

# Open-Set Crowdsourcing using Multiple-Source Transfer Learning

Guangyang Han<sup>1</sup>, Guoxian Yu<sup>2,3\*</sup>, Lei Liu<sup>2,3</sup>, Lizhen Cui<sup>2,3</sup>, Carlotta Domeniconi<sup>4</sup>, Xiangliang Zhang<sup>5</sup>

<sup>1</sup>College of Computer and Information Sciences, Southwest University, China

<sup>2</sup>School of Software, Shandong University, China

<sup>3</sup>Joint SDU-NTU Centre for Artificial Intelligence Research, Shandong University, Jinan, China

<sup>4</sup>Department of Computer Science, George Mason University, USA

<sup>5</sup>Department of Computer Science, King Abudullah University of Science and Technology, Saudi Arabia  
hanguangyang@email.swu.edu.cn, {gxyu, l.liu, clz}@sdu.edu.cn, carlotta@cs.gmu.edu, xiangliang.zhang@kaust.edu.sa

## Abstract

We raise and define a new crowdsourcing scenario, *open set crowdsourcing*, where we only know the general theme of an unfamiliar crowdsourcing project, and we don't know its label space, that is, the set of possible labels. This is still a task annotating problem, but the unfamiliarity with the tasks and the label space hampers the modelling of the task and of workers, and also the truth inference. We propose an intuitive solution, OSCrowd. First, OSCrowd integrates crowd theme related datasets into a large source domain to facilitate partial transfer learning to approximate the label space inference of these tasks. Next, it assigns weights to each source domain based on category correlation. After this, it uses multiple-source open set transfer learning to model crowd tasks and assign possible annotations. The label space and annotations given by transfer learning will be used to guide and standardize crowd workers' annotations. We validate OSCrowd in an online scenario, and prove that OSCrowd solves the open set crowdsourcing problem, works better than related crowdsourcing solutions.

## 1 Introduction

Crowdsourcing is a new computation model that employs a large number of ordinary Internet workers to work together on a large project. In recent years, with the explosion of machine learning, especially deep neural networks, the demand for labeled data has been ever-increasing. As such, how to harness crowdsourcing to provide annotated data for machine learning has become a research hot-spot [Li *et al.*, 2017]. A series of crowdsourcing platforms, such as *Amazon Mechanical Turk* and *CrowdFlower*, have emerged and serve as the workplaces for crowd workers. In the typical crowdsourcing practice, we usually assume the label space of the crowd tasks is known or fixed [Li *et al.*, 2017]; in other words, we need to predefine the labels. In this way, the crowdsourcing task is modeled as a closed set label assignment problem.

But in practice, this assumption is often violated. Suppose we have collected a set of animal-themed pictures, which include tigers, lions, horses, cats, rhinos, etc.. If we want to classify all samples in detail, due to the open set of the label space, it's hard to setup the crowd task as there are so many alternative labels for these images. This *open set* problem also hampers quality, budget and latency control of crowdsourcing, as it's difficult and costly for crowd workers to go through all the options to choose the most appropriate label, and it's also difficult for requester to infer the task truth and the worker ability model when the annotation data is sparse. On the other hand, if we let the crowd workers give out labels instead of choosing from the given options, there may be a label inconsistency problem, which requires the design of additional tasks to further determine which labels are actually the same.

Transfer learning can transfer knowledge from related fields to help the learning model achieve a better performance in the target domain [Zhuang *et al.*, 2020]. Homogeneous transfer learning assumes that the source and target domains share the same feature and label space, and their difference lies in the marginal distribution of the feature space and/or the feature-label conditional distribution [Long *et al.*, 2013]. The recently studied open set transfer learning [Panareda Busto and Gall, 2017] and partial transfer learning [Zhang *et al.*, 2018] relax the restriction of sharing the same label space for the source and target domains, thus making transfer learning more practical. Inspired by this, if we can transfer available knowledge from related domains to model target crowdsourcing tasks, and roughly infer the possible label space of target tasks, then we can better tackle an open set crowdsourcing project, and complete it with a lower budget and a lower degree of uncertainty.

To implement the above idea, we introduce a solution called Open-Set Crowdsourcing (**OSCrowd**). First, due to the lack of sufficient knowledge of target crowdsourcing tasks, OSCrowd collects as many relevant datasets as possible to ensure the coverage of the label space of target crowdsourcing tasks. Then, it integrates the datasets into a large source domain, models the problem as a partial transfer learning problem, infers the possible label space of target tasks,

\*Contact Author

and weighs each source dataset. After this, it applies open set transfer learning to predict the labels of crowdsourcing tasks, combines weight information to give the machine annotating results of target tasks. Finally, OSCrowd uses the label space information and machine annotations to assist the design and execution of crowdsourcing workflow of target tasks.

The contributions of our work are summarized as follows:

- We introduce and define the open set crowdsourcing annotating problem, discuss its relevance and challenges.
- We propose to use transfer learning to infer the label space of crowdsourcing tasks, and propose a solution to the open set crowdsourcing annotation problem.
- We simulate and verify the effectiveness of our solution in online crowdsourcing scenarios. OSCrowd solved the problem of open set crowdsourcing annotation problem, and achieved the effect of SOTA.

## 2 Related Work

Our work builds on crowdsourcing and transfer learning. A comprehensive coverage of these two areas is out of the scope of this paper. In the following, we introduce related work which is most relevant to our approach.

**Crowdsourcing** is a new computation model that coordinates the crowd (Internet workers) to do micro-tasks in order to solve computer-hard problems; as an example, the world’s largest and most popular online encyclopedia, Wikipedia, is the product of many volunteers crowd workers. The task common to the different crowdsourcing scenarios is data processing [Li *et al.*, 2017], which includes collection, classification, modification, and reorganization, with data classification (labeling) being the most typical one.

The research on crowdsourcing mainly focuses on three aspects: quality, budget, and latency control [Li *et al.*, 2017]. No worker is perfect, so the annotations provided by crowd workers inevitably contain noise. How to identify and reduce the impact of noise is a key quality control task [Zheng *et al.*, 2017]. Crowd workers are generally not free, and the scale of crowdsourcing tasks is typically huge; in addition, repeated annotations are needed to ensure quality. As such, budget control is also very important, and this is usually managed by modeling tasks/workers to reduce the number of required annotations [Ho *et al.*, 2013; Li *et al.*, 2016; Tu *et al.*, 2020]. For latency control, researchers reduce the total time required to complete crowdsourcing via improved task assignment or workflow design [Haas *et al.*, 2015].

However, most of the aforementioned worker/task modeling and truth inference methods are based on a closed label set. For example, they usually assume that the task is binary or an  $n$ -way classification problem. When the problem is open set, the quality of workers and difficulty of tasks are harder to estimate due to sparse annotated data and the label open set.

**Domain adaptation** is a common paradigm of homogeneous transfer learning. It studies the situation in which the feature and label spaces of the source and target domains are uniform, but the feature distribution (marginal distribution  $P(X)$ ) is inconsistent [Blitzer *et al.*, 2006]. Domain adaptation tries to transfer knowledge from the source domain to

the target domain by aligning the feature spaces and by finding what they have in common [Ben-David *et al.*, 2007]. It is usually assumed that the source domain has rich supervised information, while the target domain has little or no label information. In recent years, deep learning, especially adversarial learning [Goodfellow *et al.*, 2014], has achieved excellent results in domain adaptation. Deep Domain Confusion (DDC) [Tzeng *et al.*, 2014] adds a domain adaptation layer using the Maximum Mean Discrepancy (MMD) distance metric in AlexNet to help the network extract the common feature representation. Unsupervised Domain Adaptation by Backpropagation (UDAB) [Ganin and Lempitsky, 2015] sets up a domain classifier to force the network to learn domain-invariant feature representations. [Ganin *et al.*, 2016; Hoffman *et al.*, 2018] also adopt the idea of adversarial learning to bridge the differences between domains.

**Partial & Open-set transfer** relaxes the assumption that the label space is consistent, thereby making the transfer model more extensive and practical. Partial transfer assumes that the target domain label space is a subset of the source domain, while open-set transfer assumes the two label spaces partially overlap, and both have their own unique classes. [Zhang *et al.*, 2018] adopted a domain adversarial network to weight source domain samples whose classes are likely to appear in the target domain. [Cao *et al.*, 2019] further introduced two additional auxiliary discriminators to make the weights more reasonable. In open-set transfer, [Panareda Busto and Gall, 2017] assigns labels to target domain samples by calculating the sample-class center distance after the domain adaptation process is done.

In practice, it is difficult to find a source domain that completely covers the label space of the target domain, but it is easy to find multiple related source domains. Our idea is that, instead of using multiple homogeneous source domains for open set domain adaptation, we can combine multiple source domains as a whole, and then perform partial domain adaptation. In this way, we can unify multiple source domains into one framework to facilitate the transfer of knowledge.

## 3 Methods

### 3.1 Notation and Problem Formulation

In this section, the necessary definitions pertaining crowdsourcing and transfer learning are given [Zhuang *et al.*, 2020].

*Definition 1 (Task):* A task (in machine learning)  $\mathcal{T}$  contains a label space  $\mathcal{Y}$  and a decision function  $f$  which needs to be learned from the training data, that is  $\mathcal{T} = \{\mathcal{Y}, f\}$ . In crowdsourcing, the term *task* usually refer to (a batch of) sample(s) assigned to crowd workers.

*Definition 2 (Domain):* A domain consists of two parts, a feature space  $\mathcal{X}$  and a marginal distribution  $P(X)$ . In practice, a domain is often represented by an instance set with or without labels. For example, a source domain is defined as  $\mathcal{D}_s = \{(\mathbf{x}_i, y_i) | \mathbf{x}_i \in \mathcal{X}^s, y_i \in \mathcal{Y}^s, i = 1, \dots, n^s\}$ .

*Definition 3 (Transfer Learning):* Transfer learning uses knowledge implied in the source domains and source tasks ( $\{\mathcal{D}_{s_i}, \mathcal{T}_{s_i}\}$ ) to improve the performance of the decision function  $f^t$  on the target domain. Domain adaptation is

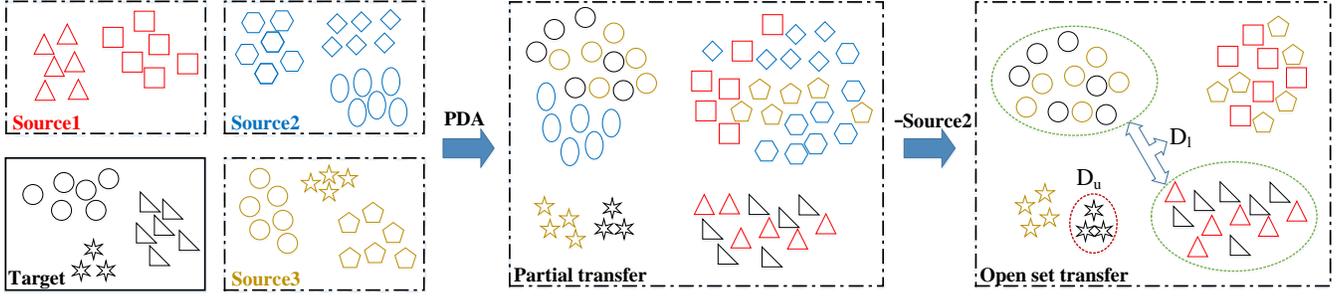


Figure 1: OScrowd first adopts partial domain adaptation (PDA) and GAN to align the feature space of the multiple-source domains and the target domain, considers the possibility of source domain classes being present in the target domain, and then excludes irrelevant source domains (Source2 in figure). Next, OScrowd uses open set transfer to assign target domain samples with one of the classes in the source domain, or to label them as unknown. The labeled samples correspond to simple tasks ( $\mathcal{D}_l$ ), while the unknown tasks are considered more difficult ( $\mathcal{D}_u$ ).  $\mathcal{D}_l$  and  $\mathcal{D}_u$  are then assigned to different quality workers for crowdsourcing. *Best viewed in color.*

a paradigm for transfer learning which focuses on aligning  $P(X)$  for the source and the target domains.

**Definition 4 (Open set Crowdsourcing):** We model Open Set Crowdsourcing as a target domain  $\mathcal{D}_t = \{\mathbf{x} | \mathbf{x}_i \in \mathcal{X}^t, i = 1, \dots, n^t\}$ , where  $n^t$  is the number of samples to be annotated by a group of  $W$  crowd workers  $\mathcal{W} = \{w_i\}_{i=1}^W$  (the annotation is represented as  $\mathbf{a}$ ). In open set crowdsourcing, we have some information about the label space  $\mathcal{Y}^t$ , but its size  $n_y^t = |\mathcal{Y}^t|$  and the specific content are unknown.

Suppose we have a crowdsourcing project  $\mathcal{D}_t$ , of which we only know the general theme, such as the fact that it pertains a set of pictures about livestock. The first thing we need to do is to gather labeled datasets about livestock from the Internet, annotating them as  $\{\mathcal{D}_{s_i}\}_{i=1}^{m'}$ . Each dataset may share ( $0 \sim n_y^t$ ) classes with  $\mathcal{D}_t$ . We merge the datasets into a large source domain  $\mathcal{D}_s$ , and the label space of  $\mathcal{D}_s$  should contain the label space of  $\mathcal{D}_t$ . We aim to: (i) discover the classes shared by  $\mathcal{D}_s$  and  $\mathcal{D}_t$ ; (ii) distinguish their differences; and (iii) transfer knowledge from  $\mathcal{D}_s$  to help model the tasks in  $\mathcal{D}_t$ .

### 3.2 Partial Domain Adaptation

To meet the above three needs, inspired by the adversarial learning based domain adaption solutions [Ganin *et al.*, 2016; Tzeng *et al.*, 2017], OSCrowd adopts a domain adversarial network built on GAN [Goodfellow *et al.*, 2014]. GAN originally contains two main parts, a generator  $G$  and a discriminator  $D$ . Here, we consider three main parts for a domain adversarial network: two feature extractors and a domain discriminator. The feature extractors  $F_s$  and  $F_t$  for  $\mathcal{D}_s$  and  $\mathcal{D}_t$  can be the same or different;  $F_s$  is trained on  $\mathcal{D}_s$  and kept unchanged, while  $F_t$  is pre-trained and acts as a generator  $G$ , which will conduct adversarial learning (minmax game) with the domain discriminator  $D$  until both reach the optimal state.

$$\min_{F_t} \max_D \mathcal{L}(D, F_t) = \mathbb{E}_{\mathbf{x} \sim p_s(\mathbf{x})} [\log D(F_s(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim p_t(\mathbf{x})} [\log (1 - D(F_t(\mathbf{x})))] \quad (1)$$

Here  $D$  is a simple binary perceptron used to distinguish source and target examples. The domain labels are set to 0 for the source samples and to 1 for the target samples.  $D$  will score each input sample and try to minimize the error on both

source and target tasks. For binary classification tasks, we choose the binary cross entropy function as the loss function, as shown in Eq. (1). On the other hand, the goal of domain adaptation is to minimize the difference in feature space between the source and target domains [Ben-David *et al.*, 2007], so that the knowledge of the source domain can be successfully applied to the target domain. A good feature extractor  $F_t$  or  $G$  should be able to extract features that can confuse  $D$ , that is, maximizing the loss of  $D$ .

Once we have learned a good  $F_s$ , we fix and copy its parameters to  $G$  as pre-trained [Tzeng *et al.*, 2017], and then iteratively train  $G$  and  $D$  until their losses no longer decrease. Following the same derivation as in GAN, when  $D$  and  $G$  are both optimal, the following relationship holds:

$$D^*(\mathbf{z}) = \frac{p_s(\mathbf{z})}{p_s(\mathbf{z}) + p_t(\mathbf{z})} \quad (2)$$

$$\mathcal{L}(F_t^*, D) = 2JS(p_s || p_t) - 2 \log 2 \geq -2 \log 2 \quad (3)$$

*Proof.* We unify the feature spaces of the source and target domains transformed by the two feature extractors to get Eq. (4). We further set the partial derivative of  $\mathcal{L}$  to  $D$  as zero, and the optimal value of  $D$  can be solved as Eq. (2), where  $p_s(\mathbf{z})$  and  $p_t(\mathbf{z})$  are the marginal distributions in the new space. The purpose of training  $G$  (domain adaptation) is to make the difference between the feature spaces as small as possible, so the ideal  $G$  should make the  $JS$  divergence close to 0, as shown in Eq. (3).

$$\begin{aligned} \mathcal{L} &= \int_{\mathbf{x}} p_s(\mathbf{x}) \log(D(F_s(\mathbf{x}))) + p_t(\mathbf{x}) \log(1 - D(F_t(\mathbf{x}))) d\mathbf{x} \\ &= \int_{\mathbf{z}} p_s(\mathbf{z}) \log(D(\mathbf{z})) + p_t(\mathbf{z}) \log(1 - D(\mathbf{z})) d\mathbf{z} \end{aligned} \quad (4)$$

Therefore, when the network training completed, for the classes shared by  $\mathcal{D}_s$  and  $\mathcal{D}_t$ , their features should be able to confuse  $D$ . In this case, when  $D$  affirmatively determines that a sample (class) comes from  $\mathcal{D}_s$ , this sample (class) should be unique to  $\mathcal{D}_s$ . In this way, we can infer the approximate classes of  $\mathcal{D}_t$ .  $D$  will score all samples in  $\mathcal{D}_s$ , then calculate the average score for each class. The scores of unique classes in  $\mathcal{D}_s$  are close to 0.

Since there may be source domains that have no intersection with  $\mathcal{D}_t$ , we remove the datasets whose classes are all unique, and re-execute the above steps to reduce negative transfer. We redefine those leftover datasets and scores as  $\mathcal{D}_s = \{\mathcal{D}_{s_i}\}_{i=1}^m$  and  $\{k_c\}_{c=1}^{n_y^s}$ , where  $m$  is the number of source datasets and  $n_y^s$  is the label space size of  $\mathcal{D}_s$ .

Through this partial transfer learning, we can exclude the irrelevant source domains, roughly guess the probability  $k_c$  of class  $c$  within the target domain, and learn a feature extraction network  $F_t \& F_s$  that extracts common features between domains. However, since some difficult samples may be difficult to be classified, so we provide an unknown class to accommodate these samples, then our next learning scenario can be modeled as an open set transfer learning, which can use  $\{k_c\}_{c=1}^{n_y^s}$  as side information.

### 3.3 Open Set Transfer

In the common feature space learned by PDA, we assign class labels to samples of the target domain based on feature vector, respectively. These samples are either classified as one of the classes in  $\mathcal{D}_s$  or as unknown. We classify by calculating the distance between the source domain sample and the center of each class in  $\mathcal{D}_s$ , similar to some metric-based meta learning work. The distance is calculated as:

$$d_{ci} = \frac{1}{k_c} \|\mathbf{s}_c - \mathbf{x}_i\|_2^2 \quad (5)$$

where  $\mathbf{s}_c$  is the mean of all samples of class  $c$  in the source domain with existence probability  $k_c$ . We assume that the sample is more likely to belong to the class with higher existence probability. We do not force each sample being assigned with a label; when the distance is larger than the max distance (Eq. (7)), we annotate the sample as unknown. We formulate the assignment problem as follows

$$\min_c d_{ci} \text{ s.t. } d_{ci} \leq \alpha d_c \quad (6)$$

$$d_c = \max\{\|\mathbf{s}_c - \mathbf{x}\|_2^2 \mid \mathbf{x} \in c\} \quad (7)$$

where hyper-parameter  $\alpha$  is the relaxation coefficient. Through transfer learning, we obtain the approximate label set of the target domain and the rough classification results of samples therein. Those information will be used to guide the subsequent crowdsourcing process.

### 3.4 Crowdsourcing

Till now, we have transformed the open set crowdsourcing problem into a close set crowdsourcing problem with a label space size of  $n_y^s$ , where  $n_y^s$  is the number of possible labels inferred by multiple source transfer learning. We divide the tasks into two pools: the task pool with machine annotations  $\mathcal{D}_l$  and the task pool marked as unknown  $\mathcal{D}_u$ .

The reason for this set up is closely related to our worker model. We believe that the characteristics of workers can be represented by two independent parameters,  $(C, A)$ . The binary parameter  $C$  is conscientiousness obtained by analyzing the behavior characteristics of workers, representing the willingness of workers to do their best and tell the truth. The probability-like parameter  $A$  represents the objective ability

of workers, and it is obtained by calculating the accuracy of workers. If the worker's subjective is insufficient, we will not accept their annotations regardless of the abilities, and we further distinguish honest workers based on their skill level.

In the crowdsourcing process, we keep three worker pools according to the worker model  $(C, A)$ :  $\mathcal{W}_u$ ,  $\mathcal{W}_r$ , and  $\mathcal{W}_e$ . For example, the four types of crowd workers defined in [Kazai *et al.*, 2011] (i.e. expert, normal, random, and spammer) can be grouped into random and spammer  $\mathcal{W}_u$ , normal  $\mathcal{W}_r$ , and expert  $\mathcal{W}_e$ . If the worker's answers are close to random guessing or always the same, the worker is considered as unreliable and put into  $\mathcal{W}_u$ . If a worker's answer is of passable quality, we add the worker into  $\mathcal{W}_r$ . Further, if the worker's answers are of high quality, and she/he has answered a sufficient number (e.g.  $\geq 10$ ) of questions, then this worker is deemed as an expert and put into  $\mathcal{W}_e$ .

To make OSCrowd practical, we design our crowdsourcing process in an online scenario. In online scenarios, task allocation algorithms such as task-centering or worker-task selection will fail due to the uncontrollable availability of workers. Following the setting in [Ho and Vaughan, 2012], workers  $\{w_i\}_{i=1}^W$  from a pool of size  $W$  arrive one at a time and the task requester must assign a batch of tasks to the newly arrived worker. We don't know whether a worker will be available again, nor the order the workers arrive.

Our crowdsourcing process can be divided into two stages based on the level of understanding of workers.

**Stage 1 Explore:** When a worker arrives for the first time, we give priority to allocate tasks in  $\mathcal{D}_l$ , and ensure that each task can receive enough ( $\geq 5$ ) annotations, so that we can infer the true label for the tasks and the worker capability.

**Stage 2 Collocate:** Case-1; to workers in  $\mathcal{W}_u$  we still assign tasks in  $\mathcal{D}_l$  that have received enough machine annotations. Case-2; to workers in  $\mathcal{W}_r$  we assign tasks in  $\mathcal{D}_l$  which do not have ascertained labels. Case-3; workers in  $\mathcal{W}_e$  are assigned hard tasks in  $\mathcal{D}_u$ .

Since estimating task truth and worker ability usually requires abundant data to ensure quality, we do not aggregate the answers immediately if workers are available. Batch update is also feasible in the middle and late stage. To address the cold start problem, we first use the machine label as the ground truth, and update the worker ability  $A^w$  of worker  $w$  according to Eq. (8) (number of correct annotations  $N_{correct}^w$  divided by number of total annotations  $N_{total}^w$  given by  $w$ ). Then Eq. (9) is used to soft  $w$ 's one-hot annotation  $\mathbf{a}^w$  according to  $A^w$ . Next, we use Eq. (10) to aggregate annotations of task  $t$  and select the label with the highest probability as the consensus label of  $t$ ;  $q^t$  is the number of total annotations received by task  $t$ . The EM algorithm is adopted to iteratively update task truth and worker model until convergence.

$$A^w = \frac{N_{correct}^w}{N_{total}^w} \times 100\% \quad (8)$$

$$\tilde{\mathbf{a}}^w = Tr(\mathbf{a}^w | A^w) = \begin{cases} \frac{1-A^w}{n_y^s} & a_i^w = 0 \\ A^w & a_i^w = 1 \end{cases} \quad (9)$$

$$\hat{\mathbf{a}}^t = \frac{1}{q^t} \sum_{i=1}^{q^t} \tilde{\mathbf{a}}_i^t \quad (10)$$

Office31			OfficeHome			
A (Target)	D (Source)	W (Source)	A (Target)	C (Source)	P (Source)	R (Source)
{0, 1, 2, 3, 4}	{0, 1}, {2, 5}	{3, 4, 6}, {7, 8}	{0, 1, 2, 3, 4}	{0, 1, 5}	{2, 3}	{4, 6}, {7, 8}

Table 1: Composition of the crowdsourcing tasks. The source and target domains share classes {0, 1, 2, 3, 4}, but {5, 6, 7, 8} are unique to the source domain. Classes in ‘{ }’ constitute a dataset (domain). In Office31, we simulated four source domains from two distributions, and in OfficeHome four source domains from three distributions. {7,8} is completely unrelated source domain.

	Score	0	1	2	3	4	5	6	7	8
O31	$k_{c-1}$	<b>.687</b>	<b>.785</b>	<b>.849</b>	<b>.712</b>	<b>.823</b>	.411	.523	.237	.341
	$k_{c-2}$	<b>.724</b>	<b>.765</b>	<b>.867</b>	<b>.723</b>	<b>.854</b>	.387	.444	-	-
OH	$k_{c-1}$	<b>.715</b>	<b>.732</b>	<b>.661</b>	<b>.849</b>	<b>.708</b>	.279	.330	.341	.458
	$k_{c-2}$	<b>.696</b>	<b>.753</b>	<b>.744</b>	<b>.802</b>	<b>.776</b>	.264	.375	-	-

Table 2: The scores of each class in  $\mathcal{D}_s$ . We remove the irrelevant source domain ({7,8}) according to  $k_c$  in round-1 ( $k_{c-1}$ ), and run PDA again to get the  $k_{c-2}$ . It can be seen that scores of shared classes are significantly higher than those of unique classes.

In the *Collocate* stage, we use Eq. (11) to measure the completion of a task. The greater the information entropy, the higher the uncertainty of the task is. So the greater  $Com(\hat{\mathbf{a}}^t)$  in Eq. (11) is, the higher the completion of task  $t$  (the denominator is the normalization of information entropy). When the completion of a task exceeds a certain threshold  $\gamma$ , or the number of annotations reaches the ceiling, we consider the task as completed.

$$Com(\hat{\mathbf{a}}^t) = 1 - \frac{H(\hat{\mathbf{a}}^t)}{\log n_y^s} = 1 + \frac{\sum_i^n p(\hat{a}_i^t) \log p(\hat{a}_i^t)}{\log n_y^s} \quad (11)$$

Each time we update task label and worker model, we sort  $\mathcal{D}_l$  according to the completion and adjust the worker pools. For workers identified in  $\mathcal{W}_u$ , their annotations do not have valid information, and we assign them completed tasks to confirm whether their evaluations need to be adjusted. If workers in  $\mathcal{W}_r$  arrive, low completion tasks are assigned to them for further annotation. Finally, if an expert comes, we let this worker complete the most difficult tasks in  $\mathcal{D}_u$ . Since an expert’s annotation information is usually sufficient, each difficult task is annotated only once.

## 4 Experiments

### 4.1 Settings

Since OSCrowd is based on domain adaptation, we simulate our crowdsourcing tasks on two benchmark datasets in domain adaptation, *Office31* (O31) [Saenko *et al.*, 2010] and *OfficeHome* (OH) [Venkateswara *et al.*, 2017]. Office31 is one of the commonly used datasets in domain adaptation. It has three sub-datasets, each containing the same 31 sub-classes collected from three different sources - Amazon (A), DSLR (D), and Webcam (W). As such, the sub-classes have different marginal distributions. Similarly, OfficeHome has four sub-datasets and 65 sub-classes extracted from four different styles - Art (A), Clipart (C), Product (P), and Real-world (R).

For Office31, we take five sub-classes in Amazon with the most samples as the annotation task. And then we select five identical sub-classes and four unrelated sub-classes from DSLR and Webcam, and combine them into four source domains. OfficeHome also uses a similar method to simulate

multi-source domain open set crowdsourcing. See Table 1 for the specific combinations. After PDA process, the scores of each class are shown in Table 2, it can be seen that the target domain classes are well recognized out.

Type	R	A	Ratio		
expert	1	[0.8,1.0]	10%	<b>20%</b>	20%
reliable	1	[0.4,0.8]	70%	<b>60%</b>	70%
unreliable	0	[0.1,0.4]	20%	<b>20%</b>	10%

Table 3: Attributes and proportions of workers. ( $R, A$ ) are the attributes of the worker’s model. The **boldfaced** ratio of workers is the baseline. The other two ratios are used to observe the model accuracy changes under different workers’ quality.

With reference to [Kazai *et al.*, 2011] and to our worker model, we simulate three different worker settings, and their models and percentage ratios are shown in Table 3. We randomly take out one worker each time to simulate the arrivals of online workers.

In order to verify the effectiveness of OSCrowd, we select four types of comparison methods: (i) crowdsourcing baseline, (ii) general crowdsourcing technology, (iii) crowdsourcing combining transfer learning and (iv) online crowdsourcing technology, with one or two representative methods for each type. Since the crowdsourcing methods below cannot handle open set annotations, we assume that the label space is known for these methods.

In crowdsourcing baseline, we assign five annotations to each task, and adopt Weighted Majority Voting and the EM algorithm to iteratively aggregate answers, which has been proved to be robust and effective in practical applications [Zheng *et al.*, 2017]. We denote it as *WMV*.

For the general crowdsourcing technology, *GLAD* [Whitehill *et al.*, 2009] uses a probabilistic approach to model the interaction between worker ability, task difficulty, and task ground truth, and jointly infer them to obtain final labels. We also fix the number of annotations to five for each task.

*ALM* [Fang *et al.*, 2014] and *CMKT* [Han *et al.*, 2020] use transfer learning to obtain better task features. ALM transfers from a single source domain and CMKT uses multi-source transfer. ALM uses active learning to select the most informative tasks while CMKT uses task assignment based on worker’s skill and task type. We use the default settings recommended by the original paper.

In an online scenario, *DTA* [Ho and Vaughan, 2012] assigns tasks with the goal of maximizing the total benefit that the requester obtains from the completed work. DTA first explores the worker’s ability, then uses the primal-dual formulation to assign tasks. *RAOC* [Welinder and Perona, 2010] ranks workers’ abilities and current tasks’ confidence, then selects the task with the lowest confidence to annotate.

Worker	Data	WMV	GLAD	ALM	CMKT	DTA	RAOC	OSCrowd
Ratio-1	O31	0.813±0.018	0.827±0.028	0.819±0.061	0.833±0.028	0.887±0.020	0.858±0.025	0.872±0.014
	OH	0.817±0.013	0.823±0.020	0.815±0.057	0.824±0.025	0.878±0.019	0.843±0.020	0.868±0.016
Ratio-2	O31	0.856±0.016	0.869±0.025	0.844±0.053	0.866±0.031	0.911±0.022	0.887±0.028	0.896±0.017
	OH	0.840±0.011	0.855±0.019	0.838±0.048	0.851±0.021	0.906±0.016	0.877±0.018	0.884±0.013
Ratio-3	O31	0.910±0.009	0.922±0.017	0.903±0.038	0.919±0.019	0.939±0.014	0.925±0.015	0.927±0.012
	OH	0.902±0.008	0.917±0.019	0.888±0.035	0.905±0.020	0.925±0.017	0.911±0.017	0.918±0.011

Table 4: Label accuracy and standard deviations of OSCrowd and compared methods on two datasets. Ratio- $i$  is related to the three worker settings in Table 3 and Ratio-2 is the baseline. Note that DTA assumes that the actual quality of the worker’s annotations is known; as such its accuracy corresponds to the ideal scenario and upper bound.

$\alpha$		0.4	0.6	0.8	1.0	<b>1.2</b>	1.4	1.6
O31	p	.841	.829	.815	.800	<b>.783</b>	.735	.659
	r	.427	.563	.668	.781	<b>.856</b>	.882	.902
	p*r	.359	.467	.544	.625	<b>.670</b>	.648	.594
OH	p	.778	.762	.739	.708	<b>.669</b>	.607	.537
	r	.416	.548	.655	.750	<b>.828</b>	.864	.884
	p*r	.324	.418	.484	.531	<b>.554</b>	.524	.475

Table 5: The relationship between the precision (p) and recall (r) of the machine label with the relaxation coefficient  $\alpha$  on two datasets.

In OSCrowd, we set a task to obtain at most five annotations; we set the relaxation coefficient  $\alpha = 1.2$  and the completion threshold  $\gamma = 0.75$ . Ten annotations are sufficient to infer a stable worker quality. Each dataset has more than  $90 \times 5 = 450$  tasks, and we simulated 30 crowd workers. The related experimental results are shown in Table 4.

## 4.2 Analysis of Results

By comparing the above methods, we have the following important observations:

(i) **OSCrowd vs. Online:** DTA achieves the best results among all methods and OSCrowd achieves sub-optimal results. However, this is because DTA assumes that the quality of workers is known without inference. As such, DTA is more like an optimal planning algorithm rather than a crowdsourcing algorithm. Therefore, the quality achieved by DTA gives an upper bound for the other methods. Although OSCrowd and RAOC have similar problem models, with the help of machine labels given by transfer learning, OSCrowd can better deal with the initial unstable stage of the model, and its overall quality is higher than that of RAOC.

(ii) **OSCrowd vs. Transfer:** Both ALM and CMKT use transfer learning to assist crowdsourcing modeling. CMKT transfers knowledge from multi-source domains, which reduces the risk of negative transfer from a single source domain, thus achieving higher accuracy than ALM. OSCrowd further calculates the correlation of multiple source domains, excludes unrelated source domains to reduce the risk of negative transfer, and thus achieves a better performance. When using transfer learning, avoiding negative transfer is an important way to improve model quality.

(iii) **OSCrowd vs. GLAD & WMV:** Both GLAD and WMV use the ‘accuracy’ to model crowd workers. GLAD further considers the correction effect of task difficulty on worker ability, thus achieving higher accuracy than WMV. However, they do not selectively assign tasks to workers. OSCrowd tries to grade workers and tasks, then assigns simple tasks to

ordinary workers and leaves difficult tasks to experts, discarding the annotations from unreliable workers, so it achieves better quality than both GLAD and WMV.

(iv) **Impact of different worker quality:** Observing Table 4 Ratio-3, when the proportion of low-quality workers decreases, almost all methods achieve better quality, and the performance gap between them is eliminated. This shows that all seven methods are effective crowdsourcing algorithms, and that the most important factor in determining the quality of real crowdsourcing is still the crowd workers. Conversely, when the quality of workers decreases (Table 4 Ratio-1), methods containing tasks assignment design are less affected (all methods except WMV and GLAD), because they can make better use of high-quality workers and reduce noise impact caused by low-quality workers.

## 4.3 Ablation Study

We studied the impact of the relaxation coefficient  $\alpha$  in Eq. (6) and the completion threshold  $\gamma$  related to Eq. (11). [Sheng *et al.*, 2008] has studied the effect of repeated annotations. The investigation of how many annotations lead to a stable worker model in section 3.4 is a purely statistical question.

In general, a smaller  $\alpha$  will bring higher machine label accuracy, but at the same time more samples will be labeled as unknown (difficult task  $D_u$ ). We borrow precision and recall to express this relationship, and adopt their product to consider both factors, the results are shown in Table 5. It can be seen that for Office31, the optimal value is between 1.2 and 1.4, and for OfficeHome, it’s between 1.0 and 1.2, therefore we choose  $\alpha = 1.2$ .

As for  $\gamma$ , there is a specific relationship between completion and the expected label quality, bigger  $\gamma$  leads to better quality and more budget. We expect the label quality to be about 0.9 under five classes. The two extreme cases of task annotation are (0.9, 0.1, 0, 0, 0) and (0.9, 0.025, 0.25, 0.025, 0.025), then the corresponding completion can be calculated out as 0.798 and 0.712 respectively, so we eclectically set  $\gamma = 0.75$ .

## 5 Conclusion

We propose and define open set crowdsourcing annotation, and analyze its difficulties and the shortcomings of current methods. We propose to use two-stage transfer learning to help solve this problem: OSCrowd first uses multi-source adversarial partial domain adaptation to infer the possible label space of crowdsourcing tasks, and excludes irrelevant source domains, and then labels crowdsourcing tasks. The label

space information and machine labels provided by transfer learning are used to assist the crowdsourcing design. For crowdsourcing, OSCrowd scores crowdsourcing workers and assigns tasks based on their abilities and commitment. Simulation experiments have confirmed that our method effectively solves the open set crowdsourcing annotation problem.

## References

- [Ben-David *et al.*, 2007] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. In *NeurIPS*, pages 137–144, 2007.
- [Blitzer *et al.*, 2006] John Blitzer, Ryan McDonald, and Fernando Pereira. Domain adaptation with structural correspondence learning. In *EMNLP*, pages 120–128, 2006.
- [Cao *et al.*, 2019] Zhangjie Cao, Kaichao You, Mingsheng Long, Jianmin Wang, and Qiang Yang. Learning to transfer examples for partial domain adaptation. In *CVPR*, pages 2985–2994, 2019.
- [Fang *et al.*, 2014] Meng Fang, Jie Yin, and Dacheng Tao. Active learning for crowdsourcing using knowledge transfer. In *AAAI*, pages 1809–1815, 2014.
- [Ganin and Lempitsky, 2015] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, pages 1180–1189, 2015.
- [Ganin *et al.*, 2016] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(1):2096–2030, 2016.
- [Goodfellow *et al.*, 2014] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, pages 2672–2680, 2014.
- [Haas *et al.*, 2015] Daniel Haas, Jiannan Wang, Eugene Wu, and Michael J Franklin. Clamshell: Speeding up crowds for low-latency data labeling. *VLDB Endowment*, 9(4):372–383, 2015.
- [Han *et al.*, 2020] Guangyang Han, Jinzheng Tu, Guoxian Yu, Jun Wang, and Carlotta Domeniconi. Crowdsourcing with multiple-source knowledge transfer. *IJCAI*, pages 2908–2914, 2020.
- [Ho and Vaughan, 2012] Chien-Ju Ho and Jennifer Vaughan. Online task assignment in crowdsourcing markets. In *AAAI*, pages 45–51, 2012.
- [Ho *et al.*, 2013] Chien-Ju Ho, Shahin Jabbari, and Jennifer Wortman Vaughan. Adaptive task assignment for crowdsourced classification. In *ICML*, pages 534–542, 2013.
- [Hoffman *et al.*, 2018] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *ICML*, pages 1989–1998, 2018.
- [Kazai *et al.*, 2011] Gabriella Kazai, Jaap Kamps, and Natasa Milic-Frayling. Worker types and personality traits in crowdsourcing relevance labels. In *CIKM*, pages 1941–1944, 2011.
- [Li *et al.*, 2016] Qi Li, Fenglong Ma, Jing Gao, Lu Su, and Christopher J Quinn. Crowdsourcing high quality labels with a tight budget. In *WSDM*, pages 237–246, 2016.
- [Li *et al.*, 2017] Guoliang Li, Yudian Zheng, Ju Fan, Jiannan Wang, and Reynold Cheng. Crowdsourced data management: Overview and challenges. In *SIGMOD*, pages 1711–1716, 2017.
- [Long *et al.*, 2013] Mingsheng Long, Jianmin Wang, Guiguang Ding, Jianguang Sun, and Philip S Yu. Transfer feature learning with joint distribution adaptation. In *ICCV*, pages 2200–2207, 2013.
- [Panareda Busto and Gall, 2017] Pau Panareda Busto and Juergen Gall. Open set domain adaptation. In *ICCV*, pages 754–763, 2017.
- [Saenko *et al.*, 2010] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *ECCV*, pages 213–226. Springer, 2010.
- [Sheng *et al.*, 2008] Victor S Sheng, Foster Provost, and Panagiotis G Ipeirotis. Get another label? improving data quality and data mining using multiple, noisy labelers. In *KDD*, pages 614–622, 2008.
- [Tu *et al.*, 2020] Jinzheng Tu, Guoxian Yu, Jun Wang, Carlotta Domeniconi, Maozu Guo, and Xiangliang Zhang. Crowdwt: Crowdsourcing via joint modeling of workers and tasks. *ACM TKDD*, 15(1):1–24, 2020.
- [Tzeng *et al.*, 2014] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- [Tzeng *et al.*, 2017] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *CVPR*, pages 7167–7176, 2017.
- [Venkateswara *et al.*, 2017] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *CVPR*, pages 5018–5027, 2017.
- [Welinder and Perona, 2010] Peter Welinder and Pietro Perona. Online crowdsourcing: rating annotators and obtaining cost-effective labels. In *CVPR*, pages 25–32, 2010.
- [Whitehill *et al.*, 2009] Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier Movellan, and Paul Ruvolo. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In *NeurIPS*, pages 2035–2043, 2009.
- [Zhang *et al.*, 2018] Jing Zhang, Zewei Ding, Wanqing Li, and Philip Ogunbona. Importance weighted adversarial nets for partial domain adaptation. In *CVPR*, pages 8156–8164, 2018.
- [Zheng *et al.*, 2017] Yudian Zheng, Guoliang Li, Yuanbing Li, Caihua Shan, and Reynold Cheng. Truth inference in

crowdsourcing: Is the problem solved? *VLDB Endowment*, 10(5):541–552, 2017.

[Zhuang *et al.*, 2020] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 99(1):43–76, 2020.