

Unsupervised Domain Adaptation by Transduction

Abstract

Abstract.

1. Introduction

Intro

2. Related Work

Related work

3. Method

3.1. Problem Definition

We are interested in the problem in which we have an unsupervised domain $\{\mathbf{x}_i\}_{i \in [N^u]}$ and a supervised domain $\{\hat{\mathbf{x}}_i, \hat{y}_i\}_{i \in [N^s]}$ such that \mathbf{x}_i is the extracted feature for point i and y_i is the corresponding label. We also assume that unsupervised and supervised features have the same dimension i.e. $\mathbf{x}, \hat{\mathbf{x}} \in \mathcal{R}^d$

We consider an asymmetric similarity metric;

$$s(\mathbf{x}_i, \hat{\mathbf{x}}_j) = \mathbf{x}_i^\top \mathbf{W} \hat{\mathbf{x}}_j \quad (1)$$

such that it is high if two points from supervised and unsupervised domains are similar to each other.

We approach to the problem from a transduction perspective; in other words, the main purpose of the method is recovering the labels y_i for each unsupervised example \mathbf{x}_i . We consider the following objective function in order to compute y_i as well as the similarity metric \mathbf{W} .

$$\begin{aligned} \min_{\mathbf{W}, y_i} \sum_{i \in [N^s]} [s(\hat{\mathbf{x}}_i, \mathbf{x}_{i^-}) - s(\hat{\mathbf{x}}_i, \mathbf{x}_{i^+}) + \alpha]_+ \\ s.t. \quad i^+ = \arg \max_{j|y_j = \hat{y}_i} s(\hat{\mathbf{x}}_i, \mathbf{x}_j) \\ i^- = \arg \max_{j|y_j \neq \hat{y}_i} s(\hat{\mathbf{x}}_i, \mathbf{x}_j) \end{aligned} \quad (2)$$

We solve this optimization problem via alternating minimization through iterating over solving for unsupervised labels y_i and the similarity metric \mathbf{W} . We explain these two steps the following sections.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

3.2. Labeling Unsupervised Points

We are using nearest neighbor rule in our labelling with an additional robustness metric. We first explain the nearest neighbor formulation then explain our extension.

Given a similarity metric \mathbf{W} , the labeling with nearest neighbor rule is;

$$(y_i)^{pred} = \hat{y}_{\arg \max_j s(\mathbf{x}_i, \hat{\mathbf{x}}_j)} \quad (3)$$

Although the nearest neighbor rule is computationally inefficient, we solve the efficiency issues through usage of batches with stochastic gradient descent as well as an efficient implementation though OpenBLAS routines. We discuss these details in the implementation section of the paper.

Although the nearest neighbor rule is expected to work with a good metric, at the initial stage of the algorithm the metric will not be accurate enough. Hence, we need a robustness measure to handle the initial stage of the algorithm. Our robustness measure comes from the consistency of labels over the unsupervised data graph.

We create a k-NN graph over the unsupervised data points such that $\mathcal{N}(\mathbf{x}_i)$ is the k unsupervised data point nearest to \mathbf{x}_i via the l_2 distance. After the k-NN graph is created, we solve the following optimization problem for labeling unsupervised data points;

$$\begin{aligned} \arg \min_{y_i} \sum_{i \in N^u} \min_{y_j = y_i} [1 - s(\hat{\mathbf{x}}_j, \mathbf{x}_i)] \\ + \lambda \sum_{i \in N^u} \sum_{j \in \mathcal{N}(\mathbf{x}_i)} [1 - \mathbf{x}_i^T \mathbf{x}_j] \mathbb{1}(y_i \neq y_j) \end{aligned} \quad (4)$$

ADD NORMALIZATION STUFF TO MAKE THIS THING NON-NEGATIVE

In its original form, this problem is multi-label energy minimization and it is submodular. However, solving this problem is still computationally inefficient. Therefore, we propose an approximate solution by minimizing the lower

bound of this algorithm.

$$\begin{aligned} & \min_{y_j=y_i} [1 - s(\hat{\mathbf{x}}_j, \mathbf{x}_i)] \\ & \geq \left(\min_{(y_j)^{pred}=y_i} [1 - s(\hat{\mathbf{x}}_j, \mathbf{x}_i)] \right) \mathbb{1}[(y_j)^{pred} = y_i] \quad (5) \\ & + \left(\min_{(y_j)^{pred} \neq y_i} [1 - s(\hat{\mathbf{x}}_j, \mathbf{x}_i)] \right) \mathbb{1}[(y_j)^{pred} \neq y_i] \end{aligned}$$

Moreover, the lower bound for the binary term can also be expressed as;

$$\begin{aligned} & [1 - \mathbf{x}_i^T \mathbf{x}_j] \mathbb{1}(y_i \neq y_j) \\ & \geq [1 - \mathbf{x}_i^T \mathbf{x}_j] \mathbb{1}[y_i = (y_i)^{pred}, y_j = (y_j)^{pred}, (y_i)^{pred} \neq (y_j)^{pred}] \\ & + [1 - \mathbf{x}_i^T \mathbf{x}_j] \mathbb{1}[y_i = (y_i)^{pred}, y_j \neq (y_j)^{pred}, (y_i)^{pred} = (y_j)^{pred}] \quad (6) \end{aligned}$$

In other words we convert the multi-label problem into binary problem with auxiliary variables $y_i^b = \mathbb{1}[y_i = (y_i)^{pred}]$ and the binary energy minimization problem becomes;

$$\begin{aligned} & \arg \min_{y_i^b} \sum_{i \in N^u} E_i(y_i^b) + \sum_{i \in N^u} \sum_{j \in \mathcal{N}(\mathbf{x}_i)} E_{i,j}(y_i^b, y_j^b) \\ & E_i(y_i^b) = \begin{cases} \min_{(y_j)^{pred}=y_i} [1 - s(\hat{\mathbf{x}}_j, \mathbf{x}_i)] & y_i^b = 1 \\ \min_{(y_j)^{pred} \neq y_i} [1 - s(\hat{\mathbf{x}}_j, \mathbf{x}_i)] & y_i^b = 0 \end{cases} \\ & E_{i,j}(y_i^b, y_j^b) = \begin{cases} 1 - \mathbf{x}_i^T \mathbf{x}_j & (y_i)^{pred} = (y_j)^{pred}, y_i^b \neq y_j^b \\ 1 - \mathbf{x}_i^T \mathbf{x}_j & (y_i)^{pred} \neq (y_j)^{pred}, y_i^b = y_j^b \\ 0 & o.w. \end{cases} \quad (7) \end{aligned}$$

Although this function is not sub-modular, it can still be very efficiently minimized via QPBO(quadratic pseudo boolean optimization)(?) internally using min-cut/max-flow. After the energy minimization, the unsupervised data points with label $y_i = 1$ keep their NN labels. For the rest, we sort them by using $\max_{j \in \mathcal{N}(\mathbf{x}_i)} \mathbf{x}_i^T \mathbf{x}_j$ and start from the maximum similarity one and greedily assign labels as;

$$y_i = \arg \max_{y_i} \max_{y_j=y_i} s(\hat{\mathbf{x}}_j, \mathbf{x}_i) + \lambda \sum_{j \in \mathcal{N}(\mathbf{x}_i) \cap P} \mathbf{x}_i^T \mathbf{x}_j \mathbb{1}(y_i \neq y_j) \quad (8)$$

where P is the set of points with label $y_i^b = 1$ and the points which are already processed via the greedy operation. We continue until all other points having the binary label $y_i^b = 0$ is labeled.

3.3. Learning Similarity Metric

Given the predicted labels y_i for unsupervised data points \mathbf{x}_i , we learn the asymmetric metric following the LMNN(Large Margin Nearest Neighbour)(?) construction. We start with finding the nearest positive and negative ex-

amples through;

$$\begin{aligned} i^+ &= \arg \max_{j|y_j=\hat{y}_i} s(\hat{\mathbf{x}}_i, \mathbf{x}_j) \\ i^- &= \arg \max_{j|y_j \neq \hat{y}_i} s(\hat{\mathbf{x}}_i, \mathbf{x}_j) \end{aligned} \quad (9)$$

Then, we construct the loss function with regularizer;

$$\min_{\mathbf{W}, y_i} \sum_{i \in [N^s]} [s(\hat{\mathbf{x}}_i, \mathbf{x}_{i^-}) - s(\hat{\mathbf{x}}_i, \mathbf{x}_{i^+}) + \alpha]_+ + r(\mathbf{W}) \quad (10)$$

which is convex in terms of the similarity metric \mathbf{W} if the regularizer is convex; and we optimize it by using Stochastic gradient descent through the subgradient $\frac{\partial \text{loss}(y_i, \mathbf{W})}{\partial \mathbf{W}} =$

$$\sum_{i \in [N^s]} \mathbb{1}[s(\hat{\mathbf{x}}_i, \mathbf{x}_{i^-}) - s(\hat{\mathbf{x}}_i, \mathbf{x}_{i^+}) > \alpha] (\hat{\mathbf{x}}_i \mathbf{x}_{i^-}^T - \hat{\mathbf{x}}_i \mathbf{x}_{i^+}^T) + \frac{\partial r(\mathbf{W})}{\partial \mathbf{W}} \quad (11)$$

As a regularizer we are using the Frobenius norm of the similarity matrix as $r(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}\|_F^2$

3.4. Learning Features

3.5. Learning with Synthetic Data

4. Convergence Analysis

It is monotonically decreasing after the each iteration

5. Experimental Results

6. Conclusion

Acknowledgements

a