



HAL
open science

Match and Reweight Strategy for Generalized Target Shift

Alain Rakotomamonjy, Rémi Flamary, Gilles Gasso, Mokhtar Z. Alaya,
Maxime Berar, Nicolas Courty

► **To cite this version:**

Alain Rakotomamonjy, Rémi Flamary, Gilles Gasso, Mokhtar Z. Alaya, Maxime Berar, et al.. Match and Reweight Strategy for Generalized Target Shift. 2020. hal-02866979v1

HAL Id: hal-02866979

<https://hal.science/hal-02866979v1>

Preprint submitted on 12 Jun 2020 (v1), last revised 15 Oct 2020 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Match and Reweight Strategy for Generalized Target Shift

Alain Rakotomamonjy
LITIS, Univ. de Rouen
Criteo AI Lab, Paris
alain.rakoto@insa-rouen.fr

Rémi Flamary
Univ Côte d'Azur, CNRS, OCA Lagrange
remi.flamary@unice.fr

Gilles Gasso
LITIS, INSA de Rouen
gilles.gasso@insa-rouen.fr

Mokhtar El Alaya
LITIS, Univ. Rouen
mokhtarzahdi.alaya@gmail.com

Maxime Bérar
LITIS, Univ. de Rouen
maxime.berar@univ-rouen.fr

Nicolas Courty
Univ. Bretagne Sud, CNRS, Irista
nicolas.courty@irisa.fr

Abstract

We address the problem of unsupervised domain adaptation under the setting of generalized target shift (both class-conditional and label shifts occur). We show that in that setting, for good generalization, it is necessary to learn with similar source and target label distributions and to match the class-conditional probabilities. For this purpose, we propose an estimation of target label proportion by blending mixture estimation and optimal transport. This estimation comes with theoretical guarantees of correctness. Based on the estimation, we learn a model by minimizing a importance weighted loss and a Wasserstein distance between weighted marginals. We prove that this minimization allows to match class-conditionals given mild assumptions on their geometry. Our experimental results show that our method performs better on average than competitors across a range domain adaptation problems including *digits*, *VisDA* and *Office*.

1 Introduction

During the last recent years, machine learning and deep learning methods managed to make significant successful breakthroughs on a large amount of difficult tasks. However, most of these models rely on the classical assumption that data from which the model has been trained and those on which it will be deployed has been sampled from the same probability distribution. In real-world applications, this assumption barely holds, due for instance to different acquisition devices, different protocols, or due to the presence of dataset bias. Unsupervised domain adaptation (UDA) methods aims at mitigating those mismatches in distributions in order to help models generalizing from labeled source domain to an unlabeled target domain.

In the context of unsupervised domain adaptation, there exists a large amount of literature addressing the DA problem under different assumptions. One of the most studied setting is based on the co-

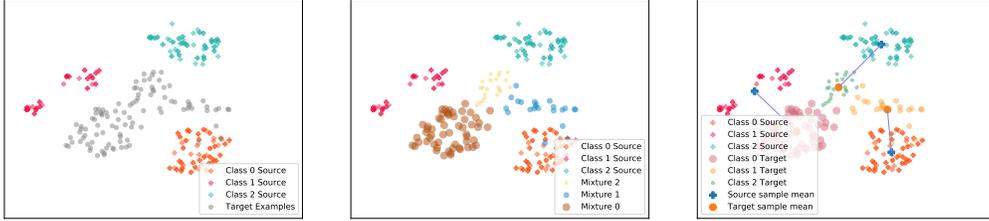


Figure 1: Illustration of how our Match and Reweight strategy works on a 3-class VisDA dataset problem. Panels represent projection into a 2D space of features of original dimension of 100. For a sake of clarity, we have plotted only a fraction of the samples. (left) After learning on the source examples, we plot the 3 classes of the source domain as well as the target examples. (middle) By learning a mixture model on the target examples, we are able to estimate the proportion (represented by the size of the markers) of the 3 classes in the target domain. At this point, we still do not know how to relate components of the mixtures to classes. (right) Based on an optimal transport hypothesis between class-conditional distributions, we match these distributions as represented by their means. Estimated label proportion is then used in a weighted Wasserstein distance suitable for adversarial domain adaptation learning.

ariate shift assumption ($p_S(x) \neq p_T(x)$ and $p_S(y|x) = p_T(y|x)$) Sugiyama et al. (2007). Indeed, leveraging on the flexibility of deep learning models to learn rich feature representations, several works have explored the idea of aligning the marginal distributions in some learned feature space while minimizing the error on the source domain Ganin & Lempitsky (2015).

When distribution mismatch comes from a shift in the distribution of the output ($p_S(y) \neq p_T(y)$) while $p_S(x|y) = p_T(x|y)$, the problem is denoted as target shift or label shift (Schölkopf et al., 2012). Such a situation is common in real-world applications where one has no control in the proportion of label categories on the test data. This is typically the case in computer-aided diagnosis system for which the frequency of a disease can not be controlled or for a computational advertising where clickthrough rate can not be predicted in advance. While less explored than covariate shift, several methods have already been proposed for addressing target shift. For instance, Lipton et al. (2018) proposed a method for estimating the ratio $p_T(y)/p_S(y)$ by analysing the confusion matrix achieved by a black box predictor. Azizzadenesheli et al. (2019) improved the stability of this importance weighting estimation in a two-step procedure, with generalization guarantees. Recently, Shrikumar et al. (2020) highlighted the importance of Expectation Maximization (EM) to efficiently correct for the difference in class proportion. Li et al. (2019) introduced an optimization scheme based on mean matching for estimating label proportion under the assumption the $p(x|y)$ are matched in a latent space. In the same flavor, Redko et al. (2019) have shown that under the target shift assumption, label proportion in the target domain can be estimated by minimizing the distance between the target marginal and the weighted source marginal.

However as most models now learn the latent representation space, in practical situations we have both a label shift ($p_S(y) \neq p_T(y)$) and class-conditional probability shift ($p_S(x|y) \neq p_T(x|y)$). For this more general DA assumption, denoted as generalized target shift, fewer works have been proposed. Zhang et al. (2013) have been the first one that proposed a methodology for handling both shifts. They used a kernel embedding of distributions for estimating importance weights or for transforming samples so as to match class-conditional distributions. For addressing the same problem Wu et al. (2019) introduced a so-called asymmetrically-relaxed distance on distribution that allows to mitigate the effect of label shift when aligning marginal distributions. Interestingly, they also show that error in the target domain is lower-bounded by the mismatch of label distributions between the two domains. Very recently, Combes et al. (2020) have presented a theoretical analysis of this problem showing that target generalization can be achieved by matching label proportion and class-conditionals in both domains. The key component of their algorithm relies on a importance weight estimation of the label distributions. Unfortunately, although relevant in practice, their label distribution estimator got theoretical guarantee only when class conditionals match across domains and empirically breaks as some class conditionals mismatch is large enough.

Our work addresses UDA with generalized target shift. As mentioned by Wu et al. (2019), in this setting, the key point is to correctly estimate the target label proportion $p_T(y)$ in an unsupervised

way and the main objective of the paper is to solve that estimation problem. More specifically, we make the following contributions. From a theoretical side, we clarify the role of the label shift and class-conditional shift in the target generalization error bound. Our theoretical analysis emphasizes the importance of learning with same label distributions in source and target domains while seeking at minimizing class-conditional shifts in a latent space. We solve the label distribution estimation problem by blending a consistent mixture proportion estimator and an optimal matching assignment problem. While conceptually simple, our strategy is supported by theoretical guarantees of correctness and consistency. Then given the estimated label proportion in the target domain, we theoretically show that finding a latent space in which the Wasserstein distance between weighted marginal distributions is minimized, guarantees that class-conditionals are also matched. Based on those analyses, we thus proposed an algorithm (named MARS from Match And Reweight Strategy) for estimating label proportion followed by minimization of weighted marginal distance. We illustrate in our experimental analyses how MARS copes with label and class-conditional shifts and show that it performs better than competitors on most classical domain adaptation problems.

2 Notations and Framework

Our goal is to address the problem of unsupervised domain adaptation. We assume a learning problem with a source and target domains and respectively note as $p_S(x, y)$ and $p_T(x, y)$ their joint distributions of features and labels. We have at our disposal a labeled source dataset $\{x_i^s, y_i^s\}_{i=1}^{n_s}$ and only unlabeled examples from the target domain $\{x_i^t\}_{i=1}^{n_t}$ with all $x_i \in \mathbb{R}^d$, sampled *i.i.d* from their respective distributions. We are interested in multi-class classification problems with C classes, so we restrict the label to be $y_i^s \in \{1, \dots, C\}$.

Domain adaptation framework Since the seminal work of Ganin & Lempitsky (2015), a common formulation of the covariate shift domain adaptation problem is to learn a mapping of the source and target samples into a latent representation space where the distance between their marginal distributions is minimized and a classifier that learns to correctly predict labels of samples in the source domain. This typically translates into the following optimization problem:

$$\min_{f, g} \frac{1}{n} \sum_{i=1}^{n_s} L(y_i^s, f(g(x_i^s))) + \lambda D(p_S^g, p_T^g) + \Omega(f, g) \quad (1)$$

where $f(\cdot)$ is the classifier, $g(\cdot)$ the feature extraction function, $L(\cdot, \cdot)$ is a continuous loss function differentiable on its second parameter and Ω a regularization term. Note that here p_S^g and p_T^g refers to the marginal distributions of the source and target domains in the latent space and $D(\cdot, \cdot)$ is a distance metric between distributions that measures discrepancy between source and target domains as mapped in a latent space induced by g . Most used distance measures are MMD Tzeng et al. (2014), Wasserstein distance Shen et al. (2018) or adversarial distance Ganin et al. (2016). In the sequel, we will consider marginal distributions in the latent space and thus drop the superscript g in p_S^g and p_T^g .

Theoretically, given the loss function L , our goal is to find a function h , here the composition of two functions f and g that minimizes an expected loss with respect to the true labelling function in the target set. The expected loss for any two functions $h, h' : \mathcal{X} \rightarrow \mathcal{Y}$ over a distribution Q , is defined as $\varepsilon_Q(h, h') = \mathbf{E}_{x \sim Q}[L(h(x), h'(x))]$. As such, our goal is to find f and g that have a small expected loss $\varepsilon_{P_T}(f(g(\cdot)), h^*)$, where h^* is the best possible classifier on the target domain in the hypothesis function class. We provide more details about the generalization in section 3.

Optimal Transport We provide here some background on optimal transport as it will be a key concept for estimating and assigning label proportion. More details can be found in Peyré et al. (2019). Optimal transport measures the distance between two distributions over a space \mathcal{X} given a transportation cost $c : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. It seeks for an optimal coupling between the two measures that minimizes a transportation cost. In a discrete case, denote the two measures as $\mu = \sum_{i=1}^n a_i \delta_{x_i}$ and $\nu = \sum_{i=1}^m b_i \delta_{x'_i}$, the so-called Kantorovitch relaxation of the OT problem seeks for a transportation coupling \mathbf{P} that minimizes the problem

$$\min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{C}, \mathbf{P} \rangle \quad (2)$$

where $\mathbf{C} \in \mathbb{R}^{n \times m}$ is the matrix of all pairwise costs, $C_{i,j} = c(x_i, x'_j)$ and $\Pi(\mathbf{a}, \mathbf{b}) = \{\mathbf{P} \in \mathbb{R}_+^{n \times m} | \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b}\}$ is the transport polytope between the two distributions. The above problem is known as the discrete optimal transport problem and in the specific case where $n = m$ and the weights \mathbf{a} and \mathbf{b} are positive and uniform then the solution of the above problem is a scaled permutation matrix (Peyré et al., 2019). One of the key features of OT that we are going to exploit for solving the domain adaptation problem is its ability to find correspondences between samples in an unsupervised way by exploiting the underlying space geometry. These features have been for instance exploited for unsupervised word translation Alvarez-Melis et al. (2019); Alaux et al. (2019).

3 Match and reweight strategy

We first discuss theoretical insights of domain adaptation for generalized target shift and propose a Match and Reweight strategy to compensate for the shifts.

3.1 Theoretical insights on generalized label shift

In this work, we are interested in a situation where both class conditional distributions in source and target mismatch (*i.e* there exists some j so that $p_S(x|y=j) \neq p_T(x|y=j)$) and target shift occurs. In the sequel, class-conditional probability and label proportion for class j will be respectively noted as $p_S^j := p_S(x|y=j)$ and $p_S^{y=j} := p_S(y=j)$.

Because we have these two sources of mismatch, the resulting domain adaptation problem is difficult and looking for latent representation that aligns the marginal distributions in source and target while minimizing source errors does not ensure small error in target domain. Indeed, as formalized by Wu et al. (2019), when target shift occurs, and both source error and distance between marginal distributions are 0 then the test error is lower bounded by the difference between the label distributions. This statement raises the necessity of adjusting the label distribution of the source domain so as to match the unknown one in the target domain. From the following proposition of Mansour et al. (2009), we derive a bound highlighting the role of both shifts.

Theorem 1. *Mansour et al. (2009) Assume that the loss function $L(\cdot, \cdot)$ is symmetric and satisfies the triangle inequality, then for any hypothesis $h \in \mathcal{H}$, the following holds,*

$$\varepsilon_T(h, f_T) \leq \varepsilon_S(h, h_S^*) + \varepsilon_T(h_T^*, f_T) + \varepsilon_S(h_S^*, h_T^*) + \text{disc}(p_S, p_T)$$

where h_T^* and h_S^* are respectively the minimizers of $\varepsilon_T(h, f_T)$ and $\varepsilon_S(h, f_S)$ over the hypothesis space \mathcal{H} , with f_S and f_T the true labelling functions, and $\text{disc}(p_S, p_T) = \max_{h, h'} |\varepsilon_S(h, h') - \varepsilon_T(h, h')|$.

The first term of the bound is the loss induced by the hypothesis h and it can be optimized over h . The second term depends only on the complexity of the hypothesis space. The third one is the average loss between h_S^* and h_T^* under p_S and the last one is the discrepancy of the marginal distributions. $\varepsilon_S(h_S^*, h_T^*)$ is expected to vanish as the discrepancy between joint distributions decrease. This term is thus related to the hardness of the adaptation problem. Now, let us analyze the last term. We can show by expanding the marginals that

$$\text{disc}(p_S, p_T) = \max_{h, h'} \left| \sum_{j=1}^C [p_S^{y=j} \varepsilon_{S^j}(h, h') - p_T^{y=j} \varepsilon_{T^j}(h, h')] \right|.$$

In this equation, $\varepsilon_{S^j}(h, h')$ and $\varepsilon_{T^j}(h, h')$ are the average losses between h and h' under the class-conditional probability of class j for the source and target domain. Denote as \mathbf{p}_S and \mathbf{p}_T the vectors in \mathbb{R}^C of label proportions and as \mathbf{e}_S and \mathbf{e}_T the vectors composed of the $\varepsilon_{S^j}(h, h')$ and $\varepsilon_{T^j}(h, h')$. Then, we have (details in the appendix)

$$\text{disc}(p_S, p_T) = \max_{h, h'} |\mathbf{p}_S^\top \mathbf{e}_S - \mathbf{p}_T^\top \mathbf{e}_T| \leq \|\mathbf{p}_S - \mathbf{p}_T\|_\infty \max_{h, h'} \|\mathbf{e}_S\|_1 + \max_{h, h'} \|\mathbf{e}_S - \mathbf{e}_T\|_\infty$$

We note that the first term of the upper bound depends on the largest difference of the label proportion in the source and target, and on the sum of the discrepancy of h and h' over the classes. The second term is the difference of losses on the source and target based on class-conditional distributions and it vanishes when distance between class-conditionals is zero. Hence, this bound suggests

that a good model should: i) seek at latent representation that reduces class-conditional probability distances, and ii) either learn from source data with similar label proportions than the target one or minimize the difference of label proportions. Note that having small discrepancies on label proportions and on all class-conditionals is a sufficient condition for having small discrepancy on the joint distribution leading thus to a small average loss $\varepsilon_S(h_S^*, h_T^*)$ in Theorem 1.

The key idea of our approach, instead of considering the learning problem and its guarantee through the marginal $p_S = \sum_{i=1}^C p_S^{y=i} p_S^i$, is to work on a reweighted version denoted as $\tilde{p}_S = \sum_{i=1}^C w_i p_S^i$, with w_i chosen so that no label shift occurs between \tilde{p}_S and p_T . A classical choice that guarantees matching label proportions is to set $w_i = p_T^{y=i} / p_S^{y=i}$ (Sugiyama et al., 2007; Combes et al., 2020), which needs an estimation of $p_T^{y=i}$ or a direct estimation of the ratio. Assuming that the ratio w is known for each sample, from an algorithmic point of view, we aim at learning a hypothesis function $h = f \circ g$ from $\tilde{p}_S(x)$ and p_T that generalizes well on the target domain by optimizing a reweighted version of Equation 1 that writes as

$$\min_{f,g} \frac{1}{n} \sum_{i=1}^{n_s} w(x_i^s) L(y_i^s, f(g(x_i^s))) + \lambda D(\tilde{p}_S^g, p_T^g) + \Omega(f, g) \quad (3)$$

where the discrepancy between marginals D is computed through the dual of Wasserstein distance

$$WD_w(\tilde{p}_s, p_t) = \sup_{\|v\|_L \leq 1} \mathbf{E}_{x \sim \tilde{p}_s} w(x)v(x) - \mathbf{E}_{x \sim p_t} v(x)$$

and the weights are defined as $w(x) = \frac{p_T^{y=i}}{p_S^{y=i}}$ if x is of class i , with $p_T^{y=i}$ being estimated using Algorithm 1 and discussed in the sequel. Note that the first term of equation (3) corresponds to the empirical loss related to the error ε_S in Theorem 1 while the distribution divergence aims at minimizing distance between class-conditional probabilities, the third term in that theorem. We employ a classical adversarial learning strategy (detailed in the supplemental) for optimizing Equation (3) and we use gradient penalty for estimating the Wasserstein distance (Gulrajani et al., 2017). Here, the choice of the Wasserstein distance is dictated by theoretical guarantee that will be made clear later. Interestingly, Combes et al. (2020) have derived a learning problem similar to Equation (3) through a different analysis.

The next subsections describe how we perform estimation of $p_T(y)$ in situations where class-conditionals do not match, which is a much more difficult problem than those investigated in Redko et al. (2019); Combes et al. (2020) that assume matching class-conditionals. First, we discuss how to estimate proportions in the target domain using mixture models. Next we aim at finding a permutation matrix that guarantees, under mild assumption, correspondence between the class-conditional probabilities in the source and estimated mixture of the target domain. This matrix allows us to properly assign the label proportions to class-conditional probabilities in the target domain leading to a proper reweighting. Our final proposition states that under the same assumption, by minimizing the Wasserstein distance of the marginals \tilde{p}_S and p_T , we also minimize the distance between class-conditionals.

3.2 Estimating mixture proportion

In practice, we observe some examples sampled from $\{p_S^j\}$ composing the marginal source distribution associated with their labels and some other instances sampled from $\{p_T^j\}$ but with unknown labels. Hence, for this first step, we assume that the target distribution is a mixture model with C components $\{p_T^j\}$ and we want to estimate the mixture proportion of each component. For this purpose, we have considered two alternative strategies coming from the literature : i) learning a Gaussian mixture model over the data in the target domain. This gives us both the estimate components $\{p_T^j\}$ and the proportion of the mixture \mathbf{p}_u . Under some conditions on its initialization and that the model is well-calibrated, Zhao et al. (2020) have shown that the sample estimator asymptotically converges towards the true mixture model. ii) applying agglomerative clustering on the target samples tells us about the membership class of each sample and thus, it provides the proportion of each component in the mixture.

Remark 1. *Although in practice, our model for each p_T^j is simple as defined by a single component of a Gaussian mixture or a cluster, more complex models can also be considered e.g., modeling p_S^j*

Algorithm 1 Label proportion estimation in the target domain

Require: $\{x_i^s, y_i^s\}, \{x_i^t\}$, number of classes C

Ensure: \mathbf{p}_T : Estimated label proportion

- 1: $\{p_T^i\}, \mathbf{p}_u \leftarrow$ Estimate C mixtures and proportions from $\{x_i^t\}$ with C modes.
 - 2: $\mathbf{D} \leftarrow$ Compute the matrix pairwise distances between all the source p_S^i and target p_T^j modes.
 - 3: $\mathbf{P}^* \leftarrow$ Solve OT problem (2) with \mathbf{D} and uniform marginals as in Proposition 1.
 - 4: $\mathbf{p}_T \leftarrow C \cdot \mathbf{P}^* \mathbf{p}_u$ Permute the mixture proportion on source ($C \cdot \mathbf{P}^*$ is a permutation matrix)
-

and p_T^j as mixture of Gaussians and then computing the OT between mixture of Gaussian mixtures Chen et al. (2018); Delon & Desolneux (2019).

After the target mixture and proportions are estimated, we use OT to recover the correspondence as illustrated in Algorithm 1. Indeed, at the end of the mixture proportion estimating step, we still need to match them with the appropriate class-conditionals in the source. We discuss next a way to retrieve this correspondence permutation matrix and hence to ensure a proper matching between source class-conditionals and the mixture components.

3.3 Matching with optimal transport

In the following, we suppose that we have an estimation of conditional-class probabilities on source domain (based on the empirical distributions) and that one class-conditional probability is associated to one mode in the target mixture model. We recall that by marginalization, $p_S(x)$ is a linear combination of C class-conditional probabilities, i.e., $p_S(x) = \sum_j^C p_S^{y=i} p_S^i(x)$. The problem of matching source and target class-conditional probabilities is difficult and depends on some conditions on their geometrical arrangement. This boils down to an optimal assignment problem with respect to the class-conditional probabilities $\{p_S^j\}$ and $\{p_T^j\}$ and under some conditions on distance between class-conditional probabilities, the assignment would achieve correct matching. More formally, denote as \mathbb{P} the set of probability distributions over \mathbb{R}^d and assume a metric over \mathbb{P} . Assume that we want to optimally assign a finite number C of probability distributions to another set of finite number C of probability distributions, in a minimizing distance sense. Based on a pairwise distance matrix \mathbf{D} between couple of class-conditional probability distributions, the assignment problems looks for the permutation that solves $\min_{\sigma} \frac{1}{C} \sum_j D_{j, \sigma(j)}$. Note that the best permutation σ^* solution to this problem can be retrieved by solving a Kantorovitch relaxed version of the optimal transport (Peyré et al., 2019) with marginals $\mathbf{a} = \mathbf{b} = \frac{1}{C} \mathbb{1}$. Hence, this OT-based formulation of the matching problem can be interpreted as an optimal transport one between discrete measures of probability distributions of the form $\frac{1}{C} \sum_{j=1}^C \delta_{P_S^j}$. In order to be able to correctly match class-conditional probabilities in source and target domain by optimal assignment, we ask ourselves:

Under which conditions on the displacement of class conditional probabilities, solving optimal assignment problem leads to a permutation matrix that achieves correct matching?

In other word, we are looking for conditions of identifiability of classes in the target domain based on their geometry with respect to the classes in source domain. Our proposition below presents an abstract sufficient condition for identifiability based on the notion of cyclical monotonicity and then we exhibit some practical situations in which this property holds.

Proposition 1. Denote as $\nu = \frac{1}{C} \sum_{j=1}^C \delta_{p_S^j}$ and $\mu = \frac{1}{C} \sum_{j=1}^C \delta_{p_T^j}$, representing respectively the class-conditional probabilities in source and target domain. Assume that for any permutation σ of C elements, the following assumption holds

$$\sum_j \mathcal{D}(p_S^j, p_T^j) \leq \sum_j \mathcal{D}(p_S^j, p_T^{\sigma(j)})$$

with \mathcal{D} a bounded distance over probability distributions then solving the optimal transport problem defined in equation (2) with uniform marginals and \mathcal{D} as the ground cost matches correctly class-conditional probabilities. The above condition is known as the \mathcal{D} -cyclical monotonicity.

Proof. The solution \mathbf{P}^* of the OT problem lies on an extremal point of Π_C . Birkhoff's theorem Birkhoff (1946) states that the set of extremal points of Π_C is the set of permutation matrices so that

there exists an optimal solution of the form $\sigma^* : [1, \dots, C] \rightarrow [1, \dots, C]$. The support of \mathbf{P}^* is \mathcal{D} -cyclically monotone (Ambrosio & Gigli, 2013; Santambrogio, 2015) (Theorem 1.38), meaning that $\sum_j^C \mathcal{D}(p_S^j, p_T^{\sigma^*(j)}) \leq \sum_j^C \mathcal{D}(p_S^j, p_T^{\sigma(j)})$, $\forall \sigma \neq \sigma^*$. Then, by hypothesis, σ^* can be identified to the identity permutation, and solving the optimal assignment problem matches correctly class-conditional probabilities. \square

While the cyclical monotonicity above can be hard to grasp, there exists a number of situations where it applies. One condition that is simple and intuitive is when class-conditionals of same source and target classes are "near" each other in the latent space. More formally, if we assume that $\forall j \mathcal{D}(p_S^j, p_T^j) \leq \mathcal{D}(p_S^j, p_T^k) \quad \forall k$, then summing over all possible j , and choosing k so that all the couples of (j, k) form a permutation, we recover the cyclical monotonicity condition $\sum_j^C \mathcal{D}(p_S^j, p_T^j) \leq \sum_j^D \mathcal{D}(p_S^j, p_T^{\sigma(j)})$, $\forall \sigma$.

Another more general condition, that does not include the above example, on the identifiability of the target class-conditional can be retrieved by exploiting the fact that, for discrete optimal transport with uniform marginals, the support of optimal transport plan satisfies the cyclical monotonicity condition (Santambrogio, 2015). This is for instance the case, when p_S^j and p_T^j are Gaussian distributions of same covariance matrices and the mean m_T^j of each p_T^j is obtained as a linear symmetric positive definite mapping of the mean m_S^j of p_S^j and the distance $\mathcal{D}(p_S^j, p_T^j)$ is $\|m_S^j - m_T^j\|_2$ (Courty et al., 2016). This situation would correspond to a linear shift of the class-conditionals of the source domain to get the target ones.

It is interesting to compare our assumptions on identifiability to other hypotheses proposed in the literature for solving (generalized) target shift problems. When handling only target shift, one common hypothesis Redko et al. (2019) is that class-conditional probabilities are equal. This in our case boils down to have a 0 distance between $\mathcal{D}(P_S^j, P_T^j)$ guaranteeing matching. When both shifts occur on labels and class-conditionals, Wu et al. (2019) assume that there exists continuity of support between the $p(x|y)$ in source and target domains. Again, this assumption may be related to the above minimum distance hypothesis if class-conditionals in source domain are far enough. Interestingly, one of the hypothesis of Zhang et al. (2013) for handling generalized target shift is that there exists a linear transformation between the class-conditional probabilities in source and target domains. This is a particular case of our Proposition 1 and subsequent discussion, where the Monge mapping T is supposed to be linear. Our conditions for correct matching and thus for identifying classes in the target domains are more general than those proposed in the current literature.

The above proposition helps us in assigning the estimated label proportions to correct class-conditional probabilities under some mild assumptions. Interestingly those assumptions give us also guarantee that minimizing the Wasserstein distance between marginals also induces minimal distances between class-conditional probabilities.

Proposition 2. *Denote as γ the optimal coupling plan for marginals with balanced class-conditionals under assumptions given in Proposition 1. Assume that the classes are ordered so that we have $\gamma = \frac{1}{C} \text{diag}(\mathbb{1})$ then $\gamma' = \text{diag}(\mathbf{a})$ is also optimal for the transportation problem with marginals $\nu' = \sum_{j=1}^C a_j \delta_{p_S^j}$ and $\mu' = \sum_{j=1}^C a_j \delta_{p_T^j}$, with $a_j > 0, \forall j$. In addition, if the Wasserstein distance between ν' and μ' is 0, it implies that the distance between class-conditionals are all 0.*

Proof. By assumption and without loss of generality, the class-conditionals are arranged so that $\gamma = \frac{1}{C} \text{diag}(\mathbb{1})$. Because the weights in the marginals are not uniform anymore, γ is not a feasible solution for the OT problem with ν' and μ' but $\gamma' = \text{diag}(\mathbf{a})$ is. Let us now show that any feasible non-diagonal plan Γ has higher cost than γ' and thus is not optimal. At first, consider any permutation σ of C elements and its corresponding permutation matrix \mathbf{P}_σ , because $\gamma = \frac{1}{C} \text{diag}(\mathbb{1})$ is optimal, the cyclical monotonicity relation $\sum_i \mathcal{D}_{i,i} \leq \sum_i \mathcal{D}_{i,\sigma(i)}$ holds true $\forall \sigma$. Starting from $\gamma' = \text{diag}(\mathbf{a})$, any direction $\Delta_\sigma = -\mathbf{I} + \mathbf{P}_\sigma$ is a feasible direction (it does not violate the marginal constraints) and due to the cyclical monotonicity, any move in this direction will increase the cost. Since any non-diagonal $\gamma_z \in \Pi(\mathbf{a}, \mathbf{a})$ can be reached with a sum of displacements Δ_σ (property of the Birkhoff polytope) it means that the transport cost induced by γ_z will always be greater or equal to the cost for the diagonal γ' implying that γ' is the solution of the OT problem with marginals \mathbf{a} .

As a corollary, it is straightforward to show that $W(\nu', \mu') = \sum_{i=1}^C \mathcal{D}_{i,i} a_i = 0 \implies \mathcal{D}_{i,i} = 0$ as $a_i > 0$ by hypothesis. \square

Interestingly, this proposition brings us the guarantee that under some geometrical assumptions on the class-conditionals in the latent space, minimizing the Wasserstein distance of the marginals also minimizes distances between class-conditionals. Hence, minimizing the divergence term in our learning problem in Equation (3) helps in reducing the upper bound on $\text{disc}(p_S, p_T)$ and on reducing the third term $\varepsilon_T(h_S^*, h_T^*)$ of the generalization bound.

4 Discussions

Most related works are the one by Wu et al. (2019) and Combes et al. (2020) that also address generalized target shift. The first approach does not seek at estimating label proportion but instead allows flexibility in the alignment by using an asymmetric relaxed distance. In the case of Wasserstein distance, their approach consists in reweighting the marginal of the target distribution and in its dual form, their distance boils to

$$WD_w(p_S, p_T) = \sup_{\|v\|_L \leq 1} \mathbf{E}_{x \sim p_S} w(x)v(x) - \mathbf{E}_{x \sim p_T} v(x)$$

where $w(\cdot)$ is actually a constant $\frac{1}{1+\beta}$. We can note that the adversarial loss we propose is a general case of this one. Indeed, in the above, the same amount of weighting applies to all the samples of the source distribution. At the contrary, our reweighting scheme depends on the class-conditional probability and their estimate target label proportion. Hence, we believe that our approach would adapt better to imbalance without the need to tune β (by validation for instance, which is hard in unsupervised domain adaptation). The work of Combes et al. (2020) and our differs only in the way the weights $w(\mathbf{x})$ are estimated. In our case, we consider a theoretically supported and consistent estimation of the target label proportion, instead they directly estimate $w(\cdot)$ by applying a technique tailored and grounded for problems without class-conditional shifts. We will show in the experimental section that their estimator in some cases lead to poor generalization.

Still in reweighting, Yan et al. (2017) proposed a weighted Maximum Mean discrepancy distance for handling target shift in UDA. However, their weights are estimated based on pseudo-labels obtained from the learned classifier and thus, it is difficult to understand whether they provide accurate estimation of label proportion even in simple setting. While their distance is MMD-transposed version of our weighted Wasserstein, our approach is more theoretically grounded as the label proportion estimation is based on sound algorithm with proven convergence guarantees (see below) and our optimal assignment hypothesis provides guarantees under which situations class-conditional probability matching is correct.

The idea of matching moment of distributions have already been proven to be an effective for handling distribution mismatch. About ten years ago, Huang et al. (2007); Gretton et al. (2009); Yu & Szepesvári (2012) already leveraged such an idea for handling covariate shift by matching means of distributions in some reproducing kernel Hilbert space. Li et al. (2019) recycled the same idea for label proportion estimation and extended the idea to distribution matching. Interestingly, our approach differs on its usage. While most above works employ mean matching for density ratio estimation or for label proportion estimation, we use it as a mean for identifying displacement of class-conditional distributions through optimal assignment/transport. Hence, it allows us to assign estimated label proportion to the appropriate class.

For estimating the label proportion, we have proposed to learn a Gaussian mixture model of the target distribution. By doing so we are actually trying to solve a harder problem than necessary. However, once the target distribution estimation has been evaluated and class-conditional probabilities being assigned from the source class, one can use that Gaussian mixture model for labelling the target samples. Note however that Gaussian mixture learned by expectation-minimization can be hard to estimate especially in high-dimension Zhao et al. (2020) and that the speed of convergence of the EM algorithm depends on smallest mixture weights Naim & Gildea (2012). Hence, in high-dimension and/or highly imbalanced situations, one may get a poor estimate of the target distribution. Nonetheless, one can consider other non-EM approach Kannan et al. (2005); Arora et al. (2005). Hence, in practice, we can expect the approach GMM estimation and OT-based match-

Table 1: Table of averaged **balanced accuracy** for the compared models and different domain adaptation problems and label proportion imbalance settings. Reported in bold are the best performances as well as other methods which achieve performance that are statistically similar according to a Wilcoxon signrank test with $p = 0.01$. Last lines present the summary of 34 experiments presented in the supplementary (including experiments on *Office*). #Win includes the statistical ties.

Setting	Source	DANN	WD $_{\beta=0}$	WD $_{\beta=1}$	WD $_{\beta=2}$	WD $_{\beta=3}$	WD $_{\beta=4}$	IW-WD	MARSg	MARS c
MNIST-USPS 10 modes										
Balanced	76.89 \pm 3.7	79.74 \pm 3.5	93.71 \pm 0.7	74.27 \pm 4.3	51.33 \pm 4.0	76.61 \pm 3.3	71.90 \pm 5.7	95.28 \pm 0.4	95.61\pm0.7	95.64\pm1.0
Mid	80.41 \pm 3.1	78.65 \pm 3.0	94.30 \pm 0.7	75.36 \pm 3.4	55.55 \pm 4.3	78.98 \pm 3.1	72.32 \pm 4.2	95.60\pm0.5	89.70 \pm 2.3	90.39 \pm 2.6
High	78.13 \pm 4.9	81.79 \pm 4.0	93.86\pm1.1	87.44 \pm 1.7	83.83 \pm 5.2	85.65 \pm 2.5	83.65 \pm 3.0	94.08\pm1.0	88.30 \pm 1.5	89.65 \pm 2.3
USPS-MNIST 10 modes										
Balanced	77.04 \pm 2.6	80.49 \pm 2.2	73.35 \pm 2.8	66.70 \pm 2.9	49.86 \pm 2.8	55.83 \pm 2.9	52.12 \pm 3.5	80.52 \pm 2.2	84.59\pm1.7	85.50\pm2.1
Mid	79.54\pm2.8	78.88\pm1.8	75.85 \pm 1.6	63.33 \pm 2.3	53.22 \pm 2.8	47.20 \pm 2.4	48.29 \pm 2.9	78.36\pm3.5	79.73\pm3.6	78.49\pm2.5
High	78.48\pm2.4	77.79\pm2.0	76.14\pm2.7	63.00 \pm 3.3	57.56 \pm 4.8	51.19 \pm 4.4	49.31 \pm 3.3	71.53 \pm 4.7	75.62 \pm 1.8	77.14\pm2.4
MNIST-MNISTM 10 modes										
Setting 1	58.34 \pm 1.3	61.22\pm1.1	57.44 \pm 1.7	50.20 \pm 4.4	47.01 \pm 2.0	57.85 \pm 1.1	55.95 \pm 1.3	63.10\pm3.1	58.08 \pm 2.3	56.58 \pm 4.6
Setting 2	59.94 \pm 1.1	61.09 \pm 1.0	58.08 \pm 1.4	53.39 \pm 3.5	48.61 \pm 2.4	59.74 \pm 0.7	58.14 \pm 0.8	65.03\pm3.5	57.69 \pm 2.3	55.64 \pm 2.1
Setting 3	58.14 \pm 1.2	60.39\pm1.4	57.68 \pm 1.2	47.72 \pm 4.9	42.15 \pm 7.3	57.09 \pm 1.0	53.52 \pm 1.1	52.46 \pm 14.8	53.68 \pm 7.2	53.72 \pm 3.3
VisDA 12 modes										
setting 1	41.90 \pm 1.5	52.79 \pm 2.1	45.81 \pm 4.3	44.23 \pm 3.0	35.45 \pm 4.6	40.96 \pm 3.0	37.59 \pm 3.4	50.35 \pm 2.3	53.31 \pm 0.9	55.05\pm1.6
setting 2	41.75 \pm 1.5	50.82 \pm 1.6	45.72 \pm 8.9	40.49 \pm 4.8	36.21 \pm 5.0	36.12 \pm 4.6	31.86 \pm 5.7	48.59 \pm 1.8	53.09 \pm 1.6	55.33\pm1.6
setting 3	40.64 \pm 4.3	49.17 \pm 1.3	47.12 \pm 1.6	42.10 \pm 3.0	36.32 \pm 4.4	37.26 \pm 3.5	34.96 \pm 5.4	46.59 \pm 1.3	50.78 \pm 1.6	52.08\pm1.2
#Wins (/34)	7	9	5	0	1	0	2	9	12	21
Aver. Rank	4.16	4.73	5.32	6.97	8.38	6.59	7.57	4.95	3.38	2.95

ing to be a strong baseline in low-dimension and well-clustered mixtures setting but to break in high-dimension one.

5 Numerical Experiments

Experimental setup We compare several domain adaptation algorithms tailored for covariate shift and two very recent methods designed for generalized target shift. As a baseline, we consider a model, denoted as Source trained for f and g on the source examples and tested without adaptation on the target examples. Two other competitors apply adversarial domain learning using approximation of the \mathcal{H} divergence and the Wasserstein distance computed in the dual as distances for measuring discrepancy between p_S and p_T , denoted as DANN and WD $_{\beta=0}$. We consider the model proposed by Wu et al. (2019) and Combes et al. (2020) as competing algorithms able to cope with generalized target shift. For this former approach, we use the asymmetrically-relaxed Wasserstein distance so as to make it similar to our approach and also reported results for different values of the relaxation β . This model is named WD $_{\beta}$ with $\beta \geq 1$. The Combes et al. (2020)’s method, named IW-WD (for importance weighted Wasserstein distance) solves the same learning problem as ours and differs only on the way the ratio $w(x_i)$ is estimated. Our approaches are denoted as MARS c or MARSg respectively when estimating proportion by clustering or by Gaussian mixture models. All methods differ only in the metric used for computing the distance between marginal distributions and most of them except DANN use a Wasserstein distance. The difference essentially relies on the reweighting strategy of the source samples. For all models, learning rate and the hyperparameter λ in Equation 3 have been chosen based on a reverse cross-validation strategy. The metric that we have used for comparison is the balanced accuracy (the average recall obtained on each class) which is better suited for imbalanced problems (Brodersen et al., 2010). All presented results have been obtained as averages over 20 runs.

Toy dataset The toy dataset is a 3-class problem in which class-conditional probabilities are Gaussian distributions. For the source distribution, we fix the mean and the covariance matrix of each of the three Gaussians and for the target, we simply shift the means (by a fixed translation). We have carried out two sets of experiments where we have fixed the shift and modified the label proportion imbalance and another one with fixed imbalance and increasing shift. For space reasons, we have deported to the supplementary the results of the latter. Figure 2 show how models perform for varying imbalance and fixed shift. The plots nicely show what we expect. DANN performs worse as the imbalance increases. WD $_{\beta}$ works well for all balancing but its parameter β needs to increase with the imbalance level. Because of the shift in class-conditional probabilities, IW-WD is not able to properly estimated the importance weights and fails. Our approaches are adaptive to the imbalance and perform very well over a large range for both a low-noise and mid-noise setting (examples of

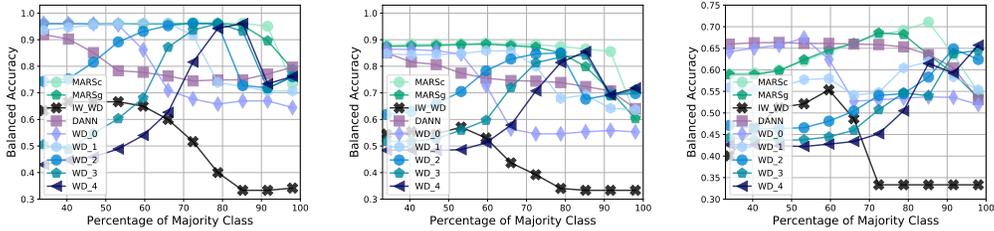


Figure 2: Performance of the compared algorithms for three different covariance matrices of the Gaussians composing the toy dataset with respect to the imbalance. The x-axis is given with respect to the percentage of majority class which is the class 1. (left) Low-error setting. (middle) mid-error setting. (right) high-error setting. Example of the source and target samples for the different cases are provided in the supplementary material.

how the Gaussians are mixed are provided in the supplementary material). For the hardest problem (most-right panel), all models have difficulties and achieve only a balanced accuracy of 0.67 over some range of imbalance. Note that for this low-dimension toy problem, as expected, the approach GMM and OT-based matching achieves the best performance as reported in the supplementary material.

Digits , VisDA (Peng et al., 2017) and Office (Venkateswara et al., 2017) We present some UDA experiments on computer vision datasets, with different imbalanced settings. Details of problem configurations as well as model architecture and training procedure can be found in the appendix. Table 1 reports the averaged balanced accuracy achieved by the different models for only for a fairly chosen subset of problems. The full table is in the supplementary. Results presented here are not comparable to results available in the literature as they mostly consider covariate shift DA (hence with balanced proportions). For these subsets of problems, our approaches yield the best average ranking. They perform better than competitors except on the MNIST-MNISTM problems. As the key issue in generalized target shift problem is the ability to estimate accurately the importance weight or the target label proportion, we believe that the learnt latent representation fairly satisfies our OT hypothesis leading to good performance.

6 Conclusion

The paper proposed a strategy for handling generalized target shift in domain adaptation. It builds upon the simple idea that if the target label proportion where known, then reweighting class-conditional probabilities in the source domain is sufficient for designing a distribution discrepancy that takes into account those shifts. In practice, our algorithm boils down to estimate the label proportion using classical methods such as Gaussian Mixture models or agglomerative clustering and then in matching source and target means of those components for allocating properly the components. Resulting label proportion is then plugged into an weighted Wasserstein distance. When employed for adversarial domain adaptation, we show that our approach outperforms competitors and is able to adapt to imbalance in target domains.

Several points are worth to be extended in future works. At the present time, we have considered simple mean-based approach for matching distributions, it is worth investigating whether higher-order moments are useful for improving the matching. Our algorithm relies mostly on our ability to estimate label proportion, we would be interested on in-depth theoretical analysis label proportion estimation and their convergence and convergence rate guarantees.

Acknowledgments

This work benefited from the support of the project OATMIL ANR-17-CE23-0012 of the French, LEAUDS ANR-18-CE23, and was performed using computing resources of CRIANN (Normandy, France).

References

- Alaux, J., Grave, E., Cuturi, M., and Joulin, A. Unsupervised hyper-alignment for multilingual word embeddings. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. URL <https://openreview.net/forum?id=HJe62s09tX>.
- Alvarez-Melis, D., Jegelka, S., and Jaakkola, T. S. Towards optimal transport with global invariances. In Chaudhuri, K. and Sugiyama, M. (eds.), *Proceedings of Machine Learning Research*, volume 89 of *Proceedings of Machine Learning Research*, pp. 1870–1879. PMLR, 16–18 Apr 2019.
- Ambrosio, L. and Gigli, N. A users guide to optimal transport. In *Modelling and optimisation of flows on networks*, pp. 1–155. Springer, 2013.
- Arora, S., Kannan, R., et al. Learning mixtures of separated nonspherical gaussians. *The Annals of Applied Probability*, 15(1A):69–92, 2005.
- Azizzadenesheli, K., Liu, A., Yang, F., and Anandkumar, A. Regularized learning for domain adaptation under label shifts. In *International Conference on Learning Representations (ICLR)*, 2019. URL <https://openreview.net/forum?id=rJ10r3R9KX>.
- Birkhoff, G. Tres observaciones sobre el algebra lineal. *Univ. Nac. Tucumán Rev. Ser. A*, 1946.
- Brodersen, K. H., Ong, C. S., Stephan, K. E., and Buhmann, J. M. The balanced accuracy and its posterior distribution. In *2010 20th International Conference on Pattern Recognition*, pp. 3121–3124. IEEE, 2010.
- Chen, Y., Georgiou, T. T., and Tannenbaum, A. Optimal transport for gaussian mixture models. *IEEE Access*, 7:6269–6278, 2018.
- Combes, R. T. d., Zhao, H., Wang, Y.-X., and Gordon, G. Domain adaptation with conditional distribution matching and generalized label shift. *arXiv preprint arXiv:2003.04475*, 2020.
- Courty, N., Flamary, R., Tuia, D., and Rakotomamonjy, A. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865, 2016.
- Delon, J. and Desolneux, A. A wasserstein-type distance in the space of gaussian mixture models. *arXiv preprint arXiv:1907.05254*, 2019.
- Ganin, Y. and Lempitsky, V. Unsupervised domain adaptation by backpropagation. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1180–1189, Lille, France, 07–09 Jul 2015. PMLR. URL <http://proceedings.mlr.press/v37/ganin15.html>.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- Gretton, A., Smola, A., Huang, J., Schmittfull, M., Borgwardt, K., and Schölkopf, B. Covariate shift by kernel mean matching. *Dataset shift in machine learning*, 3(4):5, 2009.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pp. 5767–5777, 2017.
- Huang, J., Gretton, A., Borgwardt, K., Schölkopf, B., and Smola, A. J. Correcting sample selection bias by unlabeled data. In *Advances in neural information processing systems*, pp. 601–608, 2007.
- Kannan, R., Salmasian, H., and Vempala, S. The spectral method for general mixture models. In *International Conference on Computational Learning Theory*, pp. 444–457. Springer, 2005.
- Li, Y., Murias, M., Major, S., Dawson, G., and Carlson, D. On target shift in adversarial domain adaptation. In Chaudhuri, K. and Sugiyama, M. (eds.), *Proceedings of Machine Learning Research*, volume 89 of *Proceedings of Machine Learning Research*, pp. 616–625. PMLR, 16–18 Apr 2019. URL <http://proceedings.mlr.press/v89/li19b.html>.
- Lipton, Z. C., Wang, Y.-X., and Smola, A. Detecting and correcting for label shift with black box predictors. *arXiv preprint arXiv:1802.03916*, 2018.
- Mansour, Y., Mohri, M., and Rostamizadeh, A. Domain adaptation: Learning bounds and algorithms. In *Proceedings of The 22nd Annual Conference on Learning Theory (COLT 2009)*, Montréal, Canada, 2009. URL <http://www.cs.nyu.edu/~mohri/postscript/nadap.pdf>.

- Naim, I. and Gildea, D. Convergence of the EM algorithm for gaussian mixtures with unbalanced mixing coefficients. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*, 2012. URL <http://icml.cc/2012/papers/814.pdf>.
- Peng, X., Usman, B., Kaushik, N., Hoffman, J., Wang, D., and Saenko, K. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*, 2017.
- Peyré, G., Cuturi, M., et al. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- Redko, I., Courty, N., Flamary, R., and Tuia, D. Optimal transport for multi-source domain adaptation under target shift. In Chaudhuri, K. and Sugiyama, M. (eds.), *Proceedings of Machine Learning Research*, volume 89 of *Proceedings of Machine Learning Research*, pp. 849–858. PMLR, 16–18 Apr 2019. URL <http://proceedings.mlr.press/v89/redko19a.html>.
- Santambrogio, F. Optimal transport for applied mathematicians. *Birkäuser, NY*, 55(58-63):94, 2015.
- Schölkopf, B., Janzing, D., Peters, J., Sgouritsa, E., Zhang, K., and Mooij, J. On causal and anticausal learning. In *Proceedings of the 29th International Conference on International Conference on Machine Learning, ICML12*, pp. 459466, Madison, WI, USA, 2012. Omnipress. ISBN 9781450312851.
- Shen, J., Qu, Y., Zhang, W., and Yu, Y. Wasserstein distance guided representation learning for domain adaptation. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Shrikumar, A., Alexandari, A. M., and Kundaje, A. Adapting to label shift with bias-corrected calibration, 2020. URL <https://openreview.net/forum?id=rkx-wA4YPS>.
- Sugiyama, M., Krauledat, M., and Müller, K.-R. Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research*, 8(May):985–1005, 2007.
- Tzeng, E., Hoffman, J., Zhang, N., Saenko, K., and Darrell, T. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- Venkateswara, H., Eusebio, J., Chakraborty, S., and Panchanathan, S. Deep hashing network for unsupervised domain adaptation. In *(IEEE) Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- Wu, Y., Winston, E., Kaushik, D., and Lipton, Z. Domain adaptation with asymmetrically-relaxed distribution alignment. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6872–6881, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/wu19f.html>.
- Yan, H., Ding, Y., Li, P., Wang, Q., Xu, Y., and Zuo, W. Mind the class weight bias: Weighted maximum mean discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2272–2281, 2017.
- Yu, Y. and Szepesvári, C. Analysis of kernel mean matching under covariate shift. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*, 2012. URL <http://icml.cc/2012/papers/330.pdf>.
- Zhang, K., Schölkopf, B., Muandet, K., and Wang, Z. Domain adaptation under target and conditional shift. In *International Conference on Machine Learning*, pp. 819–827, 2013.
- Zhao, R., Li, Y., Sun, Y., et al. Statistical convergence of the em algorithm on gaussian mixture models. *Electronic Journal of Statistics*, 14(1):632–660, 2020.

Supplementary material for Match and Reweight for Generalized Target Shift

This supplementary material presents some details of the theoretical and algorithmic aspects of the work as well as some additional results. They are listed as below.

1. Theoretical details on the discrepancy upper bound is given in Section 7.
2. The full algorithm for MARS is detailed and a pseudo-code is given in Algorithm 2
3. Dataset details and architecture details are given in Section 8.1 and 8.2
4. Figure 3 presents some samples of the 3-class toy data set for different configurations of covariance matrices making the problem easy, of mid-difficulty or difficult.
5. Examples of source and target class-conditionals that allow class matching through optimal transport 4 as discussed in Proposition 1.
6. Figure 5 exhibits the performances of the compared algorithms depending on the shift of the class-conditional distributions.
7. Figure 6 shows for the imbalanced toy problem, the results obtained by all competitors including a GMM.
8. Table 2 shows the performance of Source only and a simple GMM+OT on a Visda 3-class problem.
9. Table 3 depicts the different configurations of the dataset we used in our experiments
10. The full table presenting the experimental results for all competitors on different dataset settings is in Table 4.
11. Examples of label proportion error estimation is given in Figure 7.
12. Examples of *t-sne* embeddings on the VisDA-3 problem, given in Figure 8 illustrating the features obtained by DANN, $WD_\beta = 1$, IW-WD and MARSc.

7 Theoretical and algorithmic details

7.1 Bounding the discrepancy between marginals

Denote as \mathbf{p}_S and \mathbf{p}_T the vectors of label proportions and as \mathbf{e}_S and \mathbf{e}_T the vectors composed of the $\epsilon_{Sj}(h, h')$ and $\epsilon_{Tj}(h, h')$. Then, we have the following definition of the discrepancy

$$\begin{aligned}
 \text{disc}(p_S, p_T) &= \max_{h, h'} |\mathbf{p}_S^\top \mathbf{e}_S - \mathbf{p}_T^\top \mathbf{e}_T| \\
 &\leq \max_{h, h'} |\mathbf{p}_S^\top \mathbf{e}_S - \mathbf{p}_T^\top \mathbf{e}_T - \mathbf{p}_T^\top \mathbf{e}_S + \mathbf{p}_T^\top \mathbf{e}_S| \\
 &\leq \max_{h, h'} |(\mathbf{p}_S - \mathbf{p}_T)^\top \mathbf{e}_S + \mathbf{p}_T^\top (\mathbf{e}_S - \mathbf{e}_T)| \\
 &\leq \max_{h, h'} \|\mathbf{p}_S - \mathbf{p}_T\|_\infty \|\mathbf{e}_S\|_1 + \max_{h, h'} \|\mathbf{p}_T\|_1 \|\mathbf{e}_S - \mathbf{e}_T\|_\infty \\
 &\leq \|\mathbf{p}_S - \mathbf{p}_T\|_\infty \max_{h, h'} \|\epsilon_S\|_1 + \max_{h, h'} \|\mathbf{e}_S - \mathbf{e}_T\|_\infty
 \end{aligned}$$

7.2 Algorithm for training the full MARS model

We present here the algorithm we have used for training the full model. It is based on a standard backpropagation strategy using stochastic gradient descent. We estimate the label proportion using Algorithm 1 and then uses this proportion for computing the importance weights $w(\cdot)$. The first part of the algorithm consists then in computing the weighted Wasserstein distance using gradient penalty (Gulrajani et al., 2017). Once this distance is computed, we backpropagated the error through the parameters of the feature extractor g and the classifier f .

In practice, we first train the model without adaptation (hence only based on the classification loss with uniform weights, until reaching 0 training errors and then start adapting as detailed in Algorithm 2)

Algorithm 2 Training the full MARS model

Require: $\{x_i^s, y_i^s\}, \{x_i^t\}$, number of classes C , batch size B , number of critic iterations n

Ensure: \mathbf{p} : label proportion

- 1: Initialize feature extractor g , the classifier f and the domain critic $v(\cdot)$, with parameters $\theta_f, \theta_g, \theta_v$
 - 2: **repeat**
 - 3: estimate \mathbf{p}_T from $\{x_i^t\}$ using Algorithm 1 {done every 10 iterations}
 - 4: sample minibatches $\{x_B^s, y_B^s\}, \{x_B^t\}$ from $\{x_i^s, y_i^s\}$ and $\{x_i^t\}$
 - 5: compute $\{w_i\}_{i=1}^C$ based on the source proportion in the batch samples and \mathbf{p}_T
 - 6: **for** $t = 1, \dots, n$ **do**
 - 7: $x_e^s \leftarrow g(x_B^s), x_e^t \leftarrow g(x_B^t)$
 - 8: sample random points x' from the lines between x_e^s and x_e^t pairs.
 - 9: compute gradient penalty $\mathcal{L}_{\text{grad}}$ using x_e^s, x_e^t and x'
 - 10: compute empirical Wasserstein dual loss $\mathcal{L}_{wd} = \sum_i w(\mathbf{x}_i^s)v(x_i^s) - \frac{1}{B} \sum_i v(x_i^t)$
 - 11: $\theta_v \leftarrow \theta_v + \alpha_v \nabla_{\theta_v} [\mathcal{L}_{wd} - \mathcal{L}_{grad}]$
 - 12: **end for**
 - 13: compute the weighted classification loss $\mathcal{L}_w = \sum_i w(x_i^s)L(y_i^s, x_e^s)$
 - 14: $\theta_f \leftarrow \theta_f + \alpha_f \nabla_{\theta_f} \mathcal{L}_w$
 - 15: $\theta_g \leftarrow \theta_g + \alpha_g \nabla_{\theta_g} [\mathcal{L}_w + \mathcal{L}_{wd}]$
 - 16: **until** fdf
-

8 Experimental Results

8.1 Dataset details

We have considered 4 family of domain adaptation problems based on the digits, Visda, Office-31 and Office-Home dataset. For all these datasets, we have not considered the natural train/test number of examples, in order to be able to build different label distributions at constant number of examples (suppose one class has at most 800 examples, if we want that class to represent 80% of the samples, then we are limited to 1000 samples).

For the digits problem, We have used the MNIST, USPS and the MNITSM datasets. we have learned the feature extractor from scratch and considered the following train-test number of examples setting. For MNIST-USPS, USPS-MNIST and MNIST-MNITSM, we have respectively used 60000-3000, 7291-10000, 10000-10000.

The VisDA 2017 problem is a 12-class classification problem with source and target domain being simulated and real images. We have considered two sets of problem, a 3-class one (based on the classes *aeroplane*, *horse* and *truck*) and the full 12-class problem.

The Office-31 is an object categorization problem involving 31 classes with a total of 4652 samples. There exists 3 domains in the problem based on the source of the images : Amazon (A), DSLR (D) and WebCam (W). We have considered all possible pairwise source-target domains.

The Office-Home is another object categorization problem involving 65 classes with a total of 15500 samples. There exists 4 domains in the problem based on the source of the images : Art, Product, Clipart (Clip), Realworld (Real).

For the Visda and Office datasets, we have considered Imagenet pre-trained ResNet-50 features and our feature extractor (which is a fully-connected feedforward networks) aims at adapting those features. We have used pre-trained features freely available at <https://github.com/jindongwang/transferlearning/blob/master/data/dataset.md>.

8.2 Architecture details

Toy The feature extractor is a 2 layer fully connected network with 200 units and ReLU activation function. The classifier is also a 2 layer fully connected network with same number of units and activation function. Discriminators have 3 layers with same number of units.

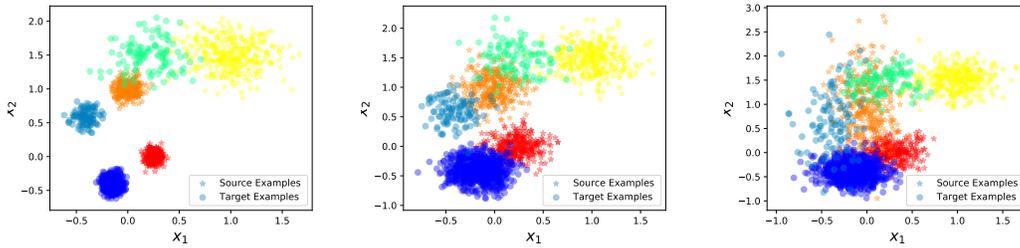


Figure 3: Examples of source and target domain examples. For each domain, data are composed of three Gaussians defining each class. In the source domain, classes are balanced whereas in the target domain, we have a ratio of 0.8, 0.1, 0.1. The three configurations presented here vary in their covariance matrices. From left to right, we have Gaussians that are larger and larger making them difficult to classify. In the most right examples, the second class of the source domain and the third one of the target domain are mixed. This region becomes indecidable for our model as the source loss want to classify it as "Class 2" while the Wasserstein distance want to match it with "Class 3" of the source domain.

Table 2: Comparing Source-Only model and GMM+OT approach on the VisDA-3-mode problems. We can note that for these problems where the latent space is of dimension 100, the GMM+OT compares poorly to Source-Only. In addition, we can note that there is very high variability in the performance.

Configuration	Source	GMM+OT
Setting 1	79.28±4.3	81.22±4.7
Setting 4	80.15±5.3	76.28±9.8
Setting 2	81.47±3.5	74.79±10.4
Setting 3	78.35±3.2	69.97±10.8
Setting 5	83.52± 3.5	76.95±10.4
Setting 6	80.84±4.2	72.86±10.2
Setting 7	79.22±3.7	69.48±9.8

Digits For the MNIST-USPS problem, the architecture of our feature extractor is composed of the two CNN layers with 32 and 20 filters of size 5×5 and 2-layer fully connected networks as discriminators with 100 and 10 units. The feature extractor uses a ReLU activation function and a max pooling. For the MNIST-MNISTM adaptation problem we have used the same feature extractor network and discriminators as in Ganin & Lempitsky (2015).

VisDA For the VisDA dataset, we have considered pre-trained 2048 features obtained from a ResNet-50 followed by 2 fully connected networks with 100 units and ReLU activations. The latent space is thus of dimension 100. Discriminators and classifiers are also a 2 layer Fully connected networks with 100 and respectively 1 and "number of class" units.

Office For the office datasets, we have considered pre-trained 2048 features obtained from a ResNet-50 followed by two fully connected networks with output of 100 and 50 units and ReLU activations. The latent space is thus of dimension 50. Discriminators and classifiers are also a 2 layer fully connected networks with 50 and respectively 1 and "number of class" units.

For Digits and VisDA and Office applications, all models have been trained using ADAM for 100 iterations with validated learning rate, while for the toy problem, we have used a SGD.

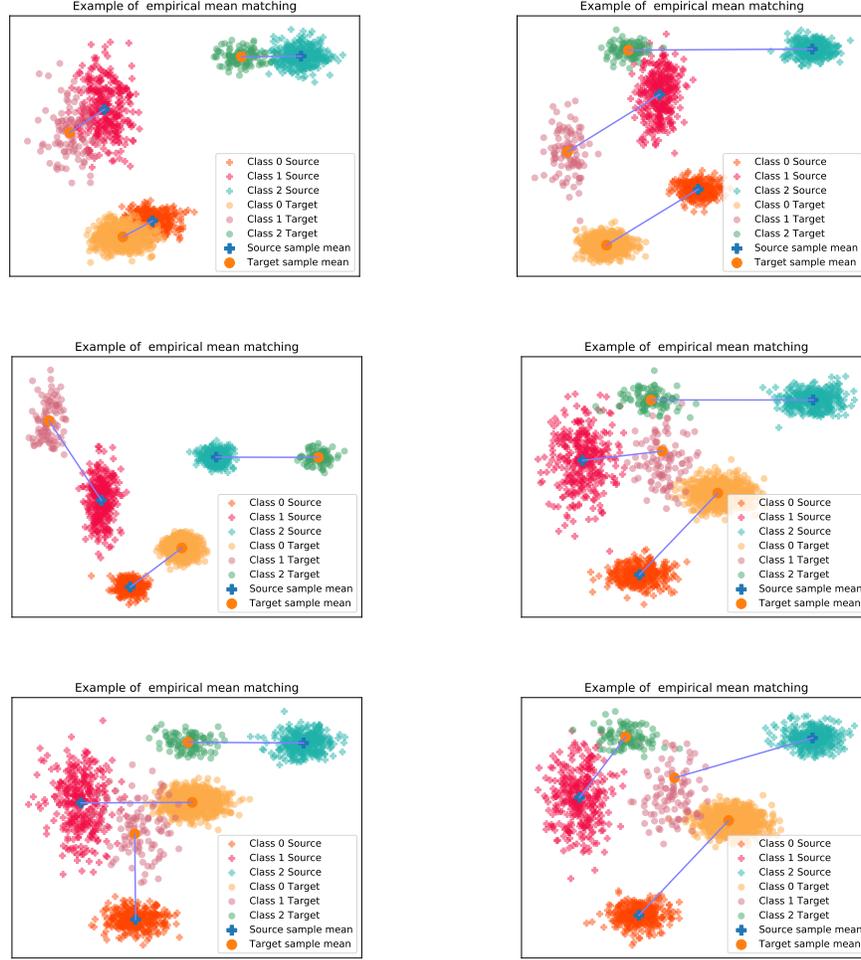


Figure 4: Example of geometrical arrangements of the source and target class-conditional distributions that allows correct and incorrect matching of classes by optimal transport of empirical means (assuming correct estimation of these means). Blue lines denote the matching. (top-left) In this setting, the displacements of each class-conditionals is so that for each class i $\|\mathbf{m}_S^i - \mathbf{m}_T^i\|_2 \leq \|\mathbf{m}_S^i - \mathbf{m}_T^j\|_2$, for all j . We are thus in the first example that we gave as satisfying Proposition 1. (top-right) Class-conditionals have been displaced such that the “nearness” hypothesis is not respected anymore. However, they have been mapped through an operator such that optimal transport allows their matchings (based on their means). (middle) We have illustrated two other examples of distribution arrangements that allow class matching. (right) Two examples that break our assumption. In both cases, one target class-conditional is “near” another source class, without the global displacements of all target class-conditionals being uniform in direction.

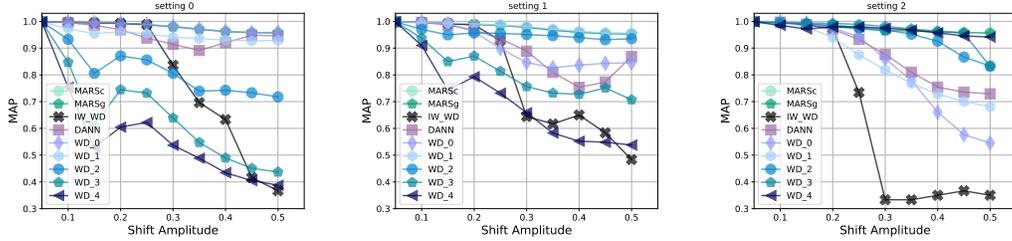


Figure 5: Performance of the compared algorithms in different label shift setting and for increasing shift between means of class-conditionals. In source domain, label distributions are uniform and shift occurs due to change only in the target domain. (left) $p_T(y = 1) = 0.33$, $p_T(y = 2) = 0.33$, $p_T(y = 3) = 0.34$. (middle) $p_T(y = 1) = 0.5$, $p_T(y = 2) = 0.2$, $p_T(y = 3) = 0.2$, (right) $p_T(y = 1) = 0.8$, $p_T(y = 2) = 0.1$, $p_T(y = 3) = 0.1$. For balanced problems, we note that best methods are $WD_{\beta=\{0,1\}}$, DANN and our approaches either using GMM or clustering for estimating label proportion. As expected, a too heavy reweighting yields to poor performance for $WD_{\beta=\{2,3,4\}}$. Then for a mild imbalance, $WD_{\beta=\{1,2\}}$ performs better than the other competitors while for higher imbalance, $WD_{\beta=\{3,4\}}$ works better. For all settings, our methods are competitive as they are adaptive to the imbalance through the estimation fo $p_T(y)$. The IW-WD of Combes et al. (2020) performs well until the distance between class-conditionals is too large. This is justified by theory as their estimator of the ratio $p_T(y)/p_S(y)$ is tailored for situations where class-conditionals are equal.

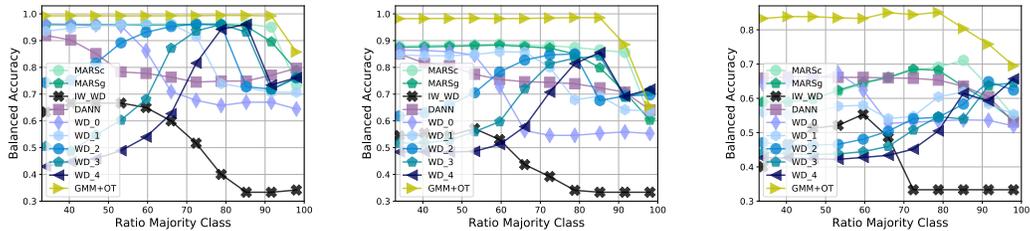


Figure 6: Performance of the compared algorithms, including **GMM+OT** for three different covariance matrices of the Gaussians composing the toy dataset with respect to the imbalance. The shift between the class-conditionals has been fixed and yields to samples similar to those presented in Figure 3. Our method is referred as **MARS**. The x-axis is given with respect to the ratio of majority class which is the class 1. (left) Low-error setting. (middle) mid-error setting. (right) high-error setting. material. We note that this toy problem can be easily solved using a GMM and a optimal transport-based label assignment. We can also remark that again as soon as the class-conditionals do not match anymore, the IW-WD of Combes et al. (2020) fails due to its inability to estimate correctly the importance weight w .

Table 3: Table of the dataset experimental settings. We have considered different domain adaptation problems and different configurations of the label shift in the source and target domain. For the digits and VisDA problem, we provide the ratio of samples of classes for each problem (*e.g.*, for the third setting of VisDA-3 problem, the second class accounts for the 70% of the samples in target domain). For Office datasets, because of large amount of classes, we have changed percent of samples of that class in the source or target (*e.g.*, the 10-class in Office Home uses respectively 20% and 100% of its sample for the source and target domain).

Configuration	Proportion Source	Proportion Target
MNIST-USPS balanced	$\{\frac{1}{10}, \dots, \frac{1}{10}\}$	$\{\frac{1}{10}, \dots, \frac{1}{10}\}$
MNIST-USPS mid	$\{\frac{1}{10}, \dots, \frac{1}{10}\}$	$\{0, \dots, 3, 6\} = 0.02, \{4, 5\} = 0.02, \{7, 8, 9\} = 0.1$
MNIST-USPS high	$\{\frac{1}{10}, \dots, \frac{1}{10}\}$	$\{0\} = 0.3665, \{1\} = 0.3651, \{2, \dots\} = 0.0335$
USPS-MNIST balanced	$\{\frac{1}{10}, \dots, \frac{1}{10}\}$	$\{\frac{1}{10}, \dots, \frac{1}{10}\}$
USPS-MNIST mid	$\{\frac{1}{10}, \dots, \frac{1}{10}\}$	$\{0, \dots, 3, 6\} = 0.02, \{4, 5\} = 0.02, \{7, 8, 9\} = 0.1$
USPS-MNIST high	$\{\frac{1}{10}, \dots, \frac{1}{10}\}$	$\{0\} = 0.3665, \{1\} = 0.3651, \{2, \dots\} = 0.0335$
MNIST-MNISTM (1)	$\{0-4\} = 0.05, \{5-9\} = 0.15$	$\{0, \dots, 3, 6\} = 0.02, \{4, 5\} = 0.02, \{7, 8, 9\} = 0.1$
MNIST-MNISTM (2)	$\{0-2\} = 0.26, \{3-9\} = 0.03$	$\{0-6\} = 0.03, \{7-9\} = 0.26$
MNIST-MNISTM (3)	$\{0-5\} = 0.05, \{6-9\} = 0.175$	$\{0-3\} = 0.175, \{4-9\} = 0.05$
VisDA-3 (1)	$\{0.33, 0.33, 0.34\}$	$\{0.33, 0.33, 0.34\}$
VisDA-3 (2)	$\{0.4, 0.2, 0.4\}$	$\{0.2, 0.6, 0.2\}$
VisDA-3 (3)	$\{0.4, 0.2, 0.4\}$	$\{0.15, 0.7, 0.15\}$
VisDA-3 (4)	$\{0.4, 0.2, 0.4\}$	$\{0.1, 0.8, 0.1\}$
VisDA-3 (5)	$\{0.6, 0.2, 0.2\}$	$\{0.2, 0.2, 0.6\}$
VisDA-3 (6)	$\{0.6, 0.2, 0.2\}$	$\{0.15, 0.2, 0.65\}$
VisDA-3 (7)	$\{0.6, 0.2, 0.2\}$	$\{0.2, 0.65, 0.15\}$
VisDA-12 (1)	$\{\frac{1}{12}, \dots, \frac{1}{12}\}$	$\{\frac{1}{12}, \dots, \frac{1}{12}\}$
VisDA-12 (2)	$\{\frac{1}{12}, \dots, \frac{1}{12}\}$	$\{0-3\} = 0.15, \{4-11\} = 0.05$
VisDA-12 (3)	$\{\frac{1}{12}, \dots, \frac{1}{12}\}$	$\{0-1\} = 0.2, \{2-5\} = 0.1, \{6-11\} = 0.03$
Office-31	$\{0-15\} : 30\% \{15-31\} : 80\%$	$\{0-15\} : 80\% \{15-31\} : 30\%$
Office-Home	$\{0-32\} : 20\% \{33-65\} : 100\%$	$\{0-32\} : 100\% \{33-65\} : 20\%$

Table 4: Table of averaged balanced accuracy for the compared models and different domain adaptation models. Number of runs used 20. Reported in bold are the best performances as well as other methods which achieves performance that are statistically similar according to a Wilcoxon signrank test with $p = 0.01$.

Setting	Source	DANN	WD $_{\beta=0}$	WD $_{\beta=1}$	WD $_{\beta=2}$	WD $_{\beta=3}$	WD $_{\beta=4}$	IW-WD	MARSG	MARSc
MNIST-USPS 10 modes										
Balanced	76.89±3.7	79.74±3.5	93.71±0.7	74.27±4.3	51.33±4.0	76.61±3.3	71.90±5.7	95.28±0.4	95.61±0.7	95.64±1.0
Mid	80.41±3.1	78.65±3.0	94.30±0.7	75.36±3.4	55.55±4.3	78.98±3.1	72.32±4.2	95.60±0.5	89.70±2.3	90.39±2.6
High	78.13±4.9	81.79±4.0	93.86±1.1	87.44±1.7	83.83±5.2	85.65±2.5	83.65±3.0	94.08±1.0	88.30±1.5	89.65±2.3
USPS-MNIST 10 modes										
Balanced	77.04±2.6	80.49±2.2	73.35±2.8	66.70±2.9	49.86±2.8	55.83±2.9	52.12±3.5	80.52±2.2	84.59±1.7	85.50±2.1
Mid	79.54±2.8	78.88±1.8	75.85±1.6	63.33±2.3	53.22±2.8	47.20±2.4	48.29±2.9	78.36±3.5	79.73±3.6	78.49±2.5
High	78.48±2.4	77.79±2.0	76.14±2.7	63.00±3.3	57.56±4.8	51.19±4.4	49.31±3.3	71.53±4.7	75.62±1.8	77.14±2.4
MNIST-MNISTM 10 modes										
Setting 1	58.34±1.3	61.22±1.1	57.44±1.7	50.20±4.4	47.01±2.0	57.85±1.1	55.95±1.3	63.10±3.1	58.08±2.3	56.58±4.6
Setting 2	59.94±1.1	61.09±1.0	58.08±1.4	53.39±3.5	48.61±2.4	59.74±0.7	58.14±0.8	65.03±3.5	57.69±2.3	55.64±2.1
Setting 3	58.14±1.2	60.39±1.4	57.68±1.2	47.72±4.9	42.15±7.3	57.09±1.0	53.52±1.1	52.46±14.8	53.68±7.2	53.72±3.3
VisdDA 3 modes										
setting 1	79.28±4.3	78.83±9.1	91.83±0.7	73.78±2.0	61.65±2.2	65.62±2.7	58.58±2.6	94.11±0.6	92.47±1.2	92.13±1.8
setting 4	80.15±5.3	75.46±9.3	72.75±1.2	86.86±7.5	86.82±1.2	80.16±6.9	75.71±2.0	85.88±5.7	87.69±3.0	91.29±4.8
setting 2	81.47±3.5	68.46±14.7	68.81±1.3	84.45±1.2	93.15±0.4	73.65±14.2	60.67±0.9	78.73±10.8	84.04±4.3	91.80±3.4
setting 3	78.35±3.2	58.93±15.9	64.13±1.9	79.17±0.8	77.12±10.3	89.93±0.5	94.38±0.3	77.96±9.3	75.68±4.1	73.81±13.2
setting 5	83.52±3.5	80.83±14.5	63.82±0.6	73.70±7.3	50.91±1.1	76.52±6.7	59.28±1.0	90.40±3.6	89.01±0.9	89.03±3.5
setting 6	80.84±4.2	54.76±19.8	45.27±2.4	63.70±5.1	67.05±6.1	42.86±10.8	62.21±1.4	94.36±1.0	93.70±0.4	93.86±1.0
setting 7	79.22±3.7	42.94±2.5	57.51±1.5	55.39±2.0	50.22±4.3	43.66±8.3	62.47±0.8	88.52±4.9	78.56±3.2	82.33±7.5
VisdDA 12 modes										
setting 1	41.90±1.5	52.79±2.1	45.81±4.3	44.23±3.0	35.45±4.6	40.96±3.0	37.59±3.4	50.35±2.3	53.31±0.9	55.05±1.6
setting 2	41.75±1.5	50.82±1.6	45.72±8.9	40.49±4.8	36.21±5.0	36.12±4.6	31.86±5.7	48.59±1.8	53.09±1.6	55.33±1.6
setting 3	40.64±4.3	49.17±1.3	47.12±1.6	42.10±3.0	36.32±4.4	37.26±3.5	34.96±5.4	46.59±1.3	50.78±1.6	52.08±1.2
Office 31										
A - D	73.73±1.4	74.26±1.8	77.22±0.7	65.10±2.0	62.65±2.6	71.47±1.2	63.89±1.1	75.74±1.6	76.07±0.9	78.20±1.3
D - W	83.64±1.1	81.89±1.5	82.61±0.6	83.53±0.8	82.80±0.7	80.10±0.5	87.09±0.9	78.93±1.5	86.32±0.6	86.20±0.8
W - A	54.05±0.9	52.16±1.0	48.94±0.4	56.81±0.4	53.02±0.5	58.83±0.4	54.93±0.5	52.23±0.7	60.68±0.8	55.18±0.8
W - D	92.76±0.9	87.64±1.4	95.07±0.3	93.13±0.5	87.60±0.9	94.69±0.6	91.18±0.6	97.04±0.9	95.14±0.8	93.80±0.6
D - A	52.51±0.9	48.06±1.2	49.78±0.4	48.75±0.5	50.13±0.4	50.28±0.7	50.75±0.5	41.39±1.8	54.65±0.9	54.95±0.9
A - W	67.45±1.5	70.15±1.0	67.07±0.6	60.62±2.1	52.92±1.4	63.98±1.3	59.73±0.8	68.76±1.6	73.09±1.5	71.90±1.2
Office Home										
Art - Clip	37.66±0.7	36.85±0.6	33.42±1.2	31.43±1.6	27.13±1.6	31.63±5.2	29.30±6.6	37.65±0.6	37.58±0.5	38.65±0.5
Art - Product	49.72±0.9	49.98±0.9	39.43±3.6	38.82±2.3	35.05±2.3	35.09±3.4	32.85±3.6	48.98±0.3	55.27±0.7	52.18±0.4
Art - Real	58.22±1.0	53.68±0.5	51.09±2.3	50.35±1.8	46.40±2.4	51.52±4.5	45.34±11.0	57.74±0.7	63.88±0.5	58.75±0.7
Clip - Art	35.29±1.4	35.70±1.5	28.92±2.9	23.13±2.0	18.37±1.5	21.95±3.1	20.44±2.3	28.74±1.2	41.15±0.6	40.73±0.8
Clip - Product	51.94±1.3	52.06±0.8	39.17±7.9	39.26±2.6	34.73±1.9	39.58±2.8	39.46±2.9	34.46±2.1	51.69±0.5	52.12±0.5
Clip - Real	50.65±1.2	51.42±1.0	43.24±2.2	40.06±2.1	32.71±1.4	39.22±2.4	35.78±2.8	35.72±1.1	53.97±0.3	56.63±0.5
Product - Art	39.59±1.6	39.47±1.5	39.17±1.0	36.11±1.0	38.77±1.1	39.50±0.6	38.24±0.6	33.95±1.4	37.77±1.1	39.31±1.3
Product - Clip	32.71±0.9	37.18±1.0	33.82±0.5	28.38±0.7	28.40±0.6	29.72±0.5	31.76±0.8	24.89±1.0	30.86±0.8	29.25±0.9
Product - Real	62.12±1.3	62.52±1.2	62.56±0.7	58.09±0.5	57.58±0.6	59.33±0.6	57.11±0.8	59.22±0.9	60.48±0.6	62.20±0.7
Real - Product	68.30±1.0	70.39±0.8	70.19±0.5	61.72±0.8	63.40±0.9	61.51±1.0	65.45±0.6	64.47±1.5	64.79±3.6	66.49±1.1
Real - Art	40.25±0.9	41.31±1.0	39.16±0.7	33.46±1.3	31.61±1.5	36.90±0.9	36.14±0.9	36.93±1.9	39.90±1.4	39.17±1.6
Real - Clip	42.74±1.1	40.86±1.0	40.42±0.5	35.59±0.8	34.90±0.9	40.42±0.5	35.64±0.8	35.60±2.0	38.69±2.1	38.82±2.5
#Wins (34)	7	9	5	0	1	0	2	9	12	21
Aver. Rank	4.16	4.73	5.32	6.97	8.38	6.59	7.57	4.95	3.38	2.95

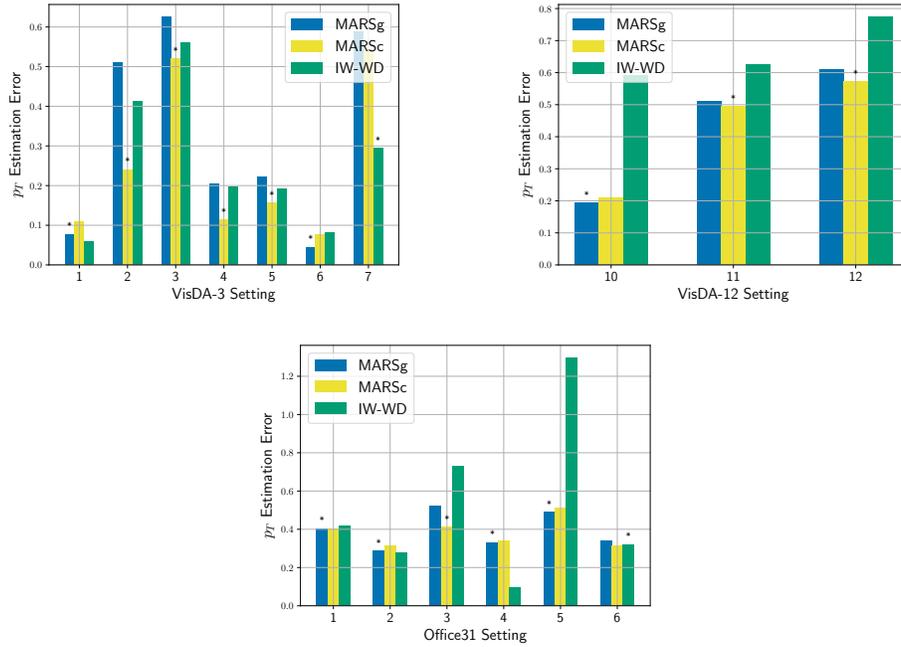


Figure 7: Examples of ℓ_1 norm error of estimated label proportion. We have reported the performance of our two methods (MARSg and MARSc) as well as the performance of IW-WD. The three panels are related to the (left) VisDA-3. (middle) VisDA-12. (right) Office 31 and the different experimental imbalance settings (see Table 3). We have also reported, with a '*' on top, among the three approaches, the best performing one in term of balanced accuracy. We note that for VisDA problems, our approaches provide better estimation than IW-WD 8 out of 10 experiments and 3 out of 6 on Office-31. We also remark the correlation between better p_T estimation and better accuracy.

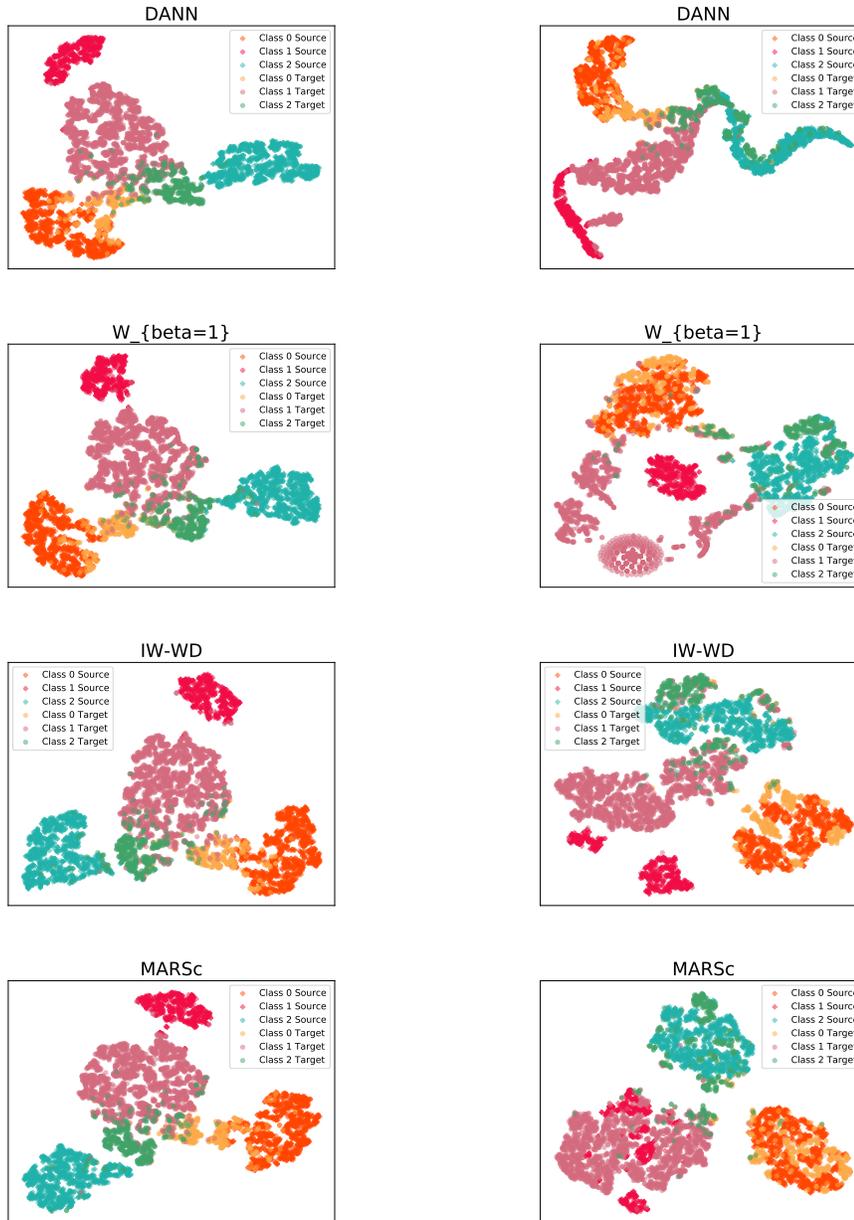


Figure 8: *t-sne* embeddings of the target sample for the VisDA-3 problem and imbalance setting 2 ($\mathbf{p}_S = [0.4, 0.2, 0.4]$ and $\mathbf{p}_T = [0.2, 0.6, 0.2]$). The columns depict the embeddings obtained (left) after training on the source data without adaptation for about 10 iterations, which is sufficient for 0 training error. (right) after adaptation by minimizing the appropriate discrepancy loss between marginal distributions. From top to bottom, we have : (first-row) DANN, (second-row) $WD_{\beta=1}$, (third-row), IW-WD (last row) MARSc. From the right column, we note how DANN and $WD_{\beta=1}$ struggles in aligning the class conditionals, especially those of Class 1, which is the class that varies the most in term of label proportion. IW-WD manages to aligns the classes “0” and “2” but is not able to correctly match the class “1”. Instead, our MARSc approach achieves high performance and correctly aligns the class conditionals, although some few examples seem to be mis-classified. Importantly, we can remark from the left column that for this example, before alignment, the embeddings seem to satisfy our Proposition 1 hypothesis. At the contrary, the assumption needed for correctly estimating \mathbf{p}_T for IW-WD is not satisfied, justifying thus the good and poor performance of those models.