

# VARIANCE MATTERS: IMPROVING DOMAIN ADAPTATION VIA STRATIFIED SAMPLING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Domain shift remains a key challenge in deploying machine learning models to the real world. Unsupervised domain adaptation (UDA) aims to address this by minimising domain discrepancy during training, but the discrepancy estimates suffer from high variance in stochastic settings, which can stifle the theoretical benefits of the method. This paper proposes Variance-Reduced Domain Adaptation via Stratified Sampling (VaRDASS), the first specialised stochastic variance reduction technique for UDA. We consider two specific discrepancy measures – correlation alignment and the maximum mean discrepancy (MMD) – and derive ad hoc stratification objectives for these terms. We then present expected and worst-case error bounds, and prove that our proposed objective for the MMD is theoretically optimal (i.e., minimises the variance) under certain assumptions. Finally, a practical k-means style optimisation algorithm is introduced and analysed. Experiments on three domain shift datasets demonstrate improved discrepancy estimation accuracy and target domain performance.

## 1 INTRODUCTION

Machine learning models often underperform when the test data distribution differs from the training distribution, a phenomenon known as domain shift. Improving robustness to domain shift has been a longstanding goal in machine learning, and is crucial to the widespread deployment of AI (Gulrajani & Lopez-Paz, 2021; Koh et al., 2021).

Unsupervised domain adaptation (UDA) is a common approach, where models learn representations invariant across source and target domains by minimising a “domain discrepancy” loss. Two widely used objectives are the correlation alignment (CORAL) loss, which matches covariance matrices (Sun & Saenko, 2016), and the maximum mean discrepancy (MMD), which aligns kernel mean embeddings (Tzeng et al., 2014; Long et al., 2015; Li et al., 2018). Although theoretically well-grounded (Ben-David et al., 2006; 2010; Redko et al., 2022), a key limitation is that empirically estimating these discrepancies is extremely noisy, especially in high dimensions and with small minibatches. This can destabilise training, and sometimes even lead to worse target domain performance than with no domain alignment at all (Dubey et al., 2021; Gao et al., 2023; Gulrajani & Lopez-Paz, 2021; Koh et al., 2021; Napoli & White, 2023; 2024; Wang et al., 2019).

Stochastic variance reduction offers a natural remedy. However, many classical such methods, including SAGA (Defazio et al., 2014), SVRG (Johnson & Zhang, 2013), and dual coordinate ascent (Shalev-Shwartz & Zhang, 2013), are incompatible with MMD and CORAL, since these losses lack finite-sum structure. Momentum-based approaches (Polyak, 1964; Polyak & Juditsky, 1992; Kingma & Ba, 2014) avoid this problem, but at the cost of biasing the gradients, since the training examples are not weighted equally (Gower et al., 2020). Estimator shrinkage (Ledoit & Wolf, 2020) may also be applied, but again this introduces bias.

A more flexible direction is to alter the sampling distribution, then correct the sampling bias using weighted losses. Importance sampling biases selection towards points with higher gradient impact (Alain et al., 2015; Johnson & Guestrin, 2018; Katharopoulos & Fleuret, 2017; 2018; Kutsuna, 2025; Loshchilov & Hutter, 2016; Zhao & Zhang, 2015). Similarly, typicality sampling upweights examples in higher-density regions (Peng et al., 2020). Diverse sampling methods reduce variance by enforcing dissimilarity within minibatches, for example using antithetic pairs (Liu & Xu, 2018),

repulsive point processes (Zhang et al., 2017; 2019; Bardenet et al., 2021; Napoli & White, 2024), or faster heuristics such as k-means++ (Napoli & White, 2024).

Finally, stratified sampling partitions the data into strata and samples independently from each. This can be considered a special case of diverse sampling (Zhang et al., 2017), but is more computationally efficient since the only significant cost is in stratifying the data, rather than drawing the samples themselves. Zhao & Zhang (2014) cluster the input space directly, which offers an iteration-independent stratification, but this assumes the gradients are Lipschitz continuous with respect to the inputs. Alternatively, Liu et al. (2020b) reconstruct the clusters periodically during training, and propose a criterion to determine when to do so. Fu & Zhang (2017) and Lu et al. (2021) extend stratified sampling to Bayesian and time-series learning respectively.

This paper proposes Variance-Reduced Domain Adaptation via Stratified Sampling (VaRDASS), the first specialised variance reduction strategy for UDA. By using stratified data samplers, VaRDASS obtains variance-reduced stochastic losses without having to trade off bias or make any other compromises with the model or learning algorithm. Moreover, we derive specific stratification objectives for two widely-used UDA losses – MMD and CORAL, present expected and worst-case error bounds, and prove that our proposed objective for the MMD is theoretically optimal (i.e., minimises the variance) under certain assumptions. We then propose a practical k-means style alternating optimisation algorithm to solve this problem in tractable time.

In experiments and simulations, we validate our choices of clustering objectives, whilst showing that VaRDASS both reduces the estimation error of the domain discrepancy for a given minibatch size, and improves target domain accuracy on 3 UDA benchmark datasets.

## 2 METHOD

### 2.1 PRELIMINARIES

In UDA, we are given labelled examples  $\mathcal{D}_s = \{(x_{s,i}, y_{s,i})\}_{i=1}^{n_s}$  from a source distribution  $P_s$ , and unlabelled examples  $\mathcal{D}_t = \{x_{t,j}\}_{j=1}^{n_t}$  from a different target distribution  $P_t$ . The goal is to produce a model  $\Theta : \mathcal{X} \rightarrow \mathcal{Y}$  such that the target risk  $\mathbb{E}_{(x_t, y_t) \sim P_t} [L_{\text{task}}(\Theta(x_t), y_t)]$  for some task loss  $L_{\text{task}}$  is minimised. We assume  $\Theta$  decomposes into a featuriser  $\Theta_F : \mathcal{X} \rightarrow \mathcal{Z}$  and prediction head  $\Theta_C : \mathcal{Z} \rightarrow \mathcal{Y}$ , and has intermediate feature representations  $z_s, z_t$ . Since the target data are unlabelled, UDA methods instead minimise  $L_{\text{task}}$  on the source domain, alongside a domain adaptation term  $L_{\text{DA}}$  which aligns the source and target feature distributions:

$$\min_{\Theta} \mathbb{E}_{(x_s, y_s) \sim P_s, x_t \sim P_t} [L_{\text{task}}(\Theta(x_s), y_s) + \lambda L_{\text{DA}}(z_s, z_t)]. \quad (1)$$

During training, minibatches  $\tilde{\mathcal{D}}_s = \{(\tilde{x}_{s,h}, \tilde{y}_{s,h})\}_{h=1}^k \subset \mathcal{D}_s$  and  $\tilde{\mathcal{D}}_t = \{\tilde{x}_{t,h}\}_{h=1}^k \subset \mathcal{D}_t$  are randomly sampled at each iteration (we assume for simplicity that  $|\tilde{\mathcal{D}}_s| = |\tilde{\mathcal{D}}_t| = k$ ), and used to compute stochastic losses  $\hat{L}_{\text{task}}(\tilde{x}_s, \tilde{y}_s)$  and  $\hat{L}_{\text{DA}}(\tilde{z}_s, \tilde{z}_t)$ .

Our aim is to reduce  $\text{Var}(\hat{L}_{\text{DA}})$  by constructing the minibatches using stratified sampling. In the following sections, we will formulate and justify stratification objectives to this end for the MMD and CORAL losses, and propose a practical clustering algorithm to solve this problem. Since  $\tilde{\mathcal{D}}_s$  and  $\tilde{\mathcal{D}}_t$  are sampled independently, they can also be stratified independently. The following sections describe the process for  $\mathcal{D}_s$ ; the same procedure can be analogously repeated for  $\mathcal{D}_t$ .

### 2.2 OPTIMAL STRATIFICATION FOR MMD

The (squared) MMD loss is

$$L_{\text{MMD}} = \|\mu_{\phi, s} - \mu_{\phi, t}\|_{\mathcal{H}}^2 \quad (2)$$

where  $\mu_{\phi, s} = \mathbb{E}[\phi(z_s)]$ ,  $\mu_{\phi, t} = \mathbb{E}[\phi(z_t)]$ ,  $\mathcal{H}$  is a reproducing kernel Hilbert space, and  $\phi : \mathcal{Z} \rightarrow \mathcal{H}$  is an implicit mapping.  $\mathcal{H}$  is associated with a unique positive-definite kernel  $\kappa : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  for which the reproducing property  $\kappa(z, z') = \langle \phi(z), \phi(z') \rangle_{\mathcal{H}}$  is satisfied (Gretton et al., 2012).

With minibatch size  $k$ , the maximum variance reduction is achieved by partitioning  $\mathcal{D}_s$  into  $k$  strata  $\mathcal{S} = \{S_h\}_{h=1}^k$  such that  $\bigcup S_h = \mathcal{D}_s$ ,  $\bigcap S_h = \emptyset$ , and drawing a single instance  $(\tilde{x}_{s,h}, \tilde{y}_{s,h}) \sim$

Uniform ( $S_h$ ) from each (Giles, 2015). Thus, a simple estimate of (2) is

$$\widehat{L}_{\text{MMD}} = \|\widehat{\mu}_{\phi,s} - \widehat{\mu}_{\phi,t}\|_{\mathcal{H}}^2, \quad (3)$$

where  $\widehat{\mu}_{\phi,s} = \frac{1}{n_s} \sum_{h=1}^k |S_h| \phi(\tilde{z}_{s,h})$  and  $\tilde{z}_{s,h} = \Theta_F(\tilde{x}_{s,h})$ , and similarly for  $\widehat{\mu}_{\phi,t}$ .

First, we show that to reduce  $\text{Var}(\widehat{L}_{\text{MMD}})$ , a good surrogate objective is to reduce  $\text{Var}(\widehat{\mu}_{\phi,s}) = \mathbb{E}[\|\widehat{\mu}_{\phi,s} - \mu_{\phi,s}\|_{\mathcal{H}}^2]$  directly. Variance expressions for estimates of the squared MMD have been derived in various prior work (Bounliphone et al., 2016; Liu et al., 2020a; Sutherland & Deka, 2022). However, to keep our analysis straightforward, we use a simpler expression that assumes the empirical mean embeddings are normally distributed. To proceed, let  $\Sigma_{\widehat{\mu}_{\phi,s}}$  and  $\Sigma_{\widehat{\mu}_{\phi,t}}$  denote the covariance matrices of the respective empirical mean embeddings.

**Lemma 1.** *Assume  $\widehat{\mu}_{\phi,s} \sim \mathcal{N}(\mu_{\phi,s}, \Sigma_{\widehat{\mu}_{\phi,s}})$  and  $\widehat{\mu}_{\phi,t} \sim \mathcal{N}(\mu_{\phi,t}, \Sigma_{\widehat{\mu}_{\phi,t}})$  are independent random vectors, and that the relevant conditions for normality under the central limit theorem are satisfied. Then,  $\widehat{\mu}_{\phi,s} - \widehat{\mu}_{\phi,t} \sim \mathcal{N}(m, \Sigma_{\widehat{\mu}_{\phi,s}} + \Sigma_{\widehat{\mu}_{\phi,t}})$ , where  $m = \mu_{\phi,s} - \mu_{\phi,t}$ , and  $\widehat{L}_{\text{MMD}}$  follows a generalised chi-squared distribution with variance (Seber, 2007)*

$$\text{Var}(\widehat{L}_{\text{MMD}}) = 2 \text{Tr}((\Sigma_{\widehat{\mu}_{\phi,s}} + \Sigma_{\widehat{\mu}_{\phi,t}})^2) + 4m^T (\Sigma_{\widehat{\mu}_{\phi,s}} + \Sigma_{\widehat{\mu}_{\phi,t}}) m. \quad (4)$$

The only quantity dependent on the stratification of  $\mathcal{D}_s$  is  $\Sigma_{\widehat{\mu}_{\phi,s}}$ . Therefore, we consider the variance as a function of  $\Sigma_{\widehat{\mu}_{\phi,s}}$  with the remaining variables constant.  $\text{Var}(\widehat{L}_{\text{MMD}})$  and  $\text{Var}(\widehat{\mu}_{\phi,s})$  will have the same minimisers if they are related by a monotonic function, which occurs when  $\Sigma_{\widehat{\mu}_{\phi,s}}$  is isotropic.

**Theorem 1.** *Assume  $\mathcal{H}$  has finite dimensionality  $d$  and  $\Sigma_{\widehat{\mu}_{\phi,s}} = \sigma^2 I$  for some positive scalar  $\sigma^2$ . Then,*

$$\arg \min_{\mathcal{S}} \text{Var}(\widehat{L}_{\text{MMD}}) = \arg \min_{\mathcal{S}} \text{Var}(\widehat{\mu}_{\phi,s}). \quad (5)$$

*Proof.* See Appendix A. □

For the anisotropic case, the relation is not monotonic and so the special case in Theorem 1 no longer holds. However, by characterising the degree of monotonicity using Spearman’s rank correlation  $\rho$ , we can derive expressions for the expected and worst-case errors of the surrogate solution, in terms of  $\rho$ ,  $n_s$ ,  $k$ , and the growth rate of  $\text{Var}(\widehat{L}_{\text{MMD}})$  with respect to the solution rankings.

To proceed, let  $f(\mathcal{S}) = \text{Var}(\widehat{L}_{\text{MMD}})$  and  $g(\mathcal{S}) = \text{Var}(\widehat{\mu}_{\phi,s})$  be the target and surrogate partitioning objectives as functions of the dataset partitioning. Define the order statistic  $\mathcal{S}_{(i)} = \arg \{f(\mathcal{S})_{(i)}\}$  for all possible partitionings (e.g.,  $\mathcal{S}_{(1)} = \arg \min_{\mathcal{S}} f(\mathcal{S})$ ) and  $p = \text{Stirling}(n_s, k)$  the number of such partitionings, given by the Stirling partition number. Given a Spearman coefficient  $\rho$  between  $f$  and  $g$ , define also a rank displacement  $\delta_i = \text{rank}(g(\mathcal{S})_{(i)}) - \text{rank}(f(\mathcal{S})_{(i)}) = \text{rank}(g(\mathcal{S})_{(i)}) - i$ , such that, from the formula for Spearman correlation,  $\sum \delta_i^2 = \frac{p(p^2-1)}{6}(1-\rho)$ . Finally, let  $\mathcal{S}_{(l)} = \arg \min_{\mathcal{S}} g(\mathcal{S})$  be the minimiser of  $g(\mathcal{S})$  which has rank  $l$  on  $f(\mathcal{S})$ .

**Theorem 2 (Worst-case error bound).** *Assume the growth rate of the sorted  $f(\mathcal{S})_{(i)}$  is bounded by some constant  $K$ , such that, for any indices  $i, j$ ,  $f(\mathcal{S})_{(i)} - f(\mathcal{S})_{(j)} \leq K(i-j)$ . Then, the error incurred by minimising  $g$  instead of  $f$  satisfies*

$$f(\mathcal{S})_{(l)} - f(\mathcal{S})_{(1)} \leq K \sqrt{\frac{p(p^2-1)}{6}}(1-\rho). \quad (6)$$

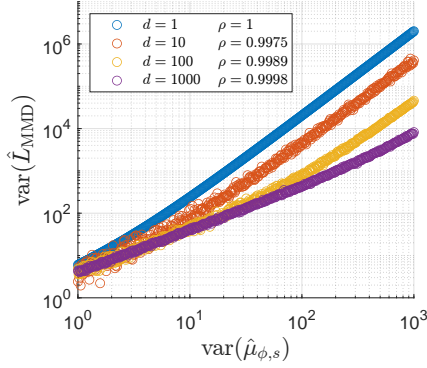


Figure 1: The relationship between the target ( $\text{Var}(\widehat{L}_{\text{MMD}})$ ) and surrogate ( $\text{Var}(\widehat{\mu}_{\phi,s})$ ) partitioning objectives, for a range of feature dimensionalities  $d$ .

*Proof.* As  $\delta_1 = l - 1$  and  $\delta_1^2 \leq \sum \delta_i^2$ ,  $l$  satisfies  $(l - 1)^2 \leq \sum \delta_i^2 = \frac{p(p^2-1)}{6}(1 - \rho)$ .  $\square$

**Theorem 3** (Expected error). *If the rank displacements are assumed to be homogeneous, then the error is*

$$f(\mathcal{S})_{(l)} - f(\mathcal{S})_{(1)} = K \sqrt{\frac{p^2-1}{6}(1-\rho)}. \quad (7)$$

*Proof.* With homogeneous displacements,  $\delta_1^2 = \sum \delta_i^2/p$ , and so  $(l - 1)^2 = \frac{p^2-1}{6}(1 - \rho)$ .  $\square$

Figure 1 shows the relationship between  $\text{Var}(\widehat{L}_{\text{MMD}})$  and  $\text{Var}(\widehat{\mu}_{\phi,s})$  for randomly generated covariances  $\Sigma_{\widehat{\mu}_{\phi,s}}$ , with  $m = 1$ ,  $\Sigma_{\widehat{\mu}_{\phi,t}} = 0$ , and 4 different values of  $d$ . Note that for  $d = 1$ , Theorem 1 holds and so  $\rho = 1$  exactly, meaning the errors defined in Theorems 2 and 3 are nil. In the remaining cases, the correlations are all  $> 0.99$ , validating our choice of surrogate clustering objective.

To solve  $\arg \min_{\mathcal{S}} \text{Var}(\widehat{\mu}_{\phi,s})$ , observe that

$$\begin{aligned} \text{Var}(\widehat{\mu}_{\phi,s}) &= \text{Var}\left(\frac{1}{n_s} \sum_{h=1}^k |S_h| \phi(\tilde{z}_{s,h})\right) = \frac{1}{n_s^2} \sum_{i=1}^k |S_h|^2 \text{Var}(\phi(\tilde{z}_{s,h})) \\ &= \frac{1}{n_s^2} \sum_{h=1}^k |S_h| \sum_{z \in S_h} \|\phi(z) - \mu_{\phi,h}\|_{\mathcal{H}}^2, \end{aligned} \quad (8)$$

where  $\mu_{\phi,h} = \frac{1}{|S_h|} \sum_{z \in S_h} \phi(z)$  is the stratum mean. Thus, we have the optimisation

$$\arg \min_{\mathcal{S}} \sum_{h=1}^k |S_h| \sum_{z \in S_h} \|\phi(z) - \mu_{\phi,h}\|_{\mathcal{H}}^2, \quad (9)$$

which can be seen as a kernel k-means-style clustering problem with dynamically weighted clusters. Compared to ordinary (kernel) k-means clustering, this objective imposes a greater penalty on larger clusters, and thus encourages (but does not strictly enforce) the clusters to be of balanced sizes.

### 2.3 OPTIMAL STRATIFICATION FOR CORAL

The CORAL loss is defined as

$$L_{\text{CORAL}} = \|R_s - R_t\|_F^2, \quad (10)$$

where  $R_s$  and  $R_t$  are the covariance matrices of  $z_s$  and  $z_t$  respectively. An estimate using stratified sampling is

$$\widehat{L}_{\text{CORAL}} = \left\| \widehat{R}_s - \widehat{R}_t \right\|_F^2, \quad (11)$$

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

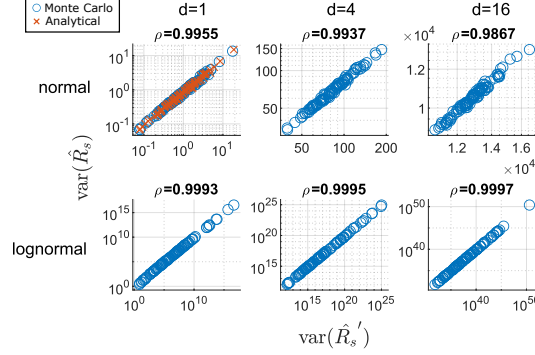


Figure 2: The relationship between the target ( $\text{Var}(\widehat{R}_s)$ ) and surrogate ( $\text{Var}(\widehat{R}'_s)$ ) partitioning objectives, for a range of feature dimensionalities  $d$ , and two distributions.

where  $\widehat{R}_s = \frac{1}{n_s-1} \sum_{h=1}^k |S_h| (\tilde{z}_{s,h} - \widehat{\mu}_s) (\tilde{z}_{s,h} - \widehat{\mu}_s)^T$ ,  $\widehat{\mu}_s = \frac{1}{n_s} \sum_{h=1}^k |S_h| \tilde{z}_{s,h}$ , and similarly for  $\widehat{R}_t$  and  $\widehat{\mu}_t$ .

As previously,  $\widehat{R}_s$  and  $\widehat{R}_t$  can be assumed to be Gaussian and so Theorem 1 implies that to minimise  $\text{Var}(\widehat{L}_{\text{CORAL}})$ , a good surrogate is to minimise  $\text{Var}(\widehat{R}_s)$  and  $\text{Var}(\widehat{R}_t)$  directly. However, unlike previously,  $\text{Var}(\widehat{R}_s)$  is not a convenient clustering objective, since  $\widehat{R}_s$  cannot be expressed as a kernel mean embedding. Therefore, we propose instead to minimise the variance of

$$\widehat{R}'_s = \frac{1}{n_s} \sum_{h=1}^k |S_h| (\tilde{z}_{s,h} - \mu_s) (\tilde{z}_{s,h} - \mu_s)^T, \quad (12)$$

which is the sample covariance of  $z_s$  in the hypothetical case where  $\mu_s$  is known.

Again, we can assess the validity of this surrogate objective by estimating its rank correlation with the  $\text{Var}(\widehat{R}_s)$ . Given population parameters  $R_s$  and  $\mu_s$ ,  $\text{Var}(\widehat{R}_s)$  and  $\text{Var}(\widehat{R}'_s)$  can be estimated for arbitrary distributions and dimensionality using Monte Carlo simulations. For the special case where  $\tilde{z}_{s,h}$ ,  $h = 1, \dots, k$  are Gaussian scalars, an exact relation can also be derived (Theorem 4).

**Theorem 4.** Assume  $\tilde{z}_{s,h} \sim \mathcal{N}(\mu_h, \Sigma_h)$ ,  $h = 1, \dots, k$ , are independent random scalars, and let  $\widehat{R}_s$  and  $\widehat{R}'_s$  be sample covariance matrices of  $z_s$  as defined. Then,

$$\text{Var}(\widehat{R}_s) = \gamma^2 (\text{Tr}((AC\Sigma)^2) + 4\mu^T AC\Sigma AC\mu) \quad (13)$$

$$\text{Var}(\widehat{R}'_s) = 2 \text{Tr}((A\Sigma)^2) + 4\mu^T AC\Sigma AC\mu \quad (14)$$

$$\text{Var}(\widehat{R}_s) = \gamma^2 (\text{Var}(\widehat{R}'_s) + 2 \text{Tr}(A^2\Sigma)^2 - 4 \text{Tr}(A^3\Sigma^2)) \quad (15)$$

where  $\mu = (\mu_1, \dots, \mu_k)^T$ ,  $\Sigma = \text{diag}((\Sigma_1, \dots, \Sigma_k))$ ,  $A = \frac{1}{n_s} \text{diag}((|S_1|, \dots, |S_k|))$ ,  $\gamma = \frac{n_s}{n_s-1}$  is a bias correction factor, and  $C = I - JA$  is the weighted centering matrix, with  $J$  the square matrix of ones.

*Proof.* See Appendix B. □

Figure 2 shows the relationship between  $\text{Var}(\widehat{R}_s)$  against  $\text{Var}(\widehat{R}'_s)$  for a range of dimensionalities of  $\mathcal{Z}$  and 2 different distributions, the normal distribution (above) and the lognormal distribution (below), to show an asymmetric non-Gaussian case. We set  $k = 8$  and the cluster size ratio  $1, \dots, k$ . The rank correlation is near-monotonic in all cases, suggesting that  $\text{Var}(\widehat{R}'_s)$  is indeed a valid surrogate objective.

Note that

$$\begin{aligned} \text{Var}(\widehat{R}'_s) &= \text{Var}\left(\frac{1}{n_s} \sum_{h=1}^k |S_h| (\tilde{z}_{s,h} - \mu_s) (\tilde{z}_{s,h} - \mu_s)^T\right) \\ &= \frac{1}{n_s^2} \sum_{h=1}^k |S_h|^2 \text{Var}\left((\tilde{z}_{s,h} - \mu_s) (\tilde{z}_{s,h} - \mu_s)^T\right), \end{aligned} \quad (16)$$

and thus we obtain the same dynamically-weighted kernel k-means problem (9) as before, with the specific mapping  $\phi_c(z) = (z - \mu_s)(z - \mu_s)^T$ , and an associated kernel  $\kappa_c(z, z') = \left((z - \mu_s)^T (z' - \mu_s)\right)^2$ .

## 2.4 OPTIMISATION ALGORITHM

Equation (9) can be solved to a local minimum by applying a Lloyd’s-style alternating optimisation algorithm (Lloyd, 1982) along with the kernel trick (Algorithm 1).

---

### Algorithm 1 Dynamically-weighted kernel k-means clustering

---

- 1: **Initialisation:** Use kernel k-means++ (Arthur & Vassilvitskii, 2007) to obtain an initial set of assignments. Then, alternate between Steps 2 and 3 until convergence or max iterations reached.
  - 2: **Distance Update:** Compute the distance matrix  $D \in \mathbb{R}^{n_s \times k}$  from each datapoint to the centroid of each cluster using the kernel trick.
  - 3: **Dynamically Weighted Assignment:** Compute the one-hot cluster assignment matrix  $U \in \{0, 1\}^{n_s \times k}$  that assigns each point to exactly one of the  $k$  clusters.
- 

$U$  is the solution to the quadratic program

$$\arg \min_U \sum_{i,j} \left[ U_{ij} D_{ij} \sum_i U_{ij} \right] \quad \text{subject to } 0 \leq U_{ij} \leq 1, \sum_j U_{ij} = 1. \quad (17)$$

Note that we are able to relax the binary constraint  $U_{ij} \in \{0, 1\}$  as Hoffman’s sufficient conditions for total unimodularity (Heller & Tompkins, 1956) are met, meaning the program will always have integer solutions without having to enforce them explicitly. Since the Hessian of (17) is indefinite in general, this problem is nonconvex and thus finding the global minimum is NP-hard. Although the problem as currently defined could be readily input to a gradient-based interior point method (to find a local minimum), these have  $O\left((kn_s)^3\right)$  complexity, and we found empirically that this is only practical up to a maximum of a few hundred data points. Instead, by considering the specific structure of the problem, we propose in this section a scalable heuristic that can find a local minimum in  $O(kn_s)$  time. The idea is to construct  $U$  incrementally row-by-row using a greedy strategy, weighting the clusters using interim values for the cluster sizes (Algorithm 2).

Greedy algorithms are susceptible to local minima and so a single run of Algorithm 2 is likely to give a poor solution. However, the algorithm can be randomised by permuting the order in which the rows of  $U$  are filled. Given the speed (each iteration simply finds the minimum value of a list), a large number of randomised trials can be performed and then the lowest-cost solution selected. We further found that by parallelising row assignments, more random trials are possible in a given time budget, and this improves the best selected solution, even though the average solution is worse (since the  $n_j$  cannot be updated between parallel assignments).

The performance characteristics of Algorithm 2 are shown in Figure 3. The input comprises Euclidian distances between randomly sampled embeddings of the Humpbacks dataset (described in the subsequent section). We use a small problem with  $n_s = 100$  and  $k = 16$  which allows us to compare performance with the interior point solution. Figure 3a compares the solutions found by Algorithm 2 (10 parallel assignments,  $10^5$  random trials), a commercial interior point solver (The MathWorks Inc., 2021), as well as an unweighted assignment which simply assigns each point to its closest centroid as in the standard k-means algorithm. Figure 3b shows the effect of varying the number of random trials between 1 and  $10^5$ , with the number of parallel assignments fixed at 10.

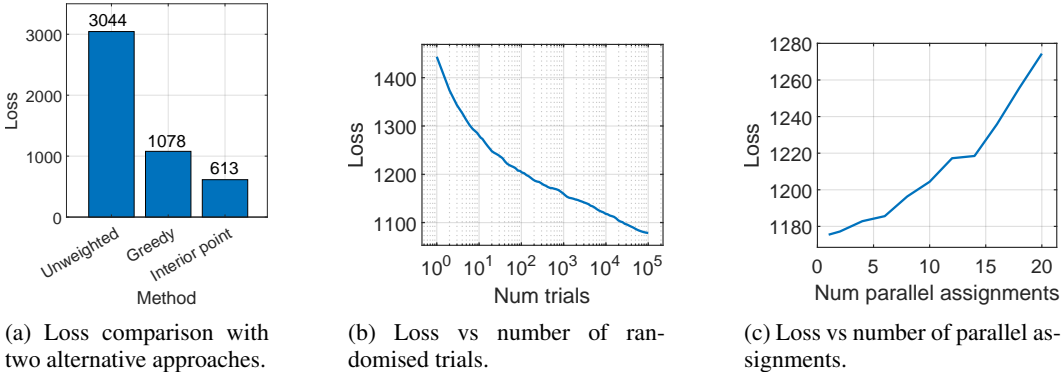


Figure 3: The performance characteristics of Algorithm 2.

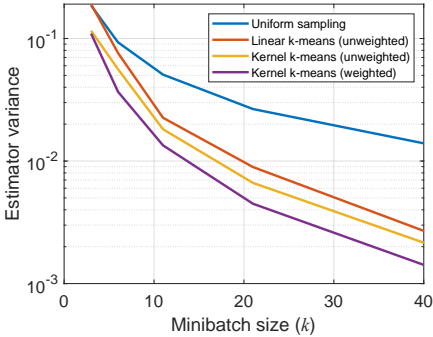


Figure 4: Estimator variance vs minibatch size for Algorithm 1 (purple) vs 3 ablations.

**Algorithm 2** Greedy incremental construction for dynamically weighted assignments

---

**Require:**  $D \in \mathbb{R}^{n_s \times k}$

**Ensure:**  $U \in \{0, 1\}^{n_s \times k}, \sum_j U_{ij} = 1$

- 1:  $U \leftarrow 0_{n_s \times k}$
- 2:  $n_1, \dots, n_k \leftarrow 0$   $\triangleright$  Interim cluster sizes
- 3: **for all**  $i \in \{1, \dots, n_s\}$  **do**
- 4:      $j \leftarrow \arg \min_j D_{ij}(n_j + 1)$
- 5:      $U_{ij} \leftarrow 1$
- 6:      $n_j \leftarrow n_j + 1$
- 7: **end for**
- 8: **return**  $U$

---

Conversely, Figure 3c shows the effect of varying the number of parallel assignments between 1 and 20, with the number of random trials fixed at 100.

An ablation study illustrating the effect of each novel component of our sampling method is shown in Figure 4. The plot shows how  $\text{Var}(\hat{\mu}_{\phi, s})$  varies with  $k$  for 4 different sampling strategies: uniform random sampling, stratified sampling with linear k-means clustering, stratified sampling with kernel k-means clustering, and finally Algorithm 1 which dynamically weights the clusters. The simulations use 1,000 feature embeddings from the Humpbacks dataset using a radial basis function kernel with unit bandwidth. It is seen that each component provides a consistent reduction in estimator variance. For  $k = 40$ , our proposed sampling method achieves around a 10-fold reduction in variance compared to uniform random sampling, and a 1.9-fold reduction compared to stratified sampling using ordinary k-means clustering.

2.5 OVERALL TRAINING STRATEGY

Although the optimal stratification depends on the network weights, it is not necessary to reconstruct the strata at every training iteration. This is because, if the network weights at consecutive training iterations are in a locally smooth region of the parameter space, the structure of the embeddings will be preserved between iterations (Liu et al., 2020b). Based on this, Liu et al. (2020b) propose a low-cost criterion to determine when to re-cluster the data. However, for simplicity, in this paper we re-cluster the data periodically every fixed  $T$  iterations, where  $T$  is thus a trade-off between the computational burden and the degree of variance reduction achieved. Based on empirical observations from Liu et al. (2020b), we choose  $T = 100$  for our experiments.

### 3 EXPERIMENTS

In this section, the proposed method is evaluated in realistic training conditions to assess whether the reduction in variance observed in Figure 4 translates to an increase in test-domain accuracy. Experiments are conducted using the DomainBed framework (Gulrajani & Lopez-Paz, 2021) on the following 3 domain shift datasets.

**Camelyon17-WILDS** (Bándi et al., 2019; Koh et al., 2021) tumour detection in microscopic tissue images across samples from different hospitals. This benchmark comprises a single training-evaluation split, 5 domains, and 14,000 examples.

**Humpbacks** (Napoli & White, 2023) detection of humpback whale vocalisations in underwater acoustic recordings across data from different acoustic monitoring programs. This benchmark comprises 4 data splits, 4 domains and 8,000 examples.

**Spawrious** (Lynch et al., 2023) classification of 4 dog breeds across images with different background environments (desert, jungle, snow etc.). This benchmark comprises 6 data splits of varying difficulty, 6 domains and 18,664 examples.

For Spawrious and Camelyon17, we use a ResNet-18 (He et al., 2015) backbone pre-trained on ImageNet. For Humpbacks, we use the same audio frontend and architecture described in Napoli & White (2023). The domain discrepancies are measured between the union of all training data and a held-out subset of the evaluation set. For the MMD, we use a radial basis function (RBF) mixture kernel (Li et al., 2018), given by  $\kappa(z, z') = \sum_{\gamma \in \mathcal{G}} e^{-\gamma \|z - z'\|^2}$ , with  $\mathcal{G} = \{0.001, 0.01, 0.1, 1, 10\}$ . With reference to Section 2.4, we run Algorithm 2 with 100 random trials and 10 parallel assignments.

Models are trained using the Adam optimiser (Kingma & Ba, 2014) for 3,000 iterations. Hyperparameters are tuned with a random search of size 10 using an in-distribution (training domain) validation set, independently for each sampler. The entire set of experiments is repeated 3 times for reproducibility, using different random seeds for hyperparameters, weight initialisations, and dataset splits. All other hyperparameter choices and training details follow the DomainBed default options.

In addition to uniform random sampling, three variants of stratified sampling are compared, with different stratification strategies: k-means clustering on the input data (Zhao & Zhang, 2014); k-means clustering on  $z_s, z_t$ ; and VaRDASS, which clusters  $z_s, z_t$  using Algorithm 1. Table 1 shows the average test accuracy and standard errors over the data splits and repeats, for each dataset. Breakdowns by data split for Humpbacks and Spawrious are given in Appendix C.

It can be seen that clustering the input data (Zhao & Zhang, 2014) is not effective in improving test domain accuracy, since the required smoothness conditions do not hold for deeper models. Performing simple k-means clustering on the features provides a decent improvement (corroborated by Figure 4, which also shows a good reduction in estimator variance), but VaRDASS performs best, with the highest accuracy in all 6 cases.

In addition to the targeted UDA algorithms, the tables also report test-domain accuracy in the case of non-adaptive training by empirical risk minimisation (ERM). The performance of ERM does not significantly change by varying the sampler, which indicates that the gains observed in CORAL and MMD are indeed due to the sampler’s effect on the adaptation component of the loss.

Finally, we investigate whether the discrepancy between the training and test domain features is lower when the variance reduced samplers are used (Appendix D). The domain discrepancy is measured between held-out datapoints which were not seen during training. The discrepancy is lower for the variance-reduced samplers in 3 out of the 6 tested scenarios. Thus, the results from this experiment are inconclusive. We note that all the discrepancy values are already extremely low (2 to 3 orders of magnitude below ERM training), so we posit that this is because any differences between samplers are likely below the noise floor and cannot be measured in this setup.

### 4 LIMITATIONS

Like all kernel methods, Algorithm 1 requires constructing a kernel matrix, which has quadratic cost in the size of the training set. Although low-rank approximation methods exist which can improve

Table 1: Average test accuracy for each dataset and training algorithm. ERM is included as an ablation but is not part of the main results comparison.

(a) Camelyon17

| Sampler                        | ERM        | CORAL             | MMD               |
|--------------------------------|------------|-------------------|-------------------|
| Uniform random                 | 89.5 ± 1.4 | 92.9 ± 0.6        | 93.7 ± 0.4        |
| Stratified: k-means (inputs)   | 84.0 ± 1.8 | 93.6 ± 0.6        | 93.5 ± 0.7        |
| Stratified: k-means (features) | 88.6 ± 0.4 | 93.8 ± 0.1        | 94.1 ± 0.8        |
| VaRDASS                        | 88.2 ± 0.5 | <b>94.5 ± 0.2</b> | <b>94.7 ± 0.3</b> |

(b) Humpbacks

| Sampler                        | ERM        | CORAL             | MMD               |
|--------------------------------|------------|-------------------|-------------------|
| Uniform random                 | 78.8 ± 1.3 | 81.5 ± 1.6        | 88.3 ± 1.4        |
| Stratified: k-means (inputs)   | 75.7 ± 1.7 | 81.7 ± 1.5        | 82.8 ± 1.4        |
| Stratified: k-means (features) | 77.8 ± 1.1 | 82.9 ± 2.0        | 83.0 ± 1.8        |
| VaRDASS                        | 78.0 ± 1.1 | <b>84.9 ± 1.1</b> | <b>90.4 ± 1.5</b> |

(c) Spawrious

| Sampler                        | ERM        | CORAL             | MMD               |
|--------------------------------|------------|-------------------|-------------------|
| Uniform random                 | 58.1 ± 0.7 | 62.0 ± 0.9        | 62.7 ± 1.0        |
| Stratified: k-means (inputs)   | 58.0 ± 0.6 | 55.8 ± 1.3        | 63.2 ± 0.9        |
| Stratified: k-means (features) | 58.0 ± 0.8 | 66.5 ± 1.1        | 71.1 ± 2.2        |
| VaRDASS                        | 58.7 ± 0.4 | <b>68.4 ± 2.0</b> | <b>72.4 ± 2.2</b> |

the scalability to linear (Chitta et al., 2014), this was not investigated in this paper, and the method will likely still be impractical for very large-scale datasets.

k-means clustering is known to be NP-hard and heuristics such as Lloyd’s (on which Algorithm 1 is based) only search for local solutions. Thus, even though our derived objective is a good surrogate for minimising variance, there is no guarantee that an improved solution can actually be found. Similarly, the runtime of Lloyd’s is known to be superpolynomial in the worst case (Arthur & Vassilvitskii, 2006), although this algorithm remains popular due to its simplicity and the fact it performs well in practice (Arthur et al., 2009; Cordeiro de Amorim & Makarenkov, 2023).

Another problem which can emerge with kernel k-means clustering, especially when  $\mathcal{H}$  is high dimensional, is large imbalance in the strata sizes or even singleton strata. This would slow down the convergence rate of the training, but can be dealt with by introducing minimum cluster size constraints to the cluster assignment step.

## 5 CONCLUSION

This paper introduced VaRDASS, a novel stochastic variance reduction technique for UDA based on stratified sampling, which specifically targets the MMD and CORAL losses. A novel clustering algorithm was proposed to solve this problem and thoroughly analysed. We showed that the method significantly reduces variance in the targeted losses, and that this results in consistent improvements in target domain accuracy. For future work, the method could adopt a strata reconstruction condition as in Liu et al. (2020b), be combined with complementary variance reduction techniques, or be extended to additional UDA or domain generalisation algorithms. Alternative clustering algorithms with superior performance to Lloyd’s could also be investigated.

## REFERENCES

Guillaume Alain, Alex Lamb, Chinnadhurai Sankar, Aaron Courville, and Yoshua Bengio. Variance Reduction in SGD by Distributed Importance Sampling. *ICLR Workshops Track*, 11 2015. URL <https://arxiv.org/abs/1511.06481v7>.

- 486 David Arthur and Sergei Vassilvitskii. How slow is the k-means method? *Proceedings of the Annual*  
487 *Symposium on Computational Geometry*, pp. 144–153, 2006. doi: 10.1145/1137856.1137880.  
488 URL <https://dl.acm.org/doi/pdf/10.1145/1137856.1137880>.  
489
- 490 David Arthur and Sergei Vassilvitskii. k-means++: The Advantages of Careful Seeding. *Proceed-*  
491 *ings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, 2007.
- 492 David Arthur, Bodo Manthey, and Heiko Röglin. k-Means has Polynomial Smoothed Complexity.  
493 *Proceedings - Annual IEEE Symposium on Foundations of Computer Science, FOCS*, pp. 405–  
494 414, 4 2009. ISSN 02725428. doi: 10.1109/FOCS.2009.14. URL [https://arxiv.org/  
495 pdf/0904.1113](https://arxiv.org/pdf/0904.1113).  
496
- 497 Péter Bándi, Oscar Geessink, Quirine Manson, Marcory Van Dijk, Maschenka Balkenhol, Meyke  
498 Hermsen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, Aoxiao Zhong,  
499 Quanzheng Li, Farhad Ghazvinian Zanjani, Svitlana Zinger, Keisuke Fukuta, Daisuke Komura,  
500 Vlado Ovtcharov, Shenghua Cheng, Shaoqun Zeng, Jeppe Thagaard, Anders B. Dahl, Huangjing  
501 Lin, Hao Chen, Ludwig Jacobsson, Martin Hedlund, Melih Çetin, Eren Halici, Hunter Jack-  
502 son, Richard Chen, Fabian Both, Jörg Franke, Heidi Kusters-Vandeveld, Willem Vreuls, Peter  
503 Bult, Bram Van Ginneken, Jeroen Van Der Laak, and Geert Litjens. From Detection of Indi-  
504 vidual Metastases to Classification of Lymph Node Status at the Patient Level: The CAME-  
505 LYON17 Challenge. *IEEE transactions on medical imaging*, 38(2):550–560, 2 2019. ISSN  
506 1558-254X. doi: 10.1109/TMI.2018.2867350. URL [https://pubmed.ncbi.nlm.nih.  
507 gov/30716025/](https://pubmed.ncbi.nlm.nih.gov/30716025/).
- 508 Remi Bardenet, Subhroshekhar Ghosh, and Meixia Lin. Determinantal point processes based on  
509 orthogonal polynomials for sampling minibatches in SGD. *NeurIPS*, 20:16226–16237, 12 2021.  
510 ISSN 10495258. URL <https://arxiv.org/abs/2112.06007v1>.
- 511 Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of Representations  
512 for Domain Adaptation. *NeurIPS*, 19, 2006.
- 513 Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wort-  
514 man Vaughan. A theory of learning from different domains. *Machine Learning*, 79(1-2):  
515 151–175, 10 2010. ISSN 15730565. doi: 10.1007/S10994-009-5152-4/METRICS. URL  
516 <https://link.springer.com/article/10.1007/s10994-009-5152-4>.  
517
- 518 Wacha Bounliphone, Eugene Belilovsky, Matthew B. Blaschko, Ioannis Antonoglou, and Arthur  
519 Gretton. A Test of Relative Similarity For Model Selection in Generative Models. *ICLR*, 11  
520 2016.
- 521 Radha Chitta, Rong Jin, Timothy C. Havens, and Anil K. Jain. Scalable Kernel Clustering: Approx-  
522 imate Kernel k-means. *arXiv*, 2 2014. URL <https://arxiv.org/abs/1402.3849v1>.  
523
- 524 Renato Cordeiro de Amorim and Vladimir Makarenkov. On k-means iterations and Gaussian clus-  
525 ters. *Neurocomputing*, 553:126547, 10 2023. ISSN 0925-2312. doi: 10.1016/J.NEUCOM.  
526 2023.126547. URL [https://www.sciencedirect.com/science/article/pii/  
527 S0925231223006707](https://www.sciencedirect.com/science/article/pii/S0925231223006707).
- 528 Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A Fast Incremental Gradient  
529 Method With Support for Non-Strongly Convex Composite Objectives. *NeurIPS*, 2(January):  
530 1646–1654, 7 2014. ISSN 10495258. URL <https://arxiv.org/abs/1407.0202v3>.  
531
- 532 Abhimanyu Dubey, Vignesh Ramanathan, Alex Pentland, and Dhruv Mahajan. Adaptive Meth-  
533 ods for Real-World Domain Generalization. *CVPR*, 2021. ISSN 10636919. doi: 10.1109/  
534 CVPR46437.2021.01411. URL <https://arxiv.org/abs/2103.15796v2>.
- 535 Tianfan Fu and Zhihua Zhang. CPSG-MCMC: Clustering-Based Preprocessing method for Stochas-  
536 tic Gradient MCMC. *AISTATS*, pp. 841–850, 4 2017. ISSN 2640-3498. URL [https:  
537 //proceedings.mlr.press/v54/ful7a.html](https://proceedings.mlr.press/v54/ful7a.html).  
538
- 539 Irena Gao, Shiori Sagawa, Pang Wei Koh, Tatsunori Hashimoto, and Percy Liang. Out-of-  
Distribution Robustness via Targeted Augmentations. *ICML*, 10 2023.

- 540 Mike Giles. Numerical Methods II Lecture 4. Technical report, Oxford University Mathematical In-  
541 stitute, 2015. URL <https://people.maths.ox.ac.uk/gilesm/mc/mc/lec4.pdf>.
- 542 Robert M. Gower, Mark Schmidt, Francis Bach, and Peter Richtarik. Variance-Reduced Methods  
543 for Machine Learning. *Proceedings of the IEEE*, 108(11):1968–1983, 11 2020. ISSN 15582256.  
544 doi: 10.1109/JPROC.2020.3028013.
- 545 Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola.  
546 A Kernel Two-Sample Test. *Journal of Machine Learning Research*, 13:723–773, 2012.
- 547 Ishaan Gulrajani and David Lopez-Paz. In Search of Lost Domain Generalization. *ICLR*, 2021.
- 548 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image  
549 Recognition. *Proceedings of the IEEE Computer Society Conference on Computer Vision and  
550 Pattern Recognition*, 2016-December:770–778, 12 2015. ISSN 10636919. doi: 10.48550/arxiv.  
551 1512.03385. URL <https://arxiv.org/abs/1512.03385v1>.
- 552 I. Heller and C.B. Tompkins. An Extension of a Theorem of Dantzig’s. *Linear Inequalities and  
553 Related Systems*, 1956.
- 554 Rie Johnson and Tong Zhang. Accelerating Stochastic Gradient Descent using Predictive Variance  
555 Reduction. *NeurIPS*, 2013.
- 556 Tyler B Johnson and Carlos Guestrin. Training Deep Models Faster with Robust, Approximate  
557 Importance Sampling. *NeurIPS*, 2018.
- 558 Angelos Katharopoulos and François Fleuret. Biased Importance Sampling for Deep Neural Net-  
559 work Training. *arXiv*, 2017.
- 560 Angelos Katharopoulos and François Fleuret. Not All Samples Are Created Equal: Deep Learning  
561 with Importance Sampling. *ICML*, 2018.
- 562 Diederik P. Kingma and Jimmy Lei Ba. Adam: A Method for Stochastic Optimization. *3rd Interna-  
563 tional Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*, 12  
564 2014. doi: 10.48550/arxiv.1412.6980. URL <https://arxiv.org/abs/1412.6980v9>.
- 565 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsub-  
566 ramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne  
567 David, Ian Stavness, Wei Guo, Berton A. Earnshaw, Imran S. Haque, Sara Beery, Jure Leskovec,  
568 Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. WILDS: A  
569 Benchmark of in-the-Wild Distribution Shifts. *ICML*, 2021.
- 570 Takuro Kutsuna. Exploring Variance Reduction in Importance Sampling for Efficient DNN Training.  
571 *arXiv*, 1 2025. URL <https://arxiv.org/abs/2501.13296v1>.
- 572 Olivier Ledoit and Michael Wolf. The Power of (Non-)Linear Shrinking: A Review and Guide to  
573 Covariance Matrix Estimation. *Journal of Financial Econometrics*, 2020. doi: 10.1093/jjfinec/  
574 nbaa007. URL [https://academic.oup.com/jfec/advance-article/doi/10.  
575 1093/jjfinec/nbaa007/5861007](https://academic.oup.com/jfec/advance-article/doi/10.1093/jjfinec/nbaa007/5861007).
- 576 Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C. Kot. Domain Generalization with  
577 Adversarial Feature Learning. *CVPR*, pp. 5400–5409, 12 2018. ISSN 10636919. doi:  
578 10.1109/CVPR.2018.00566.
- 579 Feng Liu, Wenkai Xu, Jie Lu, Guangquan Zhang, Arthur Gretton, and D. J. Sutherland. Learning  
580 Deep Kernels for Non-Parametric Two-Sample Tests. *ICML*, pp. 6272–6282, 2 2020a.
- 581 Jingchang Liu and Linli Xu. Accelerating Stochastic Gradient Descent Using Antithetic Sampling.  
582 *arXiv*, 10 2018. URL <https://arxiv.org/abs/1810.03124v1>.
- 583 Weijie Liu, Hui Qian, Chao Zhang, Zebang Shen, Jiahao Xie, and Nenggan Zheng. Accelerating  
584 Stratified Sampling SGD by Reconstructing Strata. *IJCAI*, 2020b.
- 585 Stuart P. Lloyd. Least Squares Quantization in PCM. *IEEE Transactions on Information Theory*, 28  
586 (2):129–137, 1982. ISSN 15579654. doi: 10.1109/TIT.1982.1056489.
- 587
- 588
- 589
- 590
- 591
- 592
- 593

- 594 Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning Transferable Features  
595 with Deep Adaptation Networks. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd*  
596 *International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning*  
597 *Research*, pp. 97–105, Lille, France, 2015. PMLR. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v37/long15.html)  
598 [press/v37/long15.html](https://proceedings.mlr.press/v37/long15.html).
- 599 Ilya Loshchilov and Frank Hutter. Online Batch Selection for Faster Training of Neural Networks.  
600 *ICLR workshop track*, 11 2016. URL <https://arxiv.org/pdf/1511.06343>.
- 601 Yucheng Lu, Youngsuk Park, Lifan Chen, Yuyang Wang, Christopher De Sa, and Dean Foster.  
602 Variance Reduced Training with Stratified Sampling for Forecasting Models. *ICML*, 139:7145–  
603 7155, 3 2021. ISSN 26403498. URL <https://arxiv.org/abs/2103.02062v2>.
- 604 Aengus Lynch, Gbètondji J-S Dovonon, Jean Kaddour, and Ricardo Silva. Spawrious: A Benchmark  
605 for Fine Control of Spurious Correlation Biases. *arXiv*, 3 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2303.05470v3)  
606 [abs/2303.05470v3](https://arxiv.org/abs/2303.05470v3).
- 607 Andrea Napoli and Paul White. Unsupervised Domain Adaptation for the Cross-Dataset Detection  
608 of Humpback Whale Calls. *DCASE*, 2023.
- 609 Andrea Napoli and Paul White. Improving Domain Generalisation with Diversity-based Sampling.  
610 *DCASE*, 2024. URL <http://arxiv.org/abs/2410.04235>.
- 611 Xinyu Peng, Li Li, and Fei Yue Wang. Accelerating Minibatch Stochastic Gradient Descent Using  
612 Typicality Sampling. *IEEE Transactions on Neural Networks and Learning Systems*, 31(11):  
613 4649–4659, 11 2020. ISSN 21622388. doi: 10.1109/TNNLS.2019.2957003.
- 614 B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Com-*  
615 *putational Mathematics and Mathematical Physics*, 4(5):1–17, 1964. ISSN 00415553. doi:  
616 10.1016/0041-5553(64)90137-5.
- 617 B. T. Polyak and A. B. Juditsky. Acceleration of Stochastic Approximation by Averaging. *SIAM*  
618 *Journal on Control and Optimization*, 30(4):838–855, 1992. ISSN 03630129. doi: 10.1137/  
619 0330046.
- 620 Ievgen Redko, Emilie Morvant, Amaury Habrard, Marc Sebban, and Younès Bennani. A survey on  
621 domain adaptation theory: learning bounds and theoretical guarantees. *arXiv*, 2022.
- 622 George A F Seber. *A Matrix Handbook for Statisticians*. Wiley-Interscience, USA, 1st edition,  
623 2007. ISBN 0471748692.
- 624 Shai Shalev-Shwartz and Tong Zhang. Stochastic Dual Coordinate Ascent Methods for Regularized  
625 Loss Minimization. *JMLR*, 14(1):567–599, 9 2013. ISSN 15324435. URL [https://arxiv.](https://arxiv.org/abs/1209.1873v2)  
626 [org/abs/1209.1873v2](https://arxiv.org/abs/1209.1873v2).
- 627 Baochen Sun and Kate Saenko. Deep CORAL: Correlation Alignment for Deep Domain Adaptation.  
628 *ECCV*, 9915 LNCS:443–450, 7 2016. ISSN 16113349. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1607.01719v1)  
629 [1607.01719v1](https://arxiv.org/abs/1607.01719v1).
- 630 Danica J Sutherland and Namrata Deka. Unbiased estimators for the variance of MMD estimators.  
631 *arXiv*, 6 2022. URL <https://arxiv.org/pdf/1906.02104>.
- 632 The MathWorks Inc. MATLAB, 2021.
- 633 Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep Domain Confu-  
634 sion: Maximizing for Domain Invariance. *arXiv*, 12 2014. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1412.3474v1)  
635 [1412.3474v1](https://arxiv.org/abs/1412.3474v1).
- 636 Zirui Wang, Zihang Dai, Barnabas Poczos, and Jaime Carbonell. Characterizing and Avoiding  
637 Negative Transfer. *CVPR*, 2019-June:11285–11294, 11 2019. ISSN 10636919. doi: 10.1109/  
638 CVPR.2019.01155. URL <https://arxiv.org/abs/1811.09751v4>.
- 639 Cheng Zhang, Hedvig Kjellström, and Stephan Mandt. Determinantal Point Processes for Mini-  
640 Batch Diversification. *Uncertainty in Artificial Intelligence*, 2017.

648 Cheng Zhang, Cengiz Öztireli, Stephan Mandt, and Giampiero Salvi. Active Mini-Batch Sampling  
649 Using Repulsive Point Processes. *AAAI*, 33(01):5741–5748, 7 2019. ISSN 2374-3468. doi:  
650 10.1609/AAAI.V33I01.33015741. URL [https://ojs.aaai.org/index.php/AAAI/  
651 article/view/4520](https://ojs.aaai.org/index.php/AAAI/article/view/4520).

652 Peilin Zhao and Tong Zhang. Accelerating Minibatch Stochastic Gradient Descent using Stratified  
653 Sampling. *arXiv*, 5 2014. URL <https://arxiv.org/abs/1405.3080v1>.

654 Peilin Zhao and Tong Zhang. Stochastic Optimization with Importance Sampling for Regu-  
655 larized Loss Minimization. *ICML*, pp. 1–9, 6 2015. ISSN 1938-7228. URL [https:  
656 //proceedings.mlr.press/v37/zhaoa15.html](https://proceedings.mlr.press/v37/zhaoa15.html).

657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A PROOF OF THEOREM 1

**Theorem 1.** Assume  $\mathcal{H}$  has finite dimensionality  $d$  and  $\Sigma_{\hat{\mu}_{\phi,s}} = \sigma^2 I$  for some positive scalar  $\sigma^2$ . Then,

$$\arg \min_{\mathcal{S}} \text{Var} \left( \widehat{L}_{\text{MMD}} \right) = \arg \min_{\mathcal{S}} \text{Var} \left( \widehat{\mu}_{\phi,s} \right). \quad (18)$$

*Proof.* Noting that  $\text{Var} \left( \widehat{\mu}_{\phi,s} \right) = \sigma^2 d$ , we have

$$\begin{aligned} \text{Var} \left( \widehat{L}_{\text{MMD}} \right) &= 2 \text{Tr} \left( (\sigma^2 I + \Sigma_{\widehat{\mu}_{\phi,t}})^2 \right) + 4m^T (\sigma^2 I + \Sigma_{\widehat{\mu}_{\phi,t}}) m \\ &= 2 \text{Tr} \left( \sigma^4 I + 2\sigma^2 \Sigma_{\widehat{\mu}_{\phi,t}} + \Sigma_{\widehat{\mu}_{\phi,t}}^2 \right) + 4\sigma^2 m^T m + m^T \Sigma_{\widehat{\mu}_{\phi,t}} m \\ &= 2\sigma^4 d + 4\sigma^2 (m^T m + \text{Tr} (\Sigma_{\widehat{\mu}_{\phi,t}})) + 2 \text{Tr} \left( \Sigma_{\widehat{\mu}_{\phi,t}}^2 \right) + m^T \Sigma_{\widehat{\mu}_{\phi,t}} m \\ &= \frac{2}{d} \text{Var} \left( \widehat{\mu}_{\phi,s} \right)^2 + \frac{4}{d} \text{Var} \left( \widehat{\mu}_{\phi,s} \right) (m^T m + \text{Tr} (\Sigma_{\widehat{\mu}_{\phi,t}})) \\ &\quad + 2 \text{Tr} \left( \Sigma_{\widehat{\mu}_{\phi,t}}^2 \right) + m^T \Sigma_{\widehat{\mu}_{\phi,t}} m, \end{aligned} \quad (19)$$

which is monotonically increasing in  $\text{Var} \left( \widehat{\mu}_{\phi,s} \right)$  for all fixed  $m, \Sigma_{\widehat{\mu}_{\phi,t}}, d$ .  $\square$

## B PROOF OF THEOREM 4

**Theorem 4.** Assume  $\tilde{z}_{s,h} \sim \mathcal{N}(\mu_h, \Sigma_h)$ ,  $h = 1, \dots, k$  are independent random scalars, and let  $\widehat{R}_s$  and  $\widehat{R}'_s$  be sample covariance matrices of  $z_s$  as defined. Then,

$$\text{Var} \left( \widehat{R}_s \right) = \gamma^2 \left( \text{Tr} \left( (AC\Sigma)^2 \right) + 4\mu^T AC\Sigma AC\mu \right) \quad (20)$$

$$\text{Var} \left( \widehat{R}'_s \right) = 2 \text{Tr} \left( (A\Sigma)^2 \right) + 4\mu^T AC\Sigma AC\mu \quad (21)$$

$$\text{Var} \left( \widehat{R}_s \right) = \gamma^2 \left( \text{Var} \left( \widehat{R}'_s \right) + 2 \text{Tr} \left( A^2 \Sigma \right)^2 - 4 \text{Tr} \left( A^3 \Sigma^2 \right) \right) \quad (22)$$

where  $\mu = (\mu_1, \dots, \mu_k)^T$ ,  $\Sigma = \text{diag} \left( (\Sigma_1, \dots, \Sigma_k) \right)$ ,  $A = \frac{1}{n_s} \text{diag} \left( (|S_1|, \dots, |S_k|) \right)$ ,  $\gamma = \frac{n_s}{n_s - 1}$  is a bias correction factor, and  $C = I - JA$  is the weighted centering matrix, with  $J$  the square matrix of ones.

*Proof.* For scalar  $\tilde{z}_{s,h}$ , the covariances can be written as quadratic forms

$$\widehat{R}_s = (CZ)^T A(CZ) = \gamma Z^T ACZ \quad (23)$$

$$\widehat{R}'_s = (Z - \mu_s)^T A(Z - \mu_s) \quad (24)$$

with  $Z = (\tilde{z}_{s,1}, \dots, \tilde{z}_{s,k})^T \sim \mathcal{N}(\mu, \Sigma)$ . For  $\text{Var} \left( \widehat{R}_s \right)$ , the expression is immediate from Seber (2007). For  $\text{Var} \left( \widehat{R}'_s \right)$ , one has

$$\text{Var} \left( \widehat{R}'_s \right) = 2 \text{Tr} \left( (A\Sigma)^2 \right) + 4(\mu - \mu_s)^T A\Sigma A(\mu - \mu_s) \quad (25)$$

and note that  $C\mu = \mu - \mu_s$  because  $\mu_s = \frac{1}{n_s} \sum_{h=1}^k |S_h| \mu_h$ . The result follows.

Expanding  $\text{Tr} \left( (AC\Sigma)^2 \right) = \text{Tr} \left( (A(I - JA)\Sigma)^2 \right)$  and substituting then gives

$$\begin{aligned} \frac{1}{\gamma^2} \text{Var} \left( \widehat{R}_s \right) &= 2 \left( \text{Tr} \left( (A\Sigma)^2 \right) + \text{Tr} \left( A^2 \Sigma \right)^2 - 2 \text{Tr} \left( A^3 \Sigma^2 \right) \right) + 4\mu^T AC\Sigma AC\mu \\ &= \text{Var} \left( \widehat{R}'_s \right) + 2 \text{Tr} \left( A^2 \Sigma \right)^2 - 4 \text{Tr} \left( A^3 \Sigma^2 \right). \end{aligned} \quad (26) \quad \square$$

## C PERFORMANCE BREAKDOWN BY DATA SPLIT

Table 2: Average test accuracy for the Humpbacks dataset by data split.

| Sampler                        | ERM        | CORAL             | MMD               |
|--------------------------------|------------|-------------------|-------------------|
| <i>Domain 1</i>                |            |                   |                   |
| Uniform random                 | 70.0 ± 2.7 | 72.1 ± 0.8        | <b>80.5 ± 2.4</b> |
| Stratified: k-means (inputs)   | 61.7 ± 2.3 | 67.7 ± 3.6        | 73.2 ± 2.1        |
| Stratified: k-means (features) | 73.6 ± 2.6 | 74.5 ± 4.1        | 70.2 ± 4.1        |
| VaRDASS                        | 68.0 ± 3.0 | <b>77.8 ± 2.8</b> | 80.0 ± 0.4        |
| <i>Domain 2</i>                |            |                   |                   |
| Uniform random                 | 77.7 ± 3.0 | 88.4 ± 5.9        | 92.2 ± 1.9        |
| Stratified: k-means (inputs)   | 80.1 ± 5.7 | 91.4 ± 2.4        | 88.5 ± 3.3        |
| Stratified: k-means (features) | 73.6 ± 3.1 | 89.0 ± 6.2        | 94.8 ± 1.4        |
| VaRDASS                        | 80.3 ± 2.8 | <b>95.2 ± 2.0</b> | <b>97.1 ± 1.1</b> |
| <i>Domain 3</i>                |            |                   |                   |
| Uniform random                 | 73.4 ± 2.8 | 72.4 ± 1.2        | 87.0 ± 4.6        |
| Stratified: k-means (inputs)   | 68.6 ± 2.1 | <b>77.4 ± 4.0</b> | 78.8 ± 3.9        |
| Stratified: k-means (features) | 69.8 ± 1.1 | 75.5 ± 2.0        | 73.2 ± 5.6        |
| VaRDASS                        | 68.5 ± 0.7 | 71.7 ± 1.8        | <b>89.8 ± 5.9</b> |
| <i>Domain 4</i>                |            |                   |                   |
| Uniform random                 | 94.2 ± 1.3 | 93.1 ± 2.1        | 93.5 ± 0.6        |
| Stratified: k-means (inputs)   | 92.2 ± 1.3 | 90.3 ± 0.9        | 90.8 ± 1.6        |
| Stratified: k-means (features) | 94.5 ± 0.9 | 92.7 ± 1.6        | 93.6 ± 0.9        |
| VaRDASS                        | 95.0 ± 0.7 | <b>94.8 ± 2.1</b> | <b>94.8 ± 0.8</b> |

Table 3: Average test accuracy for the Spawrious dataset by data split.

| Sampler                        | ERM        | CORAL             | MMD                |
|--------------------------------|------------|-------------------|--------------------|
| <i>O2O-Easy</i>                |            |                   |                    |
| Uniform random                 | 66.4 ± 1.6 | 69.4 ± 3.6        | 73.8 ± 2.2         |
| Stratified: k-means (inputs)   | 67.2 ± 1.2 | 74.4 ± 6.2        | 90.4 ± 0.9         |
| Stratified: k-means (features) | 69.5 ± 0.6 | <b>87.7 ± 3.2</b> | 81.0 ± 7.0         |
| VaRDASS                        | 65.2 ± 1.3 | 83.4 ± 7.1        | <b>91.7 ± 2.8</b>  |
| <i>O2O-Medium</i>              |            |                   |                    |
| Uniform random                 | 62.5 ± 1.0 | 56.0 ± 2.0        | <b>61.9 ± 1.9</b>  |
| Stratified: k-means (inputs)   | 60.7 ± 1.4 | 59.2 ± 2.3        | 54.8 ± 1.8         |
| Stratified: k-means (features) | 60.4 ± 0.7 | 60.8 ± 0.8        | 57.7 ± 2.6         |
| VaRDASS                        | 61.4 ± 0.6 | <b>61.9 ± 1.6</b> | 60.8 ± 0.4         |
| <i>O2O-Hard</i>                |            |                   |                    |
| Uniform random                 | 61.3 ± 0.9 | 64.9 ± 0.6        | 60.5 ± 1.9         |
| Stratified: k-means (inputs)   | 60.4 ± 1.7 | 51.2 ± 2.4        | 60.2 ± 2.9         |
| Stratified: k-means (features) | 63.8 ± 1.1 | 61.1 ± 2.8        | 81.4 ± 4.2         |
| VaRDASS                        | 61.5 ± 1.5 | <b>71.2 ± 8.4</b> | <b>83.4 ± 3.7</b>  |
| <i>M2M-Easy</i>                |            |                   |                    |
| Uniform random                 | 70.7 ± 2.9 | 79.1 ± 2.3        | 80.5 ± 4.0         |
| Stratified: k-means (inputs)   | 70.0 ± 1.1 | 78.8 ± 0.1        | 78.3 ± 3.2         |
| Stratified: k-means (features) | 69.6 ± 0.4 | 81.0 ± 3.1        | 80.8 ± 4.7         |
| VaRDASS                        | 71.6 ± 0.7 | <b>85.2 ± 3.6</b> | <b>84.2 ± 1.2</b>  |
| <i>M2M-Medium</i>              |            |                   |                    |
| Uniform random                 | 44.3 ± 1.4 | 54.4 ± 2.0        | 51.3 ± 2.4         |
| Stratified: k-means (inputs)   | 49.6 ± 1.3 | 52.7 ± 1.0        | 50.0 ± 1.1         |
| Stratified: k-means (features) | 47.8 ± 2.3 | 57.4 ± 2.5        | 59.8 ± 3.3         |
| VaRDASS                        | 48.4 ± 0.8 | <b>59.4 ± 1.6</b> | <b>60.0 ± 11.3</b> |
| <i>M2M-Hard</i>                |            |                   |                    |
| Uniform random                 | 43.3 ± 1.2 | 48.3 ± 0.8        | 48.1 ± 0.7         |
| Stratified: k-means (inputs)   | 40.0 ± 2.2 | 18.6 ± 3.4        | 45.5 ± 1.5         |
| Stratified: k-means (features) | 36.8 ± 4.2 | <b>51.2 ± 2.7</b> | <b>66.2 ± 8.7</b>  |
| VaRDASS                        | 44.0 ± 1.1 | 49.4 ± 2.4        | 54.1 ± 4.8         |

## D DOMAIN DISCREPANCY POST-TRAINING

Table 4: Average discrepancy between the training and test domains for each dataset.

(a) Camelyon17

| Sampler                        | CORAL                               | MMD                                 |
|--------------------------------|-------------------------------------|-------------------------------------|
| Uniform random                 | $0.151 \pm 0.074$                   | <b><math>0.135 \pm 0.021</math></b> |
| Stratified: k-means (inputs)   | <b><math>0.093 \pm 0.011</math></b> | $0.183 \pm 0.030$                   |
| Stratified: k-means (features) | $0.117 \pm 0.084$                   | $0.260 \pm 0.038$                   |
| VaRDASS                        | $0.114 \pm 0.023$                   | $0.217 \pm 0.042$                   |

(b) Humpbacks

| Sampler                        | CORAL                               | MMD                                 |
|--------------------------------|-------------------------------------|-------------------------------------|
| Uniform random                 | $5.712 \pm 1.203$                   | $6.489 \pm 0.289$                   |
| Stratified: k-means (inputs)   | <b><math>4.551 \pm 1.049</math></b> | $6.936 \pm 0.380$                   |
| Stratified: k-means (features) | $8.984 \pm 1.885$                   | $6.396 \pm 0.326$                   |
| VaRDASS                        | $8.695 \pm 1.533$                   | <b><math>6.008 \pm 0.321</math></b> |

(c) Spawrious

| Sampler                        | CORAL                               | MMD                                 |
|--------------------------------|-------------------------------------|-------------------------------------|
| Uniform random                 | $0.281 \pm 0.049$                   | <b><math>0.206 \pm 0.038</math></b> |
| Stratified: k-means (inputs)   | <b><math>0.201 \pm 0.043</math></b> | $0.229 \pm 0.041$                   |
| Stratified: k-means (features) | $0.270 \pm 0.061$                   | $0.255 \pm 0.041$                   |
| VaRDASS                        | $0.220 \pm 0.034$                   | $0.317 \pm 0.051$                   |

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917