Rethinking Memorization–Generalization Trade-Off in Generative Models

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2025

Abstract

Existing generative models exhibit a memorization–generalization trade-off, and thus, avoiding memorization is a common strategy to promote generalization. In supervised learning, this long-accepted trade-off is being challenged, as recent studies show modern overparametrized models can achieve benign overfitting; that is, they generalize well even while exactly fitting, or memorizing, the training data. This raises the question of whether overparametrized generative models can similarly bypass this trade-off and achieve superior generalization alongside memorization. We address this with an empirical risk formulation that uses presampled latent variables instead of integrating over the entire latent distribution. We then recast the generative modeling problem as a supervised learning task of learning an optimal transport map, enabling us to leverage the concept of benign overfitting. In the one-dimensional setting, we show for the first time that benign overfitting can occur in generative models. We further expand and empirically validate our approach to higher dimensions, illustrating that benign overfitting extends more broadly across generative models.

1. Introduction

Generative models, such as generative adversarial networks [11] and diffusion models [14, 35], have become increasingly influential across various domains, from image synthesis [16, 30, 37, 40] to novel material generation [2, 24, 43]. Accordingly, interest in understanding their *generalizability*, the ability to generate diverse realistic data unseen during training, has grown significantly in recent years. A conventional wisdom in generative models is that there exists a trade-off between memorization (*i.e.*, overfitting) and generalization [26, 41, 44]. Hence, techniques such as weight decay, early stopping, and underparametrization, are typically incorporated [3, 13] to suppress memorization and thereby improve generalization.

Comparably, in the context of supervised learning, the *memorization–generalization trade-off*, based on the bias–variance trade-off [12], was generally accepted as true, until just a few years ago, and techniques to prevent memorization were commonly employed [18, 39]. However, recent studies have shown that overparametrized supervised learning models trained to achieve zero empirical risk, memorizing the training data, can exhibit superior generalization performance [5, 27]. This empirically observed and theoretically supported phenomenon is now known as *benign overfitting* [4, 23].

In contrast, to the best of our knowledge, there is currently neither theoretical nor even experimental evidence showing that generative models can achieve superior generalization beyond the conventional trade-off. This motivates the following central question of our study:

> Can generative models exhibit superior generalization via benign overfitting, similar to supervised learning?

To explore this question, we first note that the empirical risk in training generative models is often the distance between the generated distribution and the *empirical* target data distribution [19, 45]. However, this strategy forces all generated outputs to match training data, resulting in poor generalization [21, 29]. Therefore, for the model to generate outputs not in the training data, we should not utilize the entire latent distribution during training. We thus consider alternatively minimizing an empirical risk that only uses a fixed subset of latent vectors presampled from the latent distribution.

While this alternative choice opens the possibility to challenge the conventional trade-off, it remains unclear whether it can lead to superior generalization performance. As a first step toward this direction, we recast generative modeling as a regression task enabling us to leverage the theory of benign overfitting in regression. This leads to the following contributions:

- 1) In Section 3, we recast generative modeling as a regression task of learning the optimal transport map from the latent distribution to the true data distribution, using presampled latent variables and training data.
- 2) In Section 4, we show that generative models can exhibit benign overfitting, based on the random feature model framework [34]. Among several relevant works [7, 15, 33], we focus on [34] as it presents the "more is better" principle; increasing the model size improves performance. This suggests that a similar phenomenon may happen in generative modeling, despite the common belief that increasing model size without proper regularization impairs performance.

2. Preliminaries

2.1. Regression: population and empirical risks

Consider a standard regression setting with a training dataset $\mathcal{D} = \{(z_i, x_i)\}_{i=1}^n$ where $\{z_i\}_{i=1}^n$ are i.i.d. samples from a probability distribution ξ over \mathbb{R}^d , and

$$x_i = f^*(\boldsymbol{z}_i) + \epsilon_i \tag{1}$$

for a target function $f^* : \mathbb{R}^d \to \mathbb{R}$ and ϵ_i denoting the additive noise component. We further consider regression with feature maps, where $f(\cdot) = \psi(\cdot)^\top \beta$ with a feature map $\psi : \mathbb{R}^d \to \mathbb{R}^p$ and model parameters $\beta \in \mathbb{R}^p$. Then (1) reduces to estimating β^* such that $f^*(\cdot) = \psi(\cdot)^\top \beta^*$. Moreover, we set ψ to be a *random* feature map, which is of the form $\psi(z) = (g(w_1, z), \dots, g(w_p, z))$ for some function g, and w_i i.i.d. samples from some distribution ρ . For later purposes, assume that g admits a singular value decomposition $g(w, z) = \sum_i \sqrt{\lambda_i} \zeta_i(w) \phi_i(z)$ with $\zeta_i \in L^2(\rho)$ and $\phi_i \in L^2(\xi)$.

Ideally, we would like to minimize the population risk

$$\mathcal{R}(\boldsymbol{\beta}) \coloneqq \mathbb{E}_{\boldsymbol{z} \sim \boldsymbol{\xi}}[\|\boldsymbol{\psi}(\boldsymbol{z})^{\top} \boldsymbol{\beta} - \boldsymbol{\psi}(\boldsymbol{z})^{\top} \boldsymbol{\beta}^{*}\|^{2}] = \|\boldsymbol{\beta} - \boldsymbol{\beta}^{*}\|_{\boldsymbol{\Sigma}}^{2},$$
(2)

where $\Sigma := \mathbb{E}_{\boldsymbol{z} \sim \xi}[\boldsymbol{\psi}(\boldsymbol{z})\boldsymbol{\psi}(\boldsymbol{z})^{\top}]$, but typically $\boldsymbol{\beta}^*$ is inaccessible. Hence, the *empirical risk*

$$\hat{\mathcal{R}}(\boldsymbol{\beta}) \coloneqq \frac{1}{n} \sum_{i=1}^{n} \left\| \boldsymbol{\psi}(\boldsymbol{z}_{i})^{\top} \boldsymbol{\beta} - \boldsymbol{x}_{i} \right\|^{2}$$
(3)

is used as a proxy to be minimized. The generalization performance of the learned model $\hat{f} = \psi^{\top} \hat{\beta}$, obtained by minimizing the empirical risk, is still formulated by the population risk $\mathcal{R}(\hat{\beta})$.

Classical statistical learning theory suggests a bias–variance trade-off [12]. However, recent studies have shown the opposite, that overparametrized models, despite being trained to perfectly memorize the training data, can still exhibit strong generalization performance [4, 5, 27, 34, 38].

2.2. Generative modeling: population and empirical risks

Now consider when the dataset $\{x_i\}_{i=1}^n$ consists of n i.i.d. samples from a distribution ν on \mathbb{R}^m . Generative models aim to learn a process which approximates the target distribution ν by transforming the latent distribution ξ . We focus on the models whose generative process is represented by a deterministic map G. In this case, the goal is to approximate ν by $G_{\sharp}\xi$, the pushforward of ξ by G.

The population risk of a generative model is defined by a divergence between the generated distribution and target data distribution. Let us consider the squared 2-Wasserstein distance as the population risk, $\mathcal{R}(G) := \mathcal{W}_2^2(G_{\sharp}\xi, \nu)$. In practice, ν is not directly accessible, so the *semi-discrete empirical risk* $\hat{\mathcal{R}}_{semi}(G) := \mathcal{W}_2^2(G_{\sharp}\xi, \hat{\nu})$ is often minimized instead, where $\hat{\nu} := \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ is the empirical target distribution. A few works, *e.g.*, [36], have also explored an alternative of minimizing

$$\hat{\mathcal{R}}_{\text{fully}}(G) \coloneqq \mathcal{W}_2^2(G_{\sharp}\hat{\xi}, \hat{\nu}) = \inf_{\pi \in S_n} \frac{1}{n} \sum_{i=1}^n \|G(\boldsymbol{z}_i) - \boldsymbol{x}_{\pi(i)}\|^2,$$
(4)

which we call the *fully-discrete empirical risk*, where $\hat{\xi} = \frac{1}{n} \sum_{i=1}^{n} \delta_{z_i}$ is the empirical latent distribution on the i.i.d. samples $z_1, \ldots, z_n \sim \xi$, and S_n is the set of all permutations of $\{1, \ldots, n\}$.

3. Recasting Generative Modeling as Regression via Optimal Transport

Generative modeling and supervised learning are typically viewed as two distinct paradigms. Here, we bridge the two by recasting generative modeling as a regression problem via optimal transport.

3.1. The first step toward bridging generative modeling and regression

Since generative modeling in general admits infinitely many mappings minimizing $\mathcal{R}(G)$, linking it with regression requires choosing a particular generator G^* . We choose G^* to be the optimal transport map. To this end, the rest of this paper considers the equal-dimensional case d = m, leaving extensions to unequal dimensions for future work.

Given a cost function c, an optimal transport (OT) map G^* from ξ to ν is an optimal solution to

$$\inf_{G:G_{\sharp}\xi=\nu} \int c(\boldsymbol{z},G(\boldsymbol{z})) d\xi(\boldsymbol{z}).$$
(5)

For the quadratic cost $c(x, y) = ||x - y||^2$, if ξ is absolutely continuous with respect to the Lebesgue measure, then Brenier's theorem [6] ensures that there uniquely exists an optimal transport map G^* .

Learning the OT map in the context of generative modeling is not new; see, *e.g.*, [17, 22, 31]. However, there were no discussions on the generalization behavior when learning from finite training data. In contrast, this paper systemically analyzes how an OT map can be learned from finite training data $\{x_i\}_{i=1}^n$ and presampled latent variables $\{z_i\}_{i=1}^n$, by casting the problem as a regression task.

3.2. Optimal transport map based regression model for generative modeling

We now formulate a regression model, analogous to (1), for generative modeling based on the OT map. This requires pairing the finite training data with presampled latent variables. We naturally choose the OT map $\tilde{\pi} \in S_n$ between them, a permutation minimizing $\min_{\pi \in S_n} \sum_{i=1}^n ||\boldsymbol{z}_i - \boldsymbol{x}_{\pi(i)}||^2$. Having paired the dataset as $\{(\boldsymbol{z}_i, \boldsymbol{x}_{\tilde{\pi}(i)})\}_{i=1}^n$, we obtain the regression model

$$\boldsymbol{x}_{\tilde{\pi}(i)} = G^*(\boldsymbol{z}_i) + \boldsymbol{\epsilon}_i, \tag{6}$$

where $\epsilon_i \in \mathbb{R}^m$ denotes noise resulting from recasting generative modeling as regression. That said, these noises are not i.i.d., and this complicates the use of existing benign overfitting theories [34].

3.2.1. Optimal transport map based regression models: One-dimensional setting

By particularly focusing on the one-dimensional setting, we can get a more refined characterization of the noises. It is well known that, for empirical distributions $\hat{\xi}$ and $\hat{\nu}$ on \mathbb{R} , the OT map pairs the sorted z_i s to the sorted x_i s [28, Remark 2.28]. Hence, the regression model (6) becomes¹

$$x_{i:n} = G^*(z_{i:n}) + \epsilon_i, \tag{7}$$

where $x_{i:n}$ denotes *i*th order statistic, *i.e.*, the *i*th smallest one among $\{x_i\}_{i=1}^n$. However, for arbitrary target distribution ν , complete characterization of order statistics remains unavailable.

Fortunately, for our subsequent benign overfitting analyses, we only need the variances of ϵ_i s to be bounded. This turns out to be achievable, under a mild assumption on the target distribution ν .

Assumption 1 The target distribution ν has finite variance; i.e., for $x \sim \nu$, we have $Var(x) \leq \sigma^2$.

Lemma 1 Under Assumption 1, for any $n \ge 1$, the noise in (7) satisfies $\frac{1}{n} \sum_{i=1}^{n} \operatorname{Var}(\epsilon_i) \le 2\sigma^2$.

4. Benign Overfitting in Generative Models

In this section, we demonstrate benign overfitting of generative models, built upon the results of [34].

4.1. Theoretical results in the one-dimensional setting

Let us first study the simpler case, of when d = m = 1. To facilitate a plausible yet theoretical analysis, we adopt the following ansatz, introduced therein.

Assumption 2 (34, Gaussian Universality Ansatz) The expected population risk remains unchanged even if we replace $\{\tilde{\phi}_i(z)\}$ with i.i.d. samples from $\mathcal{N}(0,1)$ when $z \sim \xi$, and $\{\tilde{\zeta}_i(w)\}$ with i.i.d. samples from $\mathcal{N}(0,1)$ when $w \sim \rho$.

As in kernel regression, we approximate G^* by minimizing the least squares objective,

$$\min_{G(\cdot)=\boldsymbol{\psi}(\cdot)^{\top}\boldsymbol{\beta}} \quad \frac{1}{n} \sum_{i=1}^{n} \|G(\boldsymbol{z}_{i}) - \boldsymbol{x}_{\tilde{\pi}(i)}\|^{2}.$$
(8)

We focus on its ridge regression version, in which we solve (8) with an additional term $\frac{\delta}{n} ||\beta||^2$ with $\delta > 0$ in the objective. The optimal solution is then $\hat{\beta}(\delta) = \Psi^{\top} (\Psi \Psi^{\top} + \delta I)^{-1} x$, where Ψ is a matrix whose *i*th row is $\psi(z_i)^{\top}$. Accordingly, the learned model to be studied is $\hat{G}(\cdot) = \psi(\cdot)^{\top} \hat{\beta}(\delta)$.

Let $G^*(z) = \sum_i v_i \phi_i(z)$ be the expansion of the target function with respect to $\{\phi_i\}_i$. In the random feature eigenframework [34], which deals with classical ridge regression problems under the traditional assumption that $\epsilon_1, \ldots, \epsilon_n$ are i.i.d. Gaussians, it is proposed that the expected population risk can be well approximated by

$$\mathbb{E}_{\mathcal{D}}\left[\mathcal{R}(\hat{\boldsymbol{\beta}}(\delta))\right] \approx \mathcal{E}_{\text{te}} \coloneqq \frac{1}{1 - \frac{q(p-2s) + s^2}{n(p-q)}} \left(\sum_{i} \left(\frac{\gamma}{\lambda_i + \gamma} - \frac{\kappa\lambda_i}{(\lambda_i + \gamma)^2} \frac{p}{p-q}\right) v_i^2 + \check{\sigma}^2\right), \quad (9)$$

^{1.} For simplicity, we continue to use ϵ_i , where the index *i* is assumed to have been reordered.

where, for $s \coloneqq \sum_i \frac{\lambda_i}{\lambda_i + \gamma}$ and $q \coloneqq \sum_i \left(\frac{\lambda_i}{\lambda_i + \gamma}\right)^2$, κ and γ are unique nonnegative numbers such that $n = s + \frac{\delta}{\kappa}$ and $p = s + \frac{p\kappa}{\gamma}$. It is stated in Appendix I.3 of [33], the work from which [34] is developed, that $\check{\sigma}^2$ is the term corresponding to what is "effectively noise". In their setting, that quantity is the mean squared error of $x_i - G^*(z_i)$. We argue that this reasoning also applies to our setting, and $\check{\sigma}^2$ should be set as the expected residual variance $\check{\sigma}^2 = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n (x_{i:n} - G^*(z_{i:n}))^2\right]$.

Founded on these specifications, as in [34], we can show that increasing the number of features reduces the risk, so long as we are free to choose the ridge parameter.

Theorem 2 Let $\mathcal{E}_{te}(n, p, \delta)$ denote the value of \mathcal{E}_{te} with dataset size n, number of features p, and ridge parameter δ . Suppose that $p \leq p'$. Then, under Assumptions 1 and 2, it holds that

$$\min_{\delta} \mathcal{E}_{\mathsf{te}}(n, p', \delta) \le \min_{\delta} \mathcal{E}_{\mathsf{te}}(n, p, \delta)$$

Moreover, if we further assume that p > n, then denoting $\mathcal{E}_{te, 0} = \lim_{\delta \to 0+} \mathcal{E}_{te}$, it holds that

$$\mathcal{E}_{\text{te},0}(n,p') \le \mathcal{E}_{\text{te},0}(n,p)$$

4.2. Extending the theory to higher dimensions

The theoretical results in Section 4.1 can be extended to higher dimensions, under assumptions appropriately modified. Consider where a dataset is now $\{(z_i, x_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}^m$. We wish to estimate β under the linear model generalizing (1); for $\psi : \mathbb{R}^d \to \mathbb{R}^{p \times m}$, $\beta \in \mathbb{R}^p$, and $\epsilon_i \in \mathbb{R}^m$,

$$\boldsymbol{x}_i = \boldsymbol{\psi}(\boldsymbol{z}_i)^\top \boldsymbol{\beta} + \boldsymbol{\epsilon}_i. \tag{10}$$

In particular, the change in ψ amounts to g now being an \mathbb{R}^m -valued function. Then, we can show that Theorem 2 holds almost verbatim; the exact statements are in Theorems 13 and 14. For further details, with an overview on kernel regression with vector-valued targets, see Appendices C and D.

4.3. Experiments

We present experimental results in Appendix E to validate our analyses of the random feature model, particularly the approximation of the population risk by \mathcal{E}_{te} , thereby empirically supporting the principle that using more features is better.

5. Conclusion

We showed that overparametrized generative models can generalize despite memorizing training data, contrary to conventional belief that memorization undermines generalization performance. Our approach demonstrated this through learning the optimal transport map. By reformulating generative modeling into a regression task, quantitative generalization analyses became possible, leading to theoretical demonstrations that benign overfitting can occur in generative models. Specifically, we showed that increasing the number of features improves performance. Our work hence frames a new approach for studying the interplay between overparameterization and generalization in generative models.

Our work is yet limited to the latent and target distributions defined on the same space. Future work could explore alternative mappings beyond the optimal transport map, such as those based on the Gromov–Wasserstein theory or optimal transport maps composed with embeddings, which could enable analysis when the distributions lie on different underlying spaces.

References

- [1] Barry C. Arnold, N. Balakrishnan, and H. N. Nagaraja. *A First Course in Order Statistics*. Society for Industrial and Applied Mathematics, 2008.
- [2] Pavel Avdeyev, Chenlai Shi, Yuhao Tan, Kseniia Dudnyk, and Jian Zhou. Dirichlet diffusion score model for biological sequence generation. In *International Conference on Machine Learning*, pages 1276–1301. PMLR, 2023.
- [3] Ricardo Baptista, Agnimitra Dasgupta, Nikola B. Kovachki, Assad Oberai, and Andrew M. Stuart. Memorization and regularization in generative diffusion models. *arXiv preprint* arXiv:2501.15785, 2025.
- [4] Peter L. Bartlett, Philip M. Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- [5] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machinelearning practice and the classical bias-variance trade-off. *Proceedings of the National Academy* of Sciences, 116(32):15849–15854, 2019.
- [6] Yann Brenier. Décomposition polaire et réarrangement monotone des champs de vecteurs. *Comptes Rendus de l'Académie des Sciences - Série I - Mathématique*, 305:805–808, 1987.
- [7] Abdulkadir Canatar, Brendan Bordelon, and Cengiz Pehlevan. Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks. *Nature Communications*, 12:2914, 2021.
- [8] Claudio Carmeli, Ernesto De Vito, Alessandro Toigo, and Veronica Umanità. Vector valued reproducing kernel Hilbert spaces and universality. *Analysis and Applications*, 8(1):19–61, 2010.
- [9] Daniel K. Crane. *The singular value expansion for compact and non-compact operators*. PhD thesis, Michigan Technological University, 2020.
- [10] J. D. Esary, F. Proschan, and D. W. Walkup. Association of random variables, with applications. *The Annals of Mathematical Statistics*, 38(5):1466–1474, 1967.
- [11] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. Advances in Neural Information Processing Systems, 2014.
- [12] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer, second edition, 2009.
- [13] Reinhard Heckel and Paul Hand. Deep decoder: Concise image representations from untrained non-convolutional networks. In *International Conference on Learning Representations*, 2019.
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840–6851, 2020.

- [15] Arthur Jacot, Berfin Simsek, Francesco Spadaro, Clement Hongler, and Franck Gabriel. Kernel alignment risk estimator: Risk prediction from training data. In Advances in Neural Information Processing Systems, 2020.
- [16] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *International Conference on Learning Representations*, 2018.
- [17] Alexander Korotin, Vage Egiazarian, Arip Asadulaev, Alexander Safin, and Evgeny Burnaev. Wasserstein-2 generative networks. In *International Conference on Learning Representations*, 2021.
- [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. Advances in Neural Information Processing Systems, 25, 2012.
- [19] Dohyun Kwon, Ying Fan, and Kangwook Lee. Score-based generative modeling secretly minimizes the Wasserstein distance. Advances in Neural Information Processing Systems, 35: 20205–20217, 2022.
- [20] Peter D. Lax. Functional Analysis. John Wiley and Sons, Inc., 2002.
- [21] Lorenzo Luzi, Yehuda Dar, and Richard Baraniuk. Double descent and other interpolation phenomena in GANs. *arXiv preprint arXiv:2106.04003*, 2021.
- [22] Ashok Makkuva, Amirhossein Taghvaei, Sewoong Oh, and Jason Lee. Optimal transport mapping via input convex neural networks. In *Proceedings of the 37th International Conference* on Machine Learning, 2020.
- [23] Neil Mallinar, James Simon, Amirhesam Abedsoltan, Parthe Pandit, Misha Belkin, and Preetum Nakkiran. Benign, tempered, or catastrophic: A taxonomy of overfitting. *Advances in Neural Information Processing Systems*, 35:1182–1195, 2022.
- [24] Yunwei Mao, Qi He, and Xuanhe Zhao. Designing complex architectured materials with generative adversarial networks. *Science Advances*, 6(17), 2020.
- [25] Charles A. Micchelli and Massimiliano Pontil. On learning vector-valued functions. *Neural Computation*, 17(1):177–204, 2005.
- [26] Vaishnavh Nagarajan, Colin Raffel, and Ian J. Goodfellow. Theoretical insights into memorization in GANs. In *Integration of Deep Learning Theories Workshop, NeurIPS*, 2018.
- [27] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124003, 2021.
- [28] Gabriel Peyré and Marco Cuturi. Computational Optimal Transport. Now Publishers, 2019.
- [29] Jakiw Pidstrigach. Score-based generative models detect manifolds. *Advances in Neural Information Processing Systems*, 35:35852–35865, 2022.

- [30] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Highresolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- [31] Litu Rout, Alexander Korotin, and Evgeny Burnaev. Generative modeling with optimal transport maps. In *International Conference on Learning Representations*, 2022.
- [32] Filippo Santambrogio. Optimal transport for applied mathematicians. Birkhäuser, 2015.
- [33] James B. Simon, Madeline Dickens, Dhruva Karkada, and Michael R. DeWeese. The eigenlearning framework: A conservation law perspective on kernel ridge regression and wide neural networks. *Transactions on Machine Learning Research*, 2023.
- [34] James B. Simon, Dhruva Karkada, Nikhil Ghosh, and Mikhail Belkin. More is better in modern machine learning: when infinite overparameterization is optimal and overfitting is obligatory. In *International Conference on Learning Representations*, 2024.
- [35] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265, 2015.
- [36] Ruoyu Sun, Tiantian Fang, and Alexander Schwing. Towards a better global loss landscape of GANs. *Advances in Neural Information Processing Systems*, 33:10186–10198, 2020.
- [37] Yuhta Takida, Masaaki Imaizumi, Takashi Shibuya, Chieh-Hsin Lai, Toshimitsu Uesaka, Naoki Murata, and Yuki Mitsufuji. SAN: Inducing metrizability of GAN with discriminative normalized linear layer. In *International Conference on Learning Representations*, 2024.
- [38] Alexander Tsigler and Peter L. Bartlett. Benign overfitting in ridge regression. Journal of Machine Learning Research, 24(123):1–76, 2023.
- [39] Vladimir Vapnik. Statistical learning theory. John Wiley & Sons, 1998.
- [40] Zhendong Wang, Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Diffusion-GAN: Training GANs with diffusion. In *International Conference on Learning Representations*, 2023.
- [41] TaeHo Yoon, Joo Young Choi, Sehyun Kwon, and Ernest K. Ryu. Diffusion probabilistic models generalize when they fail to memorize. In *ICML 2023 workshop on structured probabilistic inference & generative modeling*, 2023.
- [42] Kôsaku Yosida. Functional analysis. Springer, 6th edition, 1995.
- [43] Claudio Zeni, Robert Pinsler, Daniel Zügner, Andrew Fowler, Matthew Horton, Xiang Fu, Zilong Wang, Aliaksandra Shysheya, Jonathan Crabbé, Shoko Ueda, Roberto Sordillo, Lixin Sun, Jake Smith, Bichlien Nguyen, Hannes Schulz, Sarah Lewis, Chin-Wei Huang, Ziheng Lu, Yichi Zhou, Han Yang, Hongxia Hao, Jielan Li, Chunlei Yang, Wenjie Li, Ryota Tomioka, and Tian Xie. A generative model for inorganic materials design. *Nature*, 639:624–632, 2025.

- [44] Huijie Zhang, Jinfan Zhou, Yifu Lu, Minzhe Guo, Peng Wang, Liyue Shen, and Qing Qu. The emergence of reproducibility and consistency in diffusion models. In *International Conference on Machine Learning*, 2024.
- [45] Pengchuan Zhang, Qiang Liu, Dengyong Zhou, Tao Xu, and Xiaodong He. On the discrimination-generalization tradeoff in GANs. In *International Conference on Learning Representations*, 2018.

Contents

1	Introduction							
2	Preliminaries 2.1 Regression: population and empirical risks 2.2 Generative modeling: population and empirical risks	2 2 3						
3	 Recasting Generative Modeling as Regression via Optimal Transport 3.1 The first step toward bridging generative modeling and regression	3 3 3 4						
4	 Benign Overfitting in Generative Models 4.1 Theoretical results in the one-dimensional setting	4 4 5 5						
5	onclusion							
A	Missing Details for Section 3A.1Characterizing the noise model in the one-dimensional settingA.2Proof of Lemma 1	11 11 12						
B	oof for Section 4.1 1 Proof of Theorem 2							
C	Vector-Valued Kernel Regression: Prerequisites for Section 4.2C.1Vector-valued linear models	14 14 16 16 21						
D Missing Details for Section 4.2D.1 Statements on the random feature models in higher dimensions								
E	Experimental Results E.1 One dimensional examples E.2 Random feature models in higher dimensions E.3 Neural-network-like model experiments	24 24 26 27						

Appendix A. Missing Details for Section 3

For a distribution ν on \mathbb{R} , let $F_{\nu}(x) \coloneqq \nu((-\infty, x])$ be its *cumulative distribution function* (CDF), and $F_{\nu}^{[-1]}(x) \coloneqq \inf \{t \in \mathbb{R} : F_{\nu}(t) \ge x\}$ be its *quantile function*, the generalized inverse of F_{ν} . Let us begin by recalling the following well-known fact.

Theorem 3 (32, Theorem 2.9) Let ξ and ν be probability distributions on \mathbb{R} . If ξ is absolutely continuous with respect to the Lebesgue measure, then a nondecreasing function $G^* = F_{\nu}^{[-1]} \circ F_{\xi}$ is the an OT map under the quadratic cost.

A.1. Characterizing the noise model in the one-dimensional setting

We study the distributions of the noise ϵ_i in (7). Our main focus is to show that they are bounded on expectation. In proving this boundedness, the following lemma plays a pivotal role.

Lemma 4 Suppose that ξ is absolutely continuous with respect to the Lebesgue measure. Then, $x_{i:n}$ and $G^*(z_{i:n})$ are identically distributed.

Proof Showing that $G^*(z_{i:n})$ and $x_{i:n}$ are identically distributed amounts to showing that their CDFs are identical. Notice that by the definition of the generalized inverse of F_{μ} it holds that

$$F_{\mu}^{[-1]}(p) \le x \quad \Longleftrightarrow \quad p \le F_{\mu}(x).$$
 (11)

It is well-known that the CDF of the order statistic $z_{i:n}$ is

$$F_{z_{i:n}}(t) = \mathbb{P}(z_{i:n} \le t) = \sum_{j=i}^{n} \binom{n}{j} (F_{\xi}(t))^{j} (1 - F_{\xi}(t))^{n-j}.$$

By Theorem 3, we know that $G^*(t) = (F_{\nu}^{[-1]} \circ F_{\xi})(t)$. Thus, with (11), one can observe that the CDF of $G^*(z_{i:n})$ satisfies

$$F_{G^*(z_{i:n})}(t) = \mathbb{P}(G^*(z_{i:n}) \le t)$$

= $\mathbb{P}((F_{\nu}^{[-1]} \circ F_{\xi})(z_{i:n}) \le t)$
= $\mathbb{P}(F_{\xi}(z_{i:n}) \le F_{\nu}(t)).$

As we assume that ξ is absolutely continuous with respect to the Lebesgue measure, F_{ξ} has a well-defined inverse. Hence, continuing from the above, we get

$$\begin{aligned} F_{G^*(z_{i:n})}(t) &= \mathbb{P}(F_{\xi}(z_{i:n}) \leq F_{\nu}(t)) \\ &= \mathbb{P}(z_{i:n} \leq (F_{\xi}^{-1} \circ F_{\nu})(t)) \\ &= \sum_{j=i}^{n} \binom{n}{j} \left((F_{\xi} \circ F_{\xi}^{-1} \circ F_{\nu})(t) \right)^{j} \left(1 - (F_{\xi} \circ F_{\xi}^{-1} \circ F_{\nu})(t) \right)^{n-j} \\ &= \sum_{j=i}^{n} \binom{n}{j} \left(F_{\nu}(t) \right)^{j} \left(1 - F_{\nu}(t) \right)^{n-j} \\ &= \mathbb{P}(x_{i:n} \leq t) \\ &= F_{x_{i:n}}(t). \end{aligned}$$

Therefore, the distributions of $G^*(z_{i:n})$ and $x_{i:n}$ are identical.

A.2. Proof of Lemma 1

Let us begin with some simple observations. As usual, we denote $x_1, \ldots, x_n \stackrel{\text{i.i.d.}}{\sim} \nu$, and let $x_{1:n}, \ldots, x_{n:n}$ be their order statistics.

Proposition 5 Let ν be a probability distribution with $\mathbb{E}_{x \sim \nu}[|x|] < \infty$. Then for any i = 1, ..., n, *it holds that* $\mathbb{E}[x_{i:n}] \leq n \mathbb{E}[|x_1|]$.

Proof From the chain of inequalities

$$\mathbb{E}[x_{i:n}] \le \mathbb{E}[|x_{i:n}|] \le \mathbb{E}\left[\sum_{i=1}^{n} |x_{i:n}|\right] = \mathbb{E}\left[\sum_{i=1}^{n} |x_i|\right] = n \mathbb{E}[|x_1|]$$

the bound is immediate.

Proposition 6 Let ν be a probability distribution with a finite second moment. Then for any $1 \leq i, j \leq n$, it holds that $\mathbb{E}[x_{i:n}x_{j:n}] \leq n \mathbb{E}[x_1^2]$.

Proof In a similar manner to the preceding proposition, it holds for any i = 1, ..., n that

$$\mathbb{E}[x_{i:n}^2] \le \mathbb{E}\left[\sum_{i=1}^n x_{i:n}^2\right] = n \,\mathbb{E}[x_1^2]$$

Therefore, by the Cauchy-Schwarz inequality

$$\mathbb{E}[x_{i:n}x_{j:n}] \le \sqrt{\mathbb{E}[x_{i:n}^2]\mathbb{E}[x_{j:n}^2]} \le n \mathbb{E}[x_1^2]$$

and we are done.

With these results, we can prove the following lemma, which will play a key role in the proof of Lemma 1.

Lemma 7 Under Assumption 1, for any $1 \le i, j \le n$, the covariance of $x_{i:n}$ and $x_{j:n}$ is nonnegative. That is,

$$\operatorname{Cov}(x_{i:n}, x_{j:n}) \ge 0. \tag{12}$$

Proof For notational convenience, let $x = (x_1, ..., x_n)$. As $x_1, ..., x_n$ are i.i.d. samples, they are *associated* [10, Theorem 2.1], in the sense that for any two nondecreasing functions f, g with all of $\mathbb{E}[f(x)], \mathbb{E}[g(x)], \text{ and } \mathbb{E}[f(x)g(x)]$ finite, it holds that

$$\operatorname{Cov}(f(\boldsymbol{x}), g(\boldsymbol{x})) \ge 0. \tag{13}$$

Now for each $i = 1, \ldots, n$, define

$$h_i(a_1,\ldots,a_n) = (\text{the } i\text{th smallest value among } a_1,\ldots,a_n)$$

so that $h_i(x_1, ..., x_n) = x_{i:n}$. If we can take $f = h_i$ and $g = h_j$ in (13), then (12) will follow. To this end, observe that for any *i* and *j*, by Theorem 5 and Jensen's inequality,

$$\mathbb{E}[h_i(\boldsymbol{x})] = \mathbb{E}[x_{i:n}] \le n \mathbb{E}[|x_1|] \le n \sqrt{\mathbb{E}[x_1^2]} < \infty,$$

and by Theorem 6,

$$\mathbb{E}[h_i(\boldsymbol{x})h_j(\boldsymbol{x})] = \mathbb{E}[x_{i:n}x_{j:n}] \le n \mathbb{E}[x_1^2] < \infty.$$

Hence, it suffices to show that h_i is nondecreasing for each *i*. But this is immediate from the definition of h_i ; indeed, if any x_k increases (while all the other arguments are held fixed), then the *i*th order statistic cannot decrease. This completes the proof.

We are now ready to prove Lemma 1.

Lemma 8 (Lemma 1) Under Assumption 1, for any $n \ge 1$, the noise in (7) satisfies

$$\frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(\epsilon_i) \le 2\sigma^2.$$

Proof As z_i and x_j are independent for any pair of i and j, we can decompose the variances into

$$\frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(\epsilon_{i}) = \frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(\epsilon_{z_{i:n}} + \epsilon_{x_{i:n}})$$
$$= \frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(\epsilon_{z_{i:n}}) + \frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(\epsilon_{x_{i:n}}).$$

Thus, it suffices to show that both noise components induced by $z_{i:n}$ and $x_{i:n}$ have bounded variances. Bounding the variance of the noises induced by $x_{i:n}$ can be done as

$$\frac{1}{n}\sum_{i=1}^{n}\operatorname{Var}(\epsilon_{x_{i:n}}) = \frac{1}{n}\sum_{i=1}^{n}\operatorname{Var}(x_{i:n} - \mathbb{E}[x_{i:n}])$$

$$= \frac{1}{n}\sum_{i=1}^{n}\operatorname{Var}(x_{i:n})$$

$$= \frac{1}{n}\operatorname{Var}\left(\sum_{i=1}^{n}x_{i:n}\right) - \frac{2}{n}\sum_{1\leq i< j\leq n}\operatorname{Cov}(x_{i:n}, x_{j:n})$$

$$= \frac{1}{n}\operatorname{Var}\left(\sum_{i=1}^{n}x_{i}\right) - \frac{2}{n}\sum_{1\leq i< j\leq n}\operatorname{Cov}(x_{i:n}, x_{j:n})$$

$$\leq \sigma^{2} - \frac{2}{n}\sum_{1\leq i< j\leq n}\operatorname{Cov}(x_{i:n}, x_{j:n})$$

$$\leq \sigma^{2}$$

where the first inequality follows from ν being a distribution with a bounded variance (Assumption 1), and the second inequality is a direct consequence of Theorem 7.

For the other sum, observe that

$$\frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(\epsilon_{i,z}) = \frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(\mathbb{E}[G^{*}(z_{i:n})] - G^{*}(z_{i:n})) = \frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(G^{*}(z_{i:n})).$$

By Lemma 4, for each i = 1, ..., n, we know that $G^*(z_{i:n})$ and $x_{i:n}$ are identically distributed. Hence, we have

$$\sum_{i=1}^{n} \operatorname{Var}(G^{*}(z_{i:n})) = \sum_{i=1}^{n} \operatorname{Var}(x_{i:n}),$$

$$\frac{1}{n} \sum_{i=1}^{n} \operatorname{Var}(\epsilon_{i}) = \frac{2}{n} \sum_{i=1}^{n} \operatorname{Var}(x_{i:n}) \le 2\sigma^{2}$$
(14)

and therefore

which is exactly the claimed bound.

Appendix B. Proof for Section 4.1

B.1. Proof of Theorem 2

Theorem 9 (Theorem 2) Let $\mathcal{E}_{te}(n, p, \delta)$ denote the value of \mathcal{E}_{te} with dataset size n, number of features p, and ridge parameter $\delta \geq 0$. Suppose that $p \leq p'$. Then, under Assumptions 1 and 2, it holds that

$$\min_{s} \mathcal{E}_{te}(n, p', \delta) \le \min_{s} \mathcal{E}_{te}(n, p, \delta).$$

Moreover, if we further assume that p > n, then denoting $\mathcal{E}_{te, 0} = \lim_{\delta \to 0+} \mathcal{E}_{te}$, it holds that

$$\mathcal{E}_{\text{te},0}(n,p') \leq \mathcal{E}_{\text{te},0}(n,p).$$

Proof The first part of the statement is an immediate consequence of Theorem 1 in [34], which states that if σ^2 is a fixed constant then $p \mapsto \min_{\delta} \mathcal{E}_{te}(n, p, \delta)$ is nonincreasing. In our regression model, the noise may depend on the dataset size n, but it is clear that it does not depend on the number of features p. Hence, if n is kept fixed, $\check{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \operatorname{Var}(\epsilon_i)$ in (9) will also be a fixed constant, and therefore, the nonincreasing property of \mathcal{E}_{te} implies the claimed inequality.

The second part of the statement follows from Proposition 2 in [34], under a similar logic. Indeed, if n is kept fixed, $\check{\sigma}^2$ will also be fixed, so in the limit as $\delta \to 0+$, we can apply [34, Proposition 2] to conclude that $\frac{\partial \mathcal{E}_{\text{te},0}}{\partial p} \leq 0$ when p > n. The claimed inequality is then immediate.

Appendix C. Vector-Valued Kernel Regression: Prerequisites for Section 4.2

C.1. Vector-valued linear models

As a preliminary step, let us briefly review how the kernel method can be extended to when the targets, or labels, are multidimensional. For further details, see, for example, [25].

Suppose we are given a dataset $\{(z_i, x_i)\}_{i=1,...,n} \subset \mathbb{R}^d \times \mathbb{R}^m$. In finding a regressor, we consider linear models, which are functions of the form $f(z) = \psi(z)^\top \beta$ for some *matrix-valued* feature map $\psi : \mathbb{R}^d \to \mathbb{R}^{p \times m}$, so that the model is linear on the *p*-dimensional parameter $\beta \in \mathbb{R}^p$.

Such a feature map induces a *matrix-valued* kernel $k(z, z') = \psi(z)^{\top} \psi(z') \in \mathbb{R}^{m \times m}$. By mimicking the construction of the reproducing Hilbert kernel space (RKHS) associated with scalar-valued kernels, one can show the existence of the RKHS \mathcal{H} associated with k, in the sense that

(i) \mathcal{H} is a Hilbert space consisting of functions $\mathbb{R}^d \to \mathbb{R}^p$,

- (ii) there exists a function $\kappa : \mathbb{R}^d \to \mathcal{H}$ such that $\mathbf{k}(\mathbf{z}, \mathbf{z}') = \langle \kappa(\mathbf{z}), \kappa(\mathbf{z}') \rangle_{\mathcal{H}} \ \forall \mathbf{z}, \mathbf{z}' \in \mathbb{R}^d$, where $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is the inner product defined on \mathcal{H} , and
- (iii) for any $h \in \mathcal{H}$ and $z \in \mathbb{R}^d$, denoting by h_j the *j*th component function of h and by e_j the *j*th standard basis vector of \mathbb{R}^m , it holds that

$$\langle \boldsymbol{h}, \boldsymbol{k}(\cdot, \boldsymbol{z}) \boldsymbol{e}_j \rangle_{\mathcal{H}} = h_j(\boldsymbol{z}).$$
 (15)

Notice that the empirical kernel matrix (or the Gram matrix) becomes a block matrix of the form

$$\hat{oldsymbol{K}} = egin{bmatrix} oldsymbol{k}(oldsymbol{z}_1,oldsymbol{z}_1) & \cdots & oldsymbol{k}(oldsymbol{z}_1,oldsymbol{z}_n) \ dots & \ddots & dots \ oldsymbol{k}(oldsymbol{z}_n,oldsymbol{z}_1) & \cdots & oldsymbol{k}(oldsymbol{z}_n,oldsymbol{z}_n) \end{bmatrix} \in \mathbb{R}^{nm imes nm}.$$

Now we consider the empirical risk minimization problem, applied to this setting. For a loss function $\ell : \mathbb{R}^m \times \mathbb{R}^m \to \mathbb{R}$, one can either consider an empirical risk minimization problem of finding an optimal parameter,

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \ell(\boldsymbol{x}_i, \boldsymbol{\psi}(\boldsymbol{z}_i)^{\top} \boldsymbol{\beta}) + \frac{\delta}{2} \|\boldsymbol{\beta}\|^2, \qquad (16)$$

or minimizing the empirical risk directly over the RKHS,

$$\min_{\boldsymbol{f}\in\mathcal{H}} \sum_{i=1}^{n} \ell(\boldsymbol{x}_i, \boldsymbol{f}(\boldsymbol{z}_i)) + \frac{\delta}{2} \|\boldsymbol{f}\|_{\mathcal{H}}^2.$$
(17)

The representer theorem also holds for vector-valued outputs. In other words, there exist vectors $a_1, \ldots, a_n \in \mathbb{R}^m$ such that the empirical risk minimization problem (16) has an optimal solution $\hat{\beta}$ which can be written as

$$\hat{oldsymbol{eta}} = \sum_{j=1}^n \psi(oldsymbol{z}_j)oldsymbol{a}_j,$$

and moreover, for the same vectors a_1, \ldots, a_n , the function

$$\hat{f} = \sum_{j=1}^n oldsymbol{k}(\,\cdot\,,oldsymbol{z}_j)oldsymbol{a}_j$$

becomes an optimal solution \hat{f} to (17). In particular, for our purpose, we set ℓ to be the mean squared error $\ell(x, x') = \frac{1}{2} ||x - x'||^2$. Then, (16) becomes an unconstrained convex quadratic programming, for which we readily have a closed-form solution

$$\mathbf{c} = (\hat{oldsymbol{K}} + \delta oldsymbol{I})^{-1}$$
x.

Here \mathfrak{o} and \mathfrak{x} are concatenations of a_1, \ldots, a_n and x_1, \ldots, x_n , respectively, into vectors in \mathbb{R}^{nm} .

C.2. Extending the eigenlearning framework

To simplify our narration, we act as if the data points z_1, \ldots, z_n are sampled without replacement from a discrete set $\mathcal{X} \subset \mathbb{R}^d$ with $|\mathcal{X}| = M$ very large. This is a strategy also taken in [33], and the arguments supporting the validity of this choice therein still apply.

Under very mild assumptions, the kernel k admits a Mercer-type decomposition

$$\boldsymbol{k}(\boldsymbol{z}, \boldsymbol{z}') = \sum_{i} \lambda_{i} \boldsymbol{\phi}_{i}(\boldsymbol{z}) \boldsymbol{\phi}_{i}(\boldsymbol{z}')^{\top}.$$
(18)

For example, only by assuming that ψ is continuous and $z \mapsto ||k(z, z)||$ is locally bounded one can get a decomposition with each eigenvector ϕ_i being continuous [8].

The decomposition (18) of the kernel k leads to a decomposition $K = \Phi^{\top} \Lambda \Phi$ of the empirical kernel matrix, where Φ is the design matrix having a block matrix form

$$oldsymbol{\Phi} = egin{bmatrix} -- oldsymbol{\phi}_1(oldsymbol{z}_1)^ op & \cdots & -- oldsymbol{\phi}_1(oldsymbol{z}_n)^ op & -- \ dots & dots & dots & dots \ -- oldsymbol{\phi}_M(oldsymbol{z}_1)^ op & \cdots & -- oldsymbol{\phi}_M(oldsymbol{z}_n)^ op & -- \end{bmatrix}$$

and $\mathbf{\Lambda} = \operatorname{diag}(\lambda_1, \ldots, \lambda_M).$

Let $f = \sum_i v_i \phi_i$ and $\hat{f} = \sum_i \hat{v}_i \phi_i$ be the expansions of the target function f and the estimator \hat{f} , respectively, with respect to the eigenbasis. With some algebra (which we omit here, as the derivations are mostly identical to what we detail in Appendix C.3), one can verify that the identity

$$\hat{v} = \underbrace{\mathbf{\Lambda \Phi} (\mathbf{\Phi}^{ op} \mathbf{\Lambda \Phi} + \delta I)^{-1} \mathbf{\Phi}^{ op}}_{=:T} v$$

still holds mutatis mutandis, with the learning transfer matrix T.

C.3. Extending the RF eigenframework

Let us now move on to random feature models, where the feature map ψ is possibly nondeterministic. More precisely, we assume that for some function $g : \mathbb{R}^h \times \mathbb{R}^d \to \mathbb{R}^m$ and a probability distribution ρ on \mathbb{R}^h , the feature map ψ is given by

$$\boldsymbol{\psi}(\boldsymbol{z}) = \begin{bmatrix} -\boldsymbol{g}(\boldsymbol{w}_1, \boldsymbol{z})^\top & - \\ \vdots \\ -\boldsymbol{g}(\boldsymbol{w}_p, \boldsymbol{z})^\top & - \end{bmatrix}$$
(19)

where w_1, \ldots, w_p are i.i.d. samples from ρ .

Due to the targets being multidimensional, we have to deal with vector spaces of \mathbb{R}^m -valued functions. In particular, we have to work with $L^2(\xi; \mathbb{R}^m)$, the vector space of \mathbb{R}^m -valued functions that are "square-integrable with respect to ξ ". Strictly speaking, to define measurability and integrability of Banach-space-valued functions in general, one has to introduce the concept of *Bochner integrals*; see, *e.g.*, [42, Sections V.4–5]. Fortunately, when the codomain is a Euclidean space \mathbb{R}^m , the measurability (resp., integrability) of a function reduces to the measurability (resp., integrability)

of each of the component functions. Under this notion of integrability, $L^2(\xi; \mathbb{R}^m)$ becomes a Hilbert space once we give the inner product as

$$\langle \boldsymbol{f}, \boldsymbol{h} \rangle_{L^2(\xi;\mathbb{R}^m)} \coloneqq \int \boldsymbol{f}(\boldsymbol{z})^\top \boldsymbol{h}(\boldsymbol{z}) \ d\xi(\boldsymbol{z}).$$

In the RF eigenlearning framework [34] which considers the case where m = 1, the singular value decomposition (SVD) of g(w, z), a scalar-valued function in that case, plays an important role. Thus, a key step in extending this framework is to show that SVD is also possible for vector-valued functions.

Theorem 10 (Singular Value Decomposition of Vector-Valued Functions) Assume that $L^2(\rho)$ and $L^2(\xi; \mathbb{R}^m)$ are separable, and $\mathbf{g} \in L^2(\rho \times \xi; \mathbb{R}^m)$. Then, there exist positive numbers $\sigma_1 \geq \sigma_2 \geq \cdots$ and orthonormal sequences $\{\zeta_i\}_i \subset L^2(\rho)$ and $\{\phi_i\}_i \subset L^2(\xi; \mathbb{R}^m)$ such that

$$\boldsymbol{g}(\boldsymbol{w}, \boldsymbol{z}) = \sum_{i} \sigma_{i} \zeta_{i}(\boldsymbol{w}) \phi_{i}(\boldsymbol{z})$$
(20)

where the sum is either finite or converges in the L^2 norm.

Proof Consider an integral operator $\mathcal{I}: L^2(\rho) \to L^2(\xi; \mathbb{R}^m)$ defined as

$$\mathcal{I}(f) = \int \boldsymbol{g}(\boldsymbol{w},\,\cdot\,)f(\boldsymbol{w})\;d
ho(\boldsymbol{w}),$$

and denote the component operators of \mathcal{I} by $\mathcal{I}(f) = (\mathcal{I}_1(f), \ldots, \mathcal{I}_m(f))$. As g is coordinatewise square-integrable, each of $\mathcal{I}_1, \ldots, \mathcal{I}_m$ is compact [20, Theorem 22.4]. Now let f_1, f_2, \ldots be a bounded sequence in $L^2(\rho)$. By passing into a subsequence m times, with at the *j*th pass ensuring that the image of the obtained subsequence under \mathcal{I}_j converges, we can find a subsequence, say $\{f_{k_j}\}_j$, of $\{f_k\}_k$ such that all of $\{\mathcal{I}_1(f_{k_j})\}_j, \ldots, \{\mathcal{I}_m(f_{k_j})\}_j$ converge. This is equivalent to the convergence of $\{\mathcal{I}(f_{k_j})\}_j$, showing that \mathcal{I} itself is also a compact operator. Then \mathcal{I} admits a singular value decomposition [9, Theorem 1.6]; there exist positive numbers $\sigma_1 \ge \sigma_2 \ge \ldots$ and orthonormal sequences $\{\zeta_i\}_i \subset L^2(\rho)$ and $\{\phi_i\}_i \subset L^2(\xi; \mathbb{R}^m)$ such that

$$\mathcal{I}(\cdot) = \sum_{i} \sigma_i \langle \zeta_i, \cdot \rangle_{L^2(\rho)} \phi_i$$
(21)

where the sum converges in the operator norm.

Let $\{\eta_i\}_i$ an orthonormal sequence in $L^2(\rho)$ so that $\{\zeta_i\} \cup \{\eta_i\}$ is an orthonormal basis of $L^2(\rho)$. Similarly, let $\{\chi_i\}_i \subset L^2(\xi; \mathbb{R}^m)$ be so that $\{\phi_i\}_i \cup \{\chi_i\}_i$ is an orthonormal basis of $L^2(\xi; \mathbb{R}^m)$. For the moment, let us assume that the right hand side of (21) is an infinite sum. For each k = 1, 2, ..., let $g_k(w, z) = \sum_{i=1}^k \sigma_i \zeta_i(w) \phi_i(z), r_{k+1} = g - g_k$, and

$$\mathcal{J}_k(f) = \int \boldsymbol{r}_{k+1}(\boldsymbol{w},\,\cdot\,)f(\boldsymbol{w})\,d
ho(\boldsymbol{w})$$

Then by (21), we have

$$\mathcal{J}_{k}(\zeta_{i}) = \int \boldsymbol{g}(\boldsymbol{w}, \cdot)\zeta_{i}(\boldsymbol{w}) \, d\rho(\boldsymbol{w}) - \int \boldsymbol{g}_{k}(\boldsymbol{w}, \cdot)\zeta_{i}(\boldsymbol{w}) \, d\rho(\boldsymbol{w})$$

$$= \mathcal{I}(\zeta_{i}) - \sum_{j=1}^{k} \int \zeta_{i}(\boldsymbol{w})\zeta_{j}(\boldsymbol{w})\phi_{j} \, d\rho(\boldsymbol{w})$$

$$= \sigma_{i}\phi_{i} - \sigma_{i}\phi_{i}\mathbb{1}_{\{i \leq k\}}$$

$$= \sigma_{i}\phi_{i}\mathbb{1}_{\{i > k\}}.$$
 (22)

For any $i, j = 1, 2, ..., using \zeta_i \phi_j$ to denote $(w, z) \mapsto \zeta_i(w) \phi_j(z)$, it follows that

$$\langle \boldsymbol{r}_{k+1}, \zeta_i \boldsymbol{\phi}_j \rangle_{L^2(\rho \times \xi; \mathbb{R}^m)} = \iint \boldsymbol{r}_{k+1}(\boldsymbol{w}, \boldsymbol{z})^\top \left(\zeta_i(\boldsymbol{w}) \boldsymbol{\phi}_j(\boldsymbol{z})\right) \, d(\rho \times \xi)(\boldsymbol{w}, \boldsymbol{z})$$

$$= \iint \left(\int \boldsymbol{r}_{k+1}(\boldsymbol{w}, \boldsymbol{z}) \zeta_i(\boldsymbol{w}) \, d\rho(\boldsymbol{w})\right)^\top \boldsymbol{\phi}_j(\boldsymbol{z}) \, d\xi(\boldsymbol{z})$$

$$= \int \mathcal{J}_k(\zeta_i)^\top \boldsymbol{\phi}_j \, d\xi$$

$$= \langle \mathcal{J}_k(\zeta_i), \boldsymbol{\phi}_j \rangle_{L^2(\xi; \mathbb{R}^m)}$$

$$= \sigma_i \delta_{ij} \mathbb{1}_{\{i > k\}}.$$

$$(23)$$

Let us use analogous notations for $\eta_i \phi_j$, $\zeta_i \chi_j$, and $\eta_i \chi_j$. As $\langle \zeta_I, \eta_i \rangle_{L^2(\rho)} = 0$ for any pair (I, i), it is immediate from (21) that $\mathcal{I}(\eta_i) = \mathbf{0}$. Following the same steps as in (22) and (23), we also get $\mathcal{J}_k(\eta_i) = \mathbf{0}$ and hence

$$\langle \boldsymbol{r}_{k+1}, \eta_i \boldsymbol{\phi}_j \rangle_{L^2(\rho \times \xi; \mathbb{R}^m)} = 0.$$
⁽²⁴⁾

Moreover, as $\langle \phi_I, \chi_i \rangle_{L^2(\xi;\mathbb{R}^m)} = 0$ for any pair (I, i), for any $\theta \in \{\zeta_i\} \cup \{\eta_i\}$, from (21) we get $\langle \mathcal{I}(\theta), \chi_j \rangle_{L^2(\xi;\mathbb{R}^m)} = 0$, and by the construction of g_k we further have $\langle \mathcal{J}_k(\theta), \chi_j \rangle_{L^2(\xi;\mathbb{R}^m)} = 0$. So it follows that

$$\langle \boldsymbol{r}_{k+1}, \boldsymbol{\theta} \boldsymbol{\chi}_j \rangle_{L^2(\rho \times \xi; \mathbb{R}^m)} = \langle \mathcal{J}_k(\boldsymbol{\theta}), \boldsymbol{\chi}_j \rangle_{L^2(\xi; \mathbb{R}^m)} = 0.$$
 (25)

As $\{\zeta_i \phi_j\}_{i,j} \cup \{\eta_i \phi_j\}_{i,j} \cup \{\zeta_i \chi_j\}_{i,j} \cup \{\eta_i \chi_j\}_{i,j}$ is an orthonormal basis of $L^2(\rho \times \xi; \mathbb{R}^m)$ with (24) and (25), we conclude that

$$\begin{aligned} \|\boldsymbol{r}_{k+1}\|_{L^2(\rho\times\xi;\mathbb{R}^m)}^2 &= \sum_{i,j} \left| \langle \boldsymbol{r}_{k+1}, \zeta_i \boldsymbol{\phi}_j \rangle_{L^2(\rho\times\xi;\mathbb{R}^m)} \right|^2 \\ &= \sum_{i,j} \sigma_i^2 \delta_{ij} \mathbb{1}_{\{i>k\}} \\ &= \sum_{i>k} \sigma_i^2. \end{aligned}$$

In particular, by considering when k = 0, because $r_1 = g$ is square-integrable we get $\sum_i \sigma_i^2 < \infty$, and consequently, $\|g - g_k\|_{L^2(\rho \times \xi; \mathbb{R}^m)}^2 = \|r_{k+1}\|_{L^2(\rho \times \xi; \mathbb{R}^m)}^2 \to 0$ as $k \to \infty$. That is,

$$\boldsymbol{g}(\boldsymbol{w}, \boldsymbol{z}) = \sum_{i} \sigma_i \zeta_i(\boldsymbol{w}) \boldsymbol{\phi}_i(\boldsymbol{z}), \qquad (26)$$

where the sum converges in the L^2 norm. Now, if the right hand side of (21) were a finite sum, say *i* ranging from 1 to *I*, the same logic applies, but it suffices to consider *k* ranging only up to *I*, and we would have $g - g_I = 0$. Still, we do have (26) in this case also, hence the proof is complete.

To avoid unnecessary complexities, it is desirable that $\{\zeta_i\}_i$ and $\{\phi_i\}$ are orthonormal bases of their respective ambient spaces. To this end, from now on let us assume that the dimensions of $L^2(\rho)$ and $L^2(\xi; \mathbb{R}^m)$ are equal, so that even if we extend $\{\zeta_i\}_i$ and $\{\phi_i\}_i$ into orthonormal bases, we still have the expansion of g as in (20) by allowing $\sigma_i = 0$ if necessary. For example, under the mild assumption that ρ and ξ each assign nonzero measures to a countably infinite number of disjoint sets in \mathbb{R}^h and \mathbb{R}^d respectively, the dimensions of both spaces will be countably infinite.

Equipped with the SVD of vector-valued functions, we have the empirical kernel

$$\hat{\boldsymbol{k}}(\boldsymbol{z}, \boldsymbol{z}') = \frac{1}{p} \boldsymbol{\psi}(\boldsymbol{z})^{\top} \boldsymbol{\psi}(\boldsymbol{z}') = \frac{1}{p} \sum_{i=1}^{p} \boldsymbol{g}(\boldsymbol{w}_{i}, \boldsymbol{z}) \boldsymbol{g}(\boldsymbol{w}_{i}, \boldsymbol{z}')^{\top}$$
$$= \frac{1}{p} \sum_{j,j'} \sum_{i=1}^{p} \sigma_{j} \sigma_{j'} \zeta_{j}(\boldsymbol{w}_{i}) \zeta_{j'}(\boldsymbol{w}_{i}) \phi_{j}(\boldsymbol{z}) \phi_{j'}(\boldsymbol{z}')^{\top}$$
(27)

and the deterministic kernel

$$\boldsymbol{k}(\boldsymbol{z}, \boldsymbol{z}') = \mathbb{E}_{\boldsymbol{w}}[\boldsymbol{g}(\boldsymbol{w}, \boldsymbol{z})\boldsymbol{g}(\boldsymbol{w}, \boldsymbol{z}')^{\top}] = \sum_{j} \lambda_{j} \phi_{j}(\boldsymbol{z}) \phi_{j}(\boldsymbol{z}')^{\top}$$
(28)

where we set $\lambda_j = \sigma_j^2$. Thanks to the similarity of the equations (27) and (28) to their counterparts in the scalar-valued regression RF eigenframework [34], the remaining steps are mostly straightforward, resembling what we did in Appendix C.2 to extend the scalar-valued framework into the vector-valued setting. Define $\Lambda := \text{diag}(\lambda_1, \lambda_2, ...)$ and

$$oldsymbol{Z}\coloneqq egin{bmatrix} \zeta_1(oldsymbol{w}_1)&\cdots&\zeta_1(oldsymbol{w}_p)\ \zeta_2(oldsymbol{w}_1)&\cdots&\zeta_2(oldsymbol{w}_p)\ dots&dots&dots&dots\end{pmatrix}$$

so that $\tilde{\mathbf{\Lambda}} \coloneqq \frac{1}{p} \mathbf{\Lambda}^{1/2} \mathbf{Z} \mathbf{Z}^{\top} \mathbf{\Lambda}^{1/2}$ is a matrix whose entries are $\tilde{\Lambda}_{jj'} = \frac{1}{p} \sum_{i=1}^{p} \sigma_j \sigma_{j'} \zeta_j(\mathbf{w}_i) \zeta_{j'}(\mathbf{w}_i)$. Then for a design matrix

$$oldsymbol{\Phi} = egin{bmatrix} - oldsymbol{\phi}_1(oldsymbol{z}_1)^ op & \cdots & - oldsymbol{\phi}_1(oldsymbol{z}_n)^ op & \cdots \ - oldsymbol{\phi}_2(oldsymbol{z}_n)^ op & \cdots & - oldsymbol{\phi}_2(oldsymbol{z}_n)^ op & \cdots \ dots & dots &$$

it is clear from (27) that the empirical kernel matrix $\hat{K} = \left[\hat{k}(z_i, z_j)\right]_{i,j}$ can be written as

$$\hat{K} = \Phi^{+} \hat{\Lambda} \Phi_{+}$$

Again, let $f = \sum_i v_i \phi_i$ and $\hat{f} = \sum_i \hat{v}_i \phi_i$ be the expansions with respect to the eigenbasis of the target function f and the estimator \hat{f} , respectively. Then for each J = 1, 2, ..., noting that the

kernel we actually use in the regression is the empirical kernel \hat{k} , one can observe that

$$egin{aligned} \hat{v}_J &= \left\langle \phi_J, \hat{f}
ight
angle_{L^2(\xi; \mathbb{R}^d)} = \left\langle \phi_J, \sum_{s=1}^n \hat{k}(\,\cdot\,, oldsymbol{z}_s) oldsymbol{a}_s
ight
angle_{L^2(\xi; \mathbb{R}^d)} \ &= \sum_{s=1}^n rac{1}{p} \sum_{j'} \sum_{i=1}^p \sigma_J \sigma_{j'} \zeta_J(oldsymbol{w}_i) \zeta_{j'}(oldsymbol{w}_i) \phi_{j'}(oldsymbol{z}_s)^ op oldsymbol{a}_s \ &= \sum_{s=1}^n \sum_{j'} ilde{oldsymbol{\Lambda}}_{Jj'} \phi_{j'}(oldsymbol{z}_s)^ op oldsymbol{a}_s \ &= \left[ilde{oldsymbol{\Lambda}} oldsymbol{\Phi}_{\mathsf{G}}
ight]_J \end{aligned}$$

where $[\cdot]_J$ denotes the *J*th component of a vector. Therefore, when there is no noise in the targets (*i.e.*, labels), the identity

$$\hat{oldsymbol{v}} = \underbrace{ ilde{oldsymbol{\Delta} \Phi} \left(\Phi^ op ilde{oldsymbol{\Delta} \Phi} + \delta oldsymbol{I}
ight)^{-1} \Phi^ op v}_{=: ilde{oldsymbol{T}}} oldsymbol{v}$$

still holds mutatis mutandis, with the learning transfer matrix \tilde{T} .

For further analyses, it would be desirable to have an assumption similar to the Gaussian universality ansatz (Assumption 2). The orthonormality equations

$$\mathbb{E}_{\boldsymbol{w}\sim\rho}[\zeta_i(\boldsymbol{w})\zeta_j(\boldsymbol{w})] = \delta_{ij}$$
(29a)

$$\mathbb{E}_{\boldsymbol{z} \sim \boldsymbol{\xi}}[\boldsymbol{\phi}_i(\boldsymbol{z})^\top \boldsymbol{\phi}_j(\boldsymbol{z})] = \delta_{ij}$$
(29b)

suggest a natural multivariate rendition of the ansatz as follows.

Assumption M2 (Multivariate Gaussian Universality Ansatz) The expected population risk remains unchanged even if we replace $\{\phi_i\}$ and $\{\zeta_i\}$ each with random Gaussian functions $\{\tilde{\phi}_i\}$ and $\{\tilde{\zeta}_i\}$, in the sense that $\{\tilde{\zeta}_i(\boldsymbol{w})\}$ become i.i.d. samples from $\mathcal{N}(0, 1)$ when $\boldsymbol{w} \sim \mu$, and $\{\tilde{\phi}_i(\boldsymbol{z})\}$ become i.i.d. samples from $\mathcal{N}(\mathbf{0}, \frac{1}{m}\boldsymbol{I})$ when $\boldsymbol{z} \sim \xi$.

Remark 11 The covariance matrix associated to ϕ is $\frac{1}{m}I$ instead of I, as we wish to have $\mathbb{E}\left[\|\tilde{\phi}_i(z_j)\|^2\right] = 1$, based on (29b).

Under the multivariate Gaussian universality ansatz, Φ is a matrix whose entries are i.i.d. samples from $\mathcal{N}(0, \frac{1}{m})$, or equivalently, $\frac{1}{\sqrt{m}}\mathcal{N}(0, 1)$. Hence, rewriting the learning transfer matrix in terms of $\mathbf{\Omega} = \sqrt{m} \Phi$ as

$$ilde{oldsymbol{T}} = ilde{oldsymbol{\Lambda}} oldsymbol{\Omega} \left(oldsymbol{\Omega}^{ op} ilde{oldsymbol{\Lambda}} oldsymbol{\Omega} + m \delta oldsymbol{I}
ight)^{-1} oldsymbol{\Omega}^{ op}$$

shows that \tilde{T} is statistically equivalent to the learning transfer matrix appearing in the scalar-valued RF eigenframework, up to the number of samples n and the ridge parameter δ being scaled by a factor of m. That is, we can apply the results from the RF eigenframework [34], with necessary modifications, as follows.

Let $s := \sum_{i} \frac{\lambda_i}{\lambda_i + \gamma}$, $q := \sum_{i} \left(\frac{\lambda_i}{\lambda_i + \gamma}\right)^2$, and $\kappa, \gamma \ge 0$ be the unique nonnegative scalars such that

$$nm = s + \frac{m\delta}{\kappa}$$
 and $p = s + \frac{p\kappa}{\gamma}$.

The test error of the vector-valued RF regression is then given approximately by

$$\mathbb{E}_{\mathcal{D}}\left[\mathcal{R}(\hat{\boldsymbol{\beta}})\right] \approx \mathcal{E}_{\mathsf{te}} \coloneqq \frac{1}{1 - \frac{q(p-2s) + s^2}{nm(p-q)}} \left(\sum_{i} \left(\frac{\gamma}{\lambda_i + \gamma} - \frac{\kappa\lambda_i}{(\lambda_i + \gamma)^2} \frac{p}{p-q}\right) v_i^2 + \check{\sigma}^2\right) \tag{30}$$

where $\check{\sigma}^2$ is the mean squared error of the noise, or in other words, the part of the data that is effectively noise.

C.4. Neural-network-like random feature models

When the labels are scalar-valued, a noteworthy interpretation of linear models that reveals their connection to modern machine learning models is as 2-layer (fully-connected feed-forward) neural networks, with a fixed first layer and a trainable second layer. Thus, it is natural to ask if this is also the case in the vector-valued case. The answer is yes, but with a slight caveat, which we now detail.

A 2-layer network with p hidden neurons, taking d-dimensional inputs and producing m-dimensional outputs, can be modeled as $f(z) = B\sigma(Wz + c)$, for weights $W \in \mathbb{R}^{p \times d}$, $c \in \mathbb{R}^{p}$, $B \in \mathbb{R}^{m \times p}$, and an activation function σ applied elementwise. As f is linear on B, it indeed is a linear model, and this becomes more apparent from the reformulation

$$\boldsymbol{f}(\boldsymbol{z}) = \begin{bmatrix} \sigma(\boldsymbol{W}\boldsymbol{z} + \boldsymbol{b})^\top & \boldsymbol{0} & \dots & \boldsymbol{0} \\ \boldsymbol{0} & \sigma(\boldsymbol{W}\boldsymbol{z} + \boldsymbol{b})^\top & \dots & \boldsymbol{0} \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{0} & \boldsymbol{0} & \dots & \sigma(\boldsymbol{W}\boldsymbol{z} + \boldsymbol{b})^\top \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \vdots \\ \boldsymbol{\beta}_m \end{bmatrix}$$

where β_i is the *i*th row vector of **B** in the column vector form. In other words, **f** is a linear model with a block diagonal feature map.

That said, this is exactly where the caveat arises. The issue is that the feature map is not of the form of (19), as that formulation cannot produce the precise pattern of zeros we now require. This shows that a separate analysis from Appendix C.3 is required for this case.

To this end, let us consider a setting where the feature map now is a block diagonal matrix

$$\boldsymbol{\psi}(\boldsymbol{z}) = \begin{bmatrix} \boldsymbol{\psi}_1(\boldsymbol{z}) & & \\ & \ddots & \\ & & \boldsymbol{\psi}_1(\boldsymbol{z}) \end{bmatrix} \in \mathbb{R}^{mp \times m}.$$
(31)

The notation ψ_1 is used to emphasize that $\psi_1(z)$ takes the same form as feature maps (19) for when m = 1. In other words, as in the scalar-valued labels setting, for some scalar valued map $g: \mathbb{R}^h \times \mathbb{R}^d \to \mathbb{R}$ we consider when $\psi_1(z) = (g(w_1, z), \dots, g(w_p, z))$ with $w_1, \dots, w_p \stackrel{\text{i.i.d.}}{\sim} \rho$. Let us call the linear models that arise from (random) feature maps having such block diagonal form the *neural-network-like linear models*. For convenience, we also let $\Psi_1 := [\psi_1(z_1) \cdots \psi_1(z_n)]^\top$. Similar to before, we define the empirical kernel as $\hat{k}(z, z') = \frac{1}{p} \psi(z)^{\top} \psi(z')$,² then we get

$$egin{aligned} \hat{m{k}}(m{z},m{z}') &= rac{1}{p}m{\psi}_1(m{z})^{ op}m{\psi}_1(m{z}')m{I} \ &= rac{1}{p}\sum_{i=1}^p g(m{w}_i,m{z})g(m{w}_i,m{z}')m{I} \end{aligned}$$

Observe that $\frac{1}{p} \sum_{i=1}^{p} g(\boldsymbol{w}_i, \boldsymbol{z}) g(\boldsymbol{w}_i, \boldsymbol{z}')$ is in fact what we would get as the empirical kernel when m = 1 so that ψ_1 is used as the feature map. Thus, for $\hat{\boldsymbol{K}}_1$ denoting the matrix whose entries are $\hat{K}_{j,j'} = \frac{1}{p} \sum_{i=1}^{p} g(\boldsymbol{w}_i, \boldsymbol{z}_j) g(\boldsymbol{w}_i, \boldsymbol{z}_{j'})$, the empirical kernel matrix $\hat{\boldsymbol{K}} = \left[\hat{\boldsymbol{k}}(\boldsymbol{z}_i, \boldsymbol{z}_j) \right]_{i,j}$ is a block diagonal matrix whose block entries are all constant multiples of the identity matrix, that is, $\hat{\boldsymbol{K}} = \hat{\boldsymbol{K}}_1 \otimes \boldsymbol{I}$ where \otimes denotes the Kronecker product of two matrices.

Recall that $\mathbf{x} = (x_{11}, \ldots, x_{1m}, \ldots, x_{n1}, \ldots, x_{nm})$. Define a permutation matrix $\mathbf{P} \in \mathbb{R}^{nm \times nm}$ such that $\mathbf{P}\mathbf{x} = (x_{11}, \ldots, x_{n1}, \ldots, x_{1m}, \ldots, x_{nm})$. For $j = 1, \ldots, m$, denote $\mathbf{y}_j = (x_{1j}, \ldots, x_{nj})$. Then as permutation matrices satisfy $\mathbf{P}^{\top} = \mathbf{P}^{-1}$, it holds that

$$\begin{split} (\hat{K} + \delta I)^{-1} \mathbf{x} &= P^{\top} P(\hat{K} + \delta I)^{-1} P^{\top} P \mathbf{x} \\ &= P^{\top} (P \hat{K} P^{\top} + \delta I)^{-1} P \mathbf{x} \\ &= P^{\top} \begin{bmatrix} \hat{K}_1(z_1, z_1) + \delta I & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \hat{K}_1(z_1, z_1) + \delta I \end{bmatrix}^{-1} P \mathbf{x} \\ &= P^{\top} \begin{bmatrix} (\hat{K}_1(z_1, z_1) + \delta I)^{-1} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & (\hat{K}_1(z_1, z_1) + \delta I)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_m \end{bmatrix}, \end{split}$$

and thus denoting by $\hat{\beta}$ the concatenation of $\hat{\beta}_1, \ldots, \hat{\beta}_n$ into a vector in \mathbb{R}^{mp} , we obtain

$$\begin{split} \hat{\boldsymbol{\beta}} &= \begin{bmatrix} \boldsymbol{\psi}(\boldsymbol{z}_1) & \cdots & \boldsymbol{\psi}(\boldsymbol{z}_n) \end{bmatrix} (\hat{\boldsymbol{K}} + \delta \boldsymbol{I})^{-1} \boldsymbol{x} \\ &= \begin{bmatrix} \boldsymbol{\psi}(\boldsymbol{z}_1) & \cdots & \boldsymbol{\psi}(\boldsymbol{z}_n) \end{bmatrix} \boldsymbol{P}^\top \begin{bmatrix} (\hat{\boldsymbol{K}}_1(\boldsymbol{z}_1, \boldsymbol{z}_1) + \delta \boldsymbol{I})^{-1} \boldsymbol{y}_1 \\ \vdots \\ (\hat{\boldsymbol{K}}_1(\boldsymbol{z}_1, \boldsymbol{z}_1) + \delta \boldsymbol{I})^{-1} \boldsymbol{y}_m \end{bmatrix} \\ &= \begin{bmatrix} \boldsymbol{\Psi}_1^\top & \cdots & \boldsymbol{0} \\ \vdots & \ddots & \vdots \\ \boldsymbol{0} & \cdots & \boldsymbol{\Psi}_1^\top \end{bmatrix} \begin{bmatrix} (\hat{\boldsymbol{K}}_1(\boldsymbol{z}_1, \boldsymbol{z}_1) + \delta \boldsymbol{I})^{-1} \boldsymbol{y}_1 \\ \vdots \\ (\hat{\boldsymbol{K}}_1(\boldsymbol{z}_1, \boldsymbol{z}_1) + \delta \boldsymbol{I})^{-1} \boldsymbol{y}_m \end{bmatrix} \\ &= \begin{bmatrix} \boldsymbol{\Psi}_1^\top (\hat{\boldsymbol{K}}_1(\boldsymbol{z}_1, \boldsymbol{z}_1) + \delta \boldsymbol{I})^{-1} \boldsymbol{y}_1 \\ \vdots \\ \boldsymbol{\Psi}_1^\top (\hat{\boldsymbol{K}}_1(\boldsymbol{z}_1, \boldsymbol{z}_1) + \delta \boldsymbol{I})^{-1} \boldsymbol{y}_m \end{bmatrix}. \end{split}$$

^{2.} As the number of rows in ψ is now mp, one can choose the factor multiplied to the sum to be $1/m_p$ instead of 1/p to maintain consistency with Appendix C.3. However, as the main functionality of that factor is to convert the sum into a mean, and as each column of ψ now contains only p nonzero elements, we prefer 1/p over $1/m_p$. It soon turns out that such a choice also allows us to directly apply one-dimensional results to this setting.

That is, for each i = 1, ..., m, each $\hat{\beta}_i$ only depends on y_i , the collection of the *i*th coordinates from $\{x_1, ..., x_n\}$. In other words, each $\hat{\beta}_i$ is computed independently from others, as if it is obtained from a kernel regression in a scalar-valued labels setting by only looking at the *i*th coordinates in the data. On top of this, we also have

$$egin{aligned} \Sigma &= \mathbb{E}_{oldsymbol{z}\sim \xi}[oldsymbol{\psi}(oldsymbol{z})^{ op}] & oldsymbol{0} & \cdots & oldsymbol{0} \ &= egin{bmatrix} \mathbb{E}_{oldsymbol{z}\sim \xi}[oldsymbol{\psi}_1(oldsymbol{z})^{ op}] & oldsymbol{0} & \cdots & oldsymbol{0} \ &\vdots & \ddots & \vdots \ &oldsymbol{0} & oldsymbol{0} & \cdots & \mathbb{E}_{oldsymbol{z}\sim \xi}[oldsymbol{\psi}_1(oldsymbol{z})^{ op}] \end{bmatrix}, \end{aligned}$$

allowing us to conclude that the population risk is the sum of the coordinatewise population risk.

Meanwhile, notice that the mean squared error of the noise is the sum of the mean squared errors of the coordinatewise noises, no matter the noise model. Hence, estimating the population risk can be done as computing the sum of coordinatewise estimations by \mathcal{E}_{te} as in (9). It follows that the proper assumption to make when studying neural-network-like linear models is the singular value decomposition of the scalar-valued function $g(w, z) = \sum_{i} \sigma_i \zeta_i(w) \phi_i(z)$ and the Gaussian universality ansatz (Assumption 2), not their multivariate versions.

With these established, let us now consider the expansions $\boldsymbol{f} = \left(\sum_{j} v_{1j}\phi_{j}, \dots, \sum_{j} v_{mj}\phi_{j}\right)$ and $\hat{\boldsymbol{f}} = \left(\sum_{j} \hat{v}_{1j}\phi_{j}, \dots, \sum_{j} \hat{v}_{mj}\phi_{j}\right)$. Denote by $\check{\sigma}_{i}^{2}$ the mean squared noise in the *i*th coordinate. Then, from the coordinatewise approximation using (9), we have

$$\mathbb{E}_{\mathcal{D}}[\mathcal{R}(\hat{\beta})] \approx \sum_{i=1}^{m} \frac{1}{1 - \frac{q(p-2s) + s^2}{n(p-q)}} \left(\sum_{j} \left(\frac{\gamma}{\lambda_j + \gamma} - \frac{\kappa \lambda_j}{(\lambda_j + \gamma)^2} \frac{p}{p-q} \right) v_{ij}^2 + \check{\sigma}_i^2 \right)$$
(32)

where, for $s \coloneqq \sum_{j} \frac{\lambda_{j}}{\lambda_{j}+\gamma}$ and $q \coloneqq \sum_{j} \left(\frac{\lambda_{j}}{\lambda_{j}+\gamma}\right)^{2}$, κ and γ are unique nonnegative numbers such that $n = s + \frac{\delta}{\kappa}$ and $p = s + \frac{p\kappa}{\gamma}$. As n, p, q, and s do not depend on the summation index i in the above, denoting by $\check{\sigma}^{2} = \sum_{i=1}^{n} \check{\sigma}_{i}^{2}$ the overall mean squared noise, we conclude that

$$\mathbb{E}_{\mathcal{D}}[\mathcal{R}(\hat{\beta})] \approx \frac{1}{1 - \frac{q(p-2s) + s^2}{n(p-q)}} \left(\sum_{i=1}^{m} \sum_{j} \left(\frac{\gamma}{\lambda_j + \gamma} - \frac{\kappa \lambda_j}{(\lambda_j + \gamma)^2} \frac{p}{p-q} \right) v_{ij}^2 + \check{\sigma}^2 \right).$$
(33)

Remark 12 The same conclusion can be derived by repeating the logic developed in Appendix C.3 but with ψ as described in (31). In particular, the same expansions $\mathbf{f} = \left(\sum_{j} v_{1j}\phi_j, \ldots, \sum_{j} v_{mj}\phi_j\right)$ and $\hat{\mathbf{f}} = \left(\sum_{j} \hat{v}_{1j}\phi_j, \ldots, \sum_{j} \hat{v}_{mj}\phi_j\right)$ will be obtained by considering $\{\phi_j \mathbf{e}_i\}_{i,j}$ as the orthonormal basis of $L^2(\xi; \mathbb{R}^m)$, where \mathbf{e}_i denotes the *i*th standard basis vector of \mathbb{R}^m . This also shows that the multivariate Gaussian universality ansatz is not quite appropriate for neural-network-like linear models; $\phi_j \mathbf{e}_i$ is a sparse vector, whereas $\mathcal{N}(0, \frac{1}{p}\mathbf{I})$ is a full-dimensional distribution. We omit the details of the derivations in this direction.

Appendix D. Missing Details for Section 4.2

D.1. Statements on the random feature models in higher dimensions

Although our empirical results in Appendix E suggest that the bounded noise assumption analogously holds in higher dimensions, we were unable to theoretically extend the boundedness analysis of Lemma 1 beyond the one-dimensional setting. We therefore explicitly assume the following.

Assumption M1 The $\epsilon_i s$ in (10) satisfy $\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \operatorname{Var}(\epsilon_{ij}) \leq \tilde{\sigma}^2$.

Then we can prove Theorem 13, which states that using more features is better for random feature models, detailed in Appendix C.3, in high dimensions also.

Theorem 13 Consider random feature models with targets in \mathbb{R}^m . Let $\mathcal{E}_{te}(n, p, \delta)$ denote the value of \mathcal{E}_{te} in (33) with dataset size n, number of features p, and ridge parameter δ . Under Assumptions *M1* and *M2*, if $p \leq p'$ then $\min_{\delta} \mathcal{E}_{te}(n, p', \delta) \leq \min_{\delta} \mathcal{E}_{te}(n, p, \delta)$, and if moreover p > nm, then $\mathcal{E}_{te,0}(n, p') \leq \mathcal{E}_{te,0}(n, p)$.

Proof Let $\mathcal{E}'_{te}(n, p, \delta)$ denote the value of \mathcal{E}_{te} in the one-dimensional setting (9), then it holds that $\mathcal{E}_{te}(n, p, \delta) = \mathcal{E}'_{te}(nm, p, \delta m)$, given that both sides involve the same $\check{\sigma}^2$. Hence, the statements in this theorem are direct consequences of Theorem 2.

It is clear that a similar result holds for neural-network-like random feature models, in particular because, as we saw in Appendix C.4 the population risk is essentially the aggregation of the coordinatewise risks. For completeness, we formalize this result as follows.

Theorem 14 Consider neural-network-like random feature models. Let $\mathcal{E}_{te}(n, p, \delta)$ denote the value of \mathcal{E}_{te} in (33) with dataset size n, number of features p, and ridge parameter δ . Under Assumptions 2 and M1, if $p \leq p'$ then $\min_{\delta} \mathcal{E}_{te}(n, p', \delta) \leq \min_{\delta} \mathcal{E}_{te}(n, p, \delta)$, and if moreover p > n, then $\mathcal{E}_{te}, 0(n, p') \leq \mathcal{E}_{te}, 0(n, p)$.

Proof Clearly, each $\check{\sigma}_i^2$ is finite, as we assume that $\check{\sigma}^2$ is finite. Recalling (32), we know that \mathcal{E}_{te} in (33) is a sum of *m* quantities of the form specified as \mathcal{E}_{te} in the one-dimensional setting (9). Hence, applying Theorem 2, the results on the one-dimensional setting, to each of those *m* quantities, the conclusions follow.

Appendix E. Experimental Results

In this section, we discuss the claims regarding random feature models, in particular focusing on the approximation of the population risk with \mathcal{E}_{te} as stated in (9) for the one-dimensional case and (30) for higher dimensions, with empirical results.

E.1. One dimensional examples

By choosing $\nu = \text{Uniform}(0, 1)$, we can exploit the following well known fact on order statistics.

Proposition 15 (1, Examples 2.2.1 and 2.3.1) In the case of $\nu = \text{Uniform}(0, 1)$, the order statistics satisfy $x_{i:n} \sim \text{Beta}(i, n - i + 1)$ and $(x_{i:n}, x_{j:n}) \sim \text{BivariateBeta}(i, j - i, n - j + 1)$. In particular, we have $\mathbb{E}[x_{i:n}] = \frac{i}{n+1}$, $\text{Var}(x_{i:n}) = \frac{i(n-i+1)}{(n+1)^2(n+2)}$, and $\text{Cov}(x_{i:n}, x_{j:n}) = \frac{i(n-j+1)}{(n+1)^2(n+2)}$.



Figure 1: Learning the optimal transport maps with random feature models in 1D. We plot the computed theoretical and experimental population risks. For the meaning of the curves and scatter plots, see the discussions in Appendix E.

In the statement of Proposition 15, by BivariateBeta(i, j - i, n - j + 1) we are referring to the so-called bivariate Beta distribution whose probability density function is

$$f(x,y) = \frac{n!}{(i-1)!(j-i-1)!(n-j)!} x^{i-1}(y-x)^{j-i-1}(1-y)^{n-j} \mathbb{1}_{0 < x < y < 1}(x,y)$$

Proposition 15 allows us to explicitly characterize the noise model. In specific, with recalling (14), we have

$$\frac{1}{n}\sum_{i=1}^{n} \operatorname{Var}(\epsilon_i) = \frac{2}{n}\sum_{i=1}^{n} \operatorname{Var}(x_{i:n}) = \frac{2}{n}\sum_{i=1}^{n} \frac{i(n-i+1)}{(n+1)^2(n+2)} = \frac{1}{3(n+1)^2}$$

For the latent distribution ξ , we consider two options, $\mathcal{N}(0,1)$ and Uniform(0,1). As F_{ν} is the identity function, when ξ is a standard Gaussian, by Theorem 3 the optimal transport map is $G^*(x) = F_{\xi}(x) = \frac{1}{2} + \frac{1}{2} \operatorname{erf}(\frac{x}{\sqrt{2}})$. Meanwhile, when ξ is also Uniform(0,1), we have $\xi = \nu$, so the identity function itself is clearly a solution of (5), hence it is the optimal transport map.

The OT maps are learned by random feature models with $g(w, z) = \sin(wz + b)$, where w denotes the pair (w, b) of random samples $w \sim \mathcal{N}(0, 4^2)$ and $b \sim \mathcal{N}(0, 1)$.

Figure 1 shows the plots of the results. The curves drawn with solid lines, labeled theory, are the plots of the estimated population risk \mathcal{E}_{te} , computed by (9) and (30). The scatter plots with vertical error bars, labeled experiment, represent the "true" population risks computed by 200 iterations of Monte Carlo integration. The circle markers are the means, and the error bars range from the 2.5th percentile to the 97.5th percentile of the results from those 200 iterations.

As the captions indicate, Figure 1(*a*) shows the results when $\xi = \mathcal{N}(0, 1)$, and Figure 1(*b*) shows the results when $\xi = \text{Uniform}(0, 1)$. In both experiments, we set n = 64.

From the plots, we can observe that \mathcal{E}_{te} reasonably approximates the population risk $\mathbb{E}_{\mathcal{D}}[\mathcal{R}(\hat{\beta})]$, concurrently capturing the overall tendency of decreasing with respect to the number of features. These results demonstrate that \mathcal{E}_{te} serves as a reliable proxy for the expected population risk, thereby empirically supporting the results regarding "using more features is better" derived via \mathcal{E}_{te} .

E.2. Random feature models in higher dimensions

As an empirical validation of our extension of the theories to the setting of multidimensional targets, we consider the setting where both the latent and the target distributions are multivariate Gaussians, $\xi = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ and $\nu = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$.

One advantage of working with Gaussians is that we have an explicit formula for the OT map, allowing us to seamlessly apply our theories.

Theorem 16 (28, Remark 2.31) For $\xi = \mathcal{N}(\mu_0, \Sigma_0)$ and $\nu = \mathcal{N}(\mu_1, \Sigma_1)$ with Σ_0 invertible, the optimal transport map G^* from ξ to ν is an affine function of the form

$$G^{*}(\boldsymbol{z}) = \boldsymbol{\mu}_{1} + \boldsymbol{\Sigma}_{0}^{-1/2} \left(\boldsymbol{\Sigma}_{0}^{1/2} \boldsymbol{\Sigma}_{1} \boldsymbol{\Sigma}_{0}^{1/2} \right)^{1/2} \boldsymbol{\Sigma}_{0}^{-1/2} (\boldsymbol{z} - \boldsymbol{\mu}_{0}).$$

As an instantiation of the manifold hypothesis, we set Σ_1 to be degenerate, so that ν is supported on an affine subspace whose dimensionality is lower than its ambient space.

The noise part, that is, the term $\check{\sigma}^2$ in the formula (30) of \mathcal{E}_{te} , is estimated in a way described in Section 3.2. More precisely, given samples $z_1, \ldots, z_n \overset{\text{i.i.d.}}{\sim} \xi$ and $x_1, \ldots, x_n \overset{\text{i.i.d.}}{\sim} \nu$, let $\tilde{\pi} \in S_n$ be the permutation of $\{1, \ldots, n\}$ such that the optimal transport map from $\hat{\xi}$ to $\hat{\nu}$ is $z_i \mapsto x_{\tilde{\pi}(i)}$, and for such $\tilde{\pi}$ we set

$$\check{\sigma}^2 = \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n \left\|G^*(\boldsymbol{z}_i) - \boldsymbol{x}_{\tilde{\pi}(i)}\right\|^2\right].$$

In the experiments, the values of $\check{\sigma}^2$ are approximated by Monte Carlo integration using fresh samples that are not used in the regression phase. While we do not have a theoretical proof that Assumption M1 holds in this setting, the results of the Monte Carlo integration serve as empirical evidence supporting its validity.

For the function g determining the (random) feature map ψ in (19), we chose

$$\boldsymbol{g}(\boldsymbol{w}, \boldsymbol{x}) = \begin{bmatrix} \sin(\boldsymbol{\omega}_1^\top \boldsymbol{x} + b_1) \\ \vdots \\ \sin(\boldsymbol{\omega}_p^\top \boldsymbol{x} + b_p) \end{bmatrix}$$
(34)

Т

for \boldsymbol{w} here denoting the collection of all the weights $\omega_1, \ldots, \omega_p \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\boldsymbol{0}, 4\boldsymbol{I})$ and all the biases $b_1, \ldots, b_p \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. We remark that this is essentially ρ being a Gaussian distribution on \mathbb{R}^{dp+p} . Figure 2(*a*) is a result when the ambient space is \mathbb{R}^2 , with

$$\boldsymbol{\mu}_{0} = \begin{bmatrix} 0.5\\0.5 \end{bmatrix}, \quad \boldsymbol{\Sigma}_{0}^{1/2} = \begin{bmatrix} 0.16 & 0\\0 & 0.08 \end{bmatrix},$$
$$\boldsymbol{\mu}_{1} = \begin{bmatrix} -0.5\\-0.5 \end{bmatrix}, \quad \boldsymbol{\Sigma}_{1}^{1/2} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2}\\-1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 0.16 & 0\\0 & 0 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2}\\-1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

In this case, we set n = 64. In each of the Monte Carlo integrations to compute the "true" population risks, 256 samples were used.



Figure 2: Learning the optimal transport maps between multivariate Gaussians with random feature models. We plot the computed theoretical and experimental population risks. The meanings of the curves and scatter plots are the same as Figure 1.

Figure 2(b) is a result when the ambient space is \mathbb{R}^3 , with

$$\boldsymbol{\mu}_{0} = \begin{bmatrix} 0.5\\ 0.5\\ 0.5 \end{bmatrix}, \quad \boldsymbol{\Sigma}_{0}^{1/2} = \frac{1}{4} \begin{bmatrix} 1 & 0 & 0\\ 0 & 0.75 & 0\\ 0 & 0 & 0.5 \end{bmatrix},$$
$$\boldsymbol{\mu}_{1} = \begin{bmatrix} -0.3\\ -0.4\\ -0.5 \end{bmatrix}, \quad \boldsymbol{\Sigma}_{1}^{1/2} = \frac{1}{4} \boldsymbol{R} \begin{bmatrix} 0.75 & 0 & 0\\ 0 & 0.5 & 0\\ 0 & 0 & 0 \end{bmatrix} \boldsymbol{R}^{\top}$$

where \boldsymbol{R} is a rotation matrix

$$oldsymbol{R} = egin{bmatrix} 1/2 & 0 & -\sqrt{3}/2 \ 0 & 1 & 0 \ \sqrt{3}/2 & 0 & 1/2 \end{bmatrix} egin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \ -1/\sqrt{2} & 1/\sqrt{2} & 0 \ 0 & 0 & 1 \end{bmatrix}.$$

In this case, we set n = 128. In each of the Monte Carlo integrations to compute the "true" population risks, 1024 samples were used.

Both results in Figure 2 show that \mathcal{E}_{te} is a fairly accurate approximation of the population risk $\mathbb{E}_{\mathcal{D}}[\mathcal{R}(\hat{\beta})]$. This empirically validates the discussions we made in Appendix C.

E.3. Neural-network-like model experiments

We also experimented with neural-network-like models, testing the potency of the formula (33) for \mathcal{E}_{te} as the estimator of the population risk.

To exhibit a comparison to the previous section, we performed an experiment in two dimensions, using the same target and latent distributions as in the two-dimensional experiment in Appendix E.2. Meanwhile, to provide insight into more practical settings involving modern generative models where the neural networks are trained with the standard Gaussian latent distribution, we conducted the three-dimensional experiment with the same target distribution as in Appendix E.2 but changed the latent distribution to $\xi = \mathcal{N}(\mathbf{0}, \mathbf{I})$.



Figure 3: Learning the optimal transport maps with neural-network-like models. We plot the computed theoretical and experimental population risks. The meanings of the curves and scatter plots are the same as Figures 1 and 2.

The same g is used as in (34), but the random feature map follows the configuration of (31). All other details remain identical to Appendix E.2.

Figure 3 shows the plots of the results. The curves and the scatter plots are drawn in the exact same way as in the previous experiments. Again, we can observe an agreement between actual experimental results and the theoretical estimation by \mathcal{E}_{te} , this time computed using (33).