



Online Transfer Learning [☆]



Peilin Zhao ^a, Steven C.H. Hoi ^{b,*}, Jialei Wang ^c, Bin Li ^d

^a Data Analytics Department, Institute for Infocomm Research, A*STAR, Singapore

^b School of Information Systems, Singapore Management University, Singapore

^c Department of Computer Science, The University of Chicago, USA

^d Department of Finance, Economics and Management School, Wuhan University, 430072, PR China

ARTICLE INFO

Article history:

Received 19 April 2012

Received in revised form 3 June 2014

Accepted 16 June 2014

Available online 17 July 2014

Keywords:

Transfer learning

Online learning

Knowledge transfer

ABSTRACT

In this paper, we propose a novel machine learning framework called “Online Transfer Learning” (OTL), which aims to attack an online learning task on a target domain by transferring knowledge from some source domain. We do not assume data in the target domain follows the same distribution as that in the source domain, and the motivation of our work is to enhance a supervised online learning task on a target domain by exploiting the existing knowledge that had been learnt from training data in source domains. OTL is in general very challenging since data in both source and target domains not only can be different in their class distributions, but also can be diverse in their feature representations. As a first attempt to this new research problem, we investigate two different settings of OTL: (i) OTL on homogeneous domains of common feature space, and (ii) OTL across heterogeneous domains of different feature spaces. For each setting, we propose effective OTL algorithms to solve online classification tasks, and show some theoretical bounds of the algorithms. In addition, we also apply the OTL technique to attack the challenging online learning tasks with concept-drifting data streams. Finally, we conduct extensive empirical studies on a comprehensive testbed, in which encouraging results validate the efficacy of our techniques.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

Transfer learning (TL) is an emerging family of machine learning techniques and has been actively studied in machine learning and AI communities in recent years [27]. In a regular transfer learning task, we assume two datasets, one from a source domain and the other from a target domain, are available where their data distributions or representations of the two domains can be very different. TL aims to build models from the target-domain dataset by exploring information from the source-domain dataset through some knowledge transferring process. Transfer learning is important for many applications where training data in a new domain may be limited or too expensive to collect. Despite being explored actively in literature [27,26,2,12,20], most existing approaches on transfer learning often have been studied in an offline/batch learning fashion, which assumes training data in the new domain is given a priori. Such an assumption may not always hold for some real applications where training examples may arrive in an online/sequential manner.

[☆] Code and datasets are available at <http://www.stevenhoi.org/OTL/>.

* Corresponding author.

E-mail addresses: zhaop@i2r.a-star.edu.sg (P. Zhao), chhoi@smu.edu.sg (S.C.H. Hoi), jialei@cs.uchicago.edu (J. Wang), binli.whu@whu.edu.cn (B. Li).

Unlike the existing transfer learning studies, this paper investigates a new framework of *Online Transfer Learning* (OTL) [33], which addresses the transfer learning problem in an online learning framework. Specifically, OTL makes two assumptions: (i) training data in the new domain arrives sequentially; and (ii) some classifiers/models had been learnt from source domains. Online transfer learning is beneficial to many real applications. Below we give two examples to illustrate some potential applications.

The first example application is for online spam detection, such as spam email filtering. Typically, a universal classifier is trained to detect the spam as accurately as possible by a batch learning approach [25]. However, a universal classifier might not be always optimal for every individual as different persons may have different opinions on the definition of spam. This raises an open question, i.e., how to transfer useful knowledge from the universal classifier to personalize the spam detector for every individual in an online learning manner. Such a problem can be naturally attacked by applying the proposed OTL technique, in which the key challenge is that the “spam” concept in the target domain for each individual can be very different from that in the source domain. For such problems, as we assume the feature spaces of both source and target domains are the same, we thus refer to this scenario as OTL on *homogeneous domains* of common feature space.

The second example application is for climate forecast in environment and climate science [24], such as weather forecast, earthquake and tsunami prediction. For example, consider a situation where new types of instruments or sensors are introduced to improve an existing weather forecast system. In this scenario, training data with new features will be added to the forecast system while old features are still retained. Such a problem also can be formulated as an online transfer learning task, which aims to build an improved forecasting system on the new domain with the augmented features by transferring the knowledge of the old classifier in the source domain. This task can be potentially more challenging than the previous example as the feature spaces of both source and target domains are different, making it difficult to train the classifier on the new data by a simple transfer from the old classifier. We thus refer to this scenario as OTL across *heterogeneous domains* of diverse feature spaces.

As a summary, this paper addresses two challenging scenarios: (i) OTL on homogeneous domains, and (ii) OTL across heterogeneous domains. One straightforward approach to OTL is based on a continuous learning strategy, which initializes a regular online learning algorithm on the target domain with the existing classifier learnt from source domains. However, such a simple solution suffers from some critical drawbacks: (i) when studying OTL on homogeneous domains, it could suffer from negative transfer (transferred knowledge is harmful to learning target task) whenever there exists much significant difference between two conditional probabilities; and (ii) when studying OTL across heterogeneous domains, the old classifiers cannot be trained continuously with the new features because of the inconsistency of the two feature spaces.

In addition to these two challenges, we note that online transfer learning is in general more challenging than a classical batch transfer learning task. This is because in an OTL task it is very hard to directly measure the distribution difference of the two domains as only a predictive model of the source domain is provided, and the data instances on the target domain arrive on-the-fly sequentially and typically must be predicted immediately. This work aims to investigate effective and efficient OTL techniques to tackle these challenges.

In particular, to tackle the first challenge, we propose two ensemble learning based strategies for transferring knowledge from source domain by combining two sets of classifiers built on different domains. The key idea is to dynamically update the combination weights for the base classifiers according to their online performance. We propose two effective algorithms and give theoretical bounds to justify their efficacy. To tackle the second challenge, we propose a co-regularization learning strategy for knowledge transfer, which can effectively handle the learning task on diverse feature spaces. The key idea of the proposed co-regularization strategy was partially inspired by the co-training principle for batch learning tasks (semi-supervised learning or multi-view learning) [5,29], which combines classifiers co-trained from different “views” of the same training instances to boost the learning efficacy.

Last but not least, we extend the idea of the proposed OTL technique to attack a real-world open challenge in data mining and machine learning, i.e., the concept-drifting data stream mining task [21] where the underlying distributions and concepts often change over time. Despite being studied extensively in literature, it remains a critical open challenge for the existing approaches based on either batch learning or online learning techniques. In this paper, we propose an effective algorithm to attack this challenge based on a natural extension of the proposed OTL technique.

The rest of this paper is organized as follows. Section 2 reviews related work. Section 3 presents the proposed OTL framework and addresses the homogeneous and heterogeneous OTL tasks for classification. Section 4 presents the extension of OTL to address concept-drift online learning tasks. Section 5 gives our experimental results and discussions, and Section 6 concludes this work. Finally, we note that a short version of this work has appeared in the conference proceedings of ICML-2010 [33]. In contrast to the conference paper, a substantial amount of new contents and extensions have been included in this journal article.

2. Related work

Our work is mainly related to two machine learning topics: *online learning* and *transfer learning*. Below reviews some important related work.

Online learning (OL) has been extensively studied for years [28,7,9,32,34–36,30,19]. Unlike typical machine learning methods that assume training examples are available before the learning task, online learning is more appropriate for some real-world problems where training data arrives sequentially. Due to the merits of attractive efficiency and scalability, var-

Download English Version:

<https://daneshyari.com/en/article/376875>

Download Persian Version:

<https://daneshyari.com/article/376875>

[Daneshyari.com](https://daneshyari.com)