

Learning from Partially Labeled Data

Martin Szummer, Tommi Jaakkola & Tomaso Poggio

Artificial Intelligence Laboratory and
The Center for Biological and Computational Learning
Massachusetts Institute of Technology
Cambridge, Massachusetts 02139

<http://www.ai.mit.edu>



The Problem: Learning from data with both labeled training points (x, y pairs) and unlabeled training points (x alone). For the labeled points, supervised learning techniques apply, but they cannot take advantage of the unlabeled points. On the other hand, unsupervised techniques can model the unlabeled data distribution, but do not exploit the labels. Thus, this task falls between traditional supervised and unsupervised learning.

Motivation: Supervised learning performance improves with larger training data sets. Unfortunately, it is often infeasible to obtain labels for large training sets. Assigning labels can require expensive resources such as human labor or laboratory tests. In some cases ground truth labels are impossible to obtain, e.g. if the necessary measurements can no longer be made, or if the labels will be given only in the future. In contrast, unlabeled training data is frequently easy to obtain in large quantities, and can outnumber the amount of labeled data by a large factor. For example, it is expensive to collect image databases of only faces, but it is cheap to collect arbitrary imagery with occasional faces, e.g. by crawling the world wide web, or by pointing a video camera out the window.

There are also developmental motivations for studying the process of learning from partially data. Children acquire language mainly by listening and imitating, with very limited feedback from adults. Human beings also excel at other partially labeled learning tasks, such as visual discrimination with hyperacuity [2].

Previous Work: Learning from partially labeled data is not well understood from the theoretical perspective. Labeled data has been shown to be exponentially more useful than unlabeled data under certain assumptions [1]. Discriminative and generative learning architectures have been shown to take advantage of unlabeled data in different settings [7].

Several practical algorithms for partially labeled data have been proposed. The EM algorithm is the most widely used. EM iterates between estimating model parameters and inferring soft labels for unlabeled points. It assumes a generative model (typically a mixture of Gaussians) and has been applied to text classification [4] and face pose determination. Unfortunately, when the data does not match EM's generative assumptions, the algorithm goes astray, and the information from labeled data is overwhelmed by unlabeled points.

Co-training [4] classifies data that exhibits factorial structure, specifically when its attributes can be partitioned into groups that are individually sufficient for learning but mutually independent. Co-training works by feeding outputs of one learner to be fed as examples for the other and vice-versa. However, in practice the individual learners are noisy, and then co-training easily veers down the wrong path.

Transduction [6] is a technique that attempts to maximize the classification margin on both labeled and unlabeled data, while classifying the labeled data as correctly as possible (Figure 1). This discriminative method imposes fewer restrictions on the data model. However, finding the optimal decision boundary requires solving a mixed integer programming problem that is NP-complete. The maximum entropy discrimination framework [3] also optimizes the margin based on both labeled and unlabeled data. However, instead of setting the margin directly, a distribution over margins is optimized to lie as close as possible to a prior. Likewise, labels of points are optimized to lie as close as possible to another prior; priors for labeled points peak around 0 or 1, whereas priors on unlabeled points are peaked at 0.5.

Approach: I am studying the problem of learning from partially labeled data from a theoretical and a practical perspective.

1. For the theoretical part I am characterizing the principles through which unlabeled data can and cannot be

used. For example, discriminative techniques that factor the distribution as $p(x, y|\alpha) = p(x)p(y|x, \alpha)$ cannot directly use the unlabeled data, because an improved knowledge of $p(x)$ does not help in estimating the model parameters α . On the other hand, models that assume $p(x, y|\alpha) = p(x|\alpha)p(y|x, \alpha)$ can immediately take advantage of partially labeled data.

At the same time, I am also characterizing the principles through which existing algorithms use partially labeled data. An improved understanding of these algorithms enables predictions on what data sets they will work. I am also proposing new data representations and algorithms that better exploit partially labeled data [5].

- For the practical part, I am implementing the main representations (kernel expansion, markov diffusion) with matching algorithms (maximum relative entropy, EM, transduction); this requires resolving the difficulties described above. I am benchmarking the algorithms on real applications, and describing their strengths and weaknesses. Application areas include computer vision (face detection and image classification), text classification and bioinformatics.

Difficulty: Existing techniques suffer from various limitations, including restrictive probabilistic assumptions, inability to deal with noisy data, being overwhelmed by unlabeled data, and computational intractability. Moreover, the properties of these algorithms are poorly understood.

Impact: Algorithms that learn from partially labeled data will be useful in all fields where unlabeled data is abundant but labels are difficult to obtain.

Future Work: I will address outstanding difficulties resulting from challenging partially labeled data problems. Examples include learning metrics, learning from very unbalanced classes, and different distributions for labeled and unlabeled data.

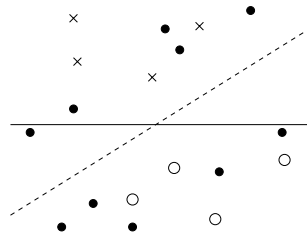


Figure 1: The decision boundary (solid line) based on labeled data (circles, crosses) has a smaller margin on all data than a boundary (dashed line) based on both labeled and unlabeled points (dots).

Research Support: Research at CBCL is sponsored by: Office of Naval Research (DARPA) under contract No. N00014-00-1-0907, National Science Foundation (ITR) under contract No. IIS-0085836, National Science Foundation (KDI) under contract No. DMS-9872936, and National Science Foundation under contract No. IIS-9800032. Additional support was provided by: Central Research Institute of Electric Power Industry, Center for e-Business (MIT), Eastman Kodak Company, DaimlerChrysler AG, Compaq, Honda R&D Co., Ltd., Komatsu Ltd., Merrill-Lynch, NEC Fund, Nippon Telegraph&Telephone, Siemens Corporate Research, Inc., Toyota Motor Corporation and The Whitaker Foundation.

References:

- [1] V. Castelli and T. Cover. The relative value of labeled and unlabeled samples in pattern recognition with an unknown mixing paramter. *IEEE T. Info. Theory*, 42:2102–2117, 1996.
- [2] M. Fahle, S. Edelman, and T. Poggio. Fast perceptual learning in visual hyperacuity. *Vision Research*, 35:3003–3013, 1995.
- [3] T. Jaakkola, M. Meila, and T. Jebara. Maximum entropy discrimination. In *nips*, volume 12, pages 470–476, 1999.
- [4] K. Nigam, A. McCallum, S. Thrun, and T. Mitchell. Text classification from labeled and unlabeled documents using em. *Machine Learning*, 39(2):103–134, 2000.

- [5] M. Szummer and T. Jaakkola. Kernel expansions with unlabeled data. In *Advances in Neural Information Processing Systems*, volume 13, pages 626–632, 2000.
- [6] V. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- [7] T. Zhang and F. Oles. A probability analysis on the value of unlabeled data for classification problems. In *Proc. of the Seventeenth International Conference on Machine Learning*, June 2000.