks. The nodes connecting to the sources are then labeled positive, and those to the sinks are labeled negative. Equivalently mincut is the *mode* of a Markov random field with binary labels (Boltzmann machine). The loss function can be viewed as a quadratic loss with infinity weight: $\infty \sum_{i \in L} (y_i - y_{i|L})^2$, so that the values on labeled data are in fact fixed at their given labels. The regularizer is

$$\frac{1}{2} \sum_{i,j} w_{ij} |y_i - y_j| = \frac{1}{2} \sum_{i,j} w_{ij} (y_i - y_j)^2 \quad (3)$$

The equality holds because the y 's take binary (0 and 1) labels. Putting the two together, mincut can be viewed to minimize the function

$$\infty \sum_{i \in L} (y_i - y_{i|L})^2 + \frac{1}{2} \sum_{i,j} w_{ij} (y_i - y_j)^2 \quad (4)$$

subject to the constraint $y_i \in \{0, 1\}, \forall i$.

One problem with mincut is that it only gives hard classification without confidence (i.e. it computes the mode, not the marginal probabilities). Blum et al. (2004) perturb the graph by adding random noise to the edge weights. Mincut is applied to multiple perturbed graphs, and the labels are determined by a majority vote. The procedure is similar to bagging, and creates a 'soft' mincut.

Pang and Lee (2004) use mincut to improve the classification of a sentence into either 'objective' or 'subjective', with the assumption that sentences close to each other tend to have the same class.

6.1.2 Discrete Markov Random Fields: Boltzmann Machines

The proper but hard way is to compute the marginal probabilities of the discrete Markov random fields. This is inherently a difficult inference problem. Zhu and Ghahramani (2002) attempted exactly this, but were limited by the MCMC sampling techniques (they used global Metropolis and Swendsen-Wang sampling).

Getz et al. (2005) computes the marginal probabilities of the discrete Markov random field at any temperature with the Multicanonical Monte-Carlo method, which seems to be able to overcome the energy trap faced by the standard Metropolis or Swendsen-Wang method. The authors discuss the relationship between temperatures and phases in such systems. They also propose a heuristic procedure to identify possible new classes.

6.1.3 Gaussian Random Fields and Harmonic Functions

The Gaussian random fields and harmonic function methods in (Zhu et al., 2003a) is a continuous relaxation to the difficult discrete Markov random fields (or Boltzmann machines). It can be viewed as having a quadratic loss function with infinity weight, so that the labeled data are clamped (fixed at given label values), and a regularizer based on the graph combinatorial Laplacian Δ :

$$\infty \sum_{i \in L} (f_i - y_i)^2 + 1/2 \sum_{i,j} w_{ij} (f_i - f_j)^2 \quad (5)$$

$$= \infty \sum_{i \in L} (f_i - y_i)^2 + f^\top \Delta f \quad (6)$$

Notice $f_i \in \mathbb{R}$, which is the key relaxation to Mincut. This allows for a simple closed-form solution for the node marginal probabilities. The mean is known as a harmonic function, which has many interesting properties (Zhu, 2005).

Recently Grady and Funka-Lea (2004) applied the harmonic function method to medical image segmentation tasks, where a user labels classes (e.g. different organs) with a few strokes. Levin et al. (2004) use the equivalent of harmonic functions for colorization of gray-scale images. Again the user specifies the desired color with only a few strokes on the image. The rest of the image is used as unlabeled data, and the labels propagation through the image. Niu et al. (2005) applied the label propagation algorithm (which is equivalent to harmonic functions) to word sense disambiguation.

6.1.4 Local and Global Consistency

The local and global consistency method (Zhou et al., 2004a) uses the loss function $\sum_{i=1}^n (f_i - y_i)^2$, and the *normalized Laplacian* $D^{-1/2} \Delta D^{-1/2} = I - D^{-1/2} W D^{-1/2}$ in the regularizer,

$$1/2 \sum_{i,j} w_{ij} (f_i / \sqrt{D_{ii}} - f_j / \sqrt{D_{jj}})^2 = f^\top D^{-1/2} \Delta D^{-1/2} f \quad (7)$$

6.1.5 Tikhonov Regularization

The Tikhonov regularization algorithm in (Belkin et al., 2004a) uses the loss function and regularizer:

$$1/k \sum_i (f_i - y_i)^2 + \gamma f^\top S f \quad (8)$$

where $S = \Delta$ or Δ^p for some integer p .

6.1.6 Manifold Regularization

The manifold regularization framework (Belkin et al., 2004b) (Belkin et al., 2005) employs two regularization terms:

$$\frac{1}{l} \sum_{i=1}^l V(x_i, y_i, f) + \gamma_A \|f\|_K^2 + \gamma_U \|f\|_I^2 \quad (9)$$

where V is an arbitrary loss function, K is a ‘base kernel’, e.g. a linear or RBF kernel. I is a regularization term induced by the labeled and unlabeled data. For example, one can use

$$\|f\|_I^2 = \frac{1}{(l+u)^2} \hat{f}^\top \Delta \hat{f} \quad (10)$$

where \hat{f} is the vector of f evaluations on $L \cup U$.

Sindhwani et al. (2005a) give perhaps the first semi-supervised kernel that is not limited to the unlabeled points, but defined over all input space. The kernel thus supports induction. Essentially the kernel is a new interpretation of the manifold regularization framework above. Starting from a base kernel K defined over the whole input space (e.g. linear kernels, RBF kernels), the authors modify the RKHS by keeping the same function space but changing the norm. Specifically a ‘point-cloud norm’ defined by $L \cup U$ is added to the original norm. The point-cloud norm corresponds to $\|f\|_I^2$. Importantly this results in a new RKHS space, with a corresponding new kernel that deforms the original one along a finite-dimensional subspace given by the data. The new kernel is defined over the whole space, yet it ‘follows the manifold’. Standard supervised kernel machines with the new kernel, trained on L only, are able to perform inductive semi-supervised learning. In fact they are equivalent to LapSVM and LapRLS (Belkin et al., 2005) with a certain parameter. Nonetheless finding the new kernel involves inverting a $n \times n$ matrix. Like many other methods it can be costly. Also notice the new kernel depends on the observed $L \cup U$ data, thus it is a random kernel.

6.1.7 Graph Kernels from the Spectrum of Laplacian

For kernel methods, the regularizer is a (typically monotonically increasing) function of the RKHS norm $\|f\|_K = f^\top K^{-1} f$ with kernel K . Such kernels are derived from the graph, e.g. the Laplacian.

Chapelle et al. (2002) and Smola and Kondor (2003) both show the spectral transformation of a Laplacian results in kernels suitable for semi-supervised learning. The diffusion kernel (Kondor & Lafferty, 2002) corresponds to a spectrum

transform of the Laplacian with

$$r(\lambda) = \exp\left(-\frac{\sigma^2}{2}\lambda\right) \quad (11)$$

The regularized Gaussian process kernel $\Delta + I/\sigma^2$ in (Zhu et al., 2003b) corresponds to

$$r(\lambda) = \frac{1}{\lambda + \sigma} \quad (12)$$

Similarly the order constrained graph kernels in (Zhu et al., 2005) are constructed from the spectrum of the Laplacian, with non-parametric convex optimization. Learning the optimal eigenvalues for a graph kernel is in fact a way to (at least partially) improve an imperfect graph. In this sense it is related to graph construction.

6.1.8 Spectral Graph Transducer

The spectral graph transducer (Joachims, 2003) can be viewed with a loss function and regularizer

$$c(f - \gamma)^\top C(f - \gamma) + f^\top Lf \quad (13)$$

where $\gamma_i = \sqrt{l_-/l_+}$ for positive labeled data, $-\sqrt{l_+/l_-}$ for negative data, l_- being the number of negative data and so on. L can be the combinatorial or normalized graph Laplacian, with a transformed spectrum.

Pham et al. (2005) perform empirical experiments on word sense disambiguation, comparing variants of co-training and spectral graph transducer. The authors notice spectral graph transducer with carefully constructed graphs ('SGT-Cotraining') produces good results.

6.1.9 Tree-Based Bayes

Kemp et al. (2003) define a probabilistic distribution $P(Y|T)$ on discrete (e.g. 0 and 1) labellings Y over an evolutionary tree T . The tree T is constructed with the labeled and unlabeled data being the leaf nodes. The labeled data is clamped. The authors assume a mutation process, where a label at the root propagates down to the leaves. The label mutates with a constant rate as it moves down along the edges. As a result the tree T (its structure and edge lengths) uniquely defines the label prior $P(Y|T)$. Under the prior if two leaf nodes are closer in the tree, they have a higher probability of sharing the same label. One can also integrate over all tree structures.

The tree-based Bayes approach can be viewed as an interesting way to incorporate structure of the domain. Notice the leaf nodes of the tree are the labeled and

unlabeled data, while the internal nodes do not correspond to physical data. This is in contrast with other graph-based methods where labeled and unlabeled data are all the nodes.

6.1.10 Some Other Methods

Szummer and Jaakkola (2001) perform a t -step Markov random walk on the graph. The influence of one example to another example is proportional to how easy the random walk goes from one to the other. It has certain resemblance to the diffusion kernel. The parameter t is important.

Chapelle and Zien (2005) use a density-sensitive connectivity distance between nodes i, j (a given path between i, j consists of several segments, one of them is the longest; now consider all paths between i, j and find the shortest ‘longest segment’). Exponentiating the negative distance gives a graph kernel.

Bousquet et al. (2004) propose ‘measure-based regularization’, the continuous counterpart of graph-based regularization. The intuition is that two points are similar if they are connected by high density regions. They define regularization based on a known density $p(x)$ and provide interesting theoretical analysis. However it seems difficult in practice to apply the theoretical results to higher ($D > 2$) dimensional tasks. Sajama and Orlitsky (2005) analyze the lower and upper bounds on estimating such density-based distance. There are two sources of error: one stems from the fact that the true density $p(x)$ is not known, the second is that for practical reasons one typically build a grid on the data points, instead of a regular grid in R^d . The authors separate these two kinds of errors (computational and estimation), and analyze them independently. It sheds light on the complexity of density-based distance, independent of the specific method one uses. It also sheds some light on approximation errors when using neighborhood graphs on data points, which is used widely in semi-supervised learning and non-linear dimensionality reduction, etc. Understanding this dichotomy is helpful when trying to improve methods for semi-supervised learning.

6.2 Graph Construction

Although the graph is at the heart of graph-based semi-supervised learning methods, its construction has not been studied extensively. The issue has been discussed in (Zhu, 2005) Chapter 3 and Chapter 7. Balcan et al. (2005a) build graphs for video surveillance using strong domain knowledge, where the graph of webcam images consists of time edges, color edges and face edges. Such graphs reflect a deep understanding of the problem structure and how unlabeled data is expected to help. Carreira-Perpinan and Zemel (2005) build robust graphs from multiple min-

imum spanning trees by perturbation and edge removal. It is possible that graph construction is domain specific because it encodes prior knowledge, and has thus far been treated on an individual basis.

6.3 Fast Computation

Many semi-supervised learning methods scale as badly as $O(n^3)$ as they were originally proposed. Because semi-supervised learning is interesting when the size of unlabeled data is large, this is clearly a problem. Many methods are also transductive (section 6.4). In 2005 several papers start to address these problems.

Fast computation of the harmonic function with conjugate gradient methods is discussed in (Argyriou, 2004). A comparison of three iterative methods: label propagation, conjugate gradient and loopy belief propagation is presented in (Zhu, 2005) Appendix F. Recently numerical methods for fast N-body problems have been applied to *dense* graphs in semi-supervised learning, reducing the computational cost from $O(n^3)$ to $O(n)$ (Mahdaviani et al., 2005). This is achieved with Krylov subspace methods and the fast Gauss transform.

The harmonic mixture models (Zhu & Lafferty, 2005) convert the original graph into a much smaller backbone graph, by using a mixture model to ‘carve up’ the original $L \cup U$ dataset. Learning on the smaller graph is much faster. Similar ideas have been used for e.g. dimensionality reduction (Teh & Roweis, 2002). The heuristics in (Delalleau et al., 2005) similarly create a small graph with a subset of the unlabeled data. They enables fast approximate computation by reducing the problem size.

Garcke and Griebel (2005) propose the use of sparse grids for semi-supervised learning. The main advantages are $O(n)$ computation complexity for sparse graphs, and the ability of induction. The authors start from the same regularization problem of (Belkin et al., 2005). The key idea is to approximate the function space with a finite basis, with sparse grids. The minimizer f in this finite dimensional subspace can be efficiently computed. As the authors point out, this method is different from the general kernel methods which rely on the representer theorem for finite representation. In practice the method is limited by data dimensionality (around 20). A potential drawback is that the method employs a regular grid, and cannot ‘zoom in’ to small interesting data regions with higher resolution.

Yu et al. (2005) solve the large scale semi-supervised learning problem by using a bipartite graph. The labeled and unlabeled points form one side of the bipartite split, while a much smaller number of ‘block-level’ nodes form the other side. The authors show that the harmonic function can be computed using the block-level nodes. The computation involves inverting a much smaller matrix on block-level nodes. It is thus cheaper and more scalable than working directly on the

LUU matrix. The authors propose two methods to construct the bipartite graph, so that it approximates the given weight matrix W on $L \cup U$. One uses Nonnegative Matrix Factorization, the other uses mixture models. The latter method has the additional benefit of induction, and is similar to the harmonic mixtures (Zhu & Lafferty, 2005). However in the latter method the mixture model is derived based on the given weight matrix W . But in harmonic mixtures W and the mixture model are independent, and the mixture model serves as a ‘second knowledge source’ in addition to W .

The original manifold regularization framework (Belkin et al., 2004b) needs to invert a $(l + u) \times (l + u)$ matrix, and is not scalable. To speed up things, Sindhwani et al. (2005b) consider *linear manifold regularization*. Effectively this is a special case when the base kernel is taken to be the linear kernel. The authors show that it is advantageous to work with the primal variables. The resulting optimization problem can be much smaller if the data dimensionality is small, or sparse.

6.4 Induction

Most graph-based semi-supervised learning algorithms are transductive, i.e. they cannot easily extend to new test points outside of $L \cup U$. Recently induction has received increasing attention. One common practice is to ‘freeze’ the graph on $L \cup U$. New points do not (although they should) alter the graph structure. This avoids expensive graph computation every time one encounters new points.

Zhu et al. (2003b) propose that new test point be classified by its nearest neighbor in LUU . This is sensible when U is sufficiently large. In (Chapelle et al., 2002) the authors approximate a new point by a linear combination of labeled and unlabeled points. Similarly in (Delalleau et al., 2005) the authors proposes an induction scheme to classify a new point x by

$$f(x) = \frac{\sum_{i \in LUU} w_{xi} f(x_i)}{\sum_{i \in LUU} w_{xi}} \quad (14)$$

This can be viewed as an application of the Nyström method (Fowlkes et al., 2004).

Yu et al. (2004) report an early attempt on semi-supervised induction using RBF basis functions in a regularization framework. In (Belkin et al., 2004b), the function f does not have to be restricted to the graph. The graph is merely used to regularize f which can have a much larger support. It is necessarily a combination of an inductive algorithm and graph regularization. The authors give the graph-regularized version of least squares and SVM. Note such an SVM is different from the graph kernels in standard SVM in (Zhu et al., 2005). The former is inductive with both a graph regularizer and an inductive kernel. The latter is transductive with only the graph regularizer. Following the work, Krishnapuram et al. (2005)

use graph regularization on logistic regression. These methods create inductive learners that naturally handle new test points.

The harmonic mixture model (Zhu & Lafferty, 2005) naturally handles new points as well. The idea is to model the labeled and unlabeled data with a mixture model, e.g. mixture of Gaussian. In standard mixture models, the class probability $p(y|i)$ for each mixture component i is optimized to maximize label likelihood. However in harmonic mixture models, $p(y|i)$ is optimized differently to minimize an underlying graph-based cost function. Under certain conditions, the harmonic mixture model converts the original graph on unlabeled data into a ‘backbone graph’, with the components being ‘super nodes’. Harmonic mixture models naturally handle induction just like standard mixture models.

Several other inductive methods have been discussed in section 6.3 together with fast computation.

6.5 Consistency

The consistency of graph-based semi-supervised learning algorithms is an open research area. By consistency we mean whether classification converges to the right solution as the number of labeled and unlabeled data grows to infinity. Recently von Luxburg et al. (2005) (von Luxburg et al., 2004) study the consistency of *spectral clustering methods*. The authors find that the normalized Laplacian is better than the unnormalized Laplacian for spectral clustering. The convergence of the eigenvectors of the unnormalized Laplacian is not clear, while the normalized Laplacian always converges under general conditions. There are examples where the top eigenvectors of the unnormalized Laplacian do not yield a sensible clustering. The corresponding problem in semi-supervised classification needs further study. One reason is that in semi-supervised learning the whole Laplacian (normalized or not) is often used for regularization, not only the top eigenvectors.

6.6 Directed Graphs

For semi-supervised learning on directed graphs, Zhou et al. (2005b) take a hub - authority approach and essentially convert a directed graph into an undirected one. Two hub nodes are connected by an undirected edge with appropriate weight if they co-link to authority nodes, and vice versa. Semi-supervised learning then proceeds on the undirected graph.

Zhou et al. (2005a) generalize the work further. The algorithm takes a transition matrix (with a unique stationary distribution) as input, and gives a closed form solution on unlabeled data. The solution parallels and generalizes the normalized Laplacian solution for undirected graphs (Zhou et al., 2004a). The previous work

(Zhou et al., 2005b) is a special case with the 2-step random walk transition matrix. In the absence of labels, the algorithm is the generalization of the normalized cut (Shi & Malik, 2000) on directed graphs.

Lu and Getoor (2003) convert the link structure in a directed graph into per-node features, and combines them with per-node object features in logistic regression. They also use an EM-like iterative algorithm.

6.7 Connection to Standard Graphical Models

The Gaussian random field formulation (Zhu et al., 2003a) is a standard undirected graphical model, with continuous random variables. Given labeled nodes (observed variables), the inference is used to obtain the mean (equivalently the mode) h_i of the remaining variables, which is the harmonic function. However the interpretation of the harmonic function as parameters for Bernoulli distributions at the nodes (i.e. each unlabeled node has label 1 with probability h_i , 0 otherwise) is non-standard.

7 Metric-Based Model Selection

Metric-based model selection (Schuurmans & Southey, 2001) is a method to detect hypotheses inconsistency with unlabeled data. We may have two hypotheses which are consistent on L , for example they all have zero training set error. However they may be inconsistent on the much larger U . If so we should reject at least one of them, e.g. the more complex one if we employ Occam’s razor.

The key observation is that a distance metric is defined in the hypothesis space H . One such metric is the number of different classifications two hypotheses make under the data distribution $p(x)$: $d_p(h_1, h_2) = E_p[h_1(x) \neq h_2(x)]$. It is easy to verify that the metric satisfies the three metric properties. Now consider the true classification function h^* and two hypotheses h_1, h_2 . Since the metric satisfies the triangle inequality (the third property), we have

$$d_p(h_1, h_2) \leq d_p(h_1, h^*) + d_p(h^*, h_2)$$

Under the premise that labels in L is noiseless, let’s assume we can approximate $d_p(h_1, h^*)$ and $d_p(h^*, h_2)$ by h_1 and h_2 ’s training set error rates $d_L(h_1, h^*)$ and $d_L(h_2, h^*)$, and approximate $d_p(h_1, h_2)$ by the difference h_1 and h_2 make on a large amount of unlabeled data U : $d_U(h_1, h_2)$. We get

$$d_U(h_1, h_2) \leq d_L(h_1, h^*) + d_L(h^*, h_2)$$

which can be verified directly. If the inequality does not hold, at least one of the assumptions is wrong. If $|U|$ is large enough and $U \stackrel{\text{iid}}{\sim} p(x)$, $d_U(h_1, h_2)$ will be a good estimate of $d_p(h_1, h_2)$. This leaves us with the conclusion that at least one of the training errors does not reflect its true error. If both training errors are close to zero, we would know that at least one model is overfitting. An Occam’s razor type of argument then can be used to select the model with less complexity. Such use of unlabeled data is very general and can be applied to almost any learning algorithms. However it only selects among hypotheses; it does not generate new hypothesis based on unlabeled data.

The co-validation method (Madani et al., 2005) also uses unlabeled data for model selection and active learning. Kaariainen (2005) uses the metric to derive a generalization error bound, see Section 8.

8 Computational Learning Theory

In this survey we have primarily focused on various semi-supervised learning algorithms. The theory of semi-supervised learning has been touched upon occasionally in the literature. However it was not until recently that the computational learning theory community began to pay more attention to this interesting problem.

Leskes (2005) presents a generalization error bound for semi-supervised learning with multiple learners, an extension to co-training. The author shows that if multiple learning algorithms are forced to produce similar hypotheses (i.e. to agree) given the same training set, and such hypotheses still have low training error, then the generalization error bound is tighter. The unlabeled data is used to assess the agreement among hypotheses. The author proposes a new Agreement-Boost algorithm to implement the procedure.

Kaariainen (2005) presents another generalization error bound for semi-supervised learning. The idea is that the target function is in the version space. If a hypothesis is in the version space (revealed by labeled data), and is close to all other hypotheses in the version space (revealed by unlabeled data), then it has to be close to the target function. Closeness is defined as classification agreement, and can be approximated using unlabeled data. This idea builds on metric-based model selection (Section 7).

Balcan and Blum (2005) propose a PAC-style model for semi-supervised learning. This is the first PAC model that explains when unlabeled data might help (notice the classic PAC model cannot incorporate unlabeled data at all). There has been previous *particular* analysis for explaining when unlabeled data helps, but they were all based on specific settings and assumptions. In contrast this PAC model is a general, unifying model. The authors define an interesting quantity:

the compatibility of a hypothesis w.r.t. the unlabeled data distribution. For example in SVM a hyperplane that cuts through high density regions would have low compatibility, while one that goes along gaps would have high compatibility. We note that the compatibility function can be defined much more generally. The intuition of the results is the following. Assuming a-priori that the target function has high compatibility with unlabeled data. Then if a hypothesis has zero training error (standard PAC style) *and* high compatibility, the theory gives the number of labeled and unlabeled data to guarantee the hypothesis is good. The number of labeled data needed can be quite small.

9 Related Areas

The focus of the survey is on classification with semi-supervised methods. There are some closely related areas with a rich literature.

9.1 Spectral Clustering

Spectral clustering is unsupervised. As such there is no labeled data to guide the process. Instead the clustering depends solely on the graph weights W . On the other hand semi-supervised learning for classification has to maintain a balance between how good the ‘clustering’ is, and how well the labeled data can be explained by it. Such balance is expressed explicitly in the regularization framework.

As we have seen in section 8.1 of (Zhu, 2005) and section 6.5 here, the top eigenvectors of the graph Laplacian can unfold the data manifold to form meaningful clusters. This is the intuition behind spectral clustering. There are several criteria on what constitutes a good clustering (Weiss, 1999).

The normalized cut (Shi & Malik, 2000) seeks to minimize

$$Ncut(A, B) = \frac{cut(A, B)}{assoc(A, V)} + \frac{cut(A, B)}{assoc(B, V)} \quad (15)$$

The *continuous relaxation* of the cluster indicator vector can be derived from the normalized Laplacian. In fact it is derived from the second smallest eigenvector of the normalized Laplacian. The continuous vector is then discretized to obtain the clusters.

The data points are mapped into a new space spanned by the first k eigenvectors of the normalized Laplacian in (Ng et al., 2001), with special normalization. Clustering is then performed with traditional methods (like k-means) in this new space. This is very similar to kernel PCA.

Fowlkes et al. (2004) use the Nyström method to reduce the computation cost for large spectral clustering problems. This is related to the method in (Zhu, 2005) Chapter 10.

Chung (1997) presents the mathematical details of spectral graph theory.

9.2 Learning with Positive and Unlabeled Data

In many real world applications, labeled data may be available from only one of the two classes. Then there is the unlabeled data, known to contain both classes. There are two ways to formulate the problem: classification or ranking.

Classification Here one builds a classifier even though there is no negative example. It is important to note that with the positive training data one can estimate the positive class conditional probability $p(x|+)$, and with the unlabeled data one can estimate $p(x)$. If the prior $p(+)$ is known or estimated from other sources, one can derive the negative class conditional as

$$p(x|-) = \frac{p(x) - p(+)p(x|+)}{1 - p(+)} \quad (16)$$

With $p(x|-)$ one can then perform classification with Bayes rule. Denis et al. (2002) use this fact for text classification with Naive Bayes models.

Another set of methods heuristically identify some ‘reliable’ negative examples in the unlabeled set, and use EM on generative (Naive Bayes) models (Liu et al., 2002) or logistic regression (Lee & Liu, 2003).

Ranking Given a large collection of items, and a few ‘query’ items, ranking orders the items according to their similarity to the queries. Information retrieval is the standard technique under this setting, and we will not attempt to include the extensive literatures on this mature field. It is worth pointing out that graph-based semi-supervised learning can be modified for such settings. Zhou et al. (2004b) treat it as semi-supervised learning with positive data on a graph, where the graph induces a similarity measure, and the queries are positive examples. Data points are ranked according to their graph similarity to the positive training set.

9.3 Clustering with Side Information

This is the ‘opposite’ of semi-supervised classification. The goal is clustering but there are some ‘labeled data’ in the form of *must-links* (two points must in the same cluster) and *cannot-links* (two points cannot in the same cluster). There is a tension between satisfying these constraints and optimizing the original clustering criterion (e.g. minimizing the sum of squared distances within clusters). Procedurally one can modify the distance metric to try to accommodate the constraints, or one can

bias the search. We refer readers to a recent short survey (Grira et al., 2004) for the literatures.

9.4 Nonlinear Dimensionality Reduction

The goal of nonlinear dimensionality reduction is to find a faithful low dimensional mapping of the high dimensional data. As such it belongs to unsupervised learning. However the way it discovers low dimensional manifold within a high dimensional space is closely related to spectral graph semi-supervised learning. Representative methods include Isomap (Tenenbaum et al., 2000), locally linear embedding (LLE) (Roweis & Saul, 2000) (Saul & Roweis, 2003), Hessian LLE (Donoho & Grimes, 2003), Laplacian eigenmaps (Belkin & Niyogi, 2003), and semidefinite embedding (SDE) (Weinberger & Saul, 2004) (Weinberger et al., 2004) (Weinberger et al., 2005).

9.5 Learning a Distance Metric

Many learning algorithms depend, either explicitly or implicitly, on a distance metric on X . We use the term metric here loosely to mean a measure of distance or (dis)similarity between two data points. The default distance in the feature space may not be optimal, especially when the data forms a lower dimensional manifold in the feature vector space. With a large amount of U , it is possible to detect such manifold structure and its associated metric. The graph-based methods above are based on this principle. We review some other methods next.

The simplest example in text classification might be Latent Semantic Indexing (LSI, a.k.a. Latent Semantic Analysis LSA, Principal Component Analysis PCA, or sometimes Singular Value Decomposition SVD). This technique defines a linear subspace, such that the variance of the data, when projected to the subspace, is maximumly preserved. LSI is widely used in text classification, where the original space for X is usually tens of thousands dimensional, while people believe meaningful text documents reside in a much lower dimensional space. Zelikovitz and Hirsh (2001) and Cristianini et al. (2001) both use U , in this case unlabeled documents, to augment the term-by-document matrix of L . LSI is performed on the augmented matrix. This representation induces a new distance metric. By the property of LSI, words that co-occur very often in the same documents are merged into a single dimension of the new space. In the extreme this allows two documents with no common words to be ‘close’ to each other, via chains of co-occur word pairs in other documents.

Oliveira et al. (2005) propose a simple procedure for semi-supervised learning: First one runs PCA on $L \cup U$ (ignoring the labels). The result is a linear subspace

that is constructed with more data points if one uses only L in PCA. In the next step, only L is mapped onto the subspace, and an SVM is learned. The method is useful when class separation is linear and along the principal component directions, and unlabeled helps by reducing the variance in estimating such directions.

Probabilistic Latent Semantic Analysis (PLSA) (Hofmann, 1999) is an important improvement over LSI. Each word in a document is generated by a ‘topic’ (a multinomial, i.e. unigram). Different words in the document may be generated by different topics. Each document in turn has a fixed topic proportion (a multinomial on a higher level). However there is no link between the topic proportions in different documents.

Latent Dirichlet Allocation (LDA) (Blei et al., 2003) is one step further. It assumes the topic proportion of each document is drawn from a Dirichlet distribution. With variational approximation, each document is represented by a posterior Dirichlet over the topics. This is a much lower dimensional representation. Griffiths et al. (2005) extend LDA model to ‘HMM-LDA’ which uses both short-term syntactic and long-term topical dependencies, as an effort to integrate semantics and syntax. Li and McCallum (2005) apply the HMM-LDA model to obtain word clusters, as a rudimentary way for semi-supervised learning on sequences.

Some algorithms derive a metric entirely from the density of U . These are motivated by unsupervised clustering and based on the intuition that data points in the same high density ‘clump’ should be close in the new metric. For instance, if U is generated from a single Gaussian, then the Mahalanobis distance induced by the covariance matrix is such a metric. Tipping (1999) generalizes the Mahalanobis distance by fitting U with a mixture of Gaussian, and define a Riemannian manifold with metric at x being the weighted average of individual component inverse covariance. The distance between x_1 and x_2 is computed along the straight line (in Euclidean space) between the two points. Rattray (2000) further generalizes the metric so that it only depends on the change in log probabilities of the density, not on a particular Gaussian mixture assumption. And the distance is computed along a curve that minimizes the distance. The new metric is invariant to linear transformation of the features, and connected regions of relatively homogeneous density in U will be close to each other. Such metric is attractive, yet it depends on the homogeneity of the initial Euclidean space. Their application in semi-supervised learning needs further investigation.

We caution the reader that the metrics proposed above are based on unsupervised techniques. They all identify a lower dimensional manifold within which the data reside. However the data manifold may or may not correlate with a particular classification task. For example, in LSI the new metric emphasizes words with prominent count variances, but ignores words with small variances. If the classification task is subtle and depends on a few words with small counts, LSI might

wipe out the salient words all together. Therefore the success of these methods is hard to guarantee without putting some restrictions on the kind of classification tasks. It would be interesting to include L into the metric learning process.

In a separate line of work, Baxter (1997) proves that there is a unique optimal metric for classification if we use 1-nearest-neighbor. The metric, named Canonical Distortion Measure (CDM), defines a distance $d(x_1, x_2)$ as the expected loss if we classify x_1 with x_2 's label. The distance measure proposed in (Yianilos, 1995) can be viewed as a special case. Yianilos assume a Gaussian mixture model has been learned from U , such that a class correspond to a component, but the correspondence is unknown. In this case CDM $d(x_1, x_2) = p(x_1, x_2 \text{ from same component})$ and can be computed analytically. Now that a metric has been learned from U , we can find within L the 1-nearest-neighbor of a new data point x , and classify x with the nearest neighbor's label. It will be interesting to compare this scheme with EM based semi-supervised learning, where L is used to label mixture components.

Weston et al. (2004) propose the neighborhood mismatch kernel and the bagged mismatch kernel. More precisely both are *kernel transformation* that modifies an input kernel. In the neighborhood method, one defines the neighborhood of a point as points close enough according to certain similarity measure (note this is *not* the measure induced by the input kernel). The output kernel between point i, j is the average of pairwise kernel entries between i 's neighbors and j 's neighbors. In bagged method, if a clustering algorithm thinks they tend to be in the same cluster (note again this is a different measure than the input kernel), the corresponding entry in the input kernel is boosted.

9.6 Inferring Label Sampling Mechanisms

Most semi-supervised learning methods assume L and U are both *i.i.d.* from the underlying distribution. However as (Rosset et al., 2005) points out that is not always the case. For example y can be the binary label whether a customer is satisfied, obtained through a survey. It is conceivable survey participation (and thus labeled data) depends on the satisfaction y .

Let s_i be the binary missing indicator for y_i . The authors model $p(s|x, y)$ with a parametric family. The goal is to estimate $p(s|x, y)$ which is the label sampling mechanism. This is done by computing the expectation of an arbitrary function $g(x)$ in two ways: on $L \cup U$ as $1/n \sum_{i=1}^n g(x_i)$, and on L only as $1/n \sum_{i \in L} g(x_i) / p(s_i = 1 | x_i, y_i)$. By equating the two $p(s|x, y)$ can be estimated. The intuition is that the expectation on L requires weighting the labeled samples inversely proportional to the labeling probability, to compensate for ignoring the unlabeled data.

Acknowledgment

I thank John Lafferty, Zoubin Ghahramani, Tommi Jaakkola, Ronald Rosenfeld, Maria Florina Balcan, Kai Yu, Sajama, Matthias Seeger, Yunpeng Xu and all other colleagues who discussed the literature with me.

References

- Argyriou, A. (2004). Efficient approximation methods for harmonic semi-supervised learning. Master's thesis, University College London.
- Balcan, M.-F., & Blum, A. (2005). A PAC-style model for learning from labeled and unlabeled data. *COLT 2005*.
- Balcan, M.-F., Blum, A., Choi, P. P., Lafferty, J., Pantano, B., Rwebangira, M. R., & Zhu, X. (2005a). Person identification in webcam images: An application of semi-supervised learning. *ICML2005 Workshop on Learning with Partially Classified Training Data*.
- Balcan, M.-F., Blum, A., & Yang, K. (2005b). Co-training and expansion: Towards bridging theory and practice. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Baluja, S. (1998). Probabilistic modeling for face orientation discrimination: Learning from labeled and unlabeled data. *Neural Information Processing Systems*.
- Baxter, J. (1997). The canonical distortion measure for vector quantization and function approximation. *Proc. 14th International Conference on Machine Learning* (pp. 39–47). Morgan Kaufmann.
- Belkin, M., Matveeva, I., & Niyogi, P. (2004a). Regularization and semi-supervised learning on large graphs. *COLT*.
- Belkin, M., & Niyogi, P. (2003). Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation*, 15, 1373–1396.
- Belkin, M., Niyogi, P., & Sindhwani, V. (2004b). *Manifold regularization: A geometric framework for learning from examples* (Technical Report TR-2004-06). University of Chicago.

- Belkin, M., Niyogi, P., & Sindhvani, V. (2005). On manifold regularization. *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics (AISTAT 2005)*.
- Bennett, K., & Demiriz, A. (1999). Semi-supervised support vector machines. *Advances in Neural Information Processing Systems, 11*, 368–374.
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research, 3*, 993–1022.
- Blum, A., & Chawla, S. (2001). Learning from labeled and unlabeled data using graph mincuts. *Proc. 18th International Conf. on Machine Learning*.
- Blum, A., Lafferty, J., Rwebangira, M., & Reddy, R. (2004). Semi-supervised learning using randomized mincuts. *ICML-04, 21st International Conference on Machine Learning*.
- Blum, A., & Mitchell, T. (1998). Combining labeled and unlabeled data with co-training. *COLT: Proceedings of the Workshop on Computational Learning Theory*.
- Bousquet, O., Chapelle, O., & Hein, M. (2004). Measure based regularization. *Advances in Neural Information Processing Systems 16*.
- Callison-Burch, C., Talbot, D., & Osborne, M. (2004). Statistical machine translation with word- and sentence-aligned parallel corpora. *Proceedings of the ACL*.
- Carreira-Perpinan, M. A., & Zemel, R. S. (2005). Proximity graphs for clustering and manifold learning. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Castelli, V., & Cover, T. (1995). The exponential value of labeled samples. *Pattern Recognition Letters, 16*, 105–111.
- Castelli, V., & Cover, T. (1996). The relative value of labeled and unlabeled samples in pattern recognition with an unknown mixing parameter. *IEEE Transactions on Information Theory, 42*, 2101–2117.
- Chapelle, O., Weston, J., & Schölkopf, B. (2002). Cluster kernels for semi-supervised learning. *Advances in Neural Information Processing Systems, 15*.
- Chapelle, O., & Zien, A. (2005). Semi-supervised classification by low density separation. *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics (AISTAT 2005)*.

- Chu, W., & Ghahramani, Z. (2004). *Gaussian processes for ordinal regression* (Technical Report). University College London.
- Chung, F. R. K. (1997). *Spectral graph theory, regional conference series in mathematics, no. 92*. American Mathematical Society.
- Corduneanu, A., & Jaakkola, T. (2001). *Stable mixing of complete and incomplete information* (Technical Report AIM-2001-030). MIT AI Memo.
- Corduneanu, A., & Jaakkola, T. (2003). On information regularization. *Nineteenth Conference on Uncertainty in Artificial Intelligence (UAI03)*.
- Corduneanu, A., & Jaakkola, T. S. (2005). Distributed information regularization on graphs. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Cozman, F., Cohen, I., & Cirelo, M. (2003). Semi-supervised learning of mixture models. *ICML-03, 20th International Conference on Machine Learning*.
- Cristianini, N., Shawe-Taylor, J., & Lodhi, H. (2001). Latent semantic kernels. *Proc. 18th International Conf. on Machine Learning*.
- Dara, R., Kremer, S., & Stacey, D. (2000). Clustering unlabeled data with SOMs improves classification of labeled real-world data. submitted.
- Delalleau, O., Bengio, Y., & Roux, N. L. (2005). Efficient non-parametric function induction in semi-supervised learning. *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics (AISTAT 2005)*.
- Demirez, A., & Bennett, K. (2000). Optimization approaches to semisupervised learning. In M. Ferris, O. Mangasarian and J. Pang (Eds.), *Applications and algorithms of complementarity*. Boston: Kluwer Academic Publishers.
- Demiriz, A., Bennett, K., & Embrechts, M. (1999). Semi-supervised clustering using genetic algorithms. *Proceedings of Artificial Neural Networks in Engineering*.
- Dempster, A., Laird, N., & Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*.
- Denis, F., Gilleron, R., & Tommasi, M. (2002). Text classification from positive and unlabeled examples. *The 9th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems (IPMU)*.

- Donoho, D. L., & Grimes, C. E. (2003). Hessian eigenmaps: locally linear embedding techniques for high-dimensional data. *Proceedings of the National Academy of Arts and Sciences*, 100, 5591–5596.
- Fowlkes, C., Belongie, S., Chung, F., & Malik, J. (2004). Spectral grouping using the Nyström method. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26, 214–225.
- Fujino, A., Ueda, N., & Saito, K. (2005). A hybrid generative/discriminative approach to semi-supervised classifier design. *AAAI-05, The Twentieth National Conference on Artificial Intelligence*.
- Fung, G., & Mangasarian, O. (1999). *Semi-supervised support vector machines for unlabeled data classification* (Technical Report 99-05). Data Mining Institute, University of Wisconsin Madison.
- Garcke, J., & Griebel, M. (2005). Semi-supervised learning with sparse grids. *Proc. of the 22nd ICML Workshop on Learning with Partially Classified Training Data*. Bonn, Germany.
- Getz, G., Shental, N., & Domany, E. (2005). Semi-supervised learning – a statistical physics approach. *Proc. of the 22nd ICML Workshop on Learning with Partially Classified Training Data*. Bonn, Germany.
- Goldman, S., & Zhou, Y. (2000). Enhancing supervised learning with unlabeled data. *Proc. 17th International Conf. on Machine Learning* (pp. 327–334). Morgan Kaufmann, San Francisco, CA.
- Grady, L., & Funka-Lea, G. (2004). Multi-label image segmentation for medical applications based on graph-theoretic electrical potentials. *ECCV 2004 workshop*.
- Grandvalet, Y., & Bengio, Y. (2005). Semi-supervised learning by entropy minimization. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Griffiths, T. L., Steyvers, M., Blei, D. M., & Tenenbaum, J. B. (2005). Integrating topics and syntax. *NIPS 17*.
- Girra, N., Crucianu, M., & Boujemaa, N. (2004). Unsupervised and semi-supervised clustering: a brief survey. in ‘A Review of Machine Learning Techniques for Processing Multimedia Content’, Report of the MUSCLE European Network of Excellence (FP6).

- Hofmann, T. (1999). Probabilistic latent semantic analysis. *Proc. of Uncertainty in Artificial Intelligence, UAI'99*. Stockholm.
- Jaakkola, T., Meila, M., & Jebara, T. (1999). Maximum entropy discrimination. *Neural Information Processing Systems, 12*, 12.
- Joachims, T. (1999). Transductive inference for text classification using support vector machines. *Proc. 16th International Conf. on Machine Learning* (pp. 200–209). Morgan Kaufmann, San Francisco, CA.
- Joachims, T. (2003). Transductive learning via spectral graph partitioning. *Proceedings of ICML-03, 20th International Conference on Machine Learning*.
- Jones, R. (2005). *Learning to extract entities from labeled and unlabeled text* (Technical Report CMU-LTI-05-191). Carnegie Mellon University. Doctoral Dissertation.
- Kaariainen, M. (2005). Generalization error bounds using unlabeled data. *COLT 2005*.
- Kemp, C., Griffiths, T., Stromsten, S., & Tenenbaum, J. (2003). Semi-supervised learning with trees. *Advances in Neural Information Processing System 16*.
- Kondor, R. I., & Lafferty, J. (2002). Diffusion kernels on graphs and other discrete input spaces. *Proc. 19th International Conf. on Machine Learning*.
- Krishnapuram, B., Williams, D., Xue, Y., Hartemink, A., Carin, L., & Figueiredo, M. (2005). On semi-supervised classification. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Lawrence, N. D., & Jordan, M. I. (2005). Semi-supervised learning via Gaussian processes. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Lee, W. S., & Liu, B. (2003). Learning with positive and unlabeled examples using weighted logistic regression. *Proceedings of the Twentieth International Conference on Machine Learning (ICML)*.
- Leskes, B. (2005). The value of agreement, a new boosting algorithm. *COLT 2005*.
- Levin, A., Lischinski, D., & Weiss, Y. (2004). Colorization using optimization. *ACM Transactions on Graphics*.

- Li, W., & McCallum, A. (2005). Semi-supervised sequence modeling with syntactic topic models. *AAAI-05, The Twentieth National Conference on Artificial Intelligence*.
- Liu, B., Lee, W. S., Yu, P. S., & Li, X. (2002). Partially supervised classification of text documents. *Proceedings of the Nineteenth International Conference on Machine Learning (ICML)*.
- Lu, Q., & Getoor, L. (2003). Link-based classification using labeled and unlabeled data. *ICML 2003 workshop on The Continuum from Labeled to Unlabeled Data in Machine Learning and Data Mining*.
- Madani, O., Pennock, D. M., & Flake, G. W. (2005). Co-validation: Using model disagreement to validate classification algorithms. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Maeireizo, B., Litman, D., & Hwa, R. (2004). Co-training for predicting emotions with spoken dialogue data. *The Companion Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Mahdaviani, M., de Freitas, N., Fraser, B., & Hamze, F. (2005). Fast computational methods for visually guided robots. *The 2005 International Conference on Robotics and Automation (ICRA)*.
- McCallum, A., & Nigam, K. (1998). A comparison of event models for naive bayes text classification. *AAAI-98 Workshop on Learning for Text Categorization*.
- Miller, D., & Uyar, H. (1997). A mixture of experts classifier with learning based on both labelled and unlabelled data. *Advances in NIPS 9* (pp. 571–577).
- Mitchell, T. (1999). The role of unlabeled data in supervised learning. *Proceedings of the Sixth International Colloquium on Cognitive Science*. San Sebastian, Spain.
- Ng, A., Jordan, M., & Weiss, Y. (2001). On spectral clustering: Analysis and an algorithm. *Advances in Neural Information Processing Systems, 14*.
- Nigam, K. (2001). *Using unlabeled data to improve text classification* (Technical Report CMU-CS-01-126). Carnegie Mellon University. Doctoral Dissertation.
- Nigam, K., & Ghani, R. (2000). Analyzing the effectiveness and applicability of co-training. *Ninth International Conference on Information and Knowledge Management* (pp. 86–93).

- Nigam, K., McCallum, A. K., Thrun, S., & Mitchell, T. (2000). Text classification from labeled and unlabeled documents using EM. *Machine Learning*, 39, 103–134.
- Niu, Z.-Y., Ji, D.-H., & Tan, C.-L. (2005). Word sense disambiguation using label propagation based semi-supervised learning. *Proceedings of the ACL*.
- Oliveira, C. S., Cozman, F. G., & Cohen, I. (2005). Splitting the unsupervised and supervised components of semi-supervised learning. *Proc. of the 22nd ICML Workshop on Learning with Partially Classified Training Data*. Bonn, Germany.
- Pang, B., & Lee, L. (2004). A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. *Proceedings of the ACL* (pp. 271–278).
- Pham, T. P., Ng, H. T., & Lee, W. S. (2005). Word sense disambiguation with semi-supervised learning. *AAAI-05, The Twentieth National Conference on Artificial Intelligence*.
- Ratsaby, J., & Venkatesh, S. (1995). Learning from a mixture of labeled and unlabeled examples with parametric side information. *Proceedings of the Eighth Annual Conference on Computational Learning Theory*, 412–417.
- Ratray, M. (2000). A model-based distance for clustering. *Proc. of International Joint Conference on Neural Networks*.
- Riloff, E., Wiebe, J., & Wilson, T. (2003). Learning subjective nouns using extraction pattern bootstrapping. *Proceedings of the Seventh Conference on Natural Language Learning (CoNLL-2003)*.
- Rosenberg, C., Hebert, M., & Schneiderman, H. (2005). Semi-supervised self-training of object detection models. *Seventh IEEE Workshop on Applications of Computer Vision*.
- Rosset, S., Zhu, J., Zou, H., & Hastie, T. (2005). A method for inferring label sampling mechanisms in semi-supervised learning. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Roweis, S. T., & Saul, L. K. (2000). Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290, 2323–2326.
- Sajama, & Orlitsky, A. (2005). Estimating and computing density based distance metrics. *ICML05, 22nd International Conference on Machine Learning*. Bonn, Germany.

- Saul, L. K., & Roweis, S. T. (2003). Think globally, fit locally: unsupervised learning of low dimensional manifolds. *Journal of Machine Learning Research*, 4, 119–155.
- Schuermans, D., & Southey, F. (2001). Metric-based methods for adaptive model selection and regularization. *Machine Learning, Special Issue on New Methods for Model Selection and Model Combination*, 48, 51–84.
- Seeger, M. (2001). *Learning with labeled and unlabeled data* (Technical Report). University of Edinburgh.
- Shahshahani, B., & Landgrebe, D. (1994). The effect of unlabeled samples in reducing the small sample size problem and mitigating the Hughes phenomenon. *IEEE Trans. On Geoscience and Remote Sensing*, 32, 1087–1095.
- Shi, J., & Malik, J. (2000). Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22, 888–905.
- Sindhwani, V., Niyogi, P., & Belkin, M. (2005a). Beyond the point cloud: from transductive to semi-supervised learning. *ICML05, 22nd International Conference on Machine Learning*. Bonn, Germany.
- Sindhwani, V., Niyogi, P., Belkin, M., & Keerthi, S. (2005b). Linear manifold regularization for large scale semi-supervised learning. *Proc. of the 22nd ICML Workshop on Learning with Partially Classified Training Data*. Bonn, Germany.
- Smola, A., & Kondor, R. (2003). Kernels and regularization on graphs. *Conference on Learning Theory, COLT/KW*.
- Szummer, M., & Jaakkola, T. (2001). Partially labeled classification with Markov random walks. *Advances in Neural Information Processing Systems*, 14.
- Szummer, M., & Jaakkola, T. (2002). Information regularization with partially labeled data. *Advances in Neural Information Processing Systems*, 15.
- Teh, Y. W., & Roweis, S. (2002). Automatic alignment of local representations. *Advances in NIPS*.
- Tenenbaum, J. B., de Silva, V., & Langford, J. C. (2000). A global geometric framework for nonlinear dimensionality reduction. *Science*, 290, 2319–2323.
- Tipping, M. (1999). Deriving cluster analytic distance functions from Gaussian mixture models.
- Vapnik, V. (1998). *Statistical learning theory*. Springer.

- von Luxburg, U., Belkin, M., & Bousquet, O. (2004). *Consistency of spectral clustering* (Technical Report TR-134). Max Planck Institute for Biological Cybernetics.
- von Luxburg, U., Bousquet, O., & Belkin, M. (2005). Limits of spectral clustering. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Weinberger, K. Q., Packer, B. D., & Saul, L. K. (2005). Nonlinear dimensionality reduction by semidefinite programming and kernel matrix factorization. *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics (AISTAT 2005)*.
- Weinberger, K. Q., & Saul, L. K. (2004). Unsupervised learning of image manifolds by semidefinite programming. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 988–995).
- Weinberger, K. Q., Sha, F., & Saul, L. K. (2004). Learning a kernel matrix for nonlinear dimensionality reduction. *Proceedings of ICML-04* (pp. 839–846).
- Weiss, Y. (1999). Segmentation using eigenvectors: A unifying view. *ICCV (2)* (pp. 975–982).
- Weston, J., Leslie, C., Zhou, D., Elisseeff, A., & Noble, W. S. (2004). Semi-supervised protein classification using cluster kernels. In S. Thrun, L. Saul and B. Schölkopf (Eds.), *Advances in neural information processing systems 16*. Cambridge, MA: MIT Press.
- Xu, L., & Schuurmans, D. (2005). Unsupervised and semi-supervised multi-class support vector machines. *AAAI-05, The Twentieth National Conference on Artificial Intelligence*.
- Yarowsky, D. (1995). Unsupervised word sense disambiguation rivaling supervised methods. *Proceedings of the 33rd Annual Meeting of the Association for Computational Linguistics* (pp. 189–196).
- Yianilos, P. (1995). *Metric learning via normal mixtures* (Technical Report). NEC Research Institute.
- Yu, K., Tresp, V., & Zhou, D. (2004). *Semi-supervised induction with basis functions* (Technical Report 141). Max Planck Institute for Biological Cybernetics, Tübingen, Germany.

- Yu, K., Yu, S., & Tresp, V. (2005). Blockwise supervised inference on large graphs. *Proc. of the 22nd ICML Workshop on Learning with Partially Classified Training Data*. Bonn, Germany.
- Zelikovitz, S., & Hirsh, H. (2001). Improving text classification with LSI using background knowledge. *IJCAI01 Workshop Notes on Text Learning: Beyond Supervision*.
- Zhang, T., & Oles, F. J. (2000). A probability analysis on the value of unlabeled data for classification problems. *Proc. 17th International Conf. on Machine Learning* (pp. 1191–1198). Morgan Kaufmann, San Francisco, CA.
- Zhou, D., Bousquet, O., Lal, T., Weston, J., & Schölkopf, B. (2004a). Learning with local and global consistency. *Advances in Neural Information Processing System 16*.
- Zhou, D., Huang, J., & Schölkopf, B. (2005a). Learning from labeled and unlabeled data on a directed graph. *ICML05, 22nd International Conference on Machine Learning*. Bonn, Germany.
- Zhou, D., Schölkopf, B., & Hofmann, T. (2005b). Semi-supervised learning on directed graphs. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.
- Zhou, D., Weston, J., Gretton, A., Bousquet, O., & Schölkopf, B. (2004b). Ranking on data manifolds. *Advances in Neural Information Processing System 16*.
- Zhu, X. (2005). *Semi-supervised learning with graphs*. Doctoral dissertation, Carnegie Mellon University. CMU-LTI-05-192.
- Zhu, X., & Ghahramani, Z. (2002). *Towards semi-supervised classification with Markov random fields* (Technical Report CMU-CALD-02-106). Carnegie Mellon University.
- Zhu, X., Ghahramani, Z., & Lafferty, J. (2003a). Semi-supervised learning using Gaussian fields and harmonic functions. *ICML-03, 20th International Conference on Machine Learning*.
- Zhu, X., Kandola, J., Ghahramani, Z., & Lafferty, J. (2005). Nonparametric transforms of graph kernels for semi-supervised learning. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*. Cambridge, MA: MIT Press.

- Zhu, X., & Lafferty, J. (2005). Harmonic mixtures: combining mixture models and graph-based methods for inductive and scalable semi-supervised learning. *ICML-05, 22nd International Conference on Machine Learning*.
- Zhu, X., Lafferty, J., & Ghahramani, Z. (2003b). *Semi-supervised learning: From Gaussian fields to Gaussian processes* (Technical Report CMU-CS-03-175). Carnegie Mellon University.