

Unsupervised Deep Domain Adaptation on People Detection

Anonymous ECCV submission

Paper ID ***

Abstract. This paper addresses the problem of unsupervised domain adaptation on the task of people detection in crowded scenes. That is, given a deep detection model well-trained on source domain, we adapt it into scene-specific detectors for any target domain on which no annotations are available. Firstly, we utilize iterative algorithm to iteratively auto-annotate target samples with high confidence on people instance as training set for scene-specific model on target domain. However, auto-annotated samples not only are lack of negative samples, but contains false positive samples. Therefore, on the one hand, we reuse negative samples from source domain to compensate for imbalance between the amount of positive samples and negative samples. On the other hand, we design an unsupervised regularizer based on deep network to mitigate influence from data error. Besides, we transform the last full connected layer into two sub-layers- element-wise layer and sum layer, on which the unsupervised regularizer can be added on. In experiments on people detection, the proposed method boosts recall by nearly 30% while precision stays almost the same. Furthermore, we perform our method on standard domain adaptation benchmarks on both supervised and unsupervised settings and our results are state of the art.

Keywords: Unsupervised Domain Adaptation, Unsupervised Regularizer, Deep Neural Network, People Detection

1 Introduction

Deep neural network has shown great power on traditional computer vision tasks, however, the labelled dataset should be large enough to train a deep model. In famous challenges such as PASCAL VOC and MS COCO, millions of labelled images are needed for training. This is also the case in surveillance applications. The annotation process for the task of people detection in crowded scenes is even more resource consuming, cause we need to label concrete locations of people instances. In modern society, there are over millions of cameras deployed for surveillance. However, these surveillance situations vary in lights, background, viewpoints, camera resolutions and so on. Directly utilizing models trained on old scenes will results in poor performance on the new situations due to data distribution changes. It is also unpractice to annotate people instances for every surveillance situation.

When there are few or even none of labelled data in target domain, domain adaptation helps to reduce the amount of labelled data needed. Most traditional works [1-5] either learn a shared representation between source and target domain, or project features into a common subspace. Recently, there are also works [6-8] proposed to learn a scene-specific detector by deep architectures. However, these approaches are heuristic either on constructing feature space or re-weighting samples. Our motivation of developing a domain adaptation architecture is to reduce heuristic methods required during adaptation process.

In this paper, we proposed a new approach of unsupervised deep domain adaptation on people detection. Using source model trained on source domain as initialization, we utilize iterative algorithm to iteratively auto-annotate target examples with high confidence as people instance on target domain for the first iteration. During each iteration, these auto-annotated data are regarded as training set to update target model, which, then, can be taken as the auto-annotation tool to auto-annotate target samples for the next iteration. However, these auto-annotated samples are defective, including lack of negative samples and existence of false positive samples, which will no doubt lead to exploration of predictions on non-people instances. Therefore, on the one hand, to compensate for the quantitative imbalance between positive and negative samples, we randomly sample negative instances from source domain and mix into training set. On the other hand, based on deep network, we design an unsupervised regularizer to mitigate influence from data error and avoid overfitting. To have better regularization effect during adaptation process, we transform the last full connected layer of deep model into two sub-layers, element-wise layer and sum layer. Thus, the unsupervised regularizer can be added on element-wise layer to adjust all parameters in the deep network and gain better performance.

Also, we further evaluate our approach on standard domain adaptation benchmark Office Dataset. The results of our adaptation approach outperform previously published works on both supervised and unsupervised scenarios, which also demonstrate the feasibility of our adaptation approach on both detection and classification tasks.

The contributions of our work are three folds.

- We proposed a feasible scheme to learn scene-specific deep detectors for target domains by unsupervised methodology, which can be easily deployed to various surveillance situations without any additional annotations.
- For better performance of unsupervised regularizer, we transform the last full connected layer of deep network into two sub-layers, element-wise layer and sum layer. Thus, all parameters contained in the deep network can be adjusted under the unsupervised regularizer. To our knowledge, this is the first attempt to transform full connected layers for the purpose of domain adaptation.
- Experiments on standard domain adaptation benchmarks for classification also demonstrate the applicability of our approach to other deep domain adaptation tasks.

The remainder of this paper is organized as follows. Section 2 reviews related works. Section 3 presents the details of our approach. Experimental results are shown in Section 4. Section 5 concludes the paper.

2 Relate Work

In many detection works, generic model trained by large amount of samples on source domain are directly utilized to detect on target domain. They assume that samples on target domain are subsets of source domain. However, when the distribution of data on target and source domain varies largely, the performance will drop significantly. Domain adaptation aims to reduce the amount of data needed for target domain.

Many domain adaptation works tried to learn a common representation space shared between source and target domain. Saenko et al. [1, 2] proposed a both linear-transform-based technique and kernel-transform-based technique to minimize domain changes. Gopalan et al. [3] projected features into Grassmann manifold instead of operating on features of raw data. Alternatively, Mesnil et al. [9] used transfer learning to obtain good representations. However, these methods are limited because scene-specific features are not learned to boost accuracy. The regularizer of our method are inspired these works.

Another group of works [4, 5, 10] on domain adaptation is to make the distribution of source and target domain more similar. Among these works, Maximum Mean Discrepancy (MMD) is used to as a metric to reselect samples from source domain in order to have similar distribution as target samples. In [11], MMD is incorporated as regularization to reduce distribution mismatch.

There are also works on deep adaptation to construct scene-specific detector. Wang et al.[6] explored context cues to compute confidence, and [7] learns distributions of target samples and proposed a cluster layer for scene-specific visual patterns. These works re-weighted auto-annotated samples for their final object function and additional context cues are needed for reliable performance. However, heuristic methods are required to re-weight samples. Alternately, Hattori el al. [8] learned scene-specific detector by generating a spatially-varying pedestrian appearance model. And Pishchulin et al. [12] used 3D shape models to generate training data. However, Synthesis for domain adaptation are also costly. Our approach minimize heuristic algorithms needed during adaptation process.

3 Our Approach

In this section, we introduce our unsupervised domain adaptation architecture on the task of people detection in crowded scenes. We denote training samples from source domain as $\mathbf{X}^S = \{x_i^S\}_{i=1}^{N^S}$. For training samples on source domain, we have corresponding annotations $Y^S = \{y_i^S\}_{i=1}^{N^S}$ with $y_i^S = (b_i^S, c_i^S)$, where $b_i^S = (x, y, w, h) \in R^4$ is the bounding box location and $c_i^S \in \{0, 1\}$ is the label

135 indicating whether x_i^S is a people instance. Also, we can denote the training
 136 samples on target domain as $\mathbf{X}^{T,n} = \{x_j^{T,n}\}_{j=1}^{N^{T,n}}$ and corresponding annotations
 137 as $Y^{T,n} = \{y_j^{T,n}\}_{j=1}^{N^{T,n}}$ with $y_j^{T,n} = (b_j^{T,n}, c_j^{T,n})$, where $n \in \{1, \dots, N^I\}$ is the
 138 index of adaptation iteration process. N^I is the maximum number of adaptation
 139 iterations. Note that different from source domain, training set on target domain
 140 are auto-annotated samples with high confidence as people instance. And during
 141 every adaptation iteration, the auto-annotation tool (also the target model) will
 142 be updated. Thus, the training samples for target domain at n^{th} adaptation
 143 iteration may differ from that at $(n+1)^{th}$ adaptation iteration.

144 The adaptation architecture of our approach consists of two streams – source
 145 stream M^S and target stream M^T , as shown in Fig x(!). Source stream takes
 146 samples from source domain as input, and target stream operates on samples
 147 from target domain. These two streams can utilize any end to end deep detection
 148 network as their model. Here we use the below mentioned network in Sec 3.1 in
 149 our experiment. In initialization stage, we firstly use abundant annotated sam-
 150 ples from source domain to train the model of source stream under a supervised
 151 loss function to regress bounding box. After its convergence, the weights of the
 152 model of source stream are used to initialize target stream. In adaptation stage,
 153 iteration algorithm is used as training method. Target model in target stream is
 154 trained and upgraded in this process while source stream stays static.

155 Both supervised loss and unsupervised regularizer are designed as loss func-
 156 tion to train the target stream for learning scene-specific detector as well as avoid
 157 overfitting. For supervised loss, the auto-annotated data can be used as train-
 158 ing labels. As the auto-annotated samples contains data error, an unsupervised
 159 loss are required to regularize the network. We take the source model in source
 160 stream as a reference for feature vector distribution. We defined the combination
 161 of supervised loss and unsupervised loss as our loss function for adaptation:

$$162 \quad L(\theta^T | \mathbf{X}^S, \mathbf{B}^S, \mathbf{X}^T, \tilde{\mathbf{B}}^T, \theta^S) = L_S + \alpha * L_U \quad (1)$$

$$164 \quad L_S = \sum_{j=1}^{N^T} \sum_{k=1}^{N^T} (r(\theta^T | x_j^T, \tilde{b}_{j,k}^T) + c(\theta^T | x_j^T, \tilde{b}_{j,k}^T)) \quad (2)$$

$$167 \quad L_U = L_{MMD}(\theta^T | \mathbf{X}^S, \mathbf{X}^T, \theta^S) \quad (3)$$

168 where L_S is supervised loss to learn the scene-specific detector and L_U is the
 169 unsupervised regularizer part. $r(\cdot)$ is a regression loss for bounding box location,
 170 like norm-1 loss, and $c(\cdot)$ is a classification loss for bounding box confidence, such
 171 as cross-entropy loss. And $L_{MMD}(\cdot)$ is the MMD-based loss for unsupervised
 172 regularization. Coefficient α balance the effect of supervised and unsupervised
 173 loss. We set $\alpha = 10$ in our experiments.

174 3.1 Detection Network

175 The generic model ¹ used in our adaptation architecture for source and tar-
 176 get stream is an end to end detection network without any precomputed region

177 ¹ Proposed by Russel et al.

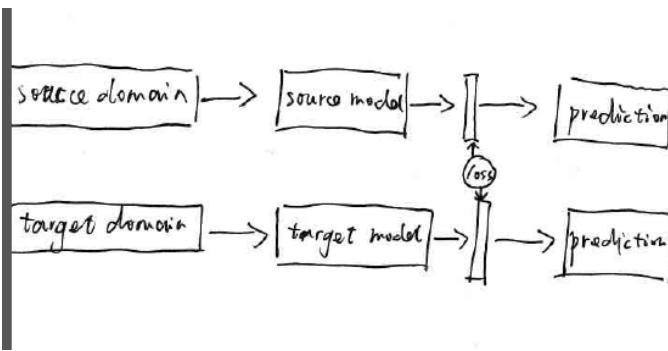


Fig. 1. One kernel at x_s (*dotted kernel*) or two kernels at x_i and x_j (*left and right*) lead to the same summed estimate at x_s . This shows a figure consisting of different types of lines. Elements of the figure described in the caption should be set in italics, in parentheses, as shown in this sample caption. The last sentence of a figure caption should generally end without a full stop

proposals needed. It consists of a GoogLeNet [13] for feature extraction and a RNN-based decoder for output of bounding box and corresponding confidence, as shown in Fig x(!). Firstly, the GoogLeNet model encode the image into a feature map (15x20x1024) of which each 1024 dimension vector are data representation of its receptive field corresponding to a subregion of the image. Then the RNN-based layers with batch size 300 will decode the data representation and sequentially predict 5 possible bounding boxes by the order of its corresponding confidence. Finally, all outputs are summarized to give final detection results. When trained with abundant samples from source domain, the obtained model has a high precision on target domain, however, its recall is low. Also, different from other detection network [14, 15], which need precomputed proposals for classification and fine regression, this generic model directly predict bounding boxes with high confidence. Thus, negative instances may contain people and non-people predictions, and cannot be employed in adaptation training.

3.2 Iterative Algorithm

In this section we introduce the iteration algorithm as training method at adaptation stage. The auto-annotating tool takes images on target domain with high confidence as training data for next iteration. To generate the auto-annotated data for the first iteration, we utilize the generic model well-trained on source domain, which results in training set on the target domain. As mentioned in Sec 3.1, the training set on target domain auto-annotated by the generic model in our experiment have low recall and high precision. At subsequent iterations, training set on target domain are auto-annotated by upgraded model resulting from training of last iteration. Among every iteration, these auto-annotated data are used as training data to upgrade the deep network.

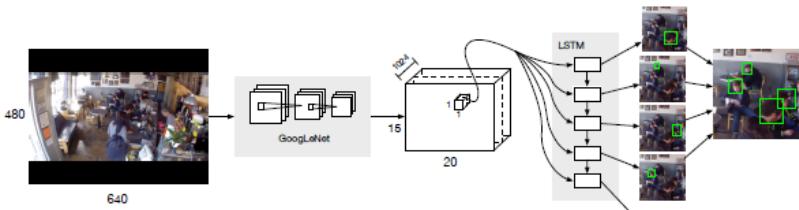


Fig. 2. One kernel at x_s (*dotted kernel*) or two kernels at x_i and x_j (*left and right*) lead to the same summed estimate at x_s . This shows a figure consisting of different types of lines. Elements of the figure described in the caption should be set in *italics*, in parentheses, as shown in this sample caption. The last sentence of a figure caption should generally end without a full stop

As the training set on target domain auto-annotated for the training of first iteration have low recall rate and people and non-people regions mixed in negative instances, we ignore back propagation of bounding boxes with low confidence during training. That is, we encourage the network to have more confidence on positive instances and stay conservative toward negative instances. This policy will no doubt resulting predictions of many non-people instance and impeding the further training when many non-people instances are regarded as people instances by the target model. To compensate for lack of true negative instances, we add annotated samples from source domain into the training data, which are human annotated and can thus provide true negative samples. In our experiment, when training on additional samples from source domain, we only do back propagation of bounding boxes with low confidence. At the same time, the unsupervised loss will also regularize the network. The complete adaptation process is illustrated in Algorithm 1. After a predetermined iteration limit N^I is reached, we obtain our final detection model on the target domain.

3.3 Unsupervised weights regularizer on Element-wise Multiply Layer

Element-wise Multiply Layer In deep neural network, the last feature vector layer are taken as an important data representation of input images. However, in this paper, we take one step further to focus on the last full connected layer which serves as an decoder to decode rich information contained in the last feature vector into final outputs. As source model are trained with abundant labelled data on source domain, the parameters of the last full connected layer are also well converged. We assume that a regularizer on the last full connected layer will achieve better results compared with the last feature vector layer. Firstly, denote the last feature vector, parameters of the last full connected layer

Algorithm 1 Deep domain adaptation algorithm (to be completed)

```

270 1: procedure DEEP DOMAIN ADAPTATION
271 2: Train source model on source stream with abundant annotated data
272 3: Use well-trained source model on source stream to initialize model on target stream
273 as  $M_0$ 
274 4: for  $i = 0:N^I$  do
275 5:    $M_i$  generate "fake ground truth"  $G_i$  of target domain
276 6:   balbal
277 7:   balbabla
278 8:    $G_i = G_i +$  random samples from source domain
279 9:   Take  $G_i$  as training data to upgrade  $M_i$  into  $M_{i+1}$ 
280 10: end for
281 11:  $M_{N^I}$ : final model.
282 12: end procedure

```

and final outputs as $\mathbf{F}_{(N^B \times N^D)}$, $\mathbf{C}_{(N^D \times N^O)}$ and $\mathbf{P}_{(N^B \times N^O)}$, respectively. The operation of full connected layer can be thus formulated as matrix multiply:

$$\mathbf{P} = \mathbf{F} * \mathbf{C} \quad (4)$$

$$P_{b,o} = \sum_d F_{b,d} * K_{d,o} \quad (5)$$

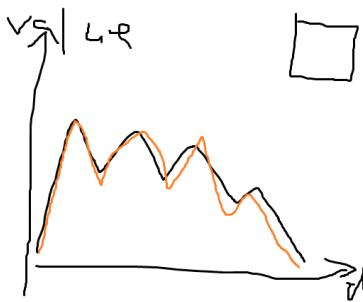
Inspired by this form, we separate the above formula into two sub-operations – element-wise multiply and sum, which can be formulated as:

$$M_{b,o,d} = F_{b,d} * C_{d,o} \quad (6)$$

$$P_{b,o} = \sum_d M_{b,o,d} \quad (7)$$

where $\mathbf{M}_{(N^B \times N^O \times N^D)} = [\mathbf{m}_{b,o}]$ is the parameter tensor of element-wise multiply operations. $\mathbf{m}_{b,o}$ is a vector with N^D dimensions, which will be the basis of unsupervised regularizer. Finally, we can equivalent-transform the last full connected layer between the last feature vector layer and final outputs layer into element-wise multiply layer and sum layer. The transformed element-wise layer is thus the last layer with parameters before output layers. Fig x (!) illustrates the transform.

Unsupervised weights regularizer on Element-wise Multiply Layer In works [decaf][], the last feature vector layer are regarded as the final representation of images. (how to introduce beyond sharing weights) In domain adaptation tasks, when generic deep model are trained with abundant data from source domain, the last element-wise layer mentioned in Sec 3.3 also includes rich information leading to the final output. For the nodes in the output layer, inputs from element-wise layers that will contribute to its value are not randomly decided. In this paper, we assume that the distribution of $\mathbf{m}_{b,o}$ of the last



319 **Fig. 3.** One kernel at x_s (*dotted kernel*) or two kernels at x_i and x_j (*left and right*)
 320 lead to the same summed estimate at x_s .

329 element-wise layer on both source and target domain should be similar. In the
 330 last element-wise multiply layer, every dimension of $\mathbf{m}_{b,o}$ may contribute to the
 331 final output. However, as the old model on source domain are trained with abun-
 332 dant images, the distribution of $\mathbf{m}_{b,o}^T$ should be consistent compared with $\mathbf{m}_{b,o}^S$ of
 333 source domain when adapting deep network on target domain. We utilize MMD
 334 (maximum mean discrepancy) to encode the similarity of element-wise multiply
 335 layer of source domain and target domain:

$$337 \quad L_{MMD}(\theta^T | \mathbf{X}^S, \mathbf{X}^T, \theta^S) = \frac{1}{N^O} \sum_{o=1}^{N^O} \text{allel} \frac{1}{N^B} \sum_{b=1}^{N^B} \mathbf{m}_{b,o}^T - \frac{1}{N^B} \sum_{b=1}^{N^B} \mathbf{m}_{b,o}^S \text{allel}^2 \quad (8)$$

340 which can also interpreted as the Euclidean distance between the center of $\mathbf{m}_{b,o}^T$
 341 and $\mathbf{m}_{b,o}^S$ across all output nodes. It's unpractical to get the distribution of
 342 the whole training set, while too few images cannot obtain a stable center for
 343 regularization. In our experiments, the $L_{MMD}(\cdot)$ loss is calculated for every
 344 batch. An example comparison of centers of $\mathbf{m}_{b_i,o}^S$ and $\mathbf{m}_{b_j,o}^S$ are shown in Fig
 345 x(1).

347 4 Experiment Results

349 In this section, we introduce experiment results on both surveillance applications
 350 and standard domain adaptation dataset. Our motivation for unsupervised do-
 351 main adaptation method is for easier deployment of We firstly evaluate our
 352 approach on video surveillance. Then we employ our approach to standard do-
 353 main adaptation benchmarks on both supervised and unsupervised settings to
 354 demonstrate the effectiveness of our method.

356 4.1 Domain Adaptation on Crowd Dataset

358 **Dataset and evaluation metrics** To show the effectiveness of our domain
 359 adaptation approach on people detection, we collected a dataset consisting of

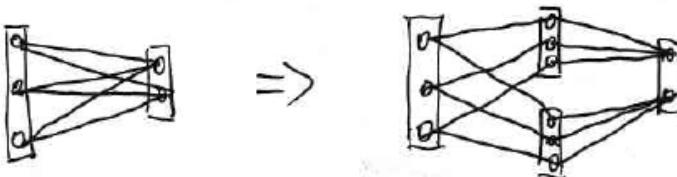


Fig. 4. One kernel at x_s (dotted kernel) or two kernels at x_i and x_j (left and right) lead to the same summed estimate at x_s .

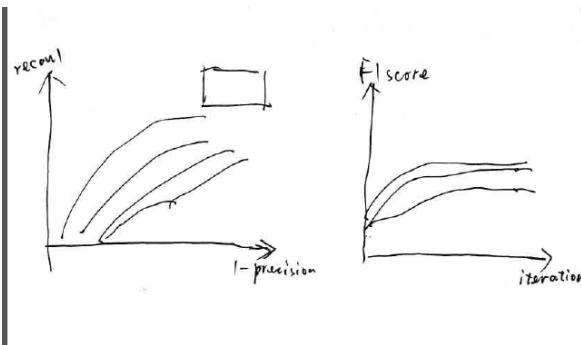


Fig. 5. One kernel at x_s (dotted kernel) or two kernels at x_i and x_j (left and right) lead to the same summed estimate at x_s . This shows a figure consisting of different types of lines. Elements of the figure described in the caption should be set in italics, in parentheses, as shown in this sample caption. The last sentence of a figure caption should generally end without a full stop

405 3 target scenes for target domain. These three scenes contain 1308, 1213 and
 406 331 un-annotated images with 0000, 0000, 0000 people instances respectively.
 407 For each scene, 100 images are annotated for evaluation. Instead of labelling the
 408 whole body of a person, we labels the head of a person as bounding box during
 409 training. The motivation for labelling only people heads comes from detection of
 410 indoor people or in crowded scenes, where the body of a person may be invisible.
 411 The dataset for source domain are Brainwash Dataset released on [http://d2.mpi-](http://d2.mpi-inf.mpg.de/datasets)
 412 inf.mpg.de/datasets. Brainwash Dataset consists of over 11917 images from 3
 413 crowded scenes. Examples of images from source and target domain are shown
 414 in Fig x(!).

415 Our evaluation metrics for detection uses the protocol defined in PASCAL
 416 VOC [16]. To judge a predicted bounding box whether correctly matches a
 417 ground truth bounding box, their intersection over their union must exceed 50%.
 418 And Multiple detections of the same ground truth bounding box are regarded as
 419 one correct prediction. We plot the precision-recall curve in Fig x(!). Also, the
 420 F1 score $F1 = 2 * precision * recall / (precision + recall)$ during the adaptation
 421 process are also shown in Fig x(!).

422
 423 **Experimental settings** We use deep learning framework Caffe [17] as the
 424 adaptation architecture of our approach. During the adaptation, we set learning
 425 rate as 0.01 and momentum as 0.5. At initialization stage, GoogLeNet weights are
 426 firstly used to initialize source model of source stream, while parameters in RNN
 427 layers are randomly initialized from a uniform distribution. For each iteration,
 428 100 auto-annotated images from target domain and 1000 annotated images from
 429 source domain are alternatively used for training. The outputs of our detection
 430 network contains bounding box locations and corresponding confidence, thus
 431 there are two full connected layers between the last feature vector layer and
 432 the final outputs. Experiments of unsupervised regularizer on element-wise layer
 433 for bounding box regression have no additional improvement on the performance
 434 when regularizer on element-wise layer for box confidence classification are added
 435 already. Our approach on domain adaptation are executed separately in the 3
 436 target scenes.

437
 438 **Comparison with baseline methods** To demonstrate the effectiveness of our
 439 approach, 4 methods are compared with method 4 as our final approach:
 440

- 441 1. Only auto-labeled samples on target domain are used for training, and with-
 out any unsupervised regularizer.
- 442 2. Only auto-labeled samples on target domain are used for training, with M-
 MD regularizer on last element-wise multiply layer as unsupervised weights
 regularizer.
- 443 3. Both auto-labeled images from target domain and labeled images from source
 domain are alternately sampled for training, with MMD regularizer on last
 feature vector as unsupervised weights regularizer.

- 450 4. Both auto-labeled images from target domain and labeled images from source
 451 domain are alternately sampled for training, with MMD regularizer on last
 452 element-wise multiply layer as unsupervised weights regularizer.

453 Fig x(!) plots the precision-recall curve of the above comparison methods in
 454 target scene 1). Also, the F1 score changes of every iteration during adaptation
 455 process are also depicted in Fig x(!). Table x(!) gives concrete precision and recall
 456 value of the 4 comparison methods on three target scenes when the F1 scores
 457 are at their highest. Examples of adaptation results are shown in Fig x(!).

	Scene 1			Scene 2			Scene 3		
	1-Pr	Re	F1	1-Pr	Re	F1	1-Pr	Re	F1
method 0	0.101	0.187	0.309	0.015	0.683	0.807	0.035	0.412	0.577
method 1	0.245	0.408	0.530	0.632	0.905	0.524	0.176	0.778	0.800
method 2	0.284	0.476	0.572	0.012	0.837	0.906	0.078	0.653	0.764
method 3	0.109	0.496	0.637	0.002	0.721	0.838	0.044	0.611	0.746
method 4	0.140	0.530	0.656	0.006	0.811	0.893	0.097	0.778	0.836

469
Performance evaluation From the Table x(!) and Fig x(1), we have the following observations:

- The recall values of method 1,2,3,4, which all utilized iteration algorithm to upgrade the target model, are larger than that of method 0, which are source model trained on source domain. This implies the effectiveness of our iteration algorithm in auto-annotation and iterative training.
- Compared with method 1, method 2 has higher F1 score. Their difference on whether a MMD regularizer are added into loss function demonstrates that our unsupervised regularizer can suppress data error and thus boost the final recall.
- Method 4 has both higher precision and higher recall than method 2, which demonstrates the effectiveness of additional samples from source domain during adaptation process.
- Compared with method 3, the recall of method 4 are further boosted. This results from the transformed element-wise layer which provided better regularization effect on the target model.

4.2 Domain Adaptation on Standard Classification Benchmark

488
Office dataset The Office dataset [1] comprises 31 categories of objects from
 489 3 domains (Amazon, DSLR, Webcam). Example images are depicted in Fig.
 490 x(!). As Amazon domain contains 2817 labelled images, which is the largest, we
 491 take it as source domain and Webcam domain as target domain. We follow the
 492 standard protocol for both supervised and unsupervised settings. Specifically,
 493 for supervised domain adaptation, we use 20 randomly sampled images with

495 labels for each category as training data for Amazon domain. When evaluate
 496 on unsupervised domain adaptation, 3 labelled images from target domain are
 497 additional selected for each class. For both settings, the rest of images on target
 498 domain are used for evaluation.



501
 502
 503
 504
 505
 506
 507
 508
Fig. 6. Some examples from three domains in the Office dataset. |

509
 510
Fig. 6. One kernel at x_s (*dotted kernel*) or two kernels at x_i and x_j (*left and right*)

511
 512
 513
 514
Experimental settings and network design On supervised setting, we
 515 reused the architecture in people detection. We utilize AlexNet [18] as the generic
 516 model of both streams. Firstly, we train the source model on source stream
 517 with provided training data from Amazon domain. Then iteration algorithm
 518 mentioned in Sec 3.2 are utilized for adaptation. The auto-labelling tool takes
 519 images on target domain with confidence exceeding 0.9 as training data for next
 520 iteration. The difference is besides auto-labelled images on target domain, 3 hu-
 521 man labelled images for each class are also included as training data for target
 522 model. For each iteration, 100 images are randomly sampled from training data.
 523 The unsupervised MMD regularizer added on the element-wise layer transformed
 524 from the last full connect layer of target model set the coefficient value α as 10
 525 on Eq 3, the same as that in people detection task.

526
 527 We use the same experimental setting for our unsupervised adaptation, ex-
 528 cept that at adaptation stage, no human labelled images can be added into the
 529 training set.

530
 531 **Performance evaluation** In Table x(!), we compare our approach with other
 532 six recently published works in both supervised and unsupervised settings. The
 533 outstanding performance on both settings confirms the effectiveness of our it-
 534 eration algorithm and MMD regularizer on the element-wise layer transformed
 535 from the last full connect layer.

536 537 5 Conclusions

538
 539 The paper ends with a conclusion.

	$A \rightarrow W$	
	Supervised	Unsupervised
GFK(PLS,PCA)[19]	46.4	15.0
SA [20]	45.0	15.3
DA-NBNN [21]	52.8	23.3
DLID [22]	51.9	26.1
DeCAF ₆ S [23]	80.7	52.2
DaNN [11]	53.6	35.0
Ours	84.3	66.3
Ours	85.4	69.3



Fig. 7. One kernel at x_s (*dotted kernel*) or two kernels at x_i and x_j (*left and right*)

585 References

- 586 1. Saenko, K., Kulis, B., Fritz, M., Darrell, T.: Adapting visual category models to
new domains. In: Computer Vision–ECCV 2010. Springer (2010) 213–226
- 587 2. Kulis, B., Saenko, K., Darrell, T.: What you saw is not what you get: Domain
adaptation using asymmetric kernel transforms. In: Computer Vision and Pattern
Recognition (CVPR), 2011 IEEE Conference on, IEEE (2011) 1785–1792
- 588 3. Gopalan, R., Li, R., Chellappa, R.: Domain adaptation for object recognition: An
unsupervised approach. In: Computer Vision (ICCV), 2011 IEEE International
Conference on, IEEE (2011) 999–1006
- 589 4. Huang, J., Gretton, A., Borgwardt, K.M., Schölkopf, B., Smola, A.J.: Correcting
sample selection bias by unlabeled data. In: Advances in neural information
processing systems. (2006) 601–608
- 590 5. Gretton, A., Smola, A., Huang, J., Schmittfull, M., Borgwardt, K., Schölkopf, B.:
Covariate shift by kernel mean matching. Dataset shift in machine learning **3**(4)
(2009) 5
- 591 6. Wang, X., Wang, M., Li, W.: Scene-specific pedestrian detection for static video
surveillance. Pattern Analysis and Machine Intelligence, IEEE Transactions on
36(2) (2014) 361–374
- 592 7. Zeng, X., Ouyang, W., Wang, M., Wang, X.: Deep learning of scene-specific clas-
sifier for pedestrian detection. In: Computer Vision–ECCV 2014. Springer (2014)
472–487
- 593 8. Hattori, H., Naresh Boddeti, V., Kitani, K.M., Kanade, T.: Learning scene-specific
pedestrian detectors without real data. In: Proceedings of the IEEE Conference
on Computer Vision and Pattern Recognition. (2015) 3819–3827
- 594 9. Mesnil, G., Dauphin, Y., Glorot, X., Rifai, S., Bengio, Y., Goodfellow, I.J., Lavoie,
E., Muller, X., Desjardins, G., Warde-Farley, D., et al.: Unsupervised and transfer
learning challenge: a deep learning approach. ICML Unsupervised and Transfer
Learning **27** (2012) 97–110
- 595 10. Gong, B., Grauman, K., Sha, F.: Connecting the dots with landmarks: Discrimi-
natively learning domain-invariant features for unsupervised domain adaptation.
In: Proceedings of The 30th International Conference on Machine Learning. (2013)
222–230
- 596 11. Ghifary, M., Kleijn, W.B., Zhang, M.: Domain adaptive neural networks for object
recognition. In: PRICAI 2014: Trends in Artificial Intelligence. Springer (2014)
898–904
- 597 12. Pishchulin, L., Jain, A., Wojek, C., Andriluka, M., Thormählen, T., Schiele, B.:
Learning people detection models from few training samples. In: Computer Vision
and Pattern Recognition (CVPR), 2011 IEEE Conference on, IEEE (2011) 1473–
1480
- 598 13. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D.,
Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: Proceedings
of the IEEE Conference on Computer Vision and Pattern Recognition. (2015) 1–9
- 599 14. Girshick, R.: Fast r-cnn. In: Proceedings of the IEEE International Conference on
Computer Vision. (2015) 1440–1448
- 600 15. Vu, T.H., Osokin, A., Laptev, I.: Context-aware cnns for person head detection.
In: Proceedings of the IEEE International Conference on Computer Vision. (2015)
2893–2901
- 601 16. Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zisser-
man, A.: The pascal visual object classes challenge: A retrospective. International
Journal of Computer Vision **111**(1) (2015) 98–136

- 630 17. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding.
631 arXiv preprint arXiv:1408.5093 (2014) 631
632 18. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep 632 convolutional neural networks. In: Advances in neural information processing systems.
633 (2012) 1097–1105 633
634 19. Gong, B., Shi, Y., Sha, F., Grauman, K.: Geodesic flow kernel for unsupervised 634 domain adaptation. In: Computer Vision and Pattern Recognition (CVPR), 2012 635 IEEE Conference on, IEEE (2012) 2066–2073 635
636 20. Fernando, B., Habrard, A., Sebban, M., Tuytelaars, T.: Unsupervised visual 636 domain adaptation using subspace alignment. In: Proceedings of the IEEE Interna- 637 tional Conference on Computer Vision. (2013) 2960–2967 637
638 21. Tommasi, T., Caputo, B.: Frustratingly easy nbnn domain adaptation. In: Proceed- 638 ings of the IEEE International Conference on Computer Vision. (2013) 897–904 639
639 22. Chopra, S., Balakrishnan, S., Gopalan, R.: Dlid: Deep learning for domain adap- 639 tation by interpolating between domains. In: ICML workshop on challenges in 640 representation learning. Volume 2. (2013) 640
641 23. Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., Darrell, T.: 641 Decaf: A deep convolutional activation feature for generic visual recognition. arXiv 642 preprint arXiv:1310.1531 (2013) 642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674