

# Separate to Adapt: Open Set Domain Adaptation via Progressive Separation\_论文笔记

---

[CVPR 2019](#)

[Blog](#)

## Abstract

---

- target domain中包含source domain中没有的类别，称之为**Open Set Domain Adaptation (OSDA)**
- 之前解决OSDA的方法，没有考虑目标域的**开放性** (openness,开放性是通过所有目标类中未知类的比例来衡量的)
- 目前的工作是将**整个target domain 和 source domain对齐**，而不排除未知类样本，这可能会由于**未知类与已知类之间的不匹配而引起负迁移**。
- 本文提出了**Separate to Adapt (STA)**，一个端到端的方法来解决OSDA问题
- 该方法采用**从粗到细的加权机制**来逐步分离未知和已知类别的样本，同时权衡其在特征分布对齐上的重要性。
- 该方法openness-robust，它可以适应目标域中的各种开放性。

## Introduction

---

- target domain 和 source domain 有完全相同的 label space，称之为 Closed Set Domain Adaptation
- 问题设置：target domain包含source domain中的所有的类，此外还有一些source domain中没有的类 (source classes 是 target classes的一个子集)

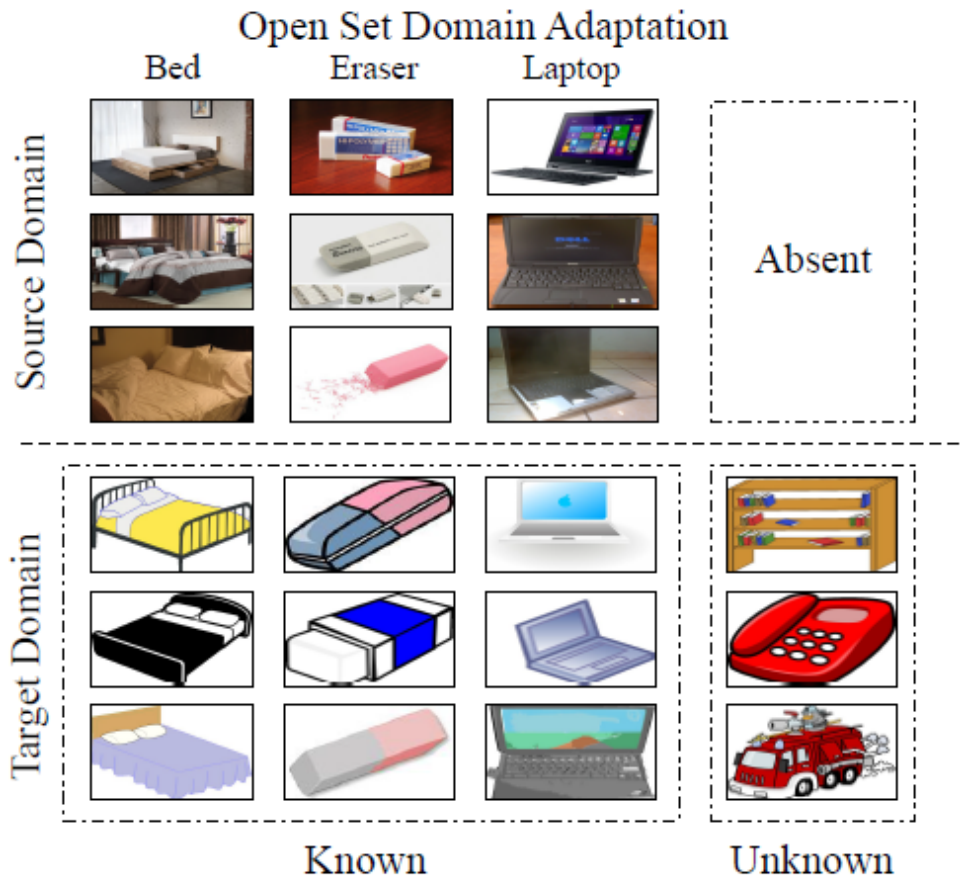


Figure 1. The open set domain adaptation problem, where the target domain contains “unknown” classes absent in the source domain.

- 目标是，将target domain中正确地~~将已知类的数据分类~~，并且将所有未知类的数据reject为“未知”
- 未知类不仅和source domain中的类存在domain gap，还存在semantic gap

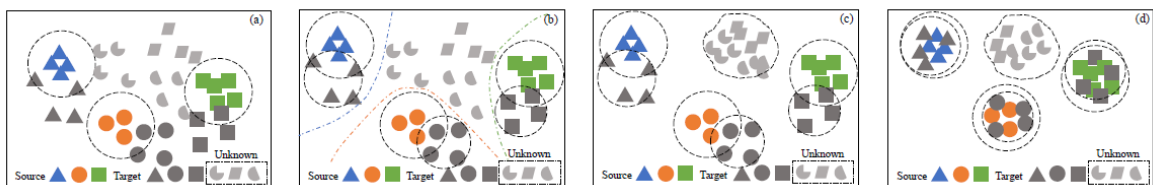


Figure 2. An overview of the proposed Separate to Adapt (STA) approach to open set domain adaptation. Gray shapes are data of target domain and shapes in color are data of source domain. Different kinds of shapes indicate different classes. (a) An example of open set domain adaptation problem, where all source classes are in the target classes and target have unknown classes. (b) The situation after training the multi-binary classifier  $G_c$  for deriving coarse weights to distinguish the unknown classes from known classes in the target domain. Dashed curve in different colors indicates the decision boundary for each binary classifier  $G_c$  for the  $c$ -th class. It forces the target data of unknown classes to move away from source data. (c) The situation after training the fine-grained binary classifier  $G_b$  for deriving more accurate weights. Target data in shared classes and in unknown classes are deviated far away. (d) The situation after the final distribution alignment, where target data in the shared classes are close to their source domain counterparts. Best viewed in color.

- 提出 a **progressive separation mechanism** consisting of a coarse-to-fine separation pipeline
  - 用source data训练一个multi-binary分类器
  - 选择具有极高和极低相似性的数据作为已知和未知类的数据，并用它们训练一个二分类器以对所有目标样本进行精细分离
  - 迭代以上两步，并使用instance-level权重来拒绝未知类的样本

## Method

### Open Set Domain Adaptation



$$L_s = \sum_{c=1}^{|\mathcal{C}_s|} \frac{1}{n_s} \sum_{i=1}^{n_s} L_{\text{bce}}(G_c(G_f(\mathbf{x}_i^s)), I(y_i^s, c)), \quad (1)$$

where  $L_{\text{bce}}$  is the binary cross-entropy loss and  $I(y_i^s, c) = 1$  if  $y_i^s = c$  and  $I(y_i^s, c) = 0$  otherwise. Each binary classifier

- $G_c$ 的output  $p_c$ 可以看作是每个样本和当前类 $c$ 的相似性
- 对于target 样本 $x_j^t$ , 找到这个样本和哪个类最相似, 计算相似性

$$s_j = \max_{c \in \mathcal{C}_s} G_c(G_f(\mathbf{x}_j^t)). \quad (2)$$

#### filtering strategy one

- 接着, 对所有的target 样本根据相似性排序, 选择具有最高和最低相似性的样本来训练二分类器  $G_b$

存在疑问, 最高和最低是只选了两个样本吗? 还是选了最高和最低的分數, 然后用这些分數对应的样本, 去训练二分类器 $G_b$ , 感觉后者比较合理一点?

- 优点:
  - 由于只使用了相似性在极限值的数据, 所以过滤相对粗糙但是拥有较高的可信度
  - 无需手动调整超参数, 鲁棒性强

#### filtering strategy two

- 分數聚类, 聚成高中低三类
- 使用高类中的均值 $s_h$ , 作为閾值,  $s_j \geq s_h$ 的是已知类, 使用低类中的均值 $s_l$ , 作为閾值,  $s_j \leq s_l$ 的是未知类

疑问, 小于 $s_h$  大于 $s_l$  的, 算啥?

$\mathbf{X}'$  to denote the set of filtered samples by the multi-binary classifier, and  $d_j$  to indicate whether a target sample  $\mathbf{x}_j \in \mathbf{X}'$  is labeled as known ( $d_j = 0$ ) or unknown ( $d_j = 1$ ). the fine-grained binary classifier  $G_c$  can be trained as follows,

$$L_b = \frac{1}{|\mathbf{X}'|} \sum_{\mathbf{x}_j \in \mathbf{X}'} L_{\text{bce}}(G_b(G_f(\mathbf{x}_j)), d_j). \quad (3)$$

## Weighted Adaptation

$$L_{\text{cls}}^s = \frac{1}{n_s} \sum_{\mathbf{x}_i \in \mathcal{D}_s} L_y \left( G_y^{1:|\mathcal{C}_s|} (G_f(\mathbf{x}_i)), y_i \right), \quad (4)$$

where  $L_y$  is cross-entropy loss,  $G_y$  is an *extended classifier* for  $|\mathcal{C}_s| + 1$  classes, i.e. the  $|\mathcal{C}_s|$  known classes in the source domain plus the additional “unknown” class in the target domain.  $G_y^{1:|\mathcal{C}_s|}$  denotes the probabilities corresponding to assigning each sample to the  $|\mathcal{C}_s|$  known classes.

- 使用 $G_b$ 的softmax输出作为instance-level权重,  $w_j = G_b(G_f(x_j))$ ,  $w_j$ 越大则属于未知类的概率越大

$$L_d = \frac{1}{n_s} \sum_{\mathbf{x}_i \in \mathcal{D}_s} L_{\text{bce}} (G_d (G_f (\mathbf{x}_i)), d_i) + \frac{1}{\sum_{\mathbf{x}_j \in \mathcal{D}_t} (1 - w_j)} \sum_{\mathbf{x}_j \in \mathcal{D}_t} (1 - w_j) L_{\text{bce}} (G_d (G_f (\mathbf{x}_j)), d_j). \quad (5)$$

- 此外, 还需要在目标域中选取未知类的样本, 以训练 $G_f$ 获得额外的unknown类

$$L_{\text{cls}}^t = \frac{1}{|\mathcal{C}_s|} \frac{1}{\sum_{\mathbf{x}_j \in \mathcal{D}_t} w_j} \sum_{\mathbf{x}_j \in \mathcal{D}_t} w_j L_y \left( G_y^{|\mathcal{C}_s|+1} (G_f (\mathbf{x}_j)), l_{\text{uk}} \right), \quad (6)$$

$l_{\text{uk}}$  是未知类, 疑问, 为什么这里也要加权呢?

- 对于已知类

$$L_e = \frac{1}{\sum_{\mathbf{x}_j \in \mathcal{D}_t} (1 - w_j)} \sum_{\mathbf{x}_j \in \mathcal{D}_t} (1 - w_j) H \left( G_y^{1:|\mathcal{C}_s|} (G_f (\mathbf{x}_j)) \right), \quad (7)$$

where  $H$  is the entropy loss and  $H(\mathbf{p}) = -\sum_k p_k \log p_k$ . It is noteworthy that we only aim to minimize the entropy of target samples estimated to be the known classes, so we use  $w_j$  as instance-level weight for the entropy minimization.

## Training Procedure

- known/unknown separation step

使用source data训练 $G_f$   $G_y$

$G_c$   $G_b$



$$(\hat{\theta}_f, \hat{\theta}_y, \hat{\theta}_b, \hat{\theta}_c |_{c=1}^{|\mathcal{C}_s|}) = \arg \min_{\theta_f, \theta_y, \theta_b, \theta_c |_{c=1}^{|\mathcal{C}_s|}} L_{\text{cls}}^s + L_s + L_b. \quad (8)$$

- **weighted adversarial adaptation step**

实现对抗性自适应，使目标域中已知类的特征分布与源域保持一致，并利用未知类中的数据为额外类训练  $G_y$

$$(\hat{\theta}_y, \hat{\theta}_d) = \arg \min_{\theta_y, \theta_d} L_{\text{cls}}^s + L_{\text{cls}}^t + L_d + \lambda L_e, \quad (9)$$

$$(\hat{\theta}_f) = \arg \min_{\theta_f} L_{\text{cls}}^s + L_{\text{cls}}^t - L_d + \lambda L_e, \quad (10)$$

## Experiment

Table 2. Classification Accuracy (%) of open set domain adaptation tasks on Office-31 (ResNet-50)

Method	A → W		A → D		D → W		W → D		D → A		W → A		Avg	
	OS	OS*	OS	OS*	OS	OS*	OS	OS*	OS	OS*	OS	OS*	OS	OS*
ResNet [9]	82.5±1.2	82.7±0.9	85.2±0.3	85.5±0.9	94.1±0.3	94.3±0.7	96.6±0.2	97.0±0.4	71.6±1.0	71.5±1.1	75.5±1.0	75.2±1.6	84.2	84.4
RTN [19]	85.6±1.2	88.1±1.0	89.5±1.4	90.1±1.6	94.8±0.3	96.2±0.7	97.1±0.2	98.7±0.9	72.3±0.9	72.8±1.5	73.5±0.6	73.9±1.4	85.4	86.8
DANN [4]	85.3±0.7	87.7±1.1	86.5±0.6	87.7±0.6	<b>97.5±0.2</b>	<b>98.3±0.5</b>	<b>99.5±0.1</b>	<b>100.0±0.0</b>	75.7±1.6	76.2±0.9	74.9±1.2	75.6±0.8	86.6	87.6
OpenMax [2]	87.4±0.5	87.5±0.3	87.1±0.9	88.4±0.9	96.1±0.4	96.2±0.3	98.4±0.3	98.5±0.3	83.4±1.0	82.1±0.6	82.8±0.9	82.8±0.6	89.0	89.3
ATI-λ [25]	87.4±1.5	88.9±1.4	84.3±1.2	86.6±1.1	93.6±1.0	95.3±1.0	96.5±0.9	98.7±0.8	78.0±1.8	79.6±1.5	80.4±1.4	81.4±1.2	86.7	88.4
OSBP [30]	86.5±2.0	87.6±2.1	88.6±1.4	89.2±1.3	97.0±1.0	96.5±0.4	97.9±0.9	98.7±0.6	88.9±2.5	90.6±2.3	85.8±2.5	84.9±1.3	90.8	91.3
<b>STA</b>	<b>89.5±0.6</b>	<b>92.1±0.5</b>	<b>93.7±1.5</b>	<b>96.1±0.4</b>	<b>97.5±0.2</b>	<b>96.5±0.5</b>	<b>99.5±0.2</b>	<b>99.6±0.1</b>	<b>89.1±0.5</b>	<b>93.5±0.8</b>	<b>87.9±0.9</b>	<b>87.4±0.6</b>	<b>92.9</b>	<b>94.1</b>

疑问，两个filtering strategy没有对比实验？

消融实验：

Table 4. Classification accuracy (%) of STA and its three variants on Office-31 (ResNet-50)

Method	A → W		A → D		D → W		W → D		D → A		W → A		Avg	
	OS	OS*	OS	OS*	OS	OS*	OS	OS*	OS	OS*	OS	OS*	OS	OS*
STA w/o w	87.5±1.4	91.4±1.1	83.0±1.2	89.6±1.2	96.2±0.9	97.3±0.4	98.1±0.7	<b>100.0±0.0</b>	80.3±1.5	79.3±1.5	71.2±1.2	74.3±1.2	86.1	88.7
STA w/o c	<b>90.4±1.7</b>	90.6±1.7	91.5±1.4	91.3±1.4	95.9±1.0	96.7±1.1	98.8±0.6	98.7±0.5	87.4±1.5	87.8±1.5	84.6±1.7	85.2±1.7	91.5	91.8
STA w/o b	85.0±1.5	89.0±1.5	90.6±1.2	91.5±1.3	94.8±1.9	<b>97.6±0.8</b>	96.2±0.6	98.2±0.5	77.7±2.2	82.5±2.4	78.9±2.6	83.6±3.5	87.2	90.4
STA w/o j	89.0±1.3	<b>92.8±1.2</b>	<b>94.8±1.5</b>	95.9±1.0	96.4±0.6	96.2±0.3	98.8±0.7	99.4±0.2	<b>89.7±1.4</b>	<b>93.6±1.4</b>	85.1±1.1	86.7±1.1	92.5	93.9
<b>STA</b>	<b>89.5±0.6</b>	<b>92.1±0.5</b>	<b>93.7±1.5</b>	<b>96.1±0.4</b>	<b>97.5±0.2</b>	<b>96.5±0.5</b>	<b>99.5±0.2</b>	<b>99.6±0.1</b>	<b>89.1±0.5</b>	<b>93.5±0.8</b>	<b>87.9±0.9</b>	<b>87.4±0.6</b>	<b>92.9</b>	<b>94.1</b>

1. w/o w (缺少对抗域训练的目标域样本的权重→对已知类和未知类的样本进行加权分离是必要的)
2. w/o c (缺少多二元分类器中的softmax分类层) →多二元分类器可以产生更好的相似度，独立地度量目标样本与每个源类之间的关系
3. w/o b (缺少二元分类器  $G_b$ ) →二元分类器可以根据多个二元分类器的结果来细化未知类和已知类样本之间的分离
4. w/o j (缺少Training Procedure中的两个steps的迭代) →联合分离和适应的有效性

## 遗留问题

- 关于 Adapting Object Detectors via Selective Cross-Domain Alignment 分组相关问题
  - 为什么要分组？

这里主要是为了找到我们感兴趣的区域，一般使用RPN出来的候选框，但是

- 候选框大小不一样，而作者希望获得**固定大小的区域**，以便于进一步处理（例如后面的patch生成等）
  - 第二个是利用K-means聚类，能刨除一些**噪声**，比如只有背景的框就不要了
- 分组之后？
- 分组之后，要确定每个区域的特征，进行feature reassignment
  - 对于第 $k$ 个分组中，里面有 $m_k$ 个proposal，每个proposal的RoI特征的维度是 $d$ ， $m_k$ 是不确定的，所以选定了一个超参 $m$
  - 对选出来的 $m$ 个proposal，进行拼接（论文里倒也没过多描述，我理解是一个cat操作？）

**Feature Reassignment.** Given the selected regions, we derive the feature representations thereof by reassigning the RoI features according to the grouping results. Specifically, each region is associated with a subset of region proposals assigned to the corresponding K-means cluster. By stacking the corresponding RoI features, we can obtain a matrix  $\Theta_k \in \mathbb{R}^{m_k \times d}$  to represent the  $k$ -th region, where  $m_k$  is the number of region proposals assigned to the  $k$ -th cluster, and  $d$  is the feature dimension.

This representation is inconvenient to work with, as the number  $m_k$  can vary. It is desirable to fix the number of features. For this purpose, we adopt a simple select-or-copy scheme. Given a pre-defined number  $m$ , if  $m_k$  is greater than  $m$ , we retain only the top- $m$  features; if  $m_k$  is less than  $m$ , we simply make copies of the assigned features until we get enough. In this way, we can derive a fixed number of features  $\tilde{\Theta}_k \in \mathbb{R}^{m \times d}$  to represent each region.

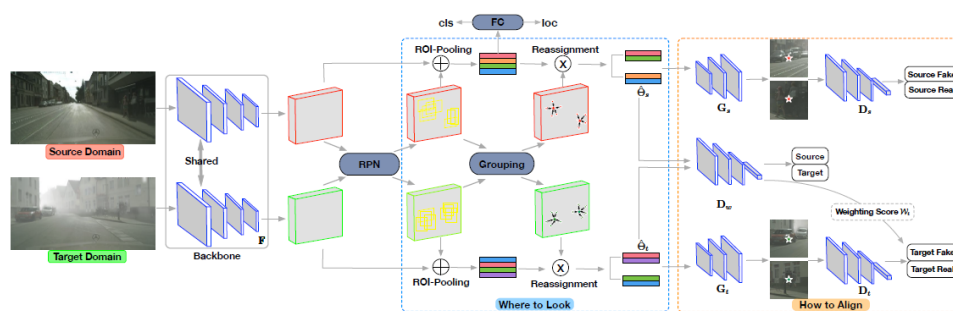


Figure 2. The pipeline of our framework. Two major components, i.e. “Where to Look” and “How to Align” are illustrated with two dashed rectangles. For the first component, an ROI-based grouping strategy is designed to mine the discriminative regions for two domains. We display the grouping procedure with cluster number = 2 (Note that  $\star$  and  $\star$  denote the centroids of clusters). For the second one, our model performs the adjusted region-level alignment using generators ( $G_s$  and  $G_t$ ), discriminators ( $D_s$  and  $D_t$ ) and weighting estimator ( $D_w$ ). We use Faster R-CNN as the detection model ( $\mathbb{F}$ ) which consists of the backbone, RPN and head part. (Best viewed in color)

