
LSGANs with Gradient Regularizers are Smooth High-dimensional Interpolators

Siddarth Asokan*

Robert Bosch Center for Cyber-Physical Systems
Indian Institute of Science
Bengaluru, India
siddartha@iisc.ac.in

Chandra Sekhar Seelamantula

Department of Electrical Engineering
Indian Institute of Science
Bengaluru, India
css@iisc.ac.in

Abstract

We consider the problem of discriminator optimization in least-squares generative adversarial networks (LSGANs) subject to higher-order gradient regularization enforced on the convex hull of all possible interpolation points between the target (real) and generated (fake) data. We analyze the proposed LSGAN cost within a *variational framework* and show that the optimal discriminator, given a generator, solves a regularized least-squares problem, and can be represented through a poly-harmonic radial basis function (RBF) interpolator. The optimal RBF discriminator can be implemented in closed-form, with the weights computed by solving a linear system of equations. We validate the proposed approach on synthetic Gaussian and standard image datasets. While the optimal LSGAN discriminator leads to a superior convergence on Gaussian data, the inherent *low-dimensional manifold* structure of images makes the optimal discriminator implementation ill-posed. Nevertheless, replacing the trainable discriminator network with a closed-form RBF interpolator results in superior convergence on 2-D Gaussian data, while overcoming pitfalls in GAN training, namely *mode dropping* and *mode collapse*.

1 Introduction

Generative adversarial networks (Goodfellow et al., 2014) are a two-player game involving a generator network G and a discriminator network D . The generator learns to transform samples from a noise distribution, typically Gaussian, into images that follow a distribution p_g . On the other hand, the discriminator accepts an input, either a *fake* sample from p_g , or a *real* sample (image from a dataset) with underlying distribution p_d , and learns to differentiate between the two.

The GAN Discriminator: Broadly speaking, the discriminator has two interpretations. In the standard GAN (SGAN) (Goodfellow et al., 2014), least-squares GAN (LSGAN) (Mao et al., 2017) or f -GANs (Nowozin et al., 2016), the discriminator has been shown to mimic a chosen divergence metric between p_d and p_g (e.g., the Jensen-Shannon or the Pearson- χ^2 divergence). Divergence-based approaches fail if the data and generator distributions are non-overlapping (Arjovsky et al., 2017). As an alternative, *integral probability metrics* (IPMs) were proposed, where the discriminator (often called the critic) approximates a distance measure between the distributions, which induces a constraint on the class of critics. For example, the Wasserstein GAN (WGAN) (Arjovsky et al., 2017), which minimizes the *earth mover's* distance, requires a Lipschitz-1 critic. Gulrajani et al. (2017); Bellemare et al. (2017); Mescheder et al. (2018); Mroueh et al. (2018) and Adler & Lunz (2018) consider regularizers on the energy of the gradients of the discriminator. Along another vertical, kernel-based discriminators with Gaussian and inverse multi-quadric kernel have been considered by Li et al. (2015, 2017a). While WGAN variants are a stable choice (Kang et al., 2022), recent

*Corresponding Author

works by Jolicoeur-Martineau et al. (2021) have been successful in employing the LSGAN loss in adversarial score-matching applications.

LSGANs and Gradient Penalties: Interpolation based regularizers have been introduced in WGAN with the gradient penalty. Roth et al. (2017); Kodali et al. (2017); Mescheder et al. (2018) observed that minimizing the norm of the discriminator’s gradient (on interpolated points) on divergence-minimizing GANs such as LSGAN results in superior performance. However, neither does the link to Lipschitz-1 discriminators hold, nor does the GAN approximate the Pearson χ^2 divergence. **What does optimizing the LSGAN with gradient penalty actually lead to?** This gap in understanding is what we seek to address in this paper.

1.1 The Proposed Approach

In this paper, we establish an explicit link between discriminator learning and high-dimensional interpolation. Intuitively, given a batch of samples in \mathbb{R}^n , the LSGAN discriminator can be seen as assigning specific *class-labels* to the real samples and the fake ones. Given an unseen sample \mathbf{x} , the ideal discriminator output $D(\mathbf{x})$ depends on the values assigned to all points in the vicinity of \mathbf{x} , which is precisely the task of kernel-based interpolation. Optimization of functions with higher-order derivatives having bounded L_2 norm to interpolate a given set of points has a unique solution (Duchon, 1977; Meinguet, 1979). This has led to successful application of higher-order gradient regularization in image processing tasks – image interpolation (Tirosh et al., 2006) and super-resolution (Ren et al., 2013), to name a few. We therefore consider the gradient-regularized LSGAN cost analyzed from a functional optimization standpoint. Our analysis shows that the optimal discriminator (given the generator) in the proposed setting involves polyharmonic radial basis functions (RBFs) for interpolation. We implement the optimum through the RBF network with predetermined weights and centers. *Essentially, we show that enforcing interpolation-based gradient regularizers on the LSGAN loss on the space of images, results in a discriminator that performs interpolation on the space of class-labels.* The proposed approach, referred to as *Poly-LSGAN*, outperforms baseline LSGANs in terms of training stability and convergence on Gaussian learning tasks. The source code for implementing Poly-LSGANs is available at <https://github.com/DarthSid95/PolyLSGANs>. However, in practice, computing the weights of the interpolator is impractical due to singularity issues arising from the *Manifold Hypothesis*, i.e., images lie in low-dimensional embeddings in high-dimensional spaces (Kelley, 2017). Nevertheless, the superior performance on synthetic data makes Poly-LSGAN a promising direction for further research.

2 LSGANs and Gradient Penalties

Mao et al. (2017) considered the GAN learning problem where the discriminator and generator networks minimize the least-squares loss. To mimic the classifier nature of the standard GAN (Goodfellow et al., 2014), an $a - b$ coding scheme is used, where a and b are the class labels of the generated samples and target data samples, respectively. On the other hand, the generator is trained to produce samples that are assigned a class label c by the discriminator. The resulting formulation is as follows:

$$\mathcal{L}_D^{\text{LS}} = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_d} [(D(\mathbf{x}) - b)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_g} [(D(\mathbf{x}) - a)^2]; \quad D^*(\mathbf{x}) = \arg \min_D \mathcal{L}_D^{\text{LS}},$$

$$\text{and } \mathcal{L}_G^{\text{LS}} = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_g} [(D^*(\mathbf{x}) - c)^2]; \quad p_g^*(\mathbf{x}) = \arg \min_{p_g} \mathcal{L}_G^{\text{LS}}.$$

While Mao et al. (2017) showed that setting $b - a = 2$ and $c - a = 1$ lead to the generator minimizing the Pearson- χ^2 divergence, a more intuitive approach is to set $c = b$, which enforces the generator to output samples that are classified as *real* by the discriminator. Rosca et al. (2020) showed that the gradient penalties in GANs have the general form $\mathbb{E}_{\mathbf{x} \sim p_r} [(\|\nabla D(\mathbf{x})\|_2 - K)^2]$, where p_r is the reference density and K is a suitable constant. Works such as Gulrajani et al. (2017); Petzka et al. (2018); Terjék (2020) consider $K > 0$ to enforce Lipschitz smoothness and use $p_r = p_{\text{int}}$, p_{int} being the interpolated distribution $\alpha p_d + (1 - \alpha) p_g$; $\alpha \in [0, 1]$. On the other hand, Kodali et al. (2017); Mescheder et al. (2018); Mroueh et al. (2018) consider p_r as limiting cases of either p_d or p_g , with $K = 0$ to promote the smoothness of the learnt discriminator. Adler & Lunz (2018) implemented m^{th} -order generalizations of the cost empirically through a Fourier representation of the cost.

2.1 Regularized LSGAN and Least-squares Interpolation

The ideal discriminator, upon convergence, will take a constant value for all values of \mathbf{x} over \mathcal{X} , the convex hull of the supports of p_d and p_g . We set p_r to be the uniform density over \mathcal{X} , which implies weighting all samples drawn from both p_d and p_g , and their possible linear combinations equally. In order to accelerate the convergence of the discriminator during training, we propose to employ higher-order gradient regularization with $K = 0$:

$$\mathcal{L}_D = \mathcal{L}_D^{\text{LS}} + \lambda_D \int_{\mathcal{X}} \|\nabla^m D(\mathbf{x})\|_2^2 d\mathbf{x}, \quad (1)$$

where $\lambda_d \geq 0$ is the Lagrange multiplier associated with the gradient penalty, and $\nabla^m D$ is the vector of m^{th} -order partial derivatives of $D(\mathbf{x})$ (cf. Appendix A).

Consider an N -sample approximation of $\mathcal{L}_D^{\text{LS}}$ in Equation (1), where N_B samples are drawn from p_d and p_g each (therefore, $N = 2N_B$), represented by the dataset batch

$$\mathcal{D} = \{(\mathbf{c}_i, y_i)\}_{i=1}^N = \{(\mathbf{x}_i, b) \mid \mathbf{x}_i \sim p_d\}_{i=1}^{N_b} \cup \{(\mathbf{x}_j, a) \mid \mathbf{x}_j \sim p_g\}_{j=1}^{N_b}.$$

The corresponding discriminator optimization problem can be formulated as follows:

$$D^* = \arg \min_D \sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N (D(\mathbf{c}_i) - y_i)^2 + \lambda_D \int_{\mathcal{X}} \|\nabla^m D(\mathbf{x})\|_2^2 d\mathbf{x}. \quad (2)$$

The above represents a regularized least-squares interpolation problem. When $\lambda_D = 0$, the optimum D^* is an interpolator that passes through the target points (\mathbf{x}_i, y_i) exactly. On the other hand, for positive values of λ_D , the minimization leads to smoother solutions, penalizing sharp transitions in the discriminator. We found out experimentally that $\lambda_D = 10$ results in superior performance. A smoother discriminator allows for more efficient training of the generator (Li et al., 2017b; Xu et al., 2018). The following theorem shows that the optimal discriminator, given the generator, that solves the above least-squares minimization problem is an interpolator between the real and fake class-labels.

Theorem 2.1. *The optimal LSGAN discriminator that minimizes the cost given in Eq. (2) is*

$$D^*(\mathbf{x}) = \sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N w_i \psi_k(\|\mathbf{x} - \mathbf{c}_i\|) + \mathcal{P}_{m-1}(\mathbf{x}; \mathbf{v}), \quad (3)$$

$$\text{where } \psi_k(r) = \begin{cases} r^k & \text{for } k = 1, 3, 5, \dots \\ r^k \ln(r) & \text{for } k = 2, 4, 6, \dots \end{cases} \quad (4)$$

is the polyharmonic radial basis function, $\mathcal{P}_{m-1}(\mathbf{x}; \mathbf{v})$ is an $(m-1)^{\text{th}}$ order polynomial parametrized by the coefficients $\mathbf{v} \in \mathbb{R}^L$, $\mathbf{x} \in \mathbb{R}^n$, $\mathcal{D} = \{(\mathbf{c}_i, y_i)\}$ is the set of real and fake centers about which the polyharmonic RBFs $\psi_k(\|\cdot\|)$ are localized, $\|\cdot\|$ denotes the ℓ_2 norm, and k denotes the order of the polyharmonic interpolator, which implicitly assumes a gradient penalty of order $m = \lceil \frac{k+n}{2} \rceil$. The N weights $\mathbf{w} = [w_1, w_2, \dots, w_N]^T$ and L polynomial coefficients $\mathbf{v} = [v_1, v_2, \dots, v_L]^T$ can be obtained by solving the linear system of equations:

$$\begin{bmatrix} \mathbf{A} + (-1)^m \lambda_D \mathbf{C}_k \mathbf{I} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix}, \quad (5)$$

$$\text{where } [\mathbf{A}]_{i,j} = \psi_k(\|\mathbf{c}_i - \mathbf{c}_j\|), \mathbf{B} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \mathbf{c}_1 & \mathbf{c}_2 & \dots & \mathbf{c}_N \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{c}_1^{m-1} & \mathbf{c}_2^{m-1} & \dots & \mathbf{c}_N^{m-1} \end{bmatrix}^T, \text{ and } \mathbf{y} = [y_1, y_2, \dots, y_N]^T,$$

\mathbf{I} is the $N \times N$ identity matrix, and \mathbf{c}_i^j is a vector of all j^{th} order monomials of \mathbf{c}_i . The matrix \mathbf{B} is required to be full-rank for invertibility, and \mathbf{C}_k is a constant that depends only on the order k .

The proof follows by applying the Euler-Lagrange equation from the *Calculus of Variations* to the cost in Eq. (2). The details are provided in Appendix B. For $k \leq 0$, the solution is non-interpolating, due to the singularity at the centers \mathbf{c}_i . Owing to the polyharmonic radial basis kernel, the proposed approach is referred to as *Poly-LSGAN*.

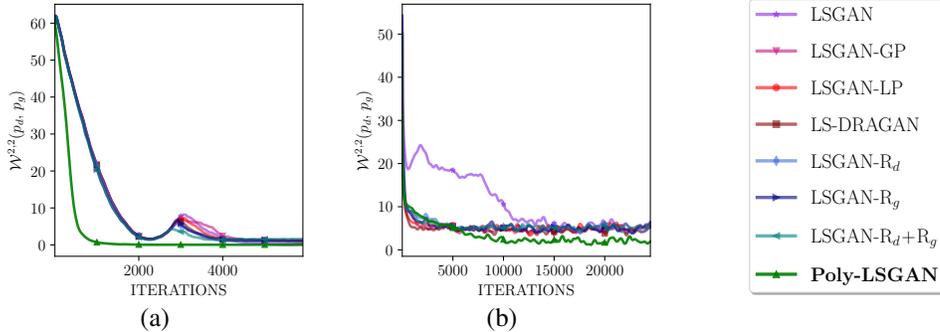


Figure 1: The Wasserstein-2 distance versus iterations on learning (a) a 2-D Gaussian; and (b) a 2-D Gaussian mixture, for various LSGAN variants. The performance of the Poly-LSGAN with the RBF discriminator is superior to the baselines in both scenarios. The convergence is also relatively smoother and stabler, unlike the baselines, which have fluctuations on the 2-D Gaussian learning task.

3 Experimental Evaluation

We now evaluate the optimal Poly-LSGAN discriminator on learning synthetic 2-D Gaussian and Gaussian mixture models (GMMs), and subsequently discuss extensions to learning images. We consider polyharmonic spline order $k = 1, 2$. For larger k , we encountered numerical instability.

3.1 Validation on Synthetic Data

We compare Poly-LSGAN against the base LSGAN (Mao et al., 2017), and LSGAN subjected to the gradient penalty (GP) (Gulrajani et al., 2017), R_d and R_g (Mescheder et al., 2018), Lipschitz penalty (LP) (Petzka et al., 2018) and the DRAGAN (Roth et al., 2017) regularizers. On the unimodal learning task, the target is $\mathcal{N}(5\mathbf{1}_2, 1.5\mathbb{I}_2)$, where $\mathbf{1}_2$ is the 2-D vector of ones, and \mathbb{I}_2 is the 2×2 identity matrix. On the GMM learning task, we consider eight components distributed uniformly about the unit circle, each having a standard deviation of 0.02. To quantify the performance, we use the Wasserstein-2 distance between the target and generator distributions ($\mathcal{W}^{2,2}(p_d, p_g)$). The network parameters are given in Appendix C. Figure 1 shows the $\mathcal{W}^{2,2}$ distance as a function of iterations on the Gaussian and Gaussian mixture learning tasks. On both datasets, we observe that using polyharmonic RBF discriminator results in superior generator performance (lower $\mathcal{W}^{2,2}$ scores). In all scenarios considered, the polyharmonic RBF discriminator learns the perfect classifier, compared to LSGAN with a trainable network discriminator. Additional comparisons are given in Appendix D.

3.2 Experiments on Image Datasets

We train on MNIST (LeCun et al., 1998), Fashion-MNIST (Xiao et al., 2017) and CelebA (Liu et al., 2015) datasets. The results are presented in Appendix D. While the underlying structure is learnt, the generated images are far from being realistic and below par compared with standard GAN results. Poly-LSGAN failed to converge as the monomial matrix \mathbf{B} (which can be viewed as a vector generalization of the Vandermonde matrix) became singular. As noted in the literature on mesh-free interpolation (Iske, 2004), \mathbf{B} must be full-rank, for the system of equations in Eq. (5) to have a unique solution. This requires the centers c_i to not lie in a subspace/manifold of \mathbb{R}^n . However, from the manifold hypothesis (Kelley, 2017; Vershynin, 2018), we know that structured image datasets lie precisely in such low-dimensional manifolds. To ascertain this behavior, we simulated a data manifold with $n = 32 \times 32 \times 3 = 3072$ dimensions, with the first M entries distributed as $\mathcal{N}(0.5\mathbf{1}_M, 0.2\mathbb{I})$ and the rest zeros, and observed that for $M \neq n$, the matrix \mathbf{B} indeed became singular.

4 Conclusions

Considering a generalization of the gradient-regularized LSGAN cost, we showed that the optimal discriminator for a given generator solves the regularized least-squares interpolation problem, whose solution can be represented using the polyharmonic RBF. The proposed discriminator results in superior performance on synthetic datasets such as 2-D Gaussians and Gaussian mixture densities. However, the rank deficiency of the monomial matrix \mathbf{B} caused convergence issues on image datasets.

Although implementing the optimal discriminator yields superior performance in idealistic settings, a straightforward extension to image generation in the ambient dimension became intractable due to singularity issues. One approach to overcoming this issue is to resort to a latent space solution as in the case of Wasserstein autoencoders (Tolstikhin et al., 2018), or to perform adversarial score matching (Jolicoeur-Martineau et al., 2021). Alternatively, one could also analyze the effect of enforcing the gradient-based regularizers on other GAN losses, such as the vanilla GAN or WGAN.

Acknowledgements

This work is supported by the Microsoft Research Ph.D. Fellowship 2018, Qualcomm Innovation Fellowship 2019, 2021 and 2022, and the Robert Bosch Centre for Cyber-Physical Systems Ph.D. Fellowships (2020-2021; 2021-2022).

References

- Abadi, M. et al. TensorFlow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint, arXiv:1603.04467*, Mar. 2016. URL <https://arxiv.org/abs/1603.04467>.
- Adler, J. and Lunz, S. Banach Wasserstein GAN. In *Advances in Neural Information Processing Systems 31*, pp. 6754–6763. 2018.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 214–223, 2017.
- Aronszajn, N., Creese, T., and Lipkin, L. *Polyharmonic Functions*. Oxford: Clarendon, 1983.
- Bellemare, M. G., Danihelka, I., Dabney, W., Mohamed, S., Lakshminarayanan, B., Hoyer, S., and Munos, R. The Cramér distance as a solution to biased Wasserstein gradients. 2017.
- Duchon, J. Splines minimizing rotation-invariant semi-norms in Sobolev spaces. In *Constructive Theory of Functions of Several Variables*, pp. 85–100. Springer, 1977.
- Fasshauer, G. E. *Meshfree Approximation Methods with MATLAB*. World Scientific, 2007.
- Flamary, R. et al. POT: Python optimal transport. *Journal of Machine Learning Research*, 22(78): 1–8, 2021. URL <http://jmlr.org/papers/v22/20-451.html>.
- Gelfand, I. M. and Fomin, S. V. *Calculus of Variations*. Prentice-Hall, 1964.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. C., and Bengio, Y. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27*, pp. 2672–2680. 2014.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. Improved training of Wasserstein GANs. In *Advances in Neural Information Processing Systems 30*. 2017.
- Iske, A. Multiresolution methods in scattered data modelling. In *Lecture Notes in Computational Science and Engineering*, 2004.
- Jolicoeur-Martineau, A., Piché-Taillefer, R., Mitliagkas, I., and des Combes, R. T. Adversarial score matching and improved sampling for image generation. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=eLfqM13z31q>.
- Kang, M., Shin, J., and Park, J. StudioGAN: A taxonomy and benchmark of GANs for image synthesis, 2022. URL <https://arxiv.org/abs/2206.09479>.
- Kelley, J. L. *General Topology*. Courier Dover Publications, Inc., 2017.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations*, 2015.
- Kodali, N., Abernethy, J. D., Hays, J., and Kira, Z. On convergence and stability of GANs. *arXiv preprint, arXiv:1705.07215*, May 2017. URL <http://arxiv.org/abs/1705.07215>.

- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Li, C. L., Chang, W. C., Cheng, Y., Yang, Y., and Póczos, B. MMD GAN: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems 30*, pp. 2203–2213. 2017a.
- Li, J., Madry, A., Peebles, J., and Schmidt, L. Towards understanding the dynamics of generative adversarial networks. *arXiv preprints, arXiv:1706.09884*, 2017b. URL <http://arxiv.org/abs/1706.09884>.
- Li, Y., Swersky, K., and Zemel, R. Generative moment matching networks. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 1718–1727, Jul 2015.
- Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision*, 2015.
- Mao, X., Li, Q., Xie, H., Lau, R. Y. K., Wang, Z., and Smolley, S. P. Least squares generative adversarial networks. In *Proceedings of International Conference on Computer Vision*, 2017.
- Meinguet, J. Multivariate interpolation at arbitrary points made simple. *Journal of Applied Mathematics and Physics (ZAMP)*, 30(2):292–304, 1979.
- Mescheder, L., Geiger, A., and Nowozin, S. Which training methods for GANs do actually converge? In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, 2018.
- Mroueh, Y., Li, C., Sercu, T., Raj, A., and Cheng, Y. Sobolev GAN. In *Proceedings of the 6th International Conference on Learning Representations*, 2018.
- Nowozin, S., Cseke, B., and Tomioka, R. f-GAN: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems 29*, pp. 271–279. 2016.
- Petzka, H., Fischer, A., and Lukovnikov, D. On the regularization of Wasserstein GANs. In *Proceedings of the 6th International Conference on Learning Representations*, 2018.
- Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. In *Proceedings of the 4th International Conference on Learning Representations*, 2016.
- Ren, Z., He, C., and Zhang, Q. Fractional-order total-variation regularization for image super-resolution. *Signal Processing*, 93(9):2408–2421, 2013.
- Rosca, M., Weber, T., Gretton, A., and Mohamed, S. A case for new neural network smoothness constraints. In *Proceedings on "I Can't Believe It's Not Better!" at NeurIPS Workshops*, Proceedings of Machine Learning Research, 12 Dec 2020.
- Roth, K., Lucchi, A., Nowozin, S., and Hofmann, T. Stabilizing training of generative adversarial networks through regularization. In *Advances in Neural Information Processing Systems 30*, 2017.
- Terjék, D. Adversarial Lipschitz regularization. In *Proceedings of the 8th International Conference on Learning Representations*, 2020.
- Tirosh, S., De Ville, D. V., and Unser, M. Polyharmonic smoothing splines and the multidimensional Wiener filtering of fractal-like signals. *IEEE Transactions on Image Processing*, 15(9), 2006.
- Tolstikhin, I. O., Bousquet, O., Gelly, S., and Schölkopf, B. Wasserstein auto-encoders. In *Proceedings of the 6th International Conference on Learning Representations*, 2018.
- Vershynin, R. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms. *arXiv preprint, arXiv:1708.07747*, Aug. 2017. URL <https://arxiv.org/abs/1708.07747>.
- Xu, Z., Li, C., and Jegelka, S. Robust GANs against dishonest adversaries, 2018. URL <https://arxiv.org/abs/1802.09700>.

A Mathematical Preliminaries

We recall the high-dimensional Euler-Lagrange conditions from the Calculus of Variations (Gelfand & Fomin, 1964), which play an important role in the optimization of the Poly-LSGAN cost in Equation 2.

Consider a vector $\mathbf{x} = [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n$ and a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. The notation $\nabla^m f(\mathbf{x})$ denotes the vector of m^{th} -order partial derivatives of f with respect to the entries of \mathbf{x} . ∇^0 is the identity operator. The elements of $\nabla^m f$ are represented using the multi-index notation $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_n]^T$, as follows:

$$\partial^{\boldsymbol{\alpha}} f = \frac{\partial^{|\boldsymbol{\alpha}|}}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \dots \partial x_n^{\alpha_n}} f, \text{ where } |\boldsymbol{\alpha}| = \sum_{i=1}^n \alpha_i,$$

and $\boldsymbol{\alpha} \in \mathbb{Z}_*^n$, the set of n -dimensional vectors with non-negative integer entries. For example, with $n = 4, m = 3$, the index $\boldsymbol{\alpha} = [2, 0, 1, 0]^T$ yields the element $\frac{\partial^3}{\partial x_1^2 \partial x_3} f(\mathbf{x})$. The square of the L_2 -norm of $\nabla^m f$ is given by the multidimensional sum:

$$\|\nabla^m f(\mathbf{x})\|_2^2 = \sum_{\boldsymbol{\alpha}: |\boldsymbol{\alpha}|=m} \binom{m!}{\boldsymbol{\alpha}!} (\partial^{\boldsymbol{\alpha}} f(\mathbf{x}))^2, \quad (6)$$

where $\boldsymbol{\alpha}! = \alpha_1! \alpha_2! \dots \alpha_n!$. The iterated Laplacian, also known as the polyharmonic operator, is defined as: $\Delta^m f(\mathbf{x}) = \Delta(\Delta^{m-1} f(\mathbf{x}))$, where $\Delta f(\mathbf{x}) = \nabla \cdot \nabla f(\mathbf{x}) = \sum_{i=1}^n \frac{\partial^2}{\partial x_i^2} f(\mathbf{x})$ is the Laplacian operator acting on $f(\mathbf{x})$. Applying the multi-index notation yields the standard form of the polyharmonic operator:

$$\Delta^m f(\mathbf{x}) = \sum_{\boldsymbol{\alpha}: |\boldsymbol{\alpha}|=m} \binom{m!}{\boldsymbol{\alpha}!} \partial^{\boldsymbol{\alpha}} (\partial^{\boldsymbol{\alpha}} f(\mathbf{x})).$$

Consider an integral cost \mathcal{L} with the integrand \mathcal{F} dependent on f and all its partial derivatives up to and including order ℓ , given by $\mathcal{L}(f(\mathbf{x}), \partial^{\boldsymbol{\alpha}} f; |\boldsymbol{\alpha}| \leq \ell) = \int_{\mathcal{X}} \mathcal{F}(f(\mathbf{x}), \partial^{\boldsymbol{\alpha}} f; |\boldsymbol{\alpha}| \leq \ell) d\mathbf{x}$, defined on a suitable domain \mathcal{X} over which f and its partial derivatives up to and including order ℓ are continuously differentiable. The optimizer f^* must satisfy the Euler-Lagrange condition:

$$\left. \frac{\partial \mathcal{F}}{\partial f} + \sum_{j=1}^{\ell} \left((-1)^j \sum_{\boldsymbol{\alpha}: |\boldsymbol{\alpha}|=j} \partial^{\boldsymbol{\alpha}} \left(\frac{\partial \mathcal{F}}{\partial (\partial^{\boldsymbol{\alpha}} f)} \right) \right) \right|_{f=f^*} = 0. \quad (7)$$

B The Optimal Poly-LSGAN Discriminator

The proof of Theorem 2.1 follows from the results in mesh-free interpolation literature (Aronszajn et al., 1983; Iske, 2004; Fasshauer, 2007) that deal with the generic polyharmonic spline interpolation problem. For completeness, we provide the proof here. While the assumption may appear to be strong, we show that this is implicitly satisfied by the optimal solution. Recall the discriminator optimization problem given in Eq. (2)

$$\arg \min_D \left\{ \sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N (D(\mathbf{c}_i) - y_i)^2 + \lambda_D \int_{\mathcal{X}} \|\nabla^m D(\mathbf{x})\|_2^2 d\mathbf{x} \right\}. \quad (8)$$

To compute the functional optimum in the Calculus of Variations setting, the above cost must be cast into an integral form. Using the Dirac delta function, we write:

$$\sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N (D(\mathbf{c}_i) - y_i)^2 = \int_{\mathcal{X}} \sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N (D(\mathbf{x}) - y_i)^2 \delta(\mathbf{x} - \mathbf{c}_i) d\mathbf{x}.$$

Then, Eq. (8) can be rewritten as an integral-cost minimization:

$$\arg \min_D \left\{ \underbrace{\int_{\mathcal{X}} \sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N (D(\mathbf{x}) - y_i)^2 \delta(\mathbf{x} - \mathbf{c}_i) + \lambda_D \|\nabla^m D(\mathbf{x})\|_2^2}_{\mathcal{F}(D, \partial^\alpha D; |\alpha|=m)} d\mathbf{x} \right\}.$$

Computing the derivatives of the integrand \mathcal{F} with respect to D and $\partial^\alpha D$ yields

$$\frac{\partial \mathcal{F}}{\partial D} = 2 \sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N (D(\mathbf{x}) - y_i) \delta(\mathbf{x} - \mathbf{c}_i), \quad \text{and} \quad \sum_{\alpha: |\alpha|=m} \partial^\alpha \left(\frac{\partial \mathcal{F}}{\partial (\partial^\alpha D)} \right) = 2\lambda_D \Delta^m D.$$

Substituting the above into the Euler-Lagrange equation (Eq. (7)) gives us the partial differential equation that the optimal discriminator $D^*(\mathbf{x})$ must satisfy:

$$\left(\sum_{i=1}^N (D(\mathbf{x}) - y_i) \delta(\mathbf{x} - \mathbf{c}_i) \right) + (-1)^m \lambda_D \Delta^m D(\mathbf{x}) \Big|_{D=D^*(\mathbf{x})} = 0.$$

While the above condition is applicable for a strong solution, a weak solution to $D(\mathbf{x})$ satisfies:

$$\int_{\mathcal{X}} \left(\left(\sum_{i=1}^N (D(\mathbf{x}) - y_i) \delta(\mathbf{x} - \mathbf{c}_i) \right) + (-1)^m \lambda_D \Delta^m D(\mathbf{x}) \right) \eta(\mathbf{x}) d\mathbf{x} \Big|_{D=D^*(\mathbf{x})} = 0, \quad (9)$$

where $\eta(\mathbf{x})$ is any test function drawn from the family of compactly-supported infinitely-differentiable functions. Aronszajn et al. (1983), an authoritative resource on polyharmonic functions, has shown that, functions of the form

$$f(\mathbf{x}) = \sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N w_i \psi_k(\|\mathbf{x} - \mathbf{c}_i\|) + \mathcal{P}_{m-1}(\mathbf{x}; \mathbf{v}), \quad \text{where } \psi_k(r) = \begin{cases} r^k & \text{for } k = 1, 3, \dots \\ r^k \ln(r) & \text{for } k = 2, 4, \dots \end{cases} \quad (10)$$

satisfy the polyharmonic PDE:

$$\Delta^m f(\mathbf{x}) = \sum_{i=1}^N C_k w_i \delta(\mathbf{x} - \mathbf{c}_i),$$

where $\mathcal{P}_{m-1}(\mathbf{x}; \mathbf{v})$ is the $(m-1)^{th}$ order polynomial parametrized by the coefficients $\mathbf{v} \in \mathbb{R}^L$. For example, with $m = 2$, we have $\mathcal{P}(\mathbf{x}; \mathbf{v}) = \langle \tilde{\mathbf{v}}, \mathbf{x} \rangle$; $\mathbf{v} \in \mathbb{R}^{n+1}$. For a polynomial $\mathcal{P}_{m-1}(\mathbf{x}; \mathbf{v})$ it holds trivially that $\Delta^m \mathcal{P}_{m-1} = \mathbf{0}$ which yields the condition. Substituting the above back into Eq. (9), we get

$$\begin{aligned} \int_{\mathcal{X}} \left(\sum_{i=1}^N ((D(\mathbf{x}) - y_i) + (-1)^m \lambda_D C_k w_i) \delta(\mathbf{x} - \mathbf{c}_i) \right) \eta(\mathbf{x}) d\mathbf{x} \Big|_{D=D^*(\mathbf{x})} &= 0, \\ \Rightarrow \sum_{i=1}^N ((D(\mathbf{c}_i) - y_i) + (-1)^m \lambda_D C_k w_i) \eta(\mathbf{c}_i) \Big|_{D=D^*(\mathbf{x})} &= 0. \end{aligned}$$

Since the above condition must hold for all possible test functions η , we have

$$D^*(\mathbf{c}_i) - y_i + (-1)^m \lambda_D C_k w_i = 0 \quad \forall i = 1, 2, \dots, N,$$

where D^* is given by Eq. (10). Substituting for D^* and stacking for all i gives the following condition that the weights and polynomial coefficients satisfy:

$$(\mathbf{A} + (-1)^m \lambda_D C_k \mathbf{I}) \mathbf{w} + \mathbf{B} \mathbf{v} = \mathbf{y}, \quad (11)$$

where $[\mathbf{A}]_{i,j} = \psi_k(\|\mathbf{c}_i - \mathbf{c}_j\|)$; $\mathbf{w} = [w_1, w_2, \dots, w_N]^T$, $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$,

$$\mathbf{B} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \mathbf{c}_1 & \mathbf{c}_2 & \dots & \mathbf{c}_N \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{c}_1^{m-1} & \mathbf{c}_2^{m-1} & \dots & \mathbf{c}_N^{m-1} \end{bmatrix}^T, \quad \text{and } \mathbf{v} = [v_0, v_1, v_2, \dots, v_L]^T.$$

The matrix \mathbf{B} corresponds to a Vandermonde matrix when $n = 1$. However, the system of linear equations is underdetermined, with fewer equations than unknowns.

To derive the second condition present in Eq. (5), $\mathbf{B}^T \mathbf{w} = \mathbf{0}$, we enforce the higher-order gradient penalty on the optimal solution D^* . To define the higher-order gradient-norm, we first consider the inner-product space associated with the higher-order gradient, where

$$\begin{aligned} \langle f, g \rangle &\triangleq \int_{\mathbb{R}^n} (\nabla^m f) \cdot (\nabla^m g) \, d\mathbf{x}. \\ &= (-1)^m \int_{\mathbb{R}^n} f \cdot (\Delta^m g) \, d\mathbf{x}, \end{aligned}$$

where in turn, the second inequality is via integration by parts. For any function D^* of the form given in Eq. (10), we have

$$\begin{aligned} \langle D^*, D^* \rangle &= (-1)^m \int_{\mathbb{R}^n} D^*(\mathbf{x}) \left(\sum_{i=1}^N (-1)^m C_k w_i \delta(\mathbf{x} - \mathbf{c}_i) \right) \, d\mathbf{x} \\ &= (-1)^m C_k \sum_{i=1}^N w_i D^*(\mathbf{c}_i). \end{aligned}$$

Substituting for D^* from Eq. (10) gives:

$$\begin{aligned} \langle D^*, D^* \rangle &= \int_{\mathbb{R}^n} \|\nabla^m D^*\|_2^2 \, d\mathbf{x} = (-1)^m C_k \sum_{i=1}^N \left(w_i \sum_{\substack{j=1 \\ (\mathbf{c}_j, y_j) \sim \mathcal{D}}}^N w_j \psi_k(\|\mathbf{c}_j - \mathbf{c}_i\|) + \mathbf{B}\mathbf{v} \right) \\ &= (-1)^m C_k \mathbf{w}^T \mathbf{A} \mathbf{w}, \end{aligned} \quad (12)$$

where the second equality holds as a result of Eq. (14). Substituting in D^* and Eq. (12) into the optimization problem in Eq. (8) yields:

$$\begin{aligned} &\arg \min_D \left\{ \sum_{\substack{i=1 \\ (\mathbf{c}_i, y_i) \sim \mathcal{D}}}^N (D(\mathbf{c}_i) - y_i)^2 + \lambda_D \int_{\mathcal{X}} \|\nabla^m D(\mathbf{x})\|_2^2 \, d\mathbf{x} \right\} \\ &= \arg \min_{\mathbf{w}, \mathbf{v}} \left\{ \underbrace{\|\mathbf{A}\mathbf{w} + \mathbf{B}\mathbf{v} - \mathbf{y}\|_2^2}_{F(\mathbf{w}, \mathbf{v})} + \lambda_D C_k \mathbf{w}^T \mathbf{A} \mathbf{w} \right\}. \end{aligned} \quad (13)$$

Minimizing the cost function in Eq. (13) with respect to \mathbf{w} and \mathbf{v} yields:

$$\begin{aligned} \frac{\partial F}{\partial \mathbf{w}} &= 2\mathbf{A}^T (\mathbf{A}\mathbf{w} + \mathbf{B}\mathbf{v} - \mathbf{y} + 2\lambda_D C_k \mathbf{w}) = 0, \quad \text{and} \quad \frac{\partial F}{\partial \mathbf{v}} = 2\mathbf{B}^T (\mathbf{A}\mathbf{w} + \mathbf{B}\mathbf{v} - \mathbf{y}) = 0 \\ &\Rightarrow \mathbf{B}^T \mathbf{w} = \mathbf{0}, \end{aligned} \quad (14)$$

which gives us the second necessary condition that the optimal weights and polynomial coefficients must satisfy. Equation (14) ensures that the solution obtained is such that the sum of the unbounded polyharmonic kernels vanish as \mathbf{x} tends to infinity. Essentially, in regions close to the centers \mathbf{c}_i , there is a large contribution in $D^*(\mathbf{x})$ from the kernel function, and when far away from the centers, the polynomial has a large contribution in $D^*(\mathbf{x})$. This ensures that the discriminator obtained by solving the system of equations does not grow to infinity. This completes the proof of Theorem 2.1.

C Training Specifications

Experiments on 2-D Gaussian Data: On the unimodal Gaussian learning task, the generator is a single layer affine transformation of the noise z , given by $x = Mz + b$, while on the GMM task, it is a three-layer neural network with Leaky ReLU activations. The discriminator in baseline LSGAN variants is a three-layer neural network with Leaky ReLU activation in both the Gaussian and GMM learning tasks. Poly-LSGAN employs the RBF discriminator while weights are computed by solving the system of equation given in Eq. (5).

The networks are trained using the Adam optimizer (Kingma & Ba, 2015) with a learning rate of $\eta_g = 0.002$ for the generator and $\eta_d = 0.0075$ for the discriminator. A batch size of 500 is employed. The models are evaluated using the Wasserstein-2 distance, $\mathcal{W}^{2,2}(p_d, p_g)$, between the target and generator distributions. The distance $\mathcal{W}^{2,2}(p_d, p_g)$ is computed in closed-form on Gaussian data, given by

$$\mathcal{W}^{2,2}(p_d, p_g) = \|\boldsymbol{\mu}_d - \boldsymbol{\mu}_g\|^2 + \text{Trace} \left(\Sigma_d + \Sigma_g - 2\sqrt{\Sigma_d \Sigma_g} \right).$$

On the Gaussian mixture data, $\mathcal{W}^{2,2}(p_d, p_g)$ is estimated empirically using the *Python optimal transport* (Flamary et al., 2021) library.

Experiments on Image Data: On image learning tasks, we employ the DCGAN (Radford et al., 2016) generator, trained using the Adam optimizer. The batch size is set to 100. The generator learning rate is set to $\eta_g = 10^{-4}$. The discriminator is the polyharmonic RBF with $k = 1$, which corresponds to enforcing a gradient penalty of order $m = \lceil \frac{n+1}{2} \rceil$. Consequently, the matrix \mathbf{B} contains monomials up to order $m - 1$.

Codebase: The TensorFlow 2.0 (Abadi et al., 2016) based implementation for Poly-LSGAN is available on GitHub at <https://github.com/DarthSid95/PolyLSGANs>.

D Additional Experimental Results

We now present additional experiments results from the authentic Gaussian and image datasets experiments conducted on Poly-LSGAN.

Experiments on Gaussian Data: Figures 2 and 3 present the generated and target data samples, superimposed on the level-sets of the discriminator, for the 2-D Gaussian, and 8-component Gaussian mixture learning tasks, respectively. For the Gaussian learning problem, we observe that Poly-LSGAN does not *mode-collapse* upon convergence to the target distribution. However, in the baseline GANs, depending on the learning rate, the generator converges to a distribution of smaller support than the target, before latching on to the desired target. Similarly, on the GMM learning task, Poly-LSGAN learns the target distribution more accurately compared to the baselines.

Experiments on Image Data: Figure 4 depicts the images generated by Poly-LSGAN when trained on the MNIST, Fashion-MNIST and CelebA datasets. In all scenarios, we observe that, although the generator is able to produce images from the target dataset, the visual quality of the images is sub-par compared to standard GAN approaches. Additional training of these models resulted in gradient explosion caused by the singularity of the monomial matrix \mathbf{B} .

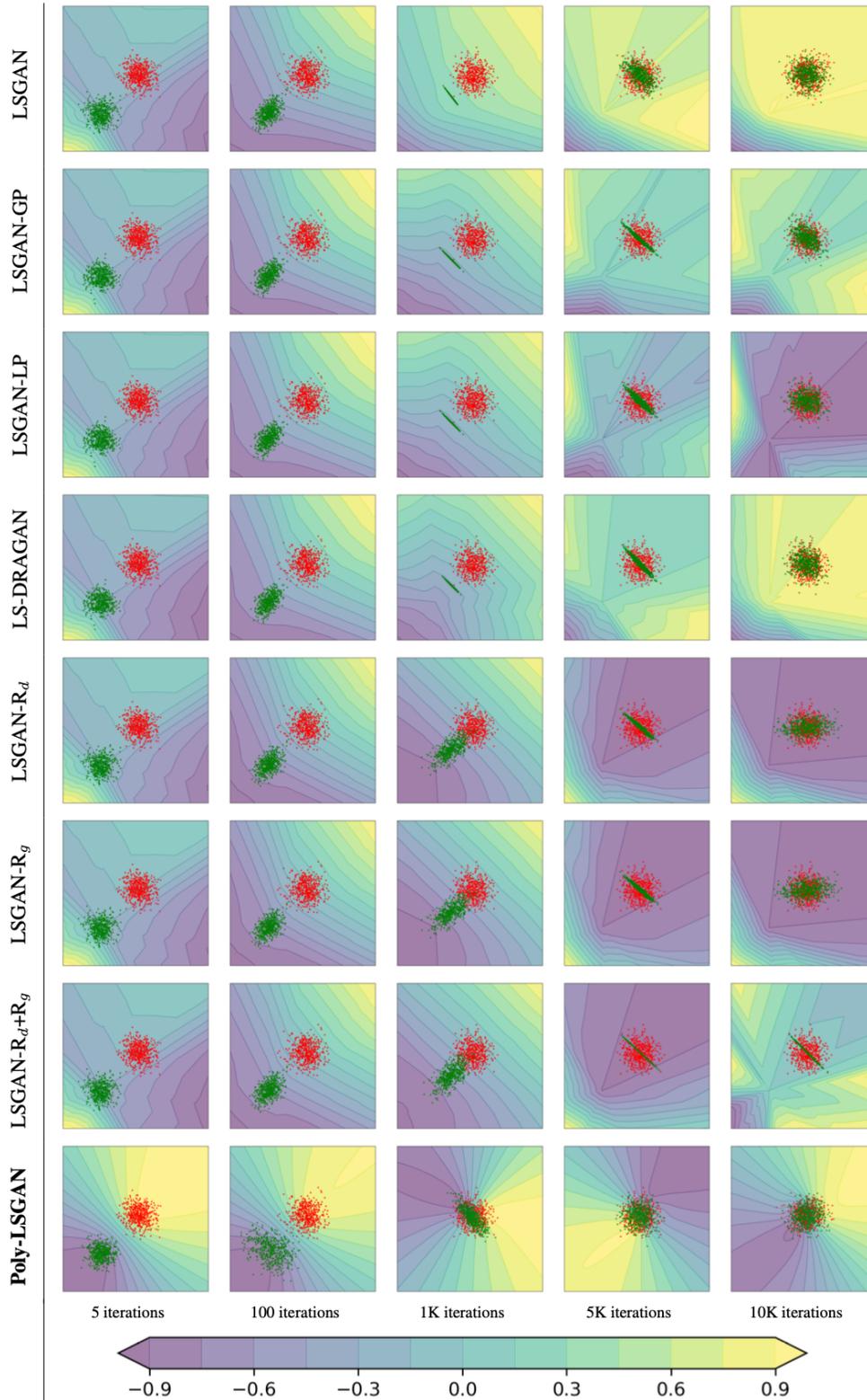


Figure 2: (Color online) Convergence of generator distribution (*green*) to the target 2-D Gaussian data (*red*) on the considered LSGAN variants. The heatmap represents the values taken by the discriminator. The Poly-LSGAN approach leads to a better representation of the discriminator function during the initial training iterations when compared to baseline approaches, leading to a faster convergence. Poly-LSGAN also does not experience *mode collapse*.

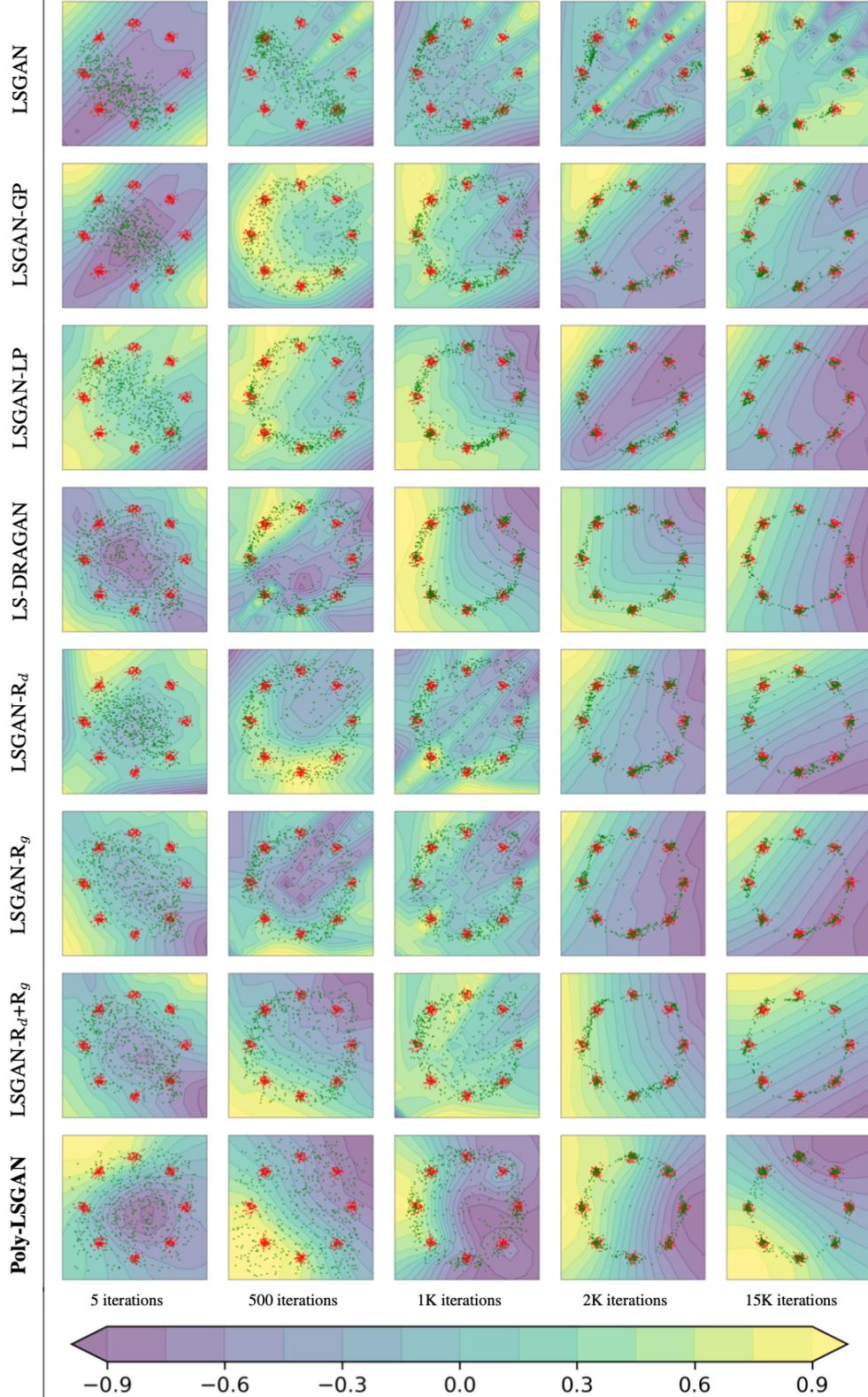


Figure 3: (Color online) Convergence of generator distribution (*green*) to the target multimodal Gaussian data (*red*) on the considered LSGAN variants, superimposed on the level-sets of the discriminator. The ideal $D(\mathbf{x})$ assigned a value of $b = 1$ to reals and $a = -1$ to fakes. Poly-LSGAN is able to identify the modes of the GMM more accurately than the baselines.



Figure 4: Images generated by training Poly-LSGAN on vectorized images drawn from (a) MNIST; (b) Fashion-MNIST; and (c) CelebA datasets. While Poly-LSGAN learns meaningful representations (although visually sub-par compared to standard GANs) on MNIST and Fashion-MNIST, the generator fails to converge in all scenarios. *The poor performance of Poly-LSGAN is attributed to training instability caused by the singularity of the monomial matrix \mathbf{B} in solving for the optimal discriminator weights.*