
Identifying the Instances Associated with Distribution Shifts using the Max-Sliced Bures Divergence

Austin J. Brockmeier*
University of Delaware
Newark, Delaware, USA
ajbrock@udel.edu

Claudio Cesar Claros-Olivares
University of Delaware
Newark, Delaware, USA
cesar@udel.edu

Matthew S. Emigh
Naval Surface Warfare Center- Panama City Division
Panama City, Florida, USA
matthew.emigh@navy.mil

Luis G. Sanchez Giraldo
University of Kentucky
Lexington, Kentucky, USA
luis.sanchez@uky.edu

Abstract

We investigate an interpretable approach to compare two distributions. The approach, max-sliced Bures divergence, approximates the max-sliced Wasserstein distance and projects the distributions into a one-dimensional subspace defined by a ‘slicing’ vector. Unlike heuristic algorithms for the max-sliced Wasserstein-2 distance that are not guaranteed to find the optimal slice, we detail a tractable algorithm that finds the global optimal slice and scales to large sample sizes, due to its expression in terms of second moments. However, it is unable to detect changes in higher-order statistics. To overcome this, we explore using a non-linear mapping provided by the internal representation of a pre-trained neural network (Inception Net). Our approach provides an interpretation of the Fréchet Inception distance by identifying the instances that are either overrepresented or underrepresented with respect to the other sample. We apply the proposed measure to detect class imbalances and underrepresentation within data sets.

1 Introduction

Divergence measures quantify the dissimilarity between probability distributions and are useful for detecting and correcting covariate shift [Shimodaira, 2000, Quionero-Candela et al., 2009, Lipton et al., 2018]. Although a variety of divergences exist, not all are easy to both interpret and estimate. We consider a divergence to be interpretable if it can be expressed in terms of a ‘witness function’ that can be used to identify subsets that contribute most to the divergence between two distributions. Additionally, rather than allow subsets to come from across the domain, the witness function should identify localized, ideally contiguous, subsets. This allows the user to pinpoint and understand specific examples from a subset of the differences between two samples or distributions.

A natural choice for a witness function is the log of the ratio of the probability density functions, as in the Kullback-Leibler divergence [Kullback and Leibler, 1951] and other f -divergences [Ali and Silvey, 1966, Rényi, 1961]. However, estimating the densities from samples is challenging [Vapnik, 2013], which motivates alternative approaches such as directly estimating the density ratio [Nguyen et al., 2008, Kanamori et al., 2009, Yamada et al., 2011, 2013, Saito et al., 2018, Lee et al., 2019] or via variational formulations [Nguyen et al., 2010, Nowozin et al., 2016].

*<https://www.eecis.udel.edu/~ajbrock>

Another class of divergences are the integral probability metrics (IPMs) [Müller, 1997]. An IPM can be expressed as the maximal difference between the first moments of witness function evaluations under each distribution for a given function family Ω , $D_\Omega(\mu, \nu) = \sup_{\omega \in \Omega} |\mathbb{E}_{X \sim \mu}[\omega(X)] - \mathbb{E}_{Y \sim \nu}[\omega(Y)]|$. Different choices of Ω yield different IPMs; examples include total variation, the dual formulation of the Wasserstein-1 distance, maximum mean discrepancy (MMD) [Gretton et al., 2007], and covariance feature matching [Mroueh et al., 2017]. Since $\omega^*(\cdot) \in \arg \max_{\omega \in \Omega} \mathbb{E}_{X \sim \mu}[\omega(X)] - \mathbb{E}_{Y \sim \nu}[\omega(Y)]$, relatively large evaluations of $\omega^*(\cdot)$ are associated to subsets with higher probability under μ . Restricting Ω to the class of functions whose large evaluations are localized would yield an interpretable divergence. An alternative option for interpretable witness functions is to consider restricted subspaces (slices) of the domain.

Before answering the question, ‘Do existing interpretable divergences exist?’ we first consider why they may not be useful for other tasks—beyond interpretability. A divergence that is restricted to localized witness functions makes inefficient use of all the instances in a sample in learning and hypothesis testing (two-sample tests). For the latter, the witness function should collect evidence across the domain: focusing on a localized subset or low-dimensional subspace limits the evidence collection. For batch sampling to train a generative model, a localized witness would only provide feedback for one specific subset at a time, wasting points that fall outside of this subset. For example, restricting the witness function to a rank-1 subspace would waste differences normal to this subspace.

Nonetheless, such a divergence exists: it is the max-sliced Wasserstein- p distance [Kolouri et al., 2019, Deshpande et al., 2019], $D_{W_p, \Omega}(\mu, \nu) = \sup_{\omega \in \Omega} W_p(\omega_\# \mu, \omega_\# \nu)$ for an appropriate function family Ω . The limitations of max-slicing for generative modeling and learning have already been noted and alternatives proposed [Nguyen et al., 2021]. While not as simple to interpret as the general case, the distributional form of the sliced Wasserstein distance [Nguyen et al., 2021], which optimizes the distribution over the slices, can also be used to investigate discrepancies across different subspaces. A difficulty with the max-sliced Wasserstein- p distance is its computation, since it is a saddle-point optimization problem that requires solving a one-dimensional optimal transport problem at each evaluation. Previous works have approximated the max-sliced Wasserstein distance using a first moment approximation [Deshpande et al., 2019] or a finite number of steps of a local optimizer [Kolouri et al., 2019]. Neither guarantees obtaining an optimal witness function.

As an alternative we propose the ‘‘max-sliced Bures divergence’’ which relaxes the Wasserstein-2 (W2) distance with a second-moment approximation that relies on the Bures distance [Dowson and Landau, 1982, Gelbrich, 1990]. The Bures distance is a distance metric between positive semidefinite matrices, which is equivalent to the W2 distance for zero-mean Gaussians [Muzellec and Cuturi, 2018, Zhang et al., 2020, Oh et al., 2020, De Meulemeester et al., 2020, Janati et al., 2020]. We derive a one-sided version of the max-sliced Bures divergence that

- provides an interpretable witness function that highlight regions of over- and under-representation in one distribution with respect to the other, especially useful for covariate shift detection;
- is expressed as an optimization problem that has a globally optimal solution, which can be computed efficiently (after initial computation of the means and covariance matrices, calculating the divergence and the optimal witness function is independent of the sample sizes, enabling it to be applied to large data sets and high-dimensional settings);
- and can be easily extended beyond second-order moments by using a nonlinear mapping, for example the hidden-layer activations of a pre-trained neural network can be used to create a task-relevant mapping.

2 Methodology

We consider the set of probability distributions over the domain \mathcal{X} and denote them $\mathcal{P}_\mathcal{X}$. Let $\mu, \nu \in \mathcal{P}_\mathcal{X}$ denote two probability measures, and let $X \sim \mu$ and $Y \sim \nu$ denote two random variables $X, Y \in \mathcal{X}$. We restrict our attention to probability distributions in the d -dimensional real-valued vector space $\mathcal{X} \subseteq \mathbb{R}^d$ with finite second moments and denote them as $\mathcal{P}_{L_2(\mathcal{X})}$. Denote by $\mathbf{m}_X = \mathbb{E}_{X \sim \mu}[X] \in \mathcal{X}$ and $\mathbf{m}_Y = \mathbb{E}_{Y \sim \nu}[Y] \in \mathcal{X}$ the first moments of the random variables. The family of probability mass functions over the discrete set $[n] = \{1, \dots, n\}$ is denoted $\mathcal{P}_{[n]}$. Two samples from μ and ν , $\{x_i = \mathbf{x}_i\}_{i=1}^m$ and $\{y_i = \mathbf{y}_i\}_{i=1}^n$, are denoted by their empirical measures $\hat{\mu}$

and $\hat{\nu}$. The max-sliced Bures divergence requires the matrices ρ_X and ρ_Y , which can be estimated from sample means $\mathbf{m}_X, \mathbf{m}_Y$ and covariance matrices Σ_X, Σ_Y : $\rho_X = \mathbf{m}_X \mathbf{m}_X^\top + \frac{m-1}{m} \Sigma_X$ and $\rho_Y = \mathbf{m}_Y \mathbf{m}_Y^\top + \frac{n-1}{n} \Sigma_Y$, assuming unbiased covariance estimates.

Let $\Omega = \{\omega(\cdot) = \langle \cdot, \mathbf{w} \rangle^2 : \mathbf{w} \in \mathcal{S}\}$, where \mathcal{S} denotes the unit hypersphere \mathbb{S}^{d-1} . Using $\mathbb{E}[\omega(X)] = \mathbb{E}[\langle X, \mathbf{w} \rangle^2] = \mathbf{w}^\top \mathbb{E}[X X^\top] \mathbf{w} = \mathbf{w}^\top \rho_X \mathbf{w}$, the max-sliced Bures divergence (MSB) is

$$D_{\text{MSB}}(\mu, \nu) = \sup_{\omega \in \Omega} |\sqrt{\mathbb{E}[\omega(X)]} - \sqrt{\mathbb{E}[\omega(Y)]}| = \sup_{\mathbf{w} \in \mathcal{S}} \left| \sqrt{\mathbf{w}^\top \rho_X \mathbf{w}} - \sqrt{\mathbf{w}^\top \rho_Y \mathbf{w}} \right| \quad (1)$$

$$= \max \left\{ \sqrt{\mathbb{E}[\omega_{\mu > \nu}(X)]} - \sqrt{\mathbb{E}[\omega_{\mu > \nu}(Y)]}, \sqrt{\mathbb{E}[\omega_{\mu < \nu}(Y)]} - \sqrt{\mathbb{E}[\omega_{\mu < \nu}(X)]} \right\}, \quad (2)$$

where $\omega_{\mu > \nu}(\cdot) = \langle \cdot, \mathbf{w}_{\mu > \nu} \rangle^2$ and $\omega_{\mu < \nu}(\cdot) = \langle \cdot, \mathbf{w}_{\mu < \nu} \rangle^2$ correspond to the optimal slices $\mathbf{w}_{\mu > \nu}$ and $\mathbf{w}_{\mu < \nu}$. In Algorithm 1, we describe how to compute the vector $\mathbf{w}_{\mu > \nu}$ defining the one-sided max sliced Bures divergence witness function $\omega_{\mu > \nu}$. It requires solving a one-dimensional bounded optimization, in which each function evaluation requires solving a maximum eigenvalue problems, which can be done efficiently using standard linear algebra packages. As an alternative, first-order

Algorithm 1: One-sided max-sliced Bures divergence

Input: $\rho_X = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top, \rho_Y = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^\top \in \mathbb{R}^{d \times d}, \epsilon > 0$

Define the function $\mathbf{v}(\cdot) : \mathbb{R} \rightarrow \mathbb{R}^d$ as $\mathbf{v}(\gamma) : \gamma \mapsto \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2 \leq 1} \mathbf{w}^\top (\gamma \rho_X - \rho_Y) \mathbf{w}$

Solve the 1D bound problem: $\gamma^* = \arg \max_{0 < \gamma \leq 1} \sqrt{\mathbf{v}(\gamma)^\top \rho_X \mathbf{v}(\gamma)} - \sqrt{\mathbf{v}(\gamma)^\top \rho_Y \mathbf{v}(\gamma)}$

Output: $\mathbf{w}_{\mu > \nu} = \mathbf{v}(\gamma^*)$

constrained algorithms can be applied after the objective of the maximization is made differentiable by smoothing the negative square root $\mathbf{w}_{\mu > \nu} \approx \arg \max_{\mathbf{w} \in \mathcal{S}} \sqrt{\mathbf{w}^\top \rho_X \mathbf{w}} - \sqrt{\mathbf{w}^\top \rho_Y \mathbf{w}} + \epsilon$.

In the Appendix, we formally state and show that the max-sliced Bures divergence is a lower bound for the max-sliced Wasserstein-2 distance, see Theorem 1. Following previous work on max-sliced Wasserstein distance by Deshpande et al. [2019] and Nguyen et al. [2021], we show that the statistical error of approximating the max-sliced Bures divergence from empirical distributions of size n is $\mathcal{O}(\sqrt{d \log(n)/n})$ (assuming a compact subset $\mathcal{X} \subset \mathbb{R}^d$). This is stated in Theorems 4 and 6.

3 Experiments

We present various examples of using the proposed max-sliced divergences to identify the discrepancies between two samples. For illustration, we use the MNIST, CIFAR10, and Office data sets, creating pseudo-real experiments by taking different subsets of real data sets. We also use the witness function to interpret samples from GANs trained for the stacked MNIST data set.

3.1 Detecting Non-Uniform Labels Distributions

We consider covariate shift on CIFAR10, created through purposeful class imbalance. $\hat{\mu}$ is a uniform sample from the training set and $\hat{\nu}$ is a sample from the test set with less instances from one class. Applying the max-sliced Bures witness function $\omega_{\mu > \nu}$ to points from $\hat{\nu}$ should assign larger values to points from the under-represented class, identifying the top-10 witness points with the largest witness function evaluations $\omega_{\mu > \nu}(y_{\hat{\nu}(1)}) \geq \dots \geq \omega_{\mu > \nu}(y_{\hat{\nu}(10)})$, where $\hat{\nu}$ is the required permutation.

For the vector representations of the samples, we use the internal representation of the Inception Network [Szegedy et al., 2016] as in the Fréchet Inception distance [Heusel et al., 2017] and Inception score [Salimans et al., 2016]. The witness function is $\omega_{\mu > \nu}(x) = \langle \mathbf{w}_{\mu > \nu}, \phi(x) \rangle^2$, where $\phi(x) \in \mathbb{R}^{2048}$ is the Inception code. The results are reported in Figure 1. When the class probabilities differ by 2% (10.2% for majority and 8.2% for minority), the mean average precision (MAP) is 0.94. This is compared to a MAP of 0.82 for the first-moment-based surrogate of the max-sliced W2 distance [Deshpande et al., 2019].

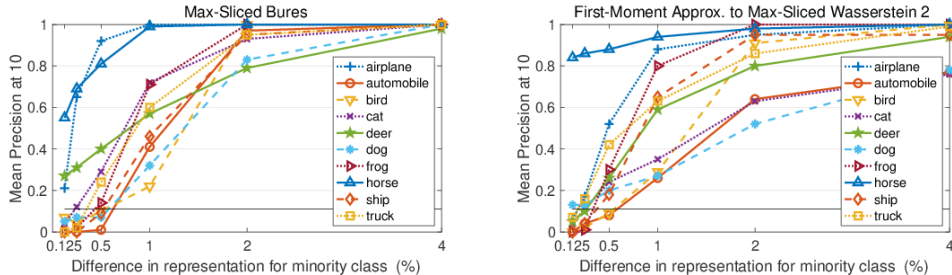


Figure 1: Max-sliced Bures divergence using the Inception Network representation are applied to samples with mismatched class distributions in CIFAR10. The first sample $\hat{\mu}$ consists of the training set (balanced classes with $m=50,000$), and the second sample $\hat{\nu}$ is an imbalanced subset of the test set with $n=10,000$, wherein the noted class is less prevalent. Each curve is the mean precision@10 (mean taken across 10 random draws) for test sets with the class noted in the legend being the minority class.

3.2 Detecting Mode Dropping

We use the one-sided max-sliced Bures divergence to monitor mode dropping during GAN training for the stacked MNIST data set. We create a 3-channel “Stacked MNIST” data set with 500,000 images from the MNIST training set to test mode dropping throughout GAN training (DCGAN architecture). The training set $\hat{\mu} = \text{Real}$ has 1000 possible modes corresponding to all 3-digit combinations. At the end of each training epoch (1000 iterations), a sample $\hat{\nu} = \text{Fake}$ of $n = 10^4$ synthetic images is generated. Figure 3 shows the top-5 fake and real images for both both witness functions after 10 epochs. To verify mode coverage we use a 4-layer convolutional network trained on single-channel MNIST; any missing combination of 3-digit labels is considered a dropped mode. Figure 2 shows the precision of the approach in detecting the dropped modes.

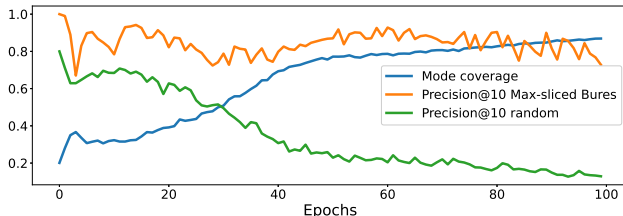


Figure 2: Detecting dropped modes using one-sided max-sliced Bures divergence. The slice $\mathbf{w}_{\text{Real} > \text{Fake}}$ is used to identify the top-10 *real* training images with the largest witness function values at each epoch. Precision@10 measures the fraction that correspond to dropped modes, as compared to random selection. The curves are the mean across 100 training trials of different GAN instances.

3.3 Interpreting the Fréchet Inception Distance

We use a witness function to identify mismatches in domain transfer tasks, specifically the Office dataset introduced by Saenko et al. [2010]. The data set contains three domains: images from Amazon.com (A), high-resolution DSLR images (D), and low-resolution webcam images (W). There are 31 classes and a varying number of images per class. The images are represented by internal activations of the Inception object classifying network [Szegedy et al., 2016]. We applied both of witness functions corresponding to the two one-sided max-sliced Bures divergences. Figure 4 shows the images identified by the witness functions. The results make sense as bicycles and notebooks with plain white backgrounds are much more common in A than in W. The closest examples are bicycles in front of white walls. Likewise binders with a wooden table and other background are common in webcam images but uncommon on Amazon.com.

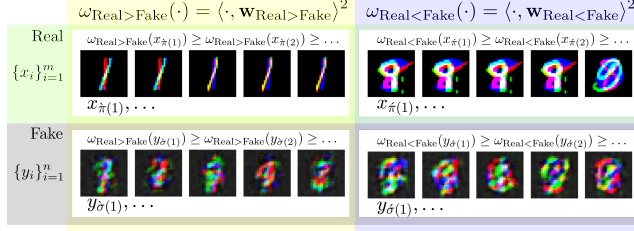


Figure 3: Example images identified by the one-sided max-sliced Bures divergence. The real images for the slice $\mathbf{w}_{\text{Real}>\text{Fake}}$ correspond to a dropped mode 1-1-1 (red-green-blue), as witnessed by the closest synthetic images which are not from the same mode. The fake images for slice $\mathbf{w}_{\text{Real}<\text{Fake}}$ highlight a subset of poorly formed synthetic images that are closest to the real mode 8-9-9.

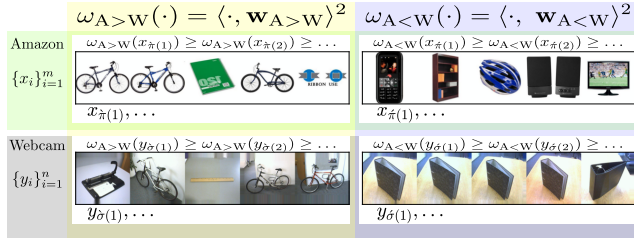


Figure 4: Max-slicing Inception codes to illustrate discrepancies in domain transfer: images from Amazon.com (A) compared to webcam images (W).

4 Limitations

While the proposed approach yields an interpretable and precise witness function, the utility of a localized witness function that assesses the differences in second moments along a 1-dimensional subspace is limited. It is ill-suited to be used in learning, as it focuses only on the worst discrepancy—ignoring the other discrepancies. Furthermore, the reliance on second moments requires a suitable non-linear learning representation, such as a pre-trained neural network or random Fourier features. We note that previous work has shown that MMD operating on top of an untrained convolutional encoding network can detect covariate shift [Rabanser et al., 2019]. Since MMD with a Gaussian kernel can be approximated using the norm of the difference of means after random Fourier feature embedding, we are hopeful that the ability to pinpoint the worst discrepancies will prove useful in practice.

On the practical side, further testing is needed to see how the method performs on high-dimensional tabular data. A specific example is to consider sparse bag-of-words representations of corpus of documents. In this case, the sparsity of the representation may be exploited when computing the largest eigenvector of the difference of the Gram matrices ρ_X and ρ_Y .

On the theoretical side, while the max-sliced Bures divergence does not suffer from the curse of dimensionality, it is only a probability metric amongst distributions completely characterized by their first and second moments. Further analysis is needed for the convergence rate of the proposed algorithm and its dependence on the eigenstructure of the Gram matrices.

5 Conclusion

We propose the max-sliced Bures divergence, a lower-bound on the max-sliced W2 distance. The underlying one-sided max-sliced Bures divergences enable direct interpretation of under- and over-representation between two samples via the witness functions, which can be used to identify discrepancies. The optimal witness functions can be found through optimization problems posed on the unit interval, as opposed to the saddle point optimization problem involved in the max-sliced Wasserstein distance. We show the potential of applying this approach to covariate shift detection.

Acknowledgments and Disclosure of Funding

Luis G. Sanchez Giraldo would like to acknowledge the support from the NVIDIA Academic Hardware Grant Program. The authors gratefully acknowledge the support of Dr. Tory Cobb and ONR321. Austin J. Brockmeier’s effort were sponsored by the Department of the Navy, Office of Naval Research under ONR award number N00014-21-1-2300. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the Office of Naval Research.

References

- Syed Mumtaz Ali and Samuel D Silvey. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society: Series B (Methodological)*, 28(1):131–142, 1966.
- Rajendra Bhatia, Tanvi Jain, and Yongdo Lim. On the Bures–Wasserstein distance between positive definite matrices. *Expositiones Mathematicae*, 37(2):165–191, 2019.
- Nicolas Bonneel, Julien Rabin, Gabriel Peyré, and Hanspeter Pfister. Sliced and Radon Wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51(1):22–45, 2015.
- Hannes De Meulemeester, Joachim Schreurs, Michaël Fanuel, Bart De Moor, and Johan AK Suykens. The Bures metric for taming mode collapse in generative adversarial networks. *arXiv preprint arXiv:2006.09096*, 2020.
- Ishan Deshpande, Ziyu Zhang, and Alexander G Schwing. Generative modeling using the sliced Wasserstein distance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3483–3491, 2018.
- Ishan Deshpande, Yuan-Ting Hu, Ruoyu Sun, Ayis Pyrros, Nasir Siddiqui, Sanmi Koyejo, Zhizhen Zhao, David Forsyth, and Alexander G Schwing. Max-sliced Wasserstein distance and its use for GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10648–10656, 2019.
- D. C. Dowson and B. V. Landau. The Fréchet distance between multivariate normal distributions. *Journal of Multivariate Analysis*, 12(3):450–455, 1982.
- Maurice Fréchet. Sur la distance de deux lois de probabilité. *C. R. Math. Acad. Sci. Paris*, 244(6):689–692, 1957.
- Christopher Fuchs and Jeroen Van De Graaf. Cryptographic distinguishability measures for quantum-mechanical states. *Information Theory, IEEE Transactions on*, 45(4):1216–1227, 1999.
- Matthias Gelbrich. On a formula for the L2 Wasserstein metric between measures on Euclidean and Hilbert spaces. *Mathematische Nachrichten*, 147(1):185–203, 1990.
- Arthur Gretton, Karsten M. Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alexander J. Smola. A kernel method for the two-sample-problem. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages 513–520. MIT Press, Cambridge, MA, 2007.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *Advances in Neural Information Processing Systems*, pages 6626–6637, 2017.
- Hicham Janati, Boris Muzellec, Gabriel Peyré, and Marco Cuturi. Entropic optimal transport between unbalanced Gaussian measures has a closed form. *Advances in Neural Information Processing Systems*, 33, 2020.
- Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. A least-squares approach to direct importance estimation. *The Journal of Machine Learning Research*, 10:1391–1445, 2009.
- Soheil Kolouri, Gustavo K Rohde, and Heiko Hoffmann. Sliced Wasserstein distance for learning Gaussian mixture models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3427–3436, 2018.
- Soheil Kolouri, Kimia Nadjahi, Umut Simsekli, Roland Badeau, and Gustavo Rohde. Generalized sliced Wasserstein distances. In *Advances in Neural Information Processing Systems*, pages 261–272, 2019.
- Solomon Kullback and Richard A Leibler. On information and sufficiency. *The Annals of Mathematical Statistics*, 22(1):79–86, 1951.

- Chen-Yu Lee, Tanmay Batra, Mohammad Haris Baig, and Daniel Ulbricht. Sliced Wasserstein discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10285–10295, 2019.
- Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift with black box predictors. In *International Conference on Machine Learning*, pages 3122–3130, 2018.
- Youssef Mroueh, Tom Sercu, and Vaibhava Goel. McGAN: Mean and covariance feature matching GAN. In *International Conference on Machine Learning*, pages 2527–2535, 2017.
- Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, pages 429–443, 1997.
- Boris Muzellec and Marco Cuturi. Generalizing point embeddings using the Wasserstein space of elliptical distributions. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 10237–10248. Curran Associates, Inc., 2018.
- Khai Nguyen, Nhat Ho, Tung Pham, and Hung Bui. Distributional sliced-Wasserstein and applications to generative modeling. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=QYj070ACDK>.
- XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals and the likelihood ratio by penalized convex risk minimization. In *Advances in Neural Information Processing Systems*, pages 1089–1096, 2008.
- XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-GAN: Training generative neural samplers using variational divergence minimization. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 271–279. Curran Associates, Inc., 2016.
- Jung Hun Oh, Maryam Pouryahya, Aditi Iyer, Aditya P Apte, Joseph O Deasy, and Allen Tannenbaum. A novel kernel Wasserstein distance on Gaussian measures: An application of identifying dental artifacts in head and neck computed tomography. *Computers in Biology and Medicine*, page 103731, 2020.
- Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset Shift in Machine Learning*. The MIT Press, 2009.
- Stephan Rabanser, Stephan Günnemann, and Zachary Lipton. Failing loudly: An empirical study of methods for detecting dataset shift. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Alfréd Rényi. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*. The Regents of the University of California, 1961.
- Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *European Conference on Computer Vision*, pages 213–226. Springer, 2010.
- Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3723–3732, 2018.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In *Advances in Neural Information Processing Systems*, pages 2234–2242, 2016.
- Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkhäuser, NY*, 55(58-63):94, 2015.
- Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90(2):227–244, 2000.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016.

Vladimir Vapnik. *The Nature of Statistical Learning Theory*. Springer Science & Business Media, 2013.

Jiqing Wu, Zhiwu Huang, Dinesh Acharya, Wen Li, Janine Thoma, Danda Pani Paudel, and Luc Van Gool. Sliced Wasserstein generative models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3713–3722, 2019.

Makoto Yamada, Taiji Suzuki, Takafumi Kanamori, Hirotaka Hachiya, and Masashi Sugiyama. Relative density-ratio estimation for robust distribution comparison. In *Advances in Neural Information Processing Systems*, pages 594–602, 2011.

Makoto Yamada, Taiji Suzuki, Takafumi Kanamori, Hirotaka Hachiya, and Masashi Sugiyama. Relative density-ratio estimation for robust distribution comparison. *Neural Computation*, 25(5):1324–1370, 2013.

Zhen Zhang, Mianzhi Wang, and Arye Nehorai. Optimal transport in reproducing kernel Hilbert spaces: theory and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(7):1741–1754, 2020.

A Appendix

These appendices detail the derivation of the proposed divergences; formal results relating them to existing divergences; and additional experimental results. First, we detail the relationship between Wasserstein-2 distance the Fréchet distance/divergence and the Bures distance/divergence. Then we detail statistical guarantees on the convergence of the estimator. Finally, we detail additional synthetic results.

A.1 Bures, Fréchet, and Wasserstein Distances

The Bures distance generalizes the Hellinger distance to positive semidefinite matrices [Fuchs and Van De Graaf, 1999, Bhatia et al., 2019]. For a pair of positive semidefinite matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$ the Bures distance is defined as

$$d_{\text{B}}(\mathbf{A}, \mathbf{B}) \triangleq \sqrt{\|\mathbf{A}\|_1 + \|\mathbf{B}\|_1 - 2 \left\| \mathbf{A}^{\frac{1}{2}} \mathbf{B}^{\frac{1}{2}} \right\|_1}, \quad (3)$$

where $\|\cdot\|_1$ denotes the trace norm (sum of the singular values), $\|\mathbf{A}^{\frac{1}{2}} \mathbf{B}^{\frac{1}{2}}\|_1 = \text{tr}(\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}} = \text{tr}(\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{\frac{1}{2}}$. $d_{\text{B}}(\mathbf{A}, \mathbf{A}) \leq \|\mathbf{A}^{\frac{1}{2}} - \mathbf{B}^{\frac{1}{2}}\|_2^2 = \text{tr}(\mathbf{A}^{\frac{1}{2}} - \mathbf{B}^{\frac{1}{2}})^2$, with equality if the matrices commute. When the matrices commute they share the same eigenbasis and the Bures distance is proportional to the Hellinger distance between the eigenvalues. While the Bures distance is a metric between positive semidefinite matrices, the corresponding divergence,

$$D_{\text{B}}(\mu, \nu) \triangleq d_{\text{B}}\left(\underbrace{\mathbb{E}_{X \sim \mu}[X X^{\top}]}_{\rho_X = \mathbf{m}_X \mathbf{m}_X^{\top} + \Sigma_X}, \underbrace{\mathbb{E}_{Y \sim \nu}[Y Y^{\top}]}_{\rho_Y = \mathbf{m}_Y \mathbf{m}_Y^{\top} + \Sigma_Y}\right), \quad (4)$$

is a semi-metric since it relies on the uncentered second moments. Thus, it is a probability metric only on the subset of distributions defined completely by their first or second moments, such as Gaussian distributions or other elliptically symmetric distributions.

To compare two Gaussian distributions, it is natural to compare the means and covariance parameters as in the Fréchet distance [Fréchet, 1957, Dowson and Landau, 1982],

$$D_{\text{F}}(\mu, \nu) \triangleq \sqrt{\|\mathbf{m}_X - \mathbf{m}_Y\|_2^2 + d_{\text{B}}^2(\Sigma_X, \Sigma_Y)}. \quad (5)$$

The Fréchet distance is also known as the Wasserstein-Bures [Janati et al., 2020] or Bures-Wasserstein distance [Bhatia et al., 2019], since it is equal to the Wasserstein-2 distance (W_2) between Gaussian measures, and generally lower bounds the W_2 distance [Gelbrich, 1990],

$$D_{\text{F}}(\mu, \nu) = D_{\text{F}}(\mathcal{N}(\mathbf{m}_X, \Sigma_X), \mathcal{N}(\mathbf{m}_Y, \Sigma_Y)) = W_2(\mathcal{N}(\mathbf{m}_X, \Sigma_X), \mathcal{N}(\mathbf{m}_Y, \Sigma_Y)) \leq W_2(\mu, \nu), \quad (6)$$

where

$$W_2(\mu, \nu) \triangleq \sqrt{\inf_{\gamma \in \Gamma(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \gamma} \|X - Y\|_2^2}, \quad (7)$$

and $\Gamma(\mu, \nu)$ defines the set of all joint distributions with marginals μ and ν .

The sliced Wasserstein distance [Wu et al., 2019, Deshpande et al., 2018, Kolouri et al., 2018] and the max-sliced Wasserstein distance [Deshpande et al., 2019], and its generalizations [Kolouri et al., 2019], evaluate discrepancies along one-dimensional linear, or non-linear, subspaces. They can be expressed by the Radon transform of integrable functions [Bonneel et al., 2015]. The linear max-sliced Wasserstein-2 distance (squared) is

$$D_{\text{MSW}_2}^2(\mu, \nu) = \sup_{\mathbf{w} \in \mathcal{S}} \inf_{\gamma \in \Gamma(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \gamma} [\langle X - Y, \mathbf{w} \rangle^2], \quad (8)$$

where \mathcal{S} can be either the unit hypersphere \mathbb{S}^{d-1} or its convex hull $\{\mathbf{w} \in \mathbb{R}^d : \|\mathbf{w}\|_2 \leq 1\}$ (in terms of the optimization problem the two are equivalent as the optimal solution necessarily lies on the unit hypersphere). Slicing is motivated by the relative ease of computing the Wasserstein- p distance in one dimension, since it has a closed form [Santambrogio, 2015, Lemma 2.8] and in the case of samples only requires a sorting procedure.

While the Bures divergence already has a closed-form solution, the Bures divergence is not easily interpretable as it cannot be written in terms of witness function. We propose a sliced version based on projecting a positive semidefinite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ to a rank-1 matrix $\mathbf{\Omega}_{\mathbf{w}} = \mathbf{w}\mathbf{w}^\top$. For compactness we denote the rank-1 projection as

$$\mathbf{A}_{\mathbf{w}} = \langle \mathbf{A}, \mathbf{\Omega}_{\mathbf{w}} \rangle \mathbf{\Omega}_{\mathbf{w}} = \langle \mathbf{A}, \mathbf{w}\mathbf{w}^\top \rangle (\mathbf{w}\mathbf{w}^\top). \quad (9)$$

When $\mathbf{w} \in \mathbb{S}^{d-1}$, $\|\mathbf{w}\|_2 = \text{tr}(\mathbf{w}\mathbf{w}^\top) = 1$, $\text{tr}(\mathbf{A}_{\mathbf{w}}) = \langle \mathbf{A}, \mathbf{w}\mathbf{w}^\top \rangle = \mathbf{w}^\top \mathbf{A} \mathbf{w}$, and

$$\begin{aligned} d_{\text{B}}(\mathbf{A}_{\mathbf{w}}, \mathbf{B}_{\mathbf{w}}) &= \sqrt{\mathbf{w}^\top (\mathbf{A} + \mathbf{B}) \mathbf{w} - 2\sqrt{\mathbf{w}^\top \mathbf{A} \mathbf{w}} \sqrt{\mathbf{w}^\top \mathbf{B} \mathbf{w}}} \\ &= |\sqrt{\mathbf{w}^\top \mathbf{A} \mathbf{w}} - \sqrt{\mathbf{w}^\top \mathbf{B} \mathbf{w}}|. \end{aligned} \quad (10)$$

Let $\rho_{X, \mathbf{w}} = \langle \mathbf{\Omega}_{\mathbf{w}}, \rho_X \rangle \mathbf{\Omega}_{\mathbf{w}}$ and $\rho_{Y, \mathbf{w}} = \langle \mathbf{\Omega}_{\mathbf{w}}, \rho_Y \rangle \mathbf{\Omega}_{\mathbf{w}}$. The max-sliced Bures divergence (MSB) is

$$D_{\text{MSB}}(\mu, \nu) \triangleq \sup_{\mathbf{w} \in \mathcal{S}} d_{\text{B}}(\rho_{X, \mathbf{w}}, \rho_{Y, \mathbf{w}}) \quad (12)$$

$$= \sup_{\mathbf{w} \in \mathcal{S}} \left| \sqrt{\mathbf{w}^\top \rho_X \mathbf{w}} - \sqrt{\mathbf{w}^\top \rho_Y \mathbf{w}} \right| \quad (13)$$

$$= \sup_{\mathbf{w} \in \mathcal{S}} \left| \sqrt{\mathbb{E}[\langle X, \mathbf{w} \rangle^2]} - \sqrt{\mathbb{E}[\langle Y, \mathbf{w} \rangle^2]} \right| \quad (14)$$

The max-sliced Bures divergence is a lower-bound on the max-sliced Wasserstein-2 distance, since the latter is lower bounded by the max-sliced Fréchet distance. This is stated in the following theorem.

Theorem 1. For $\mu, \nu \in \mathcal{P}_{L_2(\mathcal{X})}$, the max-sliced Bures divergence is a lower bound on the max-sliced Fréchet distance and the max-sliced W2 distance, $D_{\text{MSB}}(\mu, \nu) \leq D_{\text{MSF}}(\mu, \nu) \leq D_{\text{MSW}_2}(\mu, \nu)$.

To prove this theorem, we will develop some definitions and intermediate results. Recall that $\mathcal{P}_{L_2(\mathcal{X})}$ denotes the set of probability distributions in the d -dimensional real-valued vector space $\mathcal{X} \subseteq \mathbb{R}^d$, with finite second moments. The means of the random variables are denoted $\mathbf{m}_X = \mathbb{E}_{X \sim \mu}[X] \in \mathcal{X}$ and $\mathbf{m}_Y = \mathbb{E}_{Y \sim \nu}[Y] \in \mathcal{X}$. The uncentered second moments are denoted $\rho_X = \mathbb{E}[X X^\top] \in \mathcal{X} \times \mathcal{X}$ and $\rho_Y = \mathbb{E}[Y Y^\top] \in \mathcal{X} \times \mathcal{X}$. The covariances are denoted $\Sigma_X = \rho_X - \mathbf{m}_X \mathbf{m}_X^\top$ and $\Sigma_Y = \rho_Y - \mathbf{m}_Y \mathbf{m}_Y^\top$.

Let $\mathbf{w} \in \mathbb{R}^d$ denote a slicing vector. Denote the sliced means and rank-1 covariances of the random variables $\langle \mathbf{w}, X \rangle \mathbf{w}$ and $\langle \mathbf{w}, Y \rangle \mathbf{w}$ as $\mathbf{m}_{X, \mathbf{w}} = \langle \mathbf{w}, \mathbf{m}_X \rangle \mathbf{w} \in \mathbb{R}^d$ and $\mathbf{m}_{Y, \mathbf{w}} = \langle \mathbf{w}, \mathbf{m}_Y \rangle \mathbf{w} \in \mathbb{R}^d$, and $\Sigma_{X, \mathbf{w}} = \langle \mathbf{w}, \Sigma_X \mathbf{w} \rangle (\mathbf{w}\mathbf{w}^\top) \in \mathbb{R}^{d \times d}$ and $\Sigma_{Y, \mathbf{w}} = \langle \mathbf{w}, \Sigma_Y \mathbf{w} \rangle (\mathbf{w}\mathbf{w}^\top) \in \mathbb{R}^{d \times d}$, respectively.

The Fréchet difference between the sliced random variables is

$$\begin{aligned}
D_F(\mathcal{N}(\mathbf{m}_{X,\mathbf{w}}, \boldsymbol{\Sigma}_{X,\mathbf{w}}), \mathcal{N}(\mathbf{m}_{Y,\mathbf{w}}, \boldsymbol{\Sigma}_{Y,\mathbf{w}})) &= \sqrt{\|\mathbf{m}_{X,\mathbf{w}} - \mathbf{m}_{Y,\mathbf{w}}\|_2^2 + d_B^2(\boldsymbol{\Sigma}_{X,\mathbf{w}}, \boldsymbol{\Sigma}_{Y,\mathbf{w}})} \\
&= \sqrt{\|(\langle \mathbf{w}, \mathbf{m}_X \rangle - \langle \mathbf{w}, \mathbf{m}_Y \rangle) \mathbf{w}\|_2^2 + \|\boldsymbol{\Sigma}_{X,\mathbf{w}}^{\frac{1}{2}} - \boldsymbol{\Sigma}_{Y,\mathbf{w}}^{\frac{1}{2}}\|_F^2} \\
&= \sqrt{\|(\langle \mathbf{w}, \mathbf{m}_X \rangle - \langle \mathbf{w}, \mathbf{m}_Y \rangle) \mathbf{w}\|_2^2 + \|\sqrt{\langle \mathbf{w}, \boldsymbol{\Sigma}_X \mathbf{w} \rangle} (\mathbf{w} \mathbf{w}^\top)^{\frac{1}{2}} - \sqrt{\langle \mathbf{w}, \boldsymbol{\Sigma}_Y \mathbf{w} \rangle} (\mathbf{w} \mathbf{w}^\top)^{\frac{1}{2}}\|_F^2} \\
&= \sqrt{\|\mathbf{w}\|_2^2 \langle \mathbf{w}, \mathbf{m}_X - \mathbf{m}_Y \rangle^2 + \|\mathbf{w}\|_2^2 \left(\sqrt{\langle \mathbf{w}, \boldsymbol{\Sigma}_X \mathbf{w} \rangle} - \sqrt{\langle \mathbf{w}, \boldsymbol{\Sigma}_Y \mathbf{w} \rangle} \right)^2}.
\end{aligned}$$

The second equality follows from the fact that trace norm $\|\cdot\|_1$ is equal to the Frobenius norm $\|\cdot\|_F$ for rank-1 matrices. The last equality uses the fact that $\|(\mathbf{w} \mathbf{w}^\top)^{\frac{1}{2}}\|_F^2 = \text{tr} \left((\mathbf{w} \mathbf{w}^\top)^{\frac{1}{2}} (\mathbf{w} \mathbf{w}^\top)^{\frac{1}{2}} \right) = \text{tr}(\mathbf{w} \mathbf{w}^\top) = \text{tr}(\mathbf{w}^\top \mathbf{w}) = \|\mathbf{w}\|_2^2$.

We now state the max-sliced Fréchet distance and an upper bound.

$$\begin{aligned}
D_{\text{MSF}}(\mu, \nu) &= \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} D_F(\mathcal{N}(\mathbf{m}_{X,\mathbf{w}}, \boldsymbol{\Sigma}_{X,\mathbf{w}}), \mathcal{N}(\mathbf{m}_{Y,\mathbf{w}}, \boldsymbol{\Sigma}_{Y,\mathbf{w}})) \\
&= \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \underbrace{\|\mathbf{w}\|_2}_1 \sqrt{\langle \mathbf{w}, \mathbf{m}_X - \mathbf{m}_Y \rangle^2 + \left(\sqrt{\langle \mathbf{w}, \boldsymbol{\Sigma}_X \mathbf{w} \rangle} - \sqrt{\langle \mathbf{w}, \boldsymbol{\Sigma}_Y \mathbf{w} \rangle} \right)^2} \\
&= \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \sqrt{\langle \mathbf{w}, \mathbf{m}_X - \mathbf{m}_Y \rangle^2 + (\sigma_{X,\mathbf{w}} - \sigma_{Y,\mathbf{w}})^2}, \tag{15}
\end{aligned}$$

where $\sigma_{X,\mathbf{w}} = \sqrt{\mathbb{E}[(\langle \mathbf{w}, X \rangle - \langle \mathbf{w}, \mathbf{m}_X \rangle)^2]}$ is the standard deviation of $\langle \mathbf{w}, X \rangle$, and, similarly, $\sigma_{Y,\mathbf{w}} = \sqrt{\mathbb{E}[(\langle \mathbf{w}, Y \rangle - \langle \mathbf{w}, \mathbf{m}_Y \rangle)^2]}$.

The following upper bound uses a separate slice for the means and covariances:

$$\begin{aligned}
D_{\text{MSF}}^U(\mu, \nu) &= \left(\sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \langle \mathbf{w}, \mathbf{m}_X - \mathbf{m}_Y \rangle^2 + \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} d_B^2(\boldsymbol{\Sigma}_{X,\mathbf{w}}, \boldsymbol{\Sigma}_{Y,\mathbf{w}}) \right)^{\frac{1}{2}} \\
&= \sqrt{\|\mathbf{m}_X - \mathbf{m}_Y\|_2^2 + \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} (\sigma_{X,\mathbf{w}} - \sigma_{Y,\mathbf{w}})^2}. \tag{16}
\end{aligned}$$

Clearly, $D_{\text{MSF}}(\mu, \nu) \leq D_{\text{MSF}}^U(\mu, \nu)$.

We now consider the relationship between the Gauss-Wasserstein or Fréchet distance—which combines the distances between the first and second-order moments—with the max-sliced Bures when it is applied directly to the uncentered covariance matrices. For this we need the following lemma that states the reverse triangle inequality in \mathbb{R}^2 .

Lemma 2 (Reverse triangle inequality). *For two vectors in $\mathbf{u}, \mathbf{v} \in \mathbb{R}^2$, the difference between their Euclidean norms is less than or equal to the Euclidean norm of their differences $|\|\mathbf{u}\| - \|\mathbf{v}\|| \leq \|\mathbf{u} - \mathbf{v}\|$. For $\mathbf{u} = [a, b]^\top$, $\mathbf{v} = [c, d]^\top$, $|\sqrt{a^2 + b^2} - \sqrt{c^2 + d^2}| \leq \sqrt{(a - c)^2 + (b - d)^2}$.*

Theorem 3 (Sliced Fréchet and sliced Bures divergence inequality). *The max-sliced Bures divergence based on the uncentered covariance matrices is less than or equal to the max-sliced Fréchet distance, $D_{\text{MSB}}(\mu, \nu) \leq D_{\text{MSF}}(\mu, \nu)$.*

Proof. To relate the sliced Bures divergence and sliced Fréchet distances, we note that $\mathbb{E}_{X \sim \mu}[\langle X, \mathbf{w} \rangle^2] = \langle \mathbf{m}_X, \mathbf{w} \rangle^2 + \sigma_{X,\mathbf{w}}^2$ and, likewise, $\mathbb{E}_{Y \sim \nu}[\langle Y, \mathbf{w} \rangle^2] = \langle \mathbf{m}_Y, \mathbf{w} \rangle^2 + \sigma_{Y,\mathbf{w}}^2$. Then,

$$\begin{aligned}
\left| \sqrt{\mathbb{E}_{X \sim \mu}[\langle X, \mathbf{w} \rangle^2]} - \sqrt{\mathbb{E}_{Y \sim \nu}[\langle Y, \mathbf{w} \rangle^2]} \right| &= \left| \sqrt{\langle \mathbf{m}_X, \mathbf{w} \rangle^2 + \sigma_{X,\mathbf{w}}^2} - \sqrt{\langle \mathbf{m}_Y, \mathbf{w} \rangle^2 + \sigma_{Y,\mathbf{w}}^2} \right| \\
&\leq \sqrt{(\langle \mathbf{m}_X, \mathbf{w} \rangle - \langle \mathbf{m}_Y, \mathbf{w} \rangle)^2 + (\sigma_{X,\mathbf{w}} - \sigma_{Y,\mathbf{w}})^2},
\end{aligned}$$

where the inequality relies on the reverse triangle inequality (Lemma 2), with $\mathbf{u} = [a, b]^\top = [\langle \mathbf{m}_X, \mathbf{w} \rangle, \sigma_{X, \mathbf{w}}]^\top$, $\mathbf{v} = [c, d]^\top = [\langle \mathbf{m}_Y, \mathbf{w} \rangle, \sigma_{Y, \mathbf{w}}]^\top$. Taking the supremum over slices of both sides yields the desired inequality,

$$\sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \underbrace{\left| \sqrt{\mathbb{E}_{X \sim \mu}[\langle X, \mathbf{w} \rangle^2]} - \sqrt{\mathbb{E}_{Y \sim \nu}[\langle Y, \mathbf{w} \rangle^2]} \right|}_{D_{\text{MSB}}(\mu, \nu)} \leq \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \underbrace{\sqrt{(\langle \mathbf{m}_X, \mathbf{w} \rangle - \langle \mathbf{m}_Y, \mathbf{w} \rangle)^2 + (\sigma_{X, \mathbf{w}} - \sigma_{Y, \mathbf{w}})^2}}_{D_{\text{MSF}}(\mu, \nu)}.$$

□

Using these results we can now state the proof of Theorem 1.

Proof of Theorem 1. For $\mathbf{w} \in \mathbb{R}^d$, let $w(\cdot) = \langle \cdot, \mathbf{w} \rangle$. Let $X \sim \mu$ and $Y \sim \nu$, then define the zero-mean random variables $\tilde{w}(X) = \langle X - \mathbf{m}_X, \mathbf{w} \rangle$ and $\tilde{w}(Y) = \langle Y - \mathbf{m}_Y, \mathbf{w} \rangle$. By definition, $\mathbb{E}[\tilde{w}^2(X)] = \sigma_{X, \mathbf{w}}^2$ and $\mathbb{E}[\tilde{w}^2(Y)] = \sigma_{Y, \mathbf{w}}^2$. The squared Fréchet distance between random variables $w(X) = \langle \mathbf{m}_X, \mathbf{w} \rangle + \tilde{w}(X)$ and $w(Y) = \langle \mathbf{m}_Y, \mathbf{w} \rangle + \tilde{w}(Y)$ is

$$\begin{aligned} D_{\text{F}}^2(w_{\#}\mu, w_{\#}\nu) &= (\mathbb{E}[w(X)] - \mathbb{E}[w(Y)])^2 + (\sigma_{X, \mathbf{w}} - \sigma_{Y, \mathbf{w}})^2 = \langle \mathbf{m}_X - \mathbf{m}_Y, \mathbf{w} \rangle^2 + (\sigma_{X, \mathbf{w}} - \sigma_{Y, \mathbf{w}})^2 \\ &= \langle \mathbf{m}_X, \mathbf{w} \rangle^2 - 2\langle \mathbf{m}_X, \mathbf{w} \rangle \langle \mathbf{m}_Y, \mathbf{w} \rangle + \langle \mathbf{m}_Y, \mathbf{w} \rangle^2 + \sigma_{X, \mathbf{w}}^2 - 2\sigma_{X, \mathbf{w}}\sigma_{Y, \mathbf{w}} + \sigma_{Y, \mathbf{w}}^2 \\ &= \underbrace{\langle \mathbf{m}_X, \mathbf{w} \rangle^2 + \sigma_{X, \mathbf{w}}^2}_{\mathbb{E}[w^2(X)]} + \underbrace{\langle \mathbf{m}_Y, \mathbf{w} \rangle^2 + \sigma_{Y, \mathbf{w}}^2}_{\mathbb{E}[w^2(Y)]} - 2 \left(\underbrace{\langle \mathbf{m}_X, \mathbf{w} \rangle \langle \mathbf{m}_Y, \mathbf{w} \rangle}_{\mathbb{E}[w(X)]\mathbb{E}[w(Y)]} + \underbrace{\sigma_{X, \mathbf{w}}\sigma_{Y, \mathbf{w}}}_{\sqrt{\mathbb{E}[\tilde{w}^2(X)]\mathbb{E}[\tilde{w}^2(Y)]}} \right) \\ &= \mathbb{E}[w^2(X)] + \mathbb{E}[w^2(Y)] - 2 \left(\mathbb{E}[w(X)]\mathbb{E}[w(Y)] + \sqrt{\mathbb{E}[\tilde{w}^2(X)]\mathbb{E}[\tilde{w}^2(Y)]} \right) \\ &\leq \mathbb{E}[w^2(X)] + \mathbb{E}[w^2(Y)] - 2(\mathbb{E}[w(X)]\mathbb{E}[w(Y)] + \mathbb{E}[\tilde{w}(X)\tilde{w}(Y)]) \\ &= \mathbb{E}[w^2(X)] + \mathbb{E}[w^2(Y)] - 2\mathbb{E}[w(X)w(Y)] = \mathbb{E}[(w(X) - w(Y))^2], \end{aligned}$$

which follows from Hölder's inequality,

$$\mathbb{E}[\tilde{w}(X)\tilde{w}(Y)] \leq \sqrt{\mathbb{E}[\tilde{w}^2(X)]\mathbb{E}[\tilde{w}^2(Y)]} \implies -2\sqrt{\mathbb{E}[\tilde{w}^2(X)]\mathbb{E}[\tilde{w}^2(Y)]} \leq -2\mathbb{E}[\tilde{w}(X)\tilde{w}(Y)] = 0.$$

As this holds for any joint distribution, taking the infimum over all possible joint distributions $(X, Y) \sim \gamma$ that are within the coupling distribution $\gamma \in \Gamma(\mu, \nu)$, yields the Wasserstein-2 distance (squared) between the sliced distributions

$$D_{\text{F}}^2(w_{\#}\mu, w_{\#}\nu) \leq \inf_{\gamma \in \Gamma(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \gamma} [(w(X) - w(Y))^2] = D_{\text{W}_2}^2(w_{\#}\mu, w_{\#}\nu).$$

Maximizing over possible slices, yields

$$D_{\text{MSF}}^2(\mu, \nu) = \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} D_{\text{F}}^2(w_{\#}\mu, w_{\#}\nu) \leq \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} D_{\text{W}_2}^2(w_{\#}\mu, w_{\#}\nu) = D_{\text{MSW}_2}^2(\mu, \nu).$$

Combining with Theorem 3 yields the desired result. □

Theorem 4. Given a probability measure $\mu \in \mathcal{P}_{L_2(\mathcal{X})}$ supported on a compact subset $\mathcal{X} \subset \mathbb{R}^d$. Assume that $\{X_i\}_{i=1}^n$ are i.i.d. random variables with distribution μ . For the empirical measure $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$,

$$\mathbb{E}[D_{\text{MSB}}(\mu, \mu_n)] \leq \mathbb{E}[D_{\text{MSW}_2}(\mu, \mu_n)] \leq c \sqrt{\frac{d \log n}{n}}, \quad (17)$$

where $c > 0$ is some universal constant.

Proof. First we note that for all $\mu, \nu \in \mathcal{P}_{L_2(\mathcal{X})}$, $D_{\text{MSW}_2}(\mu, \nu) \leq D_{\text{MSW}_1}(\mu, \nu)$ due to the inequality between the L_1 -norm and the L_2 -norm. Then Nguyen et al. [2021] show that

$$D_{\text{MSW}_1}(\mu, \mu_n) \leq \text{diam}(\mathcal{X}) \sqrt{\frac{32}{n} [(d+1) \log(n+1) + \log(8/\delta)]}, \quad (18)$$

with probability at least $1 - \delta$, for any $\delta \in (0, 1)$, which yields $\mathbb{E}[D_{\text{MSW}_1}(\mu, \mu_n)] \leq c \sqrt{\frac{d \log n}{n}}$ for some universal constant $c > 0$. From Theorem 1, we have that $\mathbb{E}[D_{\text{MSB}}(\mu, \mu_n)] \leq \mathbb{E}[D_{\text{MSW}_2}(\mu, \mu_n)]$. Combing these two results using $\mathbb{E}[D_{\text{MSW}_2}(\mu, \nu)] \leq \mathbb{E}[D_{\text{MSW}_1}(\mu, \nu)]$ yields the desired result. □

Lemma 5. *The max-sliced Bures divergence obeys the triangle inequality, $D_{MSB}(\mu, \nu) \leq D_{MSB}(\mu, \xi) + D_{MSB}(\xi, \nu)$ for any $\mu, \nu, \xi \in \mathcal{P}_{L_2}(\mathcal{X})$.*

Proof. Using $|a - c| \leq |a - b| + |b - c|$, it is trivial to note that the sliced Bures distance obeys the triangle inequality $d_B(\mathbf{A}_w, \mathbf{C}_w) = |\sqrt{w^\top \mathbf{A} w} - \sqrt{w^\top \mathbf{C} w}| \leq |\sqrt{w^\top \mathbf{A} w} - \sqrt{w^\top \mathbf{B} w}| + |\sqrt{w^\top \mathbf{B} w} - \sqrt{w^\top \mathbf{C} w}| = d_B(\mathbf{A}_w, \mathbf{B}_w) + d_B(\mathbf{B}_w, \mathbf{C}_w)$ for $w \in \mathbb{S}^{d-1}$.

Let $\mathbf{A} = \rho_X = \mathbb{E}_{X \sim \mu}[XX^\top]$, $\mathbf{B} = \rho_Y = \mathbb{E}_{Y \sim \nu}[YY^\top]$, and $\mathbf{C} = \rho_Z = \mathbb{E}_{Z \sim \xi}[ZZ^\top]$, then

$$D_{MSB}(\mu, \nu) = \sup_{w \in \mathcal{S}} d_B(\rho_{X,w}, \rho_{Y,w}) \leq \sup_{w \in \mathcal{S}} \{d_B(\rho_{X,w}, \rho_{Z,w}) + d_B(\rho_{Z,w}, \rho_{Y,w})\} \quad (19)$$

$$\leq \sup_{w \in \mathcal{S}} d_B(\rho_{X,w}, \rho_{Z,w}) + \sup_{w \in \mathcal{S}} d_B(\rho_{Z,w}, \rho_{Y,w}) = D_{MSB}(\mu, \xi) + D_{MSB}(\xi, \nu). \quad (20)$$

□

Using the triangle inequality and Theorem 4, we have the following result directly.

Theorem 6. *Given probability measures $\mu, \nu \in \mathcal{P}_{L_2}(\mathcal{X})$ supported on a compact subset $\mathcal{X} \subset \mathbb{R}^d$. Assume that $\{X_i\}_{i=1}^n$ are i.i.d. random variables with distribution μ and $\{Y_i\}_{i=1}^n$ are i.i.d. random variables with distribution ν . For the empirical measures $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ and $\nu_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$,*

$$|D_{MSB}(\mu, \nu) - D_{MSB}(\mu_n, \nu_n)| \leq D_{MSB}(\mu, \mu_n) + D_{MSB}(\nu, \nu_n) \quad (21)$$

$$\leq 2 \text{diam}(\mathcal{X}) \sqrt{\frac{32}{n} [(d+1) \log(n+1) + \log(8/\delta)]}, \quad (22)$$

with probability at least $1 - \delta$, for any $\delta \in (0, 1)$.

A.2 Synthetic Data

We start by comparing the proposed max-sliced Bures divergence to the max-sliced W2 distance for two-dimensional data. We compare Algorithm 1 for solving each one-sided max-sliced Bures divergence with gradient-based approaches for it and the max-sliced W2 distance using ADAM with parameters `lr = 1e-3`, `beta1 = 0.9`, `beta2 = 0.999`, `epsilon = 1e-08`, capping the number of iterations at 1000 or until the change in the slice is minimal $\|w - w_{\text{old}}\|_\infty < 10^{-6}$.

In two-dimensions, a near optimal slice can be obtained by a fine grid search of the sliced Bures and sliced W2 distance as shown in Figure 5 for two zero-mean Gaussian distributions.

A.3 Pseudo-real Data: Detecting Non-Uniform Labels Distributions

We investigate the precision of the witness functions in detecting covariate shift caused by class imbalances. $\hat{\mu}$ is a uniform sample from the training set and $\hat{\nu}$ is a sample from the test set with less instances from one class, the minority class. Using the one-sided max-sliced Bures divergence we obtain the optimal witness function $\omega_{\mu > \nu}$, which should have larger evaluations for at least a subset of points in the under-represented class.

We can use this witness function to identify the minority class in the test set by examining the training set instances and their labels with the largest witness function evaluations. We apply $\omega_{\mu > \nu}$ to the sample $\hat{\mu}$, identifying the top-10 witness points with the largest witness function evaluations $\omega_{\mu > \nu}(x_{\hat{\pi}(1)}) \geq \dots \geq \omega_{\mu > \nu}(x_{\hat{\pi}(10)})$. $\{x_{\hat{\pi}(i)}\}_{i=1}^K$ are the $K = 10$ points from $\hat{\mu}$ with the largest norm after projection to this subspace defined by $w_{\mu > \nu}$. The performance is quantified by the precision of the labels of these instances. Ideally, these points should all be from the minority class.

We illustrate this simple case of covariate shift detection on the MNIST data set. Figure A.3 details a comparison of the max-sliced versions of the Bures, Fréchet, and W2 distances in terms of the divergence value estimates, the runtime, and the witness function's precision.

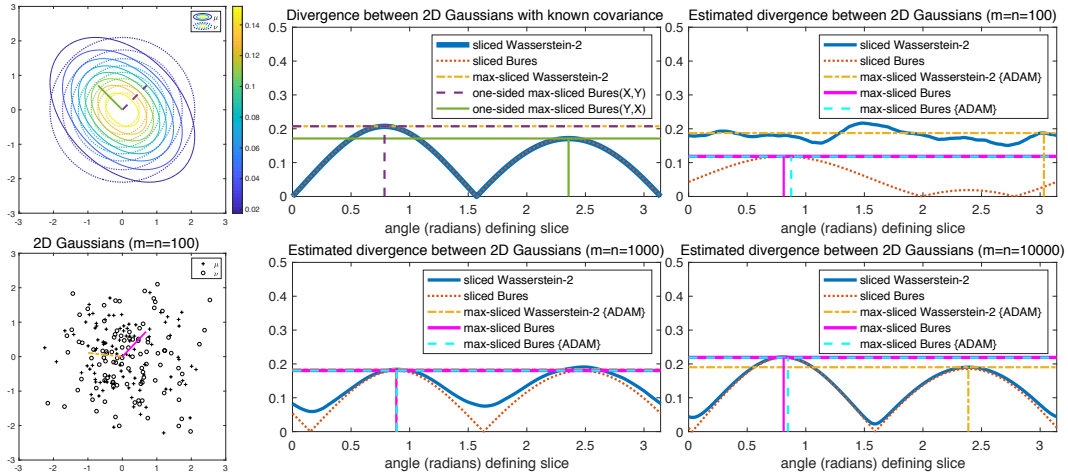


Figure 5: Sliced and max-sliced Bures and Wasserstein-2 distances are compared on population statistics and samples of varying sizes. $\mu = \mathcal{N}(\mathbf{0}, \mathbf{C})$ and $\nu = \mathcal{N}(\mathbf{0}, \mathbf{I})$, where $\mathbf{C} = \mathbf{Z}\mathbf{Z}^T$, and $\mathbf{Z} \in \mathbb{R}^{2 \times 2}$ with entries that are originally standard normals and then row normalized such that \mathbf{C} is a correlation matrix. In the population case and for zero-mean Gaussians, the Bures divergence is equivalent to the W2 distance [Gelbrich, 1990]. In the sample case, it is a lower bound. At both $m = 100$ and $m = 10^4$ the gradient optimization of the max-sliced W2 distance fails to obtain the global optimal slice (instead obtaining a local optimum).

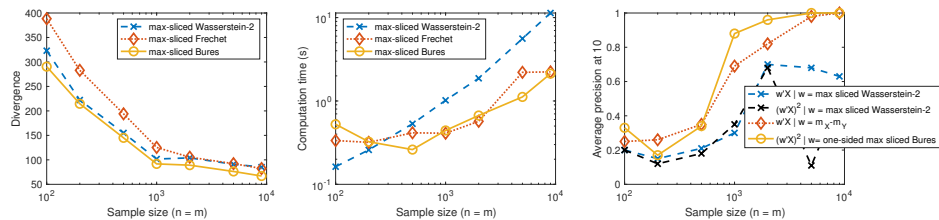


Figure 6: Max-sliced distances applied to two samples from MNIST. The first sample $\hat{\mu}$ consists of m images drawn uniformly from the training set, such that each digit has prevalence of 10%, and the second sample $\hat{\nu}$ is n images from the test set where one digit is a minority class $l \in \{0, \dots, 9\}$ with prevalence of 5%. (Left) Divergence estimates across sample size with $l = 7$. For $m < 2000$, gradient-based approach (ADAM) for the max-sliced W2 distance fails to obtain the optimal slice as it should upper bound the max-sliced Fréchet distance. (Center) Computation time. (Right) Each curve is the average precision@10 (averaged across the 10 classes). The one-sided max-sliced Bures yields the witness function $\omega_{\mu > \nu}(\cdot) = \langle \mathbf{w}, \cdot \rangle^2$, which is applied to $\hat{\mu}$ and reliably identifies points from the minority class for $m \geq 1000$. In this case the witness function corresponding to the difference of means $\omega(\cdot) = \langle \mathbf{w}, \cdot \rangle |_{\mathbf{w} = \mathbf{m}_X - \mathbf{m}_Y}$ also works, but with lower precision for $m \geq 1000$.