

A Reproducible and Realistic Evaluation of Partial Domain Adaptation Methods

Supplementary material

Outline. The supplementary material of this paper is organized as follows:

- In Section A, we give more details on our experimental protocol.
- In Section B, we provide additional results from our experiments.

A Additional details on Experimental Protocol

A.1 Implementations in BenchmarkPDA

In order to reimplement the different PDA methods, we adapted the code from the official repository associated with each of the paper. We list them in Table 1.

Method	Code Repository
PADA	https://github.com/thuml/PADA/blob/master/pytorch/src/
SAFN	https://github.com/jihanyang/AFN/blob/master/partial/OfficeHome/SAFN/code/
BA3US	https://github.com/tim-learn/BA3US/
AR	https://github.com/XJTU-XGU/Adversarial-Reweighting-for-Partial-Domain-Adaptation
JUMBOT	https://github.com/kilianFavras/JUMBOT
M-POT	https://github.com/UT-Austin-Data-Science-Group/Mini-batch-OT/tree/master/PartialDA

Table 1: Office Github code repositories for the PDA methods considered in this work.

One of our main claims regarding previous work is the use of target labels to choose the best model along training. This can be easily verified by inspecting the code. For PADA it can be seen on line 240 of the script “train_pada.py”, for BA3US in line 116 for the script “run_partial.py”, for M-POT it can be seen line 164 of the file “run_mOT.py”, for SAFN it can be seen in the “eval.py” file and finally for AR in line 149 of the script “train.py”.

For JUMBOT and M-POT which are based on optimal transport, we used the optimal transport solvers from (Flamary et al., 2021).

A.2 Model Selection

DEV requires learning a discriminative model to distinguish source samples from target samples. Its neural network architecture must be specified as well the training details. You et al. (2019) (DEV) use a multilayer perceptron, while Saito et al. (2021) (SND) use a Support Vector Machine in their reimplementation of DEV. We empirically observed the latter to yield more stable weights and so that was the one we used. In order to train the SVM discriminator, following (Saito et al., 2021), we take 3000 feature embeddings from source samples used in training and 3000 random feature embeddings from target samples, both chosen randomly. We do a 80/20 split into training and test data. The SVM is trained with a linear kernel for a maximum of 4000 iterations. Of 5 different SVM models trained with decay values spaced evenly on log space between 10^{-2} and 10^4 the one that leads to the highest accuracy (in distinguishing source from target features) on the test data split is the chosen one.

As for SND, it also requires specifying a temperature for temperature scaling component of the strategy. We used the default value of 0.05 that is suggested in (Saito et al., 2021).

Finally, we mention that the samples used for 100-RND were randomly selected and their list is made available together with the code. As for the samples used for 1-SHOT, they are the same as the ones used in semi-supervised domain adaptation.

A.3 Optimizer

In general, all methods claim to adopt Nesterov’s acceleration method as the optimization method with a momentum of 0.9 and setting the weight decay set to 5×10^{-4} . The learning rate follows the annealing strategy as in Ganin et al. (2016):

$$\mu_p = \mu_0(1 + \alpha \hat{u}p)^{-\beta},$$

where p is the training progress linearly changing from 0 to 1, $\mu_0 = 0.01$ and $\alpha = 10$ and $\beta = 0.75$.

However, inspecting the Official code repo for each PDA method, the actual learning schedule is given by

$$\mu_i = \mu_0(1 + \alpha \hat{u}i)^{-\beta},$$

where i is the iteration number in the training procedure, $\mu_0 = 0.01$ and $\alpha = 0.001$ and $\beta = 0.75$. Only when the total number of iterations is 10000 do the learning rate schedules match. In this work, we followed the latter since it is the one indeed used. For OFFICE-HOME, all methods are trained for 5000 iterations, while for VISDA they are trained for 10000 iterations, with the exception of the s. ONLY which is trained for 1000 iterations on OFFICE-HOME and 5000 iterations on VISDA.

A.4 Hyper-Parameters

In Table 2, we report the values used for each hyper-parameter in our grid search. We report in Table 3 the hyper-parameters chosen by each model selection strategy for each method on both datasets. In addition, for the reproducibility of AR with the proposed architecture in Gu et al. (2021), a feature normalization layer is added in the bottleneck which requires specifying r , the value to which the 2-norm is set. This hyper-parameter is therefore included in the hyper-parameter grid search with the possible values of $\{5, 10, 20\}$ which are the different values used in the experiments in (Gu et al., 2021).

B Additional Discussion of Results

In this section, we provide additional results that we could not add to the main paper due to the space constraints.

In Table 4, we show the accuracy per task on OFFICE-HOME averaged over three different seeds (2020, 2021, 2022) for all pairs of methods and model selection strategies.

In Table 5, we compare previously reported results with ours on VISDA. While proposed methods reported results on OFFICE-HOME, only PADA and AR results are reported in the original papers for VISDA. Gu et al. (2021) (AR) also report results for BA3US. Analysing the results, we see a 9 percentage point decrease in average task accuracy for PADA, but our experiments show that there is a significant seed dependence which we discuss in detail below. This is particularly important since Cao et al. (2018) (PADA) report results from a single run. Comparing our best seeds for PADA on the SR and RS tasks, we achieve 58.01% and 67.9% accuracy versus a reported 53.53% and 76.5%. Moreover, we point out that the official code repository for PADA does not include the details to reproduce the VISDA experiments, so it is possible that minor tweaks (e.g learning rate) are necessary. As for BA3US, our results are within the standard deviation being better on the SR task and worse on the RS task. Finally as for AR we see a decrease in performance which, as the results on OFFICE-HOME show, can be explained by the differences in the neural network architecture.

Finally in Table 6, we show all the average task accuracies from all pairs of methods and model selection strategies on the OFFICE-HOME and VISDA datasets including the 50-RND model selection strategy.

Method	HP	Values
PADA	λ	[0.1, 0.5, 1.0, 5.0, 10.0]
BA3US	λ_{wce}	[0.1, 0.5, 1, 5, 10]
	λ_{ent}	[0.01, 0.05, 0.1, 0.5, 1]
SAFN	λ	[0.005, 0.01, 0.05, 0.1, 0.5]
	Δ_r	[0.01, 0.1, 1.0]
AR	ρ_0	[2.5, 5.0, 7.5, 10.0]
	A_{up}	[5.0, 10.0]
	A_{low}	$-A_{up}$
	λ_{ent}	[0.01, 0.1, 1.0]
JUMBOT	τ	[0.001, 0.01, 0.1]
	η_1	[0.00001, 0.0001, 0.001, 0.01, 0.1]
	η_2	[0.1, 0.5, 1.]
	η_3	[5, 10, 20]
MPOT	ϵ	[0.5, 1.0, 1.5]
	η_1	[0.0001, 0.001, 0.01, 0.1, 1.0]
	η_2	[0.1, 1.0, 5.0, 10.0]
	m	[0.1, 0.2, 0.3, 0.4]

Table 2: Hyper-Parameter values for each PDA method considered in the grid search.

References

- Zhangjie Cao, Lijia Ma, Mingsheng Long, and Jianmin Wang. Partial adversarial domain adaptation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 135–150, 2018.
- Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z. Alaya, Aurélie Boisbunon, Stanislas Chambon, Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, Léo Gautheron, Nathalie T.H. Gayraud, Hicham Janati, Alain Rakotomamonjy, Ievgen Redko, Antoine Rolet, Antony Schutz, Vivien Seguy, Danica J. Sutherland, Romain Tavenard, Alexander Tong, and Titouan Vayer. Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78):1–8, 2021. URL <http://jmlr.org/papers/v22/20-451.html>.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016.
- Xiang Gu, Xi Yu, Yan Yang, Jian Sun, and Zongben Xu. Adversarial reweighting for partial domain adaptation. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=f5liPryFRoA>.
- Kuniaki Saito, Donghyun Kim, Piotr Teterwak, Stan Sclaroff, Trevor Darrell, and Kate Saenko. Tune it the right way: Unsupervised validation of domain adaptation via soft neighborhood density. *arXiv preprint arXiv:2108.10860*, 2021.
- Kaichao You, Ximei Wang, Mingsheng Long, and Michael Jordan. Towards accurate model selection in deep unsupervised domain adaptation. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 7124–7133. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/you19a.html>.

Method	Dataset	HP	ORACLE	1-SHOT	50-RND	100-RND	S-ACC	ENT	DEV	SND
PADA	OFFICE-HOME	λ	0.5	0.1	0.1	0.5	0.1	1.0	5.0	0.5
	VISDA	λ	0.5	1.0	10.0	0.5	1.0	0.5	5.0	0.1
SAFN	OFFICE-HOME	λ	0.005	0.1	0.005	0.01	0.005	0.01	0.005	0.005
		Δr	0.1	0.01	0.01	0.01	0.01	0.1	0.1	0.1
	VISDA	λ	0.005	0.005	0.05	0.05	0.005	0.05	0.005	0.05
		Δr	0.1	0.01	0.01	0.01	0.01	0.01	0.01	0.01
BA3US	OFFICE-HOME	λ_{wce}	5.0	10.0	5.0	5.0	5.0	0.1	10.0	1.0
		λ_{ent}	0.05	0.05	0.01	0.05	0.01	0.1	0.05	0.01
	VISDA	λ_{wce}	1.0	1.0	0.1	1.0	5.0	1.0	5.0	5.0
		λ_{ent}	0.5	0.5	0.5	0.5	0.05	0.5	0.05	1.0
AR	OFFICE-HOME	ρ_0	2.5	2.5	5.0	5.0	2.5	5.0	7.5	10.0
		A_{up}	5.0	5.0	10.0	5.0	5.0	10.0	10.0	10.0
		A_{low}	-5.0	-5.0	-10.0	-5.0	-5.0	-10.0	-10.0	-10.0
		λ_{ent}	0.1	0.1	1.0	1.0	0.01	1.0	0.01	1.0
	VISDA	ρ_0	2.5	2.5	2.5	2.5	2.5	7.5	2.5	10.0
		A_{up}	10.0	10.0	10.0	10.0	5.0	10.0	10.0	10.0
		A_{low}	-10.0	-10.0	-10.0	-10.0	-5.0	-10.0	-10.0	-10.0
		λ_{ent}	0.1	0.1	0.1	0.1	0.01	0.1	0.01	0.01
JUMBOT	OFFICE-HOME	τ	0.01	0.01	0.01	0.001	0.1	0.01	0.01	0.001
		η_1	0.0001	0.0001	0.001	0.0001	0.01	1e-05	0.01	1e-05
		η_2	0.5	1.0	0.5	0.1	0.1	0.5	1.0	1.0
		η_3	10.0	5.0	5.0	5.0	5.0	20.0	10.0	5.0
	VISDA	τ	0.01	0.01	0.01	0.01	0.001	0.01	0.001	0.01
		η_1	0.001	0.001	0.001	0.001	0.01	1e-05	0.01	0.0001
		η_2	1.0	1.0	0.5	1.0	0.1	0.5	1.0	1.0
		η_3	5.0	5.0	5.0	5.0	10.0	5.0	20.0	5.0
MPOT	OFFICE-HOME	ϵ	0.5	0.5	1.0	0.5	1.0	1.5	1.0	1.5
		η_1	0.01	0.01	0.01	0.01	0.001	0.0001	1.0	0.01
		η_2	10.0	1.0	1.0	1.0	1.0	10.0	0.1	1.0
		m	0.3	0.1	0.1	0.2	0.3	0.4	0.2	0.4
	VISDA	ϵ	0.5	0.5	0.5	0.5	1.0	1.0	1.0	0.5
		η_1	0.01	0.001	0.01	0.01	0.001	0.0001	0.0001	0.01
		η_2	1.0	1.0	1.0	1.0	1.0	10.0	1.0	10.0
		m	0.3	0.1	0.3	0.3	0.2	0.4	0.2	0.3

Table 3: Hyper-parameters selected for the different methods for each model selection strategy on both OFFICE-HOME and VISDA.

ALGORITHM	S2R	R2S	Avg
S. ONLY [†]	45.26	64.28	54.77
S. ONLY (Ours)	51.86	67.11	59.48
PADA [†]	53.53	76.50	65.02
PADA (Ours)	49.34	59.81	54.57
SAFN [†]	67.65	-	-
SAFN (Ours)	56.88	68.40	62.64
BA3US [†]	69.86	67.56	68.71
BA3US (Ours)	71.77	63.56	67.67
AR ^{†*}	85.30	74.82	80.06
AR (Ours)	76.33	71.36	73.85
JUMBOT [†]	-	-	-
JUMBOT (Ours)	90.55	77.46	84.01
MPOT [†]	-	-	-
MPOT (Ours)	87.23	86.67	86.95

Table 5: Comparison between reported ([†]) accuracies on partial VISDA from published methods with our implementation using the ORACLE model selection strategy. * denotes different bottleneck architectures.

DATASET	METHOD	S-ACC	ENT	DEV	SND	1-SHOT	50-RND	100-RND	ORACLE
OFFICE-HOME	S. ONLY	60.38±0.5	60.73±0.2	60.22±0.3	59.55±0.3	58.92±0.4	60.28±0.4	60.34±0.4	61.87±0.3
	PADA	63.08±0.3	59.74±0.5	52.72±2.8	62.36±0.4	62.00±0.5	63.82±0.4	63.22±0.1	63.72±0.3
	SAFN	62.09±0.2	61.37±0.3	62.03±0.4	62.59±0.1	49.30±0.7	62.00±0.2	62.36±0.2	63.30±0.2
	BA3US	68.32±1.1	73.36±0.6	62.25±7.1	75.37±0.8	65.56±7.6	73.22±0.3	75.19±0.4	75.98±0.3
	AR	65.68±0.3	70.58±0.4	64.32±0.9	70.25±0.2	70.56±0.7	70.26±0.2	70.34±0.2	72.73±0.3
	JUMBOT	62.89±0.2	74.61±0.8	61.28±0.1	72.29±0.2	74.95±0.1	64.95±0.3	75.74±0.3	77.15±0.4
	MPOT	66.24±0.1	64.46±0.1	61.37±0.2	46.92±0.4	68.28±0.2	69.90±0.5	73.06±0.3	77.31±0.5
VISDA	S. ONLY	55.15±2.4	55.24±3.2	55.07±1.2	55.02±2.9	55.72±2.2	57.90±1.1	58.16±0.6	59.48±0.4
	PADA	47.48±4.8	32.32±4.9	43.43±5.3	56.83±1.0	53.15±2.9	55.67±2.5	54.38±2.7	54.57±2.6
	SAFN	58.20±1.7	42.83±6.3	58.62±1.3	44.82±8.8	56.89±2.1	57.90±3.3	59.09±2.8	62.64±1.5
	BA3US	55.10±3.7	65.58±1.4	58.40±1.4	51.07±4.3	64.77±1.4	66.66±2.4	67.44±1.2	67.67±1.3
	AR	66.68±1.0	64.27±3.6	67.20±1.5	55.69±0.9	70.29±1.7	71.91±0.3	72.60±0.8	73.85±0.9
	JUMBOT	60.63±0.7	62.42±2.4	59.86±0.6	77.69±4.2	78.34±1.9	82.85±2.9	83.49±1.9	84.01±1.9
	MPOT	70.02±2.0	74.64±4.4	61.62±1.3	78.40±3.9	70.96±3.7	86.65±5.1	86.69±5.1	86.95±5.0

Table 6: Task accuracy average over seeds 2020, 2021, 2022 on Partial OFFICE-HOME and Partial VISDA for the PDA methods and model selection strategies. For each method, we highlight the best and worst label-free model selection strategies in green and red, respectively.