Curvature Enhanced Data Augmentation for Regression

Ilya Kaufman¹ Omri Azencot¹

Abstract

Deep learning models with a large number of parameters, often referred to as over-parameterized models, have achieved exceptional performance across various tasks. Despite concerns about overfitting, these models frequently generalize well to unseen data, thanks to effective regularization techniques, with data augmentation being among the most widely used. While data augmentation has shown great success in classification tasks using label-preserving transformations, its application in regression problems has received less attention. Recently, a novel manifold learning approach for generating synthetic data was proposed, utilizing a first-order approximation of the data manifold. Building on this foundation, we present a theoretical framework and practical tools for approximating and sampling general data manifolds. Furthermore, we introduce the Curvature-Enhanced Manifold Sampling (CEMS) method for regression tasks. CEMS leverages a second-order representation of the data manifold to enable efficient sampling and reconstruction of new data points. Extensive evaluations across multiple datasets and comparisons with state-of-the-art methods demonstrate that CEMS delivers superior performance in both in-distribution and out-of-distribution scenarios, while introducing only minimal computational overhead. Code is available at https: //github.com/azencot-group/CEMS.

1. Introduction

Deep neural networks have demonstrated remarkable performance across a wide range of applications in various fields (Krizhevsky et al., 2012; Long et al., 2015; Mnih et al., 2015; Noh et al., 2015; Vinyals et al., 2015; He et al., 2016; Nam & Han, 2016; Wu et al., 2016). Despite their success, these models are often significantly overparameterized, meaning they possess more parameters than the number of training examples. As a result, deep neural networks are prone to overfitting, whereby they "memorize" the training set rather than learning generalizable patterns, thus compromising their performance on unseen data. Regularization techniques are crucial to address this issue, as they modify the learning process to prevent overfitting by reducing the variance and increasing the bias of the underlying model (Goodfellow, 2016). Classical regularization methods, such as weight decay, dropout (Srivastava et al., 2014), and normalization techniques like Batch Normalization (Ioffe & Szegedy, 2015) and Layer Normalization (Ba et al., 2016), have been effective in many scenarios. In addition to these methods, recent research has explored the potential of data augmentation (DA) as a form of regularization. In this paper, we focus on the problem of regularizing regression models via data augmentation. That is, we explore how to artificially expand the train set (DA) for models that predict continuous values (regressors) to improve generalization and robustness.

Early work in modern computer vision revealed the effectiveness of basic image transformations such as translation and rotation (Krizhevsky et al., 2012), promoting data augmentation to become one of the key components in designing generalizable learning models (Shorten & Khoshgoftaar, 2019). In particular, classification tasks, whose goal is to predict a discrete label, benefited notably from the rapid development of DA (Simonyan & Zisserman, 2015; DeVries, 2017; Zhang et al., 2018; Zhong et al., 2020). The discrete and categorical nature of classification labels makes it easier to define label-preserving transformations and apply interpolations without compromising data integrity. In contrast, regression tasks, where the outputs are continuous, face unique challenges in ensuring that transformations produce valid input-output pairs and that interpolations maintain the underlying functional relationships. While certain regression challenges have adopted standard data augmentation approaches successfully (Redmon et al., 2016; Nochumsohn & Azencot, 2025), existing DA methods are generally less effective for regression problems (Yao et al., 2022). For this reason, developing data augmentation tools for general regression problems is an emerging field of interest

¹Department of Computer Science, Ben-Gurion University of the Negev, Beer-Sheva, Israel. Correspondence to: Ilya Kaufman <ilyakau@post.bgu.ac.il>.

Proceedings of the 42^{nd} International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

with a relatively small number of available effective techniques. One of the recent state-of-the-art (SOTA) works introduced FOMA, a data-driven and domain-independent approach based on the theory and practice of manifold learning (Kaufman & Azencot, 2024b). Our work is also inspired by manifold learning, where we consider DA as a *manifold approximation and sampling* challenge.

Manifold learning is fundamental to modern machine learning primarily through the manifold hypothesis (Belkin & Niyogi, 2003; Goodfellow, 2016), where complex and highdimensional data is assumed to lie close or on an associated low-dimensional manifold. Multiple works leveraged the relation between data and manifolds (Zhu et al., 2018; Ansuini et al., 2019; Kaufman & Azencot, 2023), and particularly, FOMA (Kaufman & Azencot, 2024b) can be viewed as a method for generating new examples by sampling from the tangent space of the data manifold, approximated using the training distribution. The tangent space at a point is a linear approximation of the manifold at that point (Lee, 2012), and thus, FOMA is a first-order approach. However, while first-order approximations work well for relatively simple or well-behaved data, they often fall short when dealing with complex, curved real-world data. We demonstrate this effect in Fig. 1B-C, where first-order approximations of points with high curvature fail to capture the structure of the manifold. While it is natural to consider higher-order approximations for improving FOMA, their computational burden may be too limiting. In this work, we advocate that second-order manifold representations offer a compelling trade-off between effectiveness and compute requirements for data augmentation for regression problems.

We propose the curvature enhanced manifold sampling (CEMS) approach, which generates new examples by drawing from a second-order representation of the data manifold. Particularly, CEMS randomly generates points in the tangent space of the manifold, whereas FOMA is different as it scales down the orthogonal complement of the tangent space which captures how the manifold deviates from the linear approximation. FOMA and CEMS are data-driven and domain-independent, i.e., their samples are based on the underlying data distribution, whose domain can be arbitrary, e.g., time series, tabular, images, and other data forms. The inclusion of curvature information allows CEMS to better capture the intrinsic geometry of the data manifold, as it accounts for the non-linearities and complex structures that first-order methods might miss. In general, second-order methods are infeasible in modern machine learning due to computational costs incurred by high-dimensional vectors. Nevertheless, our analysis shows that CEMS is governed by the *intrinsic dimension* d of the manifold, and its value is much smaller than the data dimension D, i.e., $d \ll D$. We extensively evaluate CEMS and show it is competitive in comparison to SOTA data augmentation approaches on indistribution and out-of-distribution tasks. The contributions of our work can be summarized as follows:

- We consider data augmentation for regression as a manifold learning problem, extending and generalizing prior approaches through providing the foundational theory and practice.
- We introduce CEMS, a novel fully-differentiable, datadriven and domain-independent data augmentation technique that is based on a second-order approximation of the data manifold.
- Across nine datasets, featuring several large-scale, realworld in-distribution and out-of-distribution tasks, we demonstrate that CEMS performs competitively or surpasses other augmentation strategies.

2. Related work

The theoretical foundation for data augmentation (DA) is related to the study of Empirical Risk Minimization (ERM) (Vapnik, 1991) vs. Vicinal Risk Minimization (VRM) (Chapelle et al., 2000). In VRM, one considers an extended distribution to train on, in comparison to ERM, where the train distribution is used. Although related to generative modeling (Kingma & Welling, 2014; Goodfellow et al., 2014; Sohl-Dickstein et al., 2015; Naiman et al., 2024b;a), data augmentation is a distinct approach commonly used to expand data distributions by generating synthetic samples. With the increased dependence of deep models on large volumes of data, DA has become a cornerstone in enhancing the performance and generalization of neural networks (DeVries, 2017; Chen et al., 2020b; Feng et al., 2021; Yang et al., 2022). Early work focused on domain-dependent augmentations for image, audio, and natural language data (Krizhevsky et al., 2012; He et al., 2016; Huang et al., 2017; Kobayashi, 2018; Park et al., 2019; Zhong et al., 2020). Later, automatic augmentation tools have been proposed (Cubuk et al., 2019; Lim et al., 2019), including domain-dependent search spaces for transformations. Still, adapting these methods to new data domains remains a challenge. This has sparked interest in developing domain-independent approaches that make minimal assumptions about the data domain, and it is the focus of our research.

DA for Classification. Zhang et al. (2018) introduced mixup, a well-known domain-independent DA method that convexly combines pairs of input samples and their one-hot label representations during training. Following their work, a plethora of mixup-based techniques have been suggested such as manifold mixup (Verma et al., 2019) which extends the idea of mixing examples to the latent space.



Figure 1. We demonstrate the effect of sampling from a one-dimensional manifold embedded in a two-dimensional space. A) The original data representing a sine wave where the color of each point represents the curvature at that point (brighter means higher curvature). B) Sampling using FOMA. C) Sampling using a first-order approximation. D) Sampling using CEMS (our approach).

CutMix (Yun et al., 2019) implants a random rectangular region of the input into another and many others (Guo et al., 2019; Hendrycks et al., 2020; Berthelot et al., 2019; Greenewald et al., 2023; Lim et al., 2022). Recently, (Erichson et al., 2024) extended the work (Lim et al., 2022) via stable training and further noise injections. While the family of mixup techniques have been shown to consistently improve classification learning systems (Cao et al., 2022), its efficacy is inconsistent on regression tasks (Yao et al., 2022).

DA for Regression. Unfortunately, there has been considerably less focus on developing data augmentation methods for regression tasks in comparison to the classification setting. Due to the simplicity and effectiveness of mixup-based tools in classification, a growing body of literature is drawn to adapting and extending the mixing process for regression. For instance, RegMix (Hwang & Whang, 2021) learns the optimal number of nearest neighbors to mix per sample. C-mixup (Yao et al., 2022) employs a Gaussian kernel to create a sampling probability distribution for each sample, taking label distances into account, and selecting samples for mixing according to this distribution. Anchor Data Augmentation (Schneider et al., 2023) clusters data points and adjusts the original points either towards or away from the cluster centroids. R-Mixup (Kan et al., 2023) focuses on enhancing model performance specifically for biological networks, whereas RC-Mixup (Hwang et al., 2024) extends C-mixup to be more robust against noise. Perhaps closest to our work is the recent FOMA method (Kaufman & Azencot, 2024b) that does not rely on mixing samples, but rather, it samples from a first-order approximation of the data manifold. Still, to the best of our knowledge, our work is first in suggesting fundamental manifold learning theory and tools for DA, accompanied by an effective second-order augmentation technique.

Manifold Learning. Manifold learning has been a fundamental research area in machine learning, aiming to discover the intrinsic low-dimensional structure of high-dimensional data. While early work focused on dimensionality reduction of points and preserving their geometric features (Tenenbaum et al., 2000; Roweis & Saul, 2000; Belkin & Niyogi, 2003; Weinberger & Saul, 2004; Zhang & Zha, 2005; Coifman & Lafon, 2006), modern approaches also considered regularization (Ma et al., 2018; Zhu et al., 2018), explainable artificial intelligence (Ansuini et al., 2019; Kaufman & Azencot, 2023; 2024a), and autoencoding (Chen et al., 2020a), among other applications.

Second-order methods have also been investigated in this domain. Notably, Donoho & Grimes (2003) introduced Hessian Eigenmaps, a technique that incorporates secondorder information to more accurately preserve local curvature during nonlinear dimensionality reduction. Although their approach is designed for unsupervised embedding, our method extends the use of second-order local approximations to the supervised setting. Specifically, we utilize this geometric insight to guide data augmentation, enabling the generation of new training samples in a manner that is both geometry-aware and differentiable.

Recent work from the statistical learning community has also proposed more sophisticated local chart models. These include Gaussian process-based manifold inference (Dunson & Wu, 2021), spherelet-based approximations (Li et al., 2022), and manifold denoising via generalized L_1 medians (Faigenbaum-Golovin & Levin, 2020). While these methods aim to infer or reconstruct manifold structure with statistical guarantees, our focus is on utilizing local geometric information specifically for data augmentation in supervised learning.

Recent advancements in manifold learning have enhanced

anomaly detection and out-of-distribution (OOD) recognition. Li et al. (2024) leveraged submanifold geometry, estimating tangent spaces and curvatures to define indistribution regions for OOD detection. Gao et al. (2022) proposed a hyperbolic feature augmentation method, using the Poincaré ball model for distribution estimation and infinite sampling, improving few-shot learning performance. Humayun et al. (2022) proposed MaGNET, a framework that enables uniform sampling on data manifolds derived from generative adversarial networks providing a retrainingfree solution for data augmentation. Similarly, Chadebec & Allassonnière (2021) introduced a geometry-aware variational autoencoder that leverages second-order Runge-Kutta schemes for effective data generation in low-sample-size scenarios. Extending these ideas, Cui et al. (2023) presented a trajectory-aware principal manifold framework for image generation and data augmentation, which aligns sampled data with learned projection indices to improve representation and synthesis quality.

3. Background

3.1. Manifold learning

A manifold $\mathcal{M} \subset \mathbb{R}^d$ is a mathematical structure that locally resembles an Euclidean space near each of its points (Lee, 2012). A ubiquitous assumption in machine learning states that high-dimensional point clouds $Z \subset \mathbb{R}^D$ satisfy the manifold hypothesis. Namely, the data Z lie on a manifold \mathcal{M} whose intrinsic dimension d is significantly lower than the extrinsic dimension of the ambient space D, i.e., $d \ll D$ (Goodfellow, 2016). Manifold learning is a field in machine learning that develops theory and tools for analyzing and processing high-dimensional data under the lens of geometric manifolds.

3.2. Curvature-aware manifold learning

Our curvature enhanced manifold sampling (CEMS) data augmentation method is based on a second-order approximation of the data manifold. There are several existing practical approaches for parameterizing a manifold given a collection of data points $Z = \{z^1, z^2, \dots, z^N\} \subset \mathbb{R}^D$. Here, we focus on the curvature-aware manifold learning (CAML) algorithm (Li, 2018), since it scales to high-dimensional problems, it is numerically stable, and it is easy to code. Below, we include the necessary details for describing our approach, and we refer the reader to (Lee, 2012; 2018; Li, 2018) for additional details on Riemannian geometry and its realization in machine learning.

Following our discussion above, we assume Z satisfies the manifold hypothesis. Formally, it means that there exists an embedding map $f : \mathcal{M} \to \mathbb{R}^D$ such that

$$z^{i} = f(u^{i}), \quad i = 1, \dots, N$$
, (1)

where $U = \{u^1, u^2, \cdots, u^N\} \subset \mathbb{R}^d$ are low-dimensional representations of Z. In practice, CAML parameterizes f by projecting $z \in Z$ to its tangent and normal spaces at $u \in \mathcal{M}$, where the tangent space is obtained by a linear transformation and the normal space is provided via a second-order local approximation. Specifically, given a point $z \in Z$, we find close points $N_z = \{z_j\}_{j=1}^k$, forming the neighborhood of z. Next, we construct an orthonormal basis $B_u := [B_{\mathcal{T}_u} \in \mathbb{R}^{D \times d}, B_{\mathcal{N}_u} \in \mathbb{R}^{D \times (D-d)}] \in \mathbb{R}^{D \times D}$ for the tangent space $\mathcal{T}_u \mathcal{M}$ and the normal space $\mathcal{N}_u \mathcal{M}$ at u, where the operation $[\cdot, \cdot]$ denotes column-wise concatenation. We then project N_z and z onto $B_{\mathcal{T}_u}$ and $B_{\mathcal{N}_u}$, yielding $U_z = \{u_j\}_{j=1}^k$, u and $G_z = \{g_j\}_{j=1}^k$, g respectively. To allow arbitrary sampling from \mathcal{M} , we assume that g is a map from the tangent space to the normal space, i.e., $g: \mathcal{T}_u \mathcal{M} \to \mathcal{N}_u \mathcal{M}$. The second-order Taylor expansion of g around a point $u \in \mathcal{M}$ is given by

$$g(u_j) = g(u) + (u_j - u)^T \nabla g(u) + \frac{1}{2} (u_j - u)^T H(u) (u_j - u) + \mathcal{O}(|u_j|_2^3), \quad (2)$$

where the linear (gradient) and quadratic (Hessian) terms are unknown. To compute $\nabla g(u)$ and H(u), one needs to collect the linear coefficients and constants arising from Eq. 2 into matrices Ψ and G, respectively, solve a linear system of equations (Eq. 6), and extract the numerical estimates of the gradient and Hessian. Finally, we can map the pair (u, g) back into its original space $z \in Z$ by computing $z := f(u) = B_u[u, g(u)]$. See also App. A.

4. Curvature Enhanced Manifold Sampling

We assume to be given a regression training set \mathcal{D} := $\{(x^i,y^i)\}_{i=1}^N,$ where $x^i \in \mathbb{R}^{k_1}$ is the data sample and $y^i \in \mathbb{R}^{k_2}$ is its corresponding prediction. Following the approach used in FOMA (Kaufman & Azencot, 2024b), we denote by $z^i = [x^i, y^i] \in \mathbb{R}^D$ the concatenation of the input x^i and its corresponding label y^i along the columns, treating the pair as a point on a joint input-output manifold. During training, given a mini-batch $Z = \{z^1, \dots z^b\}$, where b is the batch size, we perform the following procedure for each $z \in Z$ to create a new sample \tilde{z} , omitting the superscript to simplify notation. To account for the discrepancies in scale between X and Y, we normalize Y to match the range of X, i.e., [0, 1]. Our curvature enhanced manifold sampling (CEMS) augmentation approach consists of four main steps: 1) Extract a neighborhood N_z from \mathcal{D} for every point z; 2) Construct a basis $B_u := [B_{\mathcal{T}_u}, B_{\mathcal{N}_u}]$ for the tangent space $\mathcal{T}_u \mathcal{M}$ and the normal space $\mathcal{N}_u \mathcal{M}$ and project the neighborhood onto it; 3) Form and solve the linear system of equations in Eq. 6 to obtain the parameterization q; 4) Sample a new point from \mathcal{T}_u , evaluate its g via Eq. 2, and un-project it onto the ambient space using f. Below, we detail how we perform each step, and we motivate our design

Curvature Enhanced Data Augmentation for Regression



Figure 2. CEMS forms a neighborhood for every point z (left), computes a basis $B_{\mathcal{T}_u}$ for the tangent space via SVD (middle) while obtaining an estimate for the embedding f, samples a new point η close to u (right), and un-projects it back to \mathbb{R}^D using f (red arrow).

choices. Pseudo-code and illustration of CEMS are given in Alg. 1 and Fig. 2.

Neighborhood extraction. The approach we consider in this work is *local* in the sense that we represent the manifold structure in the vicinity of a specific point $z \in \mathcal{D}$ or within its neighborhood N_z . High-quality approximations of local properties of the manifold $\nabla g(u)$ and H(u) depend directly on the proximity of the elements in N_z to z. A straightforward approach to extracting N_z is to compute the k-nearest neighbors (kNN) (Cover & Hart, 1967) of z. Specifically, we construct neighborhoods in the joint input-output space $\mathcal{X} \times \mathcal{Y}$ to align with the manifold hypothesis, preserving the local continuity of the data and avoiding the artificial separations introduced by clustering-based methods. For each point $z_i \in \mathcal{Z}$, N_{z_i} is defined as the k closest points in feature space, which ensures the local geometry is captured reliably while minimizing diversity within the neighborhoods. Furthermore, our reliance on the local-Euclidean prior assumes that the manifold is sufficiently smooth at small scales, justifying the validity of linear approximations such as those employed by our method. This assumption underpins the extraction of neighborhood sets and their utility in manifold analysis, as it guarantees that the neighborhoods respect the underlying manifold structure.

Our analysis below and in App. D shows that the computational complexity of CEMS is governed by singular value decomposition (SVD) calculations, required for basis construction. To reduce runtime, \mathcal{D} can be pre-processed, storing $\nabla g(u)$ and H(u) for every z := f(u) on the disk, as these properties remain unchanged during training. While this pre-computation reduces runtime significantly, it incurs high memory complexity, $\mathcal{O}(2d(D-d))$, and the choice of neighbors is fixed during training. To address these limitations, we use the *same neighborhood* for all points in N_z , re-using neighborhoods and basis computations for every $z_j \in N_z$. This improves efficiency, though at the cost of accuracy, since every $z_j \in N_z$ is assumed to share the same neighborhood. Finally, the batch size determines the number of neighbors for each point, providing a balance between computational efficiency and accuracy.

Basis construction and projection. To find an orthonormal basis for the tangent space $\mathcal{T}_u \mathcal{M}$ and the normal space $\mathcal{N}_{u}\mathcal{M}$, we follow standard approaches (Singer & Wu, 2012; Li, 2018) that utilize the singular value decomposition (SVD). Specifically, we compute SVD on the centered points $\{z_j - z\}_{j=1}^k = USV^T$, while keeping our pipeline to be fully differentiable (Ionescu et al., 2015). Importantly, SVD is calculated once for every batch, as discussed above. The first d columns of U determine the basis for the tangent space, i.e., $B_{\mathcal{T}_u} := U[:, 1:d]$ and the last D - d columns determine the basis for the normal space $B_{\mathcal{N}_u} := U[:, d+1:D]$ such that $B_u := [B_{\mathcal{T}_u}, B_{\mathcal{N}_u}]$ is the concatenation of the bases. While the intrinsic dimension d can be viewed as a hyper-parameter of CEMS, we estimate it in practice using a robust estimator (Facco et al., 2017). The centered neighbors $z_i - z$ are projected to the tangent space and the normal space via $u_j := B_{\mathcal{T}_u}^T \cdot (z_j - z), g(u_j) := B_{\mathcal{N}_u}^T \cdot (z_j - z),$ respectively, where A^T is the transposed matrix of A, and $A \cdot v$ is a matrix-vector multiplication. Centering the points around z map the point u to the zero vector, and thus, Eq. 2 is transformed to the following approximation:

$$g(u_j) = u_j^T \nabla g(u) + \frac{1}{2} u_j^T H(u) u_j$$
. (3)

Linear system of equations. Under the change of basis B_u , we form the matrices Ψ and G as described in App. A, containing $\{u_j\}_{j=1}^k$ and $\{g(u_j)\}_{j=1}^k$, respectively. We then solve Eq. 6 via differentiable least squares, and we obtain an estimation of $\nabla g(u)$ and H(u), allowing to map new points in the vicinity of u by computing g. Note that while N_z

Algorithm 1 CEMS $_p$: Sample generation **Require:** Training data $Z = \{z^i = [x^i, y^i]\}_{i=1}^N$ **Require:** A sample $z \in Z$ 1: Find k-nearest neighbors N_z of z 2: Construct an orthonormal basis B_u spanning $N_z - z$ 3: Project $z_j - z$ to local orthonormal coordinates: $u_i \leftarrow B_u^T(z_i - z)$ 4: Construct G and Ψ as in Eq. (6) 5: Solve $\Psi A = G$ 6: Extract $\nabla q(z)$ and H(z) from A 7: Sample noise $\eta \sim \mathcal{N}(0, \sigma I_d)$ 8: Calculate $g(\eta)$ from: $g(\eta) \leftarrow \eta^T \nabla g + \frac{1}{2} \eta^T H \eta.$ 9: Un-project η back to the original coordinates: $z_\eta \leftarrow f(\eta) = B_u[\eta, g(\eta)] + z$ Return z_{η}

and B_u are shared across the batch, the linear solve is still computed separately per point.

Sampling and un-projecting. To generate new examples using the above machinery, we need to sample a point $\eta \in \mathbb{R}^d$ from the neighborhood of u, and un-project it to the ambient space \mathbb{R}^D (through the parameterization g and map f). While various sophisticated sampling techniques could be devised, we opted for a simple sampler with a single hyperparameter. In practice, we draw $\eta \sim \mathcal{N}(0, \sigma I_d)$. To un-project η back to the original space, we compute

$$z_{\eta} := f(\eta) = B_u \cdot [\eta, g(\eta)] + z . \tag{4}$$

Adaptation to batches. For completeness, we also describe briefly the adaptation of CEMS to the training setting where we utilize mini-batches. We provide a comparison between the two methods in App. G. In contrast to CEMS_{p} , where a neighborhood of close points is constructed for each individual sample to generate a corresponding synthetic point, CEMS requires a batch of mutually close points to estimate a shared tangent space and produce a new batch of samples. As mentioned above, given z and its neighborhood N_z , we re-use the same neighborhood and subsequent basis computations for every $z_j \in N_z$. This adaptation requires a small modification to the method. 1) We include the point $z := z_0$ in the neighborhood $N_z = \{z_j\}_{j=0}^k$. 2) We find an orthonormal basis B_u that spans $N_z - \mu$, where μ is the mean of N_z . 3) After projecting to the coordinates of B_u , we get $U_z = \{u_j\}_{j=0}^k$ and $G_z = \{g_j\}_{j=0}^k$. For every point u_i , we gather a set of close points and their embeddings via g. 4) In contrast to the point-wise basis estimation, where u served as the origin (zero vector), in the batch-wise computation we need to account for u_i and $q(u_i)$ in Eq. 2. While Steps 5-7 in Alg. 1 remain unchanged, at Step 7 we

sample a point η near u_l : $\eta \sim \mathcal{N}(u_l, \sigma I_d)$. Step 8 changes to $g(\eta_l) = g(u_l) + (\eta_l - u_l)^T \nabla g + \frac{1}{2} (\eta_l - u_l)^T H(\eta_l - u_l)$ and Step 9 changes to $z_{\eta_l} := f(\eta_l) = B_u \cdot [\eta_l, g(\eta_l)] + \mu_{N_z}$. A full description of the algorithm appears in Alg. 2.

Complexity analysis. There are two computationally demanding calculations used by CEMS, SVD and least squares. Given a data mini-batch $Z \in \mathbb{R}^{b \times D}$, where *b* is the batch size. Then, SVD requires $\mathcal{O}(\min(bD^2, Db^2))$ operations, whereas the solution of under-determined least squares costs $\mathcal{O}(b^2d^2)$. Using the manifold hypothesis, we assume that $d \ll D$ therefore $d^2 \in \mathcal{O}(D)$ and thus, the overall time complexity of CEMS is given by $\mathcal{O}(b^2D)$ which is proportional to the ambient dimension *D*. See also App. D for a more detailed analysis.

Memory analysis. The memory requirements of CEMS are primarily dictated by the computation of the SVD. Notably, the SVD is computed independently for each batch rather than for the entire dataset. In our PyTorch implementation, we leverage the economy/reduced SVD variant, which significantly reduces memory usage compared to the full SVD. For a batch matrix of size $b \times D$ (where *b* is the batch size and *D* is the ambient dimension), the space complexity is $O(bD + \min(b, D)(b + D))$. This is substantially more efficient than the full SVD, which requires $O(bD+b^2+D^2)$ memory. In practice, CEMS is particularly effective in scenarios where the batch size *b* is much smaller than the ambient dimension *D* (common in deep learning), resulting in a memory complexity that is proportional to *D*.

Comparison with FOMA. FOMA (Kaufman & Azencot, 2024b) can be interpreted as a special case of CEMS. Specifically, we can describe FOMA using our notations as follows: given a sample z and its neighborhood N_z , FOMA constructs a basis $B_u := [B_{\mathcal{T}_u}, B_{\mathcal{N}_u}]$ for the tangent space $\mathcal{T}_u \mathcal{M}$ and the normal space $\mathcal{N}_u \mathcal{M}$ and projects the neighborhood onto it, yielding $U_z = \{u_j\}_{j=1}^k$ and $G_z = \{g_j\}_{j=1}^k$, respectively. Rather than estimating the gradient $\nabla g(u_i)$ and Hessian $H(u_i)$ at each point $u_i \in U_z$ and then sampling using the Taylor expansion as performed in CEMS, FOMA generates new samples by scaling down G_z . That is, for each $z_i \in N_z$, the corresponding $\tilde{g}_i = \lambda g_i$ is scaled where $\lambda \in (0, 1)$. To complete the sampling process, every u_i is un-projected back to the original coordinates by computing $\tilde{z} := f(u_i) = B_u \cdot [u_i, \lambda g(u_i)]$. Therefore, FOMA does not use the embedding map q as detailed in Eq. 2, but it samples random points instead. Unfortunately, this sampling technique may yield new points that are not on the data manifold, especially on highly curved locations, as is also illustrated in Fig. 1.

	Airfoil		NO2		Exchange-Rate		Electricity	
	RMSE	MAPE	RMSE	MAPE	RMSE	MAPE	RMSE	MAPE
ERM	2.901	1.753	0.537	13.615	0.024	2.423	0.058	13.861
Mixup	3.730	2.327	0.528	13.534	0.024	2.441	0.058	14.306
Mani Mixup	3.063	1.842	0.522	13.382	0.024	2.475	0.058	14.556
C-Mixup	2.717	1.610	0.509	12.998	0.020	2.041	0.057	13.372
ADA	2.360	1.373	0.515	13.128	0.021	2.116	0.059	13.464
FOMA	<u>1.471</u>	<u>0.816</u>	0.512	<u>12.894</u>	0.013	1.262	<u>0.058</u>	14.614
CEMS	1.455	0.809	0.507	12.807	0.014	1.293	0.058	13.353

Table 1. Results for in-distribution generalization. Values in bold indicate the best results, while underlined values represent the second best. We present the average RMSE and MAPE across three seeds. Detailed results, including standard deviation, are available in App K.

5. Experiments

5.1. Sine example

Real-world data is often complex and curved, exhibiting intricate patterns that cannot be adequately captured by linear or simplistic models. By employing higher-order approximations of the manifold, we can generate samples that align with the true nature of real-world data. In Fig. 1, we demonstrate a toy sine example, highlighting the differences between first-order and second-order approaches. Specifically, we generated a two-dimensional point cloud of a sine wave whose intrinsic dimension is one (Fig. 1, left). Then, we sampled points from this distribution using minibatches from the train set and various data augmentation techniques. The first-order method, FOMA (Kaufman & Azencot, 2024b), struggles to adhere to the curvature of the manifold in highly-curved points, as can be seen in Fig. 1, middle left. Similarly, restricting CEMS to a first-order approximation presents a similar behavior (Fig. 1, middle right). Finally, our second-order CEMS method samples the manifold well, even near high curvature areas (Fig. 1, right). See App. A.1 for a theoretical justification for the sampling error of CEMS in comparison to first-order approaches.

5.2. In-Distribution Generalization

In what follows, we consider the in-distribution benchmark that was introduced in (Yao et al., 2022). This benchmark evaluates the performance of various data augmentation techniques in the setting of training on a train set and its augmentations, while testing on a test set that was sampled from the same distribution as the train set. Thus, a strong performance in this benchmark implies that the underlying DA method mimics the train distribution well. Below, we compare CEMS to other recent SOTA approaches, while using the same datasets that were studied in (Yao et al., 2022) and closely replicating their experimental setup. **Datasets.** We evaluate in-distribution generalization using four datasets. Two of these are tabular datasets: Airfoil Self-Noise (Airfoil) (Brooks et al., 2014), containing aerodynamic and acoustic measurements of airfoil blade sections, and NO2 (Aldrin, 2004), which predicts air pollution levels at specific locations. We also use two time series datasets: Exchange-Rate and Electricity (Lai et al., 2018), where Exchange-Rate includes daily exchange rates of several currencies and Electricity contains measurements of electric power consumption in private households. For a detailed description of these datasets, see App. J.

Experimental Settings. We perform a comparative analysis between our method, CEMS, and several established baseline approaches, including the standard empirical risk minimization (ERM) training, Mixup (Zhang et al., 2018), Manifold-Mixup (Verma et al., 2019), C-Mixup (Yao et al., 2022), Anchor Data Augmentation (ADA) (Schneider et al., 2023), and FOMA (Kaufman & Azencot, 2024b). The neural networks we trained are the same as considered in (Yao et al., 2022), where a fully connected three layer model was used for tabular datasets, and an LST-Attn (Lai et al., 2018) is utilized for time series data. The evaluation metrics include the root mean square error (RMSE) and mean absolute percentage error (MAPE). Additional details and hyperparameters are available in App. I.

Results. We present the in-distribution generalization benchmark results in Tab. 1. The results of all previous methods are reported as they appear in the corresponding original papers. Lower values are preferred either in RMSE or in MAPE. Boldface and underline denote the best and second best approaches, respectively. Remarkably, across all datasets and metrics, CEMS attains the best or second best error measures. In particular, CEMS outperforms other data augmentation strategies on Airfoil and NO2 while being comparable with FOMA on Electricity and Exchange-Rate. We also note that when CEMS is second best, its result is relatively close to the best result. We present the full results

Table 2. Comparison of out-of-distribution robustness. Bold values indicate the best results, while underlined values represent the second
best. We present the average RMSE across domains as well as the "worst within-domain" RMSE from three different seeds. For the DTI
and Poverty datasets, we provide the average R and the "worst within-domain" R. Complete results, including standard deviation, can be
found in App K.

	RCF (RMSE)	Crimes (RMSE) SkillCraft (ft (RMSE)	DTI(R)		Poverty (R)		
	Avg.↓	Avg. \downarrow	Worst \downarrow	Avg.↓	Worst \downarrow	Avg. ↑	Worst ↑	Avg. ↑	Worst \uparrow
ERM	0.164	0.136	0.170	6.147	7.906	0.483	0.439	<u>0.80</u>	0.50
Mixup	<u>0.159</u>	0.134	0.168	6.460	9.834	0.459	0.424	0.81	0.46
ManiMixup	0.157	0.128	0.155	5.908	9.264	0.474	0.431	-	-
C-Mixup	0.146	0.123	0.146	<u>5.201</u>	7.362	0.498	0.458	0.81	0.53
ADA	0.175	0.130	0.156	5.301	<u>6.877</u>	0.493	0.448	0.79	<u>0.52</u>
FOMA	<u>0.159</u>	0.128	0.158	-	-	<u>0.503</u>	<u>0.459</u>	0.78	0.49
CEMS	0.146	0.128	0.159	5.142	6.322	0.511	0.465	0.81	0.50

including standard deviation measures in App. K.

5.3. Out-of-distribution

To extend our in-distribution evaluation, we also consider an out-of-distribution benchmark, as was proposed in (Yao et al., 2022). Unlike the in-distribution case, here the test set is sampled from a distribution different from that of the train set. Therefore, excelling in this scenario provides valuable information regarding the generalization capabilities of data augmentation tools. In what follows, we perform a comparison between CEMS and several SOTA methods, while using the same datasets that were studied in (Yao et al., 2022) and closely replicating their experimental setup.

Datasets. We leverage five datasets to evaluate the performance of out-of-distribution robustness. 1) RCFashionM-NIST (RCF) (Yao et al., 2022) is a synthetic variation of Fashion-MNIST, designed to model sub-population distribution shifts, with the aim of predicting the rotation angle for each object. 2) Communities and Crime (Crime) (Redmond, 2009) is a tabular dataset focused on predicting the total number of violent crimes per 100,000 population, aiming to create a model that generalizes to states not included in the training data. 3) SkillCraft1 Master Table (Skill-Craft) (Blair et al., 2013) is a tabular dataset designed to predict the average latency in milliseconds from the onset of perception-action cycles to the first action where "LeagueIndex" is considered as domain information. 4) Drug-Target Interactions (DTI) (Huang et al., 2021) seeks to predict drugtarget interactions that are out-of-distribution, using the year as domain data. 5) PovertyMap (Poverty) (Koh et al., 2021) is a satellite image regression dataset created to estimate asset wealth in countries that were not part of the training set. For more details, please refer to App. J.

Experimental Settings. Similar to Sec. 5.2, we consider the same baseline DA approaches. For metrics, we report the RMSE (lower values are preferable) for RCF, Crimes, and SkillCraft. In addition, we use R (higher values are preferable) as the evaluation metric for Poverty and DTI, as was originally proposed in their corresponding papers (Koh et al., 2021; Huang et al., 2021). For a fair comparison, we follow the methodology in (Yao et al., 2022), and we train a ResNet-18 on the RCF and Poverty datasets, three-layer fully connected networks on Crimes and SkillCraft, and DeepDTA (Öztürk et al., 2018) on DTI. We provide further details on hyperparameters and experiments in App. I.

Results. We detail our out-of-distribution benchmark results in Tab. 2. Similarly to the in-distribution setting, the error measures of previous SOTA approaches were taken from the related original papers. We include both the average (Avg.) and worst domain performance metrics. Lower values are preferred in RMSE, and higher values are opted for R. We denote in bold and underline the best and second best results, respectively. Our results indicate that CEMS attains strong performance measures, achieving the best results in 6/9 tests. Further, the rest of the error measures of CEMS are either second best or very close to the second best. We particularly note the SkillCraft test where CEMS improves the second best results by a relative 1% and 8%for the average and worst metrics. The relative improvement is computed via $e_{\rm rel} \cdot 100$, where $e_{\rm rel} = (e - e_{\rm CEMS})/e$, with e_{CEMS} and e denoting the errors of CEMS and the second best approach, respectively. We present the full results including standard deviation measures in App. K.

6. Conclusions

This work introduces CEMS, a novel data augmentation method tailored specifically for regression problems, framed within the context of manifold learning. By leveraging second-order manifold approximations, CEMS captures the underlying curvature and structure of the data more accurately than previous first-order methods. Our extensive evaluation across nine diverse benchmark datasets, spanning both in-distribution and out-of-distribution tasks, shows that CEMS achieves competitive performance compared to SOTA techniques, often surpassing them in challenging settings. The main contributions of this work are threefold: (1) extend the view of DA for regression as a manifold learning problem, thereby providing a principled foundation and practical tools; (2) proposing CEMS, a fully differentiable, data-driven, and domain-independent second-order augmentation method; and (3) empirically validating CEMS across a variety of regression scenarios, showing its potential to serve as a robust and effective regularization technique for models predicting continuous values. Our results suggest that higher-order manifold sampling approaches hold promise for improving the generalization of regression models, especially in scenarios with limited data.

One limitation of CEMS is that the linear system in Eq. 6 might be underdetermined for data sets with a large intrinsic dimension d. In practice, the number of neighbors has to be $\mathcal{O}(d^2)$, for an overdetermined system. Our implementation sets the number of neighbors to be a constant size and thus it is independent of d. While this requirement is reasonable for low d values, it can become expensive for large d. This can be resolved by regularizing the linear system via, e.g., ridge regression. Another limitation is related to the SVD computation, where CEMS needs at least d columns. This may require a full SVD calculation, demanding $\mathcal{O}(bD^2)$ memory, where b is the batch size and D is the extrinsic dimension, which may be impractical for datasets with many features. A potential solution is to consider a different intrinsic dimension d, such that d < d. In future work, we plan to investigate these ideas, and, in addition, we plan to explore extensions of CEMS to include adaptive strategies for dynamically selecting the appropriate order of approximation based on local data properties. By pushing the boundaries of data augmentation for regression, we hope to pave the way for more robust and versatile learning systems capable of tackling complex, real-world prediction tasks.

Acknowledgments

This research was partially supported by the Lynn and William Frankel Center of the Computer Science Department, Ben-Gurion University of the Negev, an ISF grant 668/21, an ISF equipment grant, and by the Israeli Council for Higher Education (CHE) via the Data Science Research Center, Ben-Gurion University of the Negev, Israel.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Aldrin, M. CMU statlib dataset. http://lib.stat.cmu.edu/datasets/, 2004.
- Ansuini, A., Laio, A., Macke, J. H., and Zoccolan, D. Intrinsic dimension of data representations in deep neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- Ba, J. L., Kiros, J. R., and Hinton, G. E. Layer normalization. In *NIPS*, 2016.
- Belkin, M. and Niyogi, P. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
- Berthelot, D., Carlini, N., Goodfellow, I., Papernot, N., Oliver, A., and Raffel, C. A. MixMatch: A holistic approach to semi-supervised learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Birdal, T., Lou, A., Guibas, L. J., and Simsekli, U. Intrinsic dimension, persistent homology and generalization in neural networks. *Advances in neural information processing systems*, 34:6776–6789, 2021.
- Blair, M., Thompson, J., Henrey, A., and Chen, B. Skillcraft1 master table dataset. UCI Machine Learning Repository, 2013.
- Brooks, T., Pope, D., and Marcolini, M. Airfoil Self-Noise. UCI Machine Learning Repository, 2014.
- Cao, C., Zhou, F., Dai, Y., and Wang, J. A survey of mix-based data augmentation: Taxonomy, methods, applications, and explainability. *arXiv preprint arXiv:2212.10888*, 2022.
- Chadebec, C. and Allassonnière, S. Data generation in low sample size setting using manifold sampling and a geometry-aware vae. *CoRR*, 2021.
- Chapelle, O., Weston, J., Bottou, L., and Vapnik, V. Vicinal risk minimization. Advances in neural information processing systems, 13, 2000.
- Chen, N., Klushyn, A., Ferroni, F., Bayer, J., and van der Smagt, P. Learning flat latent manifolds with VAEs. In Proceedings of the 37th International Conference on Machine Learning, ICML, volume 119 of Proceedings of Machine Learning Research, pp. 1587–1596. PMLR, 2020a.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020b.
- Coifman, R. R. and Lafon, S. Diffusion maps. *Applied and Computational Harmonic Analysis*, 21:5–30, July 2006.

- Cover, T. and Hart, P. Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1):21–27, 1967.
- Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. AutoAugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 113– 123, 2019.
- Cui, E., Li, B., Li, Y., Wong, W., and Wang, D. Trajectory-aware principal manifold framework for data augmentation and image generation. *arXiv preprint arXiv:2310.07801*, 2023.
- DeVries, T. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- Donoho, D. L. and Grimes, C. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *Proceedings of the National Academy of Sciences*, 100 (10):5591–5596, 2003.
- Dunson, D. B. and Wu, N. Inferring manifolds from noisy data using gaussian processes. arXiv preprint arXiv:2110.07478, 2021.
- Erichson, B., Lim, S. H., Xu, W., Utrera, F., Cao, Z., and Mahoney, M. NoisyMix: Boosting model robustness to common corruptions. In *International Conference* on Artificial Intelligence and Statistics, pp. 4033–4041. PMLR, 2024.
- Facco, E., d'Errico, M., Rodriguez, A., and Laio, A. Estimating the intrinsic dimension of datasets by a minimal neighborhood information. *Scientific reports*, 7(1):12140, 2017.
- Faigenbaum-Golovin, S. and Levin, D. Manifold reconstruction and denoising from scattered data in high dimension via a generalization of 11-median. arXiv preprint arXiv:2012.12546, 2020.
- Feng, S. Y., Gangal, V., Wei, J., Chandar, S., Vosoughi, S., Mitamura, T., and Hovy, E. H. A survey of data augmentation approaches for NLP. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pp. 968– 988. Association for Computational Linguistics, 2021.
- Fowkes, J. M., Gould, N. I., and Farmer, C. L. A branch and bound algorithm for the global optimization of hessian lipschitz continuous functions. *Journal of Global Optimization*, 56:1791–1815, 2013.
- Gao, Z., Wu, Y., Jia, Y., and Harandi, M. Hyperbolic feature augmentation via distribution estimation and infinite

sampling on manifolds. *Advances in neural information* processing systems, 35:34421–34435, 2022.

Goodfellow, I. Deep learning, 2016.

- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- Grassberger, P. and Procaccia, I. Characterization of strange attractors. *Physical review letters*, 50(5):346, 1983.
- Greenewald, K. H., Gu, A., Yurochkin, M., Solomon, J., and Chien, E. k-Mixup regularization for deep learning via optimal transport. *Trans. Mach. Learn. Res.*, 2023.
- Guo, H., Mao, Y., and Zhang, R. Mixup as locally linear out-of-manifold regularization. In *Proceedings of the* AAAI Conference on Artificial Intelligence, volume 33, pp. 3714–3722, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *CVPR*, 2016.
- Hendrycks, D., Mu, N., Cubuk, E. D., Zoph, B., Gilmer, J., and Lakshminarayanan, B. AugMix: A simple data processing method to improve robustness and uncertainty. In 8th International Conference on Learning Representations, ICLR, 2020.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4700–4708, 2017.
- Huang, K., Fu, T., Gao, W., Zhao, Y., Roohani, Y., Leskovec, J., Coley, C. W., Xiao, C., Sun, J., and Zitnik, M. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. In *Proceedings of the Neural Information Processing Systems Track* on Datasets and Benchmarks 1, 2021.
- Humayun, A. I., Balestriero, R., and Baraniuk, R. Magnet: Uniform sampling from deep generative network manifolds without retraining. In *The International Conference* on Learning Representations (ICLR) 2022, 2022.
- Hwang, S.-H. and Whang, S. E. Regmix: Data mixing augmentation for regression. arXiv preprint arXiv:2106.03374, 2021.
- Hwang, S.-H., Kim, M., and Whang, S. E. RC-Mixup: A data augmentation strategy against noisy data for regression tasks. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1155–1165, 2024.

- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- Ionescu, C., Vantzos, O., and Sminchisescu, C. Matrix backpropagation for deep networks with structured layers. In *Proceedings of the IEEE international conference on computer vision*, pp. 2965–2973, 2015.
- Kan, X., Li, Z., Cui, H., Yu, Y., Xu, R., Yu, S., Zhang, Z., Guo, Y., and Yang, C. R-mixup: Riemannian mixup for biological networks. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1073–1085, 2023.
- Kaufman, I. and Azencot, O. Data representations' study of latent image manifolds. In *International Conference on Machine Learning*, pp. 15928–15945. PMLR, 2023.
- Kaufman, I. and Azencot, O. Analyzing deep transformer models for time series forecasting via manifold learning. *Trans. Mach. Learn. Res.*, 2024a.
- Kaufman, I. and Azencot, O. First-order manifold data augmentation for regression learning. In *Forty-first Inter*national Conference on Machine Learning, ICML, 2024b.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. In 2nd International Conference on Learning Representations, ICLR, 2014.
- Kobayashi, S. Contextual augmentation: Data augmentation by words with paradigmatic relations. In *Proceedings of* the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pp. 452–457, 2018.
- Koh, P. W., Sagawa, S., Marklund, H., Xie, S. M., Zhang, M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips, R. L., Gao, I., et al. Wilds: A benchmark of in-thewild distribution shifts. In *International Conference on Machine Learning*, pp. 5637–5664. PMLR, 2021.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- Lai, G., Chang, W.-C., Yang, Y., and Liu, H. Modeling long-and short-term temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pp. 95–104, 2018.
- Lee, J. M. Smooth manifolds. Springer, 2012.
- Lee, J. M. *Introduction to Riemannian manifolds*, volume 2. Springer, 2018.

- Levina, E. and Bickel, P. Maximum likelihood estimation of intrinsic dimension. *Advances in neural information processing systems*, 17, 2004.
- Li, D., Mukhopadhyay, M., and Dunson, D. B. Efficient manifold approximation with spherelets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 84(4):1129–1149, 2022.
- Li, X., Fang, Z., Zhang, Y., Ma, N., Bu, J., Han, B., and Wang, H. Characterizing submanifold region for out-ofdistribution detection. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- Li, Y. Curvature-aware manifold learning. Pattern Recognition, 83:273–286, 2018.
- Lim, S., Kim, I., Kim, T., Kim, C., and Kim, S. Fast AutoAugment. Advances in neural information processing systems, 32, 2019.
- Lim, S. H., Erichson, N. B., Utrera, F., Xu, W., and Mahoney, M. W. Noisy feature mixup. In *The Tenth International Conference on Learning Representations, ICLR*, 2022.
- Long, J., Shelhamer, E., and Darrell, T. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- Ma, X., Wang, Y., Houle, M. E., Zhou, S., Erfani, S., Xia, S., Wijewickrema, S., and Bailey, J. Dimensionality-driven learning with noisy labels. In *International Conference* on Machine Learning, pp. 3355–3364. PMLR, 2018.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M. A., Fidjeland, A., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. Human-level control through deep reinforcement learning. *nature*, 518:529–533, 2015.
- Naiman, I., Berman, N., Pemper, I., Arbiv, I., Fadlon, G., and Azencot, O. Utilizing image transforms and diffusion models for generative modeling of short and long time series. In Advances in Neural Information Processing Systems 38: NeurIPS, 2024a.
- Naiman, I., Erichson, N. B., Ren, P., Mahoney, M. W., and Azencot, O. Generative modeling of regular and irregular time series data via Koopman VAEs. In *The Twelfth International Conference on Learning Representations*, *ICLR*, 2024b.
- Nam, H. and Han, B. Learning multi-domain convolutional neural networks for visual tracking. In *CVPR*, 2016.
- Nochumsohn, L. and Azencot, O. Data augmentation policy search for long-term forecasting. *Trans. Mach. Learn. Res.*, 2025.

- Noh, H., Hong, S., and Han, B. Learning deconvolution network for semantic segmentation. In *ICCV*, 2015.
- Öztürk, H., Özgür, A., and Ozkirimli, E. DeepDTA: deep drug–target binding affinity prediction. *Bioinformatics*, 34(17):i821–i829, 2018.
- Park, D. S., Chan, W., Zhang, Y., Chiu, C., Zoph, B., Cubuk, E. D., and Le, Q. V. SpecAugment: A simple data augmentation method for automatic speech recognition. In 20th Annual Conference of the International Speech Communication Association, Interspeech, pp. 2613–2617. ISCA, 2019.
- Redmon, J., Divvala, S. K., Girshick, R. B., and Farhadi, A. You only look once: Unified, real-time object detection. In 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR, pp. 779–788, 2016.
- Redmond, M. Communities and Crime. UCI Machine Learning Repository, 2009.
- Roweis, S. and Saul, L. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323– 2326, 2000.
- Schneider, N., Goshtasbpour, S., and Perez-Cruz, F. Anchor data augmentation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Shorten, C. and Khoshgoftaar, T. M. A survey on image data augmentation for deep learning. *Journal of big data*, 6(1):1–48, 2019.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. In 3rd International Conference on Learning Representations, ICLR, 2015.
- Singer, A. and Wu, H.-T. Vector diffusion maps and the connection laplacian. *Communications on pure and applied mathematics*, 65(8):1067–1144, 2012.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pp. 2256–2265. pmlr, 2015.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- Tenenbaum, J., de Silva, V., and Langford, J. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.

- Vapnik, V. Principles of risk minimization for learning theory. Advances in neural information processing systems, 4, 1991.
- Verma, V., Lamb, A., Beckham, C., Najafi, A., Mitliagkas, I., Lopez-Paz, D., and Bengio, Y. Manifold mixup: Better representations by interpolating hidden states. In *International Conference on Machine Learning (ICML)*, 2019.
- Vinyals, O., Toshev, A., Bengio, S., and Erhan, D. Show and tell: A neural image caption generator. In *CVPR*, 2015.
- Weinberger, K. and Saul, L. Unsupervised learning of image manifolds by semidefinite programming. In *Proceedings* of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), volume 2, pp. 988–995, 2004.
- Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- Yang, S., Xiao, W., Zhang, M., Guo, S., Zhao, J., and Shen, F. Image data augmentation for deep learning: A survey. arXiv preprint arXiv:2204.08610, 2022.
- Yao, H., Wang, Y., Zhang, L., Zou, J., and Finn, C. C-mixup: Improving generalization in regression. In *Proceeding* of the Thirty-Sixth Conference on Neural Information Processing Systems, 2022.
- Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., and Yoo, Y. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6023–6032, 2019.
- Zhang, H., Cissé, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization. In 6th International Conference on Learning Representations, ICLR, 2018.
- Zhang, Z. and Zha, H. Principal manifolds and nonlinear dimension reduction via local tangent space alignment. *SIAM Journal on Scientific Computing*, 26(1):313–338, 2005.
- Zhong, Z., Zheng, L., Kang, G., Li, S., and Yang, Y. Random erasing data augmentation. In *Proceedings of the* AAAI conference on artificial intelligence, volume 34, pp. 13001–13008, 2020.
- Zhu, W., Qiu, Q., Huang, J., Calderbank, R., Sapiro, G., and Daubechies, I. LDMNet: Low dimensional manifold regularized neural networks. In *Proceedings of the IEEE*

conference on computer vision and pattern recognition, pp. 2743–2751, 2018.

A. Curvature-aware manifold learning

Given a train set $Z = \{z^1, \dots, z^N\}$, we parameterize the data manifold around a point $z \in Z$ that we map to u = 0, resulting in the following truncated Maclaurin series for a nearby point u_i :

$$g^{\alpha}(u_j) = u_j^T \nabla g^{\alpha} + \frac{1}{2} u_j^T H^{\alpha} u_j + \mathcal{O}(|u_j|_2^3) , \quad \alpha = 1, \dots, D - d .$$
 (5)

In order to estimate the gradient and Hessian of the embedding mapping g^{α} , we build a set of linear equations that solves Eq. 5. Particularly, we approximate g^{α} by solving the system $G = \Psi X$, where X holds the unknown elements of the gradients ∇g^{α} and the Hessians H^{α} , for every α . We define $g^{\alpha} = [g^{\alpha}(u_1), \dots, g^{\alpha}(u_k)]^T \in \mathbb{R}^k$, where u_j are points in the neighborhood of z := f(u), and $G = [g^1, \dots, g^{D-d}]$. The point u and points $\{u_j\}$ are associated with the train set Z under a natural orthogonal transformation. The local natural orthogonal coordinates are a set of coordinates that are defined at a specific point u of the manifold. They are constructed by finding a basis for the tangent space and normal space at a point u by applying principal component analysis on the neighborhood $N_z = \{z_j\}_{j=1}^k$. Namely, the first d coordinates (associated with the most significant modes, i.e., largest singular values) represent the tangent space, and the rest represent the normal space. Then, we define $\Psi = [\Psi_1, \dots, \Psi_k]$, stacking Ψ_j in rows, where Ψ_j is given via

$$\Psi_j = \left[u_j^1, \cdots, u_j^d, \left(u_j^1\right)^2, \cdots, \left(u_j^d\right)^2, \left(u_j^1 \times u_j^2\right), \cdots, \left(u_j^{d-1} \times u_j^d\right)\right],$$

and

$$X^{\alpha} = \left[\nabla g^{\alpha 1}, \cdots, \nabla g^{\alpha d}, H^{1,1}, \cdots, H^{d,d}, H^{1,2}, \cdots, H^{\alpha d-1,d}\right]^{T}$$

with $X = [X^1, \cdots, X^{D-d}]$. The set of linear equations

$$G = \Psi X$$
, xwhere $G \in \mathbb{R}^{k \times (D-d)}$ and $\Psi \in \mathbb{R}^{k \times \left(d + \frac{d(d+1)}{2}\right)}$ (6)

is solved by using the least square estimation given $X = \Psi^{\dagger}G$. In practice, we estimate only the upper triangular part of H^{α} since it is a symmetric matrix. We refer the reader for a more comprehensive and detailed treatment in (Li, 2018).

A.1. Approximation Error Bounds

In what follows, we provide a theoretical justification for the sampling error of CEMS in comparison to first-order approaches. Let $f : \mathbb{R}^d \to \mathbb{R}^D$ be a twice-differentiable function, we can express the Taylor expansion around a point u_0 up to first and second order as follows,

$$f^{(1)}(u) = f(u_0) + \nabla f(u_0)^T (u - u_0) ,$$

$$f^{(2)}(u) = f(u_0) + \nabla f(u_0)^T (u - u_0)$$
(7)

$$+\frac{1}{2}(u-u_0)^T H_f(u_0)(u-u_0).$$
(8)

Under standard smoothness assumptions, the approximation errors can be bounded as follows:

Theorem A.1 (Error Bounds). (Fowkes et al., 2013) For a twice-differentiable function f with Lipschitz continuous Hessian in a neighborhood of u_0 , we have that

$$\|f(u) - f^{(1)}(u)\| \le \frac{M}{2} \|u - u_0\|^2,$$
(9)

$$\|f(u) - f^{(2)}(u)\| \le \frac{L}{6} \|u - u_0\|^3,$$
(10)

where M bounds the spectral norm of $H_f(u)$ and L is the Lipschitz constant of $H_f(u)$ in the neighborhood of u_0 .

The second-order error decreases as $O(||u - u_0||^3)$ compared to $O(||u - u_0||^2)$ for first-order methods. This faster convergence rate ensures more accurate sampling in the vicinity of training points.

B. Intrinsic Dimension

Many techniques have been proposed for estimating a dataset's intrinsic dimension (ID). Global linear methods like principal component analysis (PCA) provide a simple baseline: one examines the eigenvalue spectrum to find the number of principal components needed to explain most variance. This works well if the data lie near a linear subspace, but it can overestimate ID for nonlinear manifolds or noisy data. In contrast, fractal geometric approaches such as the correlation dimension (Grassberger & Procaccia, 1983) measure how the number of point pairs grows with distance. The classic Grassberger–Procaccia algorithm computes the count of pairs within a radius r and uses the slope of the log–log plot of this count versus r as the estimated dimension. While the correlation dimension can capture non-integer (fractal) structure, it typically requires large samples and is sensitive to noise, limiting its robustness in practice. Overall, these early global methods can struggle with high curvature or sparse, high-dimensional data.

Local neighbor statistics have inspired more robust ID estimators. A prominent example is the maximum likelihood estimator (MLE) (Levina & Bickel, 2004), which uses distances to each point's k-nearest neighbors. By modeling the distribution of neighbor distances (under a Poisson process assumption), this method finds the ID that maximizes the likelihood of the observed spacing of points. The MLE approach showed improved accuracy over earlier techniques and is relatively robust across different data distributions. However, it requires choosing the neighborhood size k and can be biased if k is too small or large. Building on this idea, Facco et al. introduced the Two-Nearest Neighbors (TwoNN) estimator (Facco et al., 2017), which uses only the first and second nearest-neighbor distances for each point. TwoNN analytically derives an ID estimate from the ratio of these two distances, a minimalistic approach that greatly reduces bias from curved manifolds or non-uniform densities. This method is computationally lightweight and has been shown to produce consistent ID estimates even when data lie on a twisted manifold or are embedded in high-dimensional noise, highlighting its robustness.

Recently, topological data analysis methods, specifically persistent homology, have emerged as robust tools for ID estimation. Birdal et al. (Birdal et al., 2021) introduced a persistent homology-based estimator that leverages stable topological features, such as connected components and loops, which persist across multiple scales. The intuition behind persistent homology is that true geometric and topological features of a dataset tend to persist as one examines data at varying resolutions, while noise-induced features disappear quickly. This persistence-based approach provides reliable dimension estimates even in highly noisy, nonlinear, or sparse settings.

In summary, contemporary ID estimation techniques such as TwoNN, persistent homology-based measures, and MLE represent the state of the art, offering significantly improved robustness to noise, curvature, and sampling heterogeneity compared to classical PCA or correlation dimension methods.

C. Batch Selection

In our framework, for any data point $(x, y) \in \mathcal{X} \times \mathcal{Y}$, we seek to generate training samples that lie near the manifold \mathcal{M} , which represents the true data distribution \mathcal{P} . Since directly sampling from \mathcal{M} is impossible without knowing its structure, we instead approximate it locally using the tangent plane $\mathcal{T}(x, y)$ at point (x, y).

To construct this tangent plane approximation, we require a neighborhood set N_z around z. We compute the tangent plane using Singular Value Decomposition (SVD) on these neighboring points. The accuracy of this approximation heavily depends on how close the points in N_z are to z.

To implement this approach, we explored three distinct batch construction methods:

- 1. Random batch selection (random)
- 2. k-Nearest Neighbors batch selection (knn), where points are grouped based on their Euclidean distances in Z
- 3. k-Nearest Neighbors Probability batch selection (knnp), where points are sampled with probabilities inversely proportional to their distance from the original point, ensuring closer points have higher sampling probabilities

Our hypothesis was that both proximity-based and probability-based batch selection methods would generate samples closer to the data manifold, thereby improving model performance. We first trained models using the proximity-based batch selection method and selected the best-performing model based on the validation set. Using the optimal parameters found from this model, we then trained two additional models using random batch selection and knnp methods for fair comparison. The experimental results, presented in Table 3, demonstrate the relative performance of these three approaches

under identical parameter settings. The results reveal interesting patterns across different batch selection methods. While no single method dominates across all datasets, each approach shows strengths in specific scenarios. The knnp method demonstrates superior performance on the SkillCraft dataset and matches the best performance on NO2. The knn approach excels in both RCF and DTI datasets, showing particular strength in structured data scenarios. Interestingly, random batch selection remains competitive, achieving the best performance on Airfoil and matching the best result on NO2. These results suggest that the effectiveness of batch selection methods may be dataset-dependent, highlighting the importance of considering data characteristics when choosing a batch selection strategy.

C.1. Neighborhood Construction

The validity of CEMS relies on three key theoretical foundations. First, our neighborhood construction approach balances computational efficiency with geometric fidelity by sharing neighborhoods across points in close proximity. While this might appear to reduce sampling diversity, it actually preserves manifold structure because points that are close in the normalized input-output space typically share similar geometric properties. The shared neighborhood assumption is particularly valid because we normalize both input X and output Y features to the same range, ensuring that proximity in the combined space meaningfully reflects similarity in both domains.

The reliability of our construction is maintained even in regions of high output diversity through our careful treatment of the input-output space. For a point $z_i = [x_i, y_i]$, its neighborhood \mathcal{N}_z is constructed considering distances in both input and output spaces simultaneously, naturally limiting the diversity of outputs within each neighborhood. This approach ensures that points sharing a neighborhood basis have similar geometric properties, maintaining the validity of our second-order approximation.

The stochastic nature of mini-batch training provides an additional beneficial property for CEMS. As different batches are sampled each epoch, the method naturally explores varying neighborhoods and their associated tangent spaces. This dynamic sampling process enables CEMS to build a comprehensive representation of the manifold's local geometry, adapting to variations in data density across different regions. The continuous exploration of diverse local structures throughout training enhances the method's ability to capture the full geometric complexity of the underlying manifold.

The local-Euclidean structure of CEMS is supported by two complementary mechanisms. First, the second-order approximation naturally captures local curvature through the Hessian term, enabling accurate representation of nonlinear geometries. Second, our stochastic batch sampling strategy ensures exposure to different neighborhoods and tangent spaces throughout training. The projection of points onto the tangent space at μ (the neighborhood mean) maintains validity through the second-order terms in our Taylor expansion, which account for the primary nonlinearities within each neighborhood. This approach is particularly robust because the neighborhood size adapts with the batch size, preserving accurate local approximations even in regions of high curvature.

The empirical success of this construction is demonstrated in our ablation studies (Table 5), where we compare point-wise basis computation (CEMS_p) with our more efficient shared neighborhood approach (CEMS). The comparable performance across multiple datasets validates our theoretical assumptions about the effectiveness of shared geometric information within local neighborhoods.

C.2. Local vs. Global Sampling

It is important to note that CEMS focuses on local sampling within neighborhoods where the second-order approximation is valid. While alternative approaches based on geodesics can enable global sampling along the manifold, they typically incur significantly higher computational costs. Our local approach strikes a balance between sampling accuracy and computational efficiency, making it particularly well-suited for data augmentation during training.

Dataset	Airfoil↓	$ $ NO2 \downarrow	SkillCraft↓	RCF↓	DTI↑
CEMS - knn	1.455	0.507	5.142	0.146	0.511
CEMS - Knnp CEMS - random	1.441 1.435	0.506	4.941 5.155	0.175	0.309

Table 3. Results for different batch selection methods CEMS.

The theoretical guarantees provided by Theorem A.1 hold within the neighborhood where our Taylor approximations are valid. This aligns with the manifold hypothesis, which posits that real data typically lies on or near a lower-dimensional manifold with locally Euclidean structure (Goodfellow, 2016; Belkin & Niyogi, 2003). By focusing on accurate local sampling, CEMS can effectively augment the training data while maintaining the essential geometric structure of the underlying manifold.

D. Computational cost

General analysis of *n***-th order embedding maps:** To estimate the *n*-th order Taylor approximation of the embedding *g*, we need to find all the partial derivatives up to and including order *n*. The gradient ∇ vector is composed of *d* first-order partial derivatives, the Hessian matrix is composed of d^2 second-order partial derivatives. Third-order and higher partial derivatives are represented using mathematical objects called tensors. In general, the number of partial derivatives of a multi-input function grows exponentially with the order.

The order n of the partial derivatives determines the number of unknowns X and the size of the matrix Ψ which used to solve Eq. 6 via least squares. The number of columns in the matrix Ψ is $\sum_{i=1}^{n} d^{i} = \mathcal{O}(d^{n})$, and thus, the dimensions of the matrix Ψ are $k \times d^{n}$. It is preferred to solve an over-determined set of equations such that the number of neighbors $k > d^{n}$, which is memory and computationally expensive. For small values of d and n it is feasible to solve Eq. 6 in a run time complexity of $\mathcal{O}(kd^{2n})$. For larger values of n, it becomes extremely computationally and memory expensive to achieve $k \ge d^{n}$, therefore, we can assume that $k < d^{n}$.

There are two computationally expensive operations used to estimate the *n*-th order approximation of the manifold, SVD and least squares. The matrix on which we perform SVD is of shape $N_z \in \mathbb{R}^{k \times D}$ where *k* represents the number of neighbors and *D* is the extrinsic dimension. Thus, the run time complexity of the SVD operation per batch is $\mathcal{O}(\min(k^2D, kD^2))$. The complexity of least squares is determined by the dimensions of the matrix $\Psi \in \mathbb{R}^{k \times \mathcal{O}(d^n)}$ resulting in $\mathcal{O}(kd^{2n})$. If we wish to solve an over-determined system of equations, we need to set $k > d^n$ e.g., $k = 2d^n$ resulting in a run time of $\mathcal{O}(\min(d^{2n}D, d^nD^2))$ for SVD and $\mathcal{O}(d^nd^{2n}) = \mathcal{O}(d^{3n})$ for least squares.

Computational analysis of CEMS. Given a batch of data $A \in \mathbb{R}^{b \times D}$ where b is the batch size and D is the ambient dimension, our analysis considers the point-wise and batch-wise settings.

In the point-wise case, we construct a matrix $N_a \in \mathbb{R}^{b \times D}$ for every sample $a \in \mathbb{R}^D$ in the batch, containing its *b* closest neighbors, where the number of neighbors is fixed as the batch size. This practical choice decouples the computational complexity from the intrinsic dimension *d*. On each matrix N_a , we perform SVD at the complexity of $\mathcal{O}(\min(bD^2, Db^2))$. We then solve the set of *b* equations with $l = d \times (d + 1)/2$ variables (representing the unknowns in the gradient ∇ and Hessian *H*) at a complexity of $\mathcal{O}(b \times l^2) = \mathcal{O}(b \times d^4)$. In practice, the batch size is small and constant which leads to an underdetermined system that can be solved used Ridgr-Regression at a complexity of $\mathcal{O}(b^2 \times l) = \mathcal{O}(b^2 \times d^2)$. The total complexity per point is therefore $\mathcal{O}(\min(bD^2, Db^2)) + \mathcal{O}(b^2 \times d^2)$. Since the batch size *b* is usually smaller than the ambient dimension *D*, the total complexity can be revised as $\mathcal{O}(b^2(D + d^2))$. Under the manifold hypothesis, we assume that d << D and thus $d \in \mathcal{O}(D^2)$, resulting in the following complexity $\mathcal{O}(b^2D)$ for a single point and $\mathcal{O}(b^3D)$ for the entire batch. Our analysis reveals that under our assumptions, the complexity is proportional to the ambient dimension *D*.

In the batch-wise setting, the entire batch $A \in \mathbb{R}^{b \times D}$ is processed collectively. We compute the SVD of A at a complexity of $\mathcal{O}(\min(bD^2, Db^2))$. The subsequent step involves solving b equations with $l = d \times (d+1)/2$ variables at a complexity of $\mathcal{O}(b \times l^2)$. As in the point-wise case The total computational complexity of CEMS for the batch-wise setting for a single batch is $\mathcal{O}(b^2D)$

Run time comparison In Table 4, we compare the total run time of the training process in seconds to provide an estimate for the empirical computational cost of CEMS and competing methods. The results are obtained with a single RTX3090 GPU. For each data set, all the methods were estimated using the same parameters (e.g., batch size, number of epochs) for a fair comparison. It is evident from the results that the empirical run time of CEMS is o par with competing methods and does not require a large overhead.

E. Adaptation to batches

Below we provide the algorithm for the batched version as described in Sec. 4:

	AIRFOIL	NO2	RCF	DTI
ERM	3.84	1.01	172	653
C-MIXUP ADA FOMA CEMS	11.64 8.72 7.11 12.2	2.04 3.22 1.85 3.04	1700 465 364 445	1064 3519 1095 1317

Table 4. Training times comparison (in seconds).

Algorithm 2 CEMS: Batch generation

Require: Training data $Z = \{z^i = [x^i, y^i]\}_{i=1}^N$ **Require:** A sample $z \in Z$, batch size B. Set k = BFind K-nearest neighbors $N_z \leftarrow \{z_j\}_{j=1}^k \cup \{z = z_0\}$ of z Find an orthonormal basis B_u that spans $N_z - \mu_{N_z}$ 1: 2: Project every $z_j - \mu_{N_z}$ to the local orthonormal coordinates: $u_j \leftarrow B_{T_u}^T \cdot (z_j - \mu_{N_z}), \ g_j \leftarrow B_{N_u}^T \cdot (z_j - \mu_{N_z})$ For each $l = 1, \dots, k$ construct: $U_z^l \leftarrow \{u_j - u_l\}_{j \neq l}^k$ and $G_z^l \leftarrow \{g_j - g_l\}_{j \neq l}^k$ Construct G and Ψ as in Eq. 6 3: 4: 5: 6: 7: 8: Solve $\Psi A = G$ Extract $\nabla q(z)$ and H(z) from A 9: Sample a point η near u_l : $\eta \sim \mathcal{N}(u_l, \sigma I_d)$ 10: Calculate $g(\eta_l) \leftarrow g(u_l) + (\eta_l - u_l)^T \nabla g + \frac{1}{2} (\eta_l - u_l)^T H(\eta_l - u_l)$ Un-project η_l back to the original coordinates: 11: 12: $z_{\eta_l} := f(\eta_l) = B_u \cdot [\eta_l, g(\eta_l)] + \mu_{N_z}$ Return $z_\eta = \{z_{\eta_l}\}_{l=1}^k$ 13:

F. Application in Data Space vs. Latent Space

Our method can be applied in either the *data space* (i.e., raw input features) or a learned *latent space* (e.g., the output of a hidden layer in a neural network). Each setting presents different trade-offs and use cases:

- **Data Space:** Applying CEMS directly in the input space offers interpretability and preserves direct relationships between augmented samples and original features. This approach is model-agnostic and particularly suitable when the input space has meaningful geometric structure (e.g., tabular data). However, high-dimensional or noisy input spaces may not exhibit a well-defined manifold, limiting the effectiveness of local approximations.
- Latent Space: In contrast, applying CEMS in a learned latent space—typically a lower-dimensional and semantically structured representation—can lead to more compact and smoother manifolds. Latent representations often disentangle underlying factors of variation, making them more amenable to local geometric modeling. Moreover, when CEMS is implemented as a differentiable module in latent space, gradient signals can flow through the augmentation step. This enables the network to *shape the latent space with respect to the augmented samples* during training, potentially improving the alignment between the learned representation and the underlying manifold geometry.

Trade-offs: While latent space augmentation benefits from smoother geometry and end-to-end optimization, it relies on the quality of the learned representations and is tied to a specific model architecture. Data space augmentation, on the other hand, is more general and interpretable but may struggle in high-dimensional or non-smooth input spaces.

In our experiments, we found both variants to be viable, and the choice between them depends on the dataset characteristics and model architecture. We view a deeper investigation of this trade-off as an important direction for future work.

	Airfoil		NO2		Crimes (RMSE)		SkillCraft (RMSE)	
	RMSE	MAPE	RMSE	MAPE	Avg. \downarrow	Worst \downarrow	Avg. \downarrow	Worst ↓
CEMS_p	1.462	0.783	0.503	12.759	0.130	0.157	5.026	8.063
CEMS	1.455	0.809	0.507	12.807	0.128	0.159	5.142	6.322

Table 5. Ablation results for estimating the basis per point in the neighborhood (CEMS_p) vs. estimating it once and re-using the basis for every point N_z (CEMS).

G. Ablation Study

G.1. Ablation: basis computation per point vs. per neighborhood

As mentioned in Sec. 4, given a data point z, we construct a neighborhood N_z which can be used to 1) sample a point \tilde{z} near z 2) sample points \tilde{N}_z near N_z . The first option requires estimating a basis for the tangent space for each point in the dataset, whereas the second option estimates a single basis for the the entire batch of points N_z . In practice, the first option requires significantly more SVD calculations, determined by the batch size b. For large datasets, using the first option becomes very time consuming. In Tab. 5, we compare between Option 1 (CEMS_p), where p stands for point and Option 2 (CEMS), considering specifically the smaller datasets. Based on these results, we find that CEMS_p achieves error measures similar to CEMS. However, the computational complexity of CEMS_p is much higher, and thus we advocate the batch-wise computation as suggested in CEMS.

G.2. Ablation: Intrinsic dimension sensitivity

To evaluate the sensitivity of our method to intrinsic dimension estimation, we conduct a systematic analysis by perturbing the estimated values obtained from the TwoNN estimator. Table 6 presents results for in-distribution generalization (Airfoil and NO2) and out-of-distribution generalization (Crimes and SkillCraft), where the parameter id denotes the baseline dimension estimated by TwoNN, and id±1, id±2 correspond to configurations where the dimension is manually adjusted up or down. We find that performance is relatively stable under moderate changes in the intrinsic dimension. In many cases, the baseline value (id) yields the best or second-best results. For example, in the NO2 dataset, the lowest RMSE is achieved exactly at the TwoNN-estimated dimension (id), and nearby offsets produce slightly worse but comparable performance. Notably, in this dataset, the results for id+1 and id+2 are identical because the intrinsic dimension is capped at the extrinsic dimension minus one (i.e., 7), given the total number of features is 8. Similar robustness patterns are observed across other datasets. For instance, SkillCraft achieves the best average RMSE at id+1, while Crimes shows minimal variation across offsets, with the best worst-case RMSE appearing at id-2 and id+2. Additionally, Table 6 (bottom) compares intrinsic dimension estimates produced by different estimators (TwoNN, PH (Birdal et al., 2021), MLE (Levina & Bickel, 2004)). The strong agreement among estimators e.g., all three return 10 for Crimes supports the reliability of TwoNN and validates its use as our default choice. Overall, these findings suggest that our method is not overly sensitive to the intrinsic dimension, and that TwoNN provides a stable and empirically effective estimate.

Table 6. Comparison of intrinsic dimension estimates across different estimators: TwoNN, PH, and MLE. The TwoNN values are used in our method as the baseline intrinsic dimensions. Consistency between methods (e.g., all return 10 for Crimes) supports the robustness of the chosen estimates.

	Airfoil	NO2	Crimes	SkillCraft
TwoNN	3	6	10	12
PH	3	1	10	11
MLE	3	4	10	11

G.3. Ablation: Noise sensitivity

To assess the robustness of our method to the choice of the noise scale parameter σ , we conducted a sensitivity analysis by varying σ across a wide range of values, from $\sigma/10$ to $10 \cdot \sigma$. Table 7 summarizes the performance across four datasets in terms of average RMSE and MAPE (or worst-case RMSE) over three seeds. We observe that performance remains stable

Table 7. Sensitivity analysis of performance to changes in the estimated intrinsic dimension. The row id corresponds to the baseline
intrinsic dimension estimated by the TwoNN method. Rows id ± 1 and id ± 2 represent perturbations of the estimated dimension.
Results show average RMSE and MAPE across three seeds. Bold and underlined values indicate the best and second-best results per
column, respectively.

	Airfoil		NO2		Crimes		SkillCraft	
	RMSE	MAPE	RMSE	MAPE	Avg.	Worst	Avg.	Worst
id -2	1.466	0.818	<u>0.510</u>	12.688	0.129	0.158	5.183	6.541
id -1	<u>1.444</u>	0.810	0.511	<u>12.762</u>	0.128	0.161	5.211	6.705
id	1.454	<u>0.809</u>	0.507	12.807	0.128	<u>0.159</u>	<u>5.142</u>	<u>6.322</u>
id +1	1.473	0.820	0.523	13.375	0.133	0.172	5.070	6.192
id +2	1.433	0.798	0.523	13.375	<u>0.129</u>	0.158	5.196	6.599

Table 8. Sensitivity of performance to variations in the σ parameter. Each row corresponds to a different scaling factor applied to the default σ . Results show average RMSE and MAPE across three seeds. Bold and underlined values indicate the best and second-best results per column, respectively.

	Airfoil		NO2		Crimes		SkillCraft	
	RMSE	MAPE	RMSE	MAPE	Avg.	Worst	Avg.	Worst
σ/10	1.534	0.843	0.521	13.280	0.133	0.176	5.640	8.730
σ / 5	1.510	0.836	0.515	13.049	0.132	0.174	5.367	8.034
σ / 2	1.518	0.865	<u>0.513</u>	12.927	0.128	0.156	5.297	6.948
σ	1.454	0.809	0.507	12.807	0.128	<u>0.159</u>	5.142	<u>6.322</u>
$\sigma \times 2$	1.526	0.853	0.507	<u>12.824</u>	0.133	<u>0.159</u>	<u>5.173</u>	6.244
$\sigma \times 5$	<u>1.494</u>	0.831	0.520	12.961	0.137	0.164	5.849	7.613
$\sigma \times 10$	<u>1.494</u>	<u>0.826</u>	0.524	12.953	0.138	0.168	5.971	7.370

within a reasonable range around the default value σ , with minimal degradation at extreme low or high values. Notably, the best performance on NO2 is achieved exactly at the default σ , while SkillCraft achieves optimal or near-optimal results at both σ and $2 \cdot \sigma$. The variation in metrics is relatively small across the range, indicating that the method is not overly sensitive to the exact value of σ . These findings suggest that our approach is robust to noise scaling and does not require fine-grained tuning of this parameter.

G.4. Ablation: Batch sensitivity (neighborhood)

To assess the robustness of our method to the neighborhood size parameter B, which determines the size of the local batch used in tangent space estimation, we conduct a sensitivity analysis by scaling B across a wide range—from B/4 to 4B. As shown in Table 9, our method demonstrates stable performance across all datasets for both in-distribution (Airfoil, NO2) and out-of-distribution (Crimes, SkillCraft) settings. The default value B achieves the best or second-best results in most metrics, notably yielding the best RMSE and worst-case RMSE in SkillCraft, and the lowest RMSE in NO2. Interestingly, while smaller neighborhoods (e.g., B/2 or B/4) occasionally improve MAPE in Airfoil or average RMSE in Crimes, overly small or large values (e.g., B/4 or 4B) tend to degrade performance in other metrics. These results indicate that the method is not overly sensitive to the exact batch size, and that the default setting strikes a good balance between local geometric fidelity and stability. Overall, this analysis provides empirical support for the robustness of our hyperparameter selection strategy and clarifies the trade-offs involved in adjusting B, thereby addressing concerns regarding parameter sensitivity and general applicability of the method.

H. Geometric Properties Effect on CEMS

To evaluate the impact of curvature on CEMS for regression tasks, we generated a synthetic dataset where features lie on a manifold of constant scalar curvature. The manifold is defined as a hypersphere with a constant scalar curvature. Data points are sampled uniformly on the hypersphere using normalized random directions and embedded into a higher-dimensional

В	Airfoil		NO2		Crimes		SkillCraft	
	RMSE	MAPE	RMSE	MAPE	Avg.	Worst	Avg.	Worst
<i>B</i> / 4	<u>1.435</u>	0.792	0.544	14.021	<u>0.131</u>	0.166	<u>5.148</u>	<u>6.667</u>
<i>B</i> / 2	1.427	<u>0.800</u>	0.517	13.120	<u>0.131</u>	<u>0.162</u>	5.178	6.829
B	1.454	0.809	0.507	12.807	0.128	0.159	5.142	6.322
$B \times 2$	1.499	0.831	0.513	12.714	0.134	0.167	5.571	7.502
$B \times 4$	1.597	0.895	0.518	13.080	0.132	0.162	5.871	8.420

Table 9. Sensitivity analysis of performance to the batch size B. Each row corresponds to a different scale of the batch size. Results show average RMSE and MAPE (or worst RMSE) across three seeds. Bold and underlined values indicate the best and second-best results per column, respectively.

ambient space through a deterministic projection matrix to preserve the manifold structure. The regression target Y is computed as a non-linear function of the intrinsic coordinates, $Y = \sin (\sum X_{\text{intrinsic}})$, introducing a smooth dependency on the features. The features and targets are normalized to lie within the range [0, 1] using min-max scaling. This setup enables the systematic study of the effects of curvature of CEMS on regression model performance. In Fig. 3 the graph illustrate the relative improvement in RMSE between the CEMS and baseline ERM for regression tasks across varying scalar curvatures of the data manifold. The graph plots relative improvement against scalar curvature, highlighting that CEMS provides minimal advantage for nearly flat manifolds with low curvature but exhibits increasing improvement as the curvature grows. At higher curvatures (e.g., 16–64), CEMS demonstrates substantial gains, reflecting its ability to exploit geometric information in highly curved spaces. The hyperparameters of CEMS were not fine-tuned for this experiment and remained consistent across all intrinsic dimensions and curvature values. This likely explains why, in some cases, CEMS does not achieve better performance than ERM.

I. Hyperparameters

We present the hyperparameters for each dataset in Table 10. In our main results, we apply our method to the input space or the latent space, and we report the configuration with the best performance. All hyperparameters were selected through cross-validation and evaluated on the validation set. Some hyperparameters, such as architecture and optimizer, are not included in the tables since they remained unchanged and were used as specified in previous works (Yao et al., 2022; Schneider et al., 2023).



Figure 3. Illustration of the relative improvement in RMSE of CEMS over ERM. The graph demonstrates the effect of curvature, indicating minimal gains for nearly flat manifolds but substantial improvements for highly curved manifolds. These results emphasize the influence of data geometry on the performance of CEMS relative to ERM.

Curvature	Enhanced	Data	Augmentation	for	Regression
Cur / urur v					

Dataset	Airfoil	NO2	Exchange-Rate	Electricity	RCF	Crimes	SkillCraft	PovertyMap	DTI
Learning rate	$1e^{-3}$	$1e^{-3}$	$1e^{-3}$	$5e^{-4}$	$1e^{-4}$	$1e^{-3}$	$1e^{-3}$	$5e^{-3}$	$1e^{-4}$
Batch size	16	32	32	32	64	16	16	32	32
Input/Manifold	manifold	input	input	manifold	manifold	manifold	input	input	input
Epochs	700	100	200	100	50	100	100	50	60
σ	$1e^{-4}$	0.2	0.1	$1e^{-3}$	0.01	0.3	0.2	0.1	$1e^{-3}$

Table 10. Hyperparameter choices for the experiments using CEMS.

J. Dataset Description

This section provides detailed descriptions of the datasets used in our experiments.

Airfoil Self-Noise (Brooks et al., 2014). This dataset contains aerodynamic and acoustic test data for various NACA 0012 airfoils, recorded at different wind tunnel speeds and angles of attack. Each data point includes five features: frequency, angle of attack, chord length, free-stream velocity, and suction side displacement thickness, with the label representing the scaled sound pressure level. Input features are normalized using min-max normalization. The dataset is divided into 1003 training examples, 300 validation examples, and 200 test examples as noted in (Hwang & Whang, 2021).

NO2 (Aldrin, 2004). The NO2 dataset examines the relationship between air pollution near a road and traffic volume along with meteorological variables. Each input consists of seven features: the logarithm of the number of cars per hour, temperature at 2 meters above ground, wind speed, temperature difference between 25 and 2 meters above ground, wind direction, hour of the day, and the day number since October 1, 2001. The response variable is the hourly logarithm of NO2 concentration measured in Oslo from October 2001 to August 2003. The dataset is split into 200 training examples, 200 validation examples, and 100 test examples as in (Hwang & Whang, 2021).

Exchange-Rate (Lai et al., 2018). This time series dataset includes daily exchange rates for eight countries (Australia, Britain, Canada, Switzerland, China, Japan, New Zealand, and Singapore) from 1990 to 2016, totaling 7,588 observations with daily frequency. A sliding window size of 168 days is applied, resulting in an input dimension of 168×8 and a label dimension of 1×8 data points. The dataset is partitioned into training (60%), validation (20%), and test (20%) sets in chronological order as described in (Lai et al., 2018).

Electricity (Lai et al., 2018). This dataset contains hourly electricity consumption data from 321 clients, recorded every 15 minutes from 2012 to 2014, totaling 26,304 observations. A sliding window size of 168 is used, resulting in an input dimension of 168×321 and a label dimension of 1×321 . The dataset is divided into training, validation, and test sets following a methodology similar to that used for the Exchange-Rate dataset.

RCF (Yao et al., 2022). The RCF-MNIST (Rotated-Colored-Fashion) dataset features images with specific color and rotation attributes. Images are colored using RGB vectors based on the rotation angle $g \in [0, 1]$. In the training set, 80% of images are colored with [g, 0, 1 - g], and 20% with [1 - g, 0, g], creating a spurious correlation between color and label.

PovertyMap (Koh et al., 2021). Part of the WILDS benchmark (Koh et al., 2021), this dataset consists of satellite images from 23 African countries used to predict village-level asset wealth. Each input is a 224×22 multispectral LandSat image with 8 channels, and the label is the real-valued asset wealth index. The dataset is divided into 5 cross-validation folds with disjoint countries to facilitate the out-of-distribution setting, following the methodology in (Koh et al., 2021).

Crime (Redmond, 2009). The Communities And Crimes dataset merges socio-economic data from the 1990 US Census, law enforcement data from the 1990 US LEMAS survey, and crime data from the 1995 FBI UCR. It includes 122 attributes related to crime, such as median family income and percentage of officers in drug units. The target is the per capita violent crime rate. Numeric features are normalized to a range of 0.00 to 1.00, and missing values are imputed. The dataset is divided into training (1,390), validation (231), and test (373) sets, with 31, 6, and 9 disjoint domains, respectively.

SkillCraft (Blair et al., 2013). The SkillCraft dataset from UCI consists of video game telemetry data from real-time strategy (RTS) games, focusing on player expertise development. Each input includes 17 player-related parameters, such as

cognition-action-cycle variables and hotkey usage, while the label is the action latency. Missing data are filled by mean padding. The dataset is divided into training (1,878), validation (806), and test (711) sets with 4, 1, and 3 disjoint domains, respectively.

DTI (Huang et al., 2021). The Drug-Target Interactions dataset aims to predict the binding activity score between small molecules and target proteins. Input features include one-hot vectors for drugs and target proteins, and the output is the binding activity score. Training and validation data are from 2013 to 2018, while the test data spans 2019 to 2020. The "Year" attribute serves as domain information.

K. Results with Standard Deviation

In Table 11 we report the full results of in-distribution generalization and in Table 12 we report the full results of out-ofdistribution robustness.

Table 11. Full results for in-distribution generalization. Standard deviations are calculated over 3 seeds.

	Air	foil	NO2		
	RMSE	MAPE	RMSE	MAPE	
ERM	2.901 ± 0.067	1.753 ± 0.078	0.537 ± 0.005	13.615 ± 0.165	
Mixup	3.730 ± 0.190	2.327 ± 0.159	0.528 ± 0.005	13.534 ± 0.125	
Mani Mixup	3.063 ± 0.113	1.842 ± 0.114	0.522 ± 0.008	13.357 ± 0.214	
C-Mixup	2.717 ± 0.067	1.610 ± 0.085	$\underline{0.509 \pm 0.006}$	12.998 ± 0.271	
ADA	2.360 ± 0.133	1.373 ± 0.056	0.515 ± 0.007	13.128 ± 0.147	
FOMA	$\underline{1.471\pm0.047}$	$\underline{0.816\pm0.008}$	0.512 ± 0.008	$\underline{12.894\pm0.217}$	
CEMS	$\textbf{1.455} \pm \textbf{0.119}$	$\textbf{0.809} \pm \textbf{0.050}$	$\textbf{0.507} \pm \textbf{0.003}$	$\textbf{12.807} \pm \textbf{0.044}$	

	Exchan	ge-Rate	Electricity		
	RMSE	MAPE	RMSE	MAPE	
ERM	0.023 ± 0.003	2.423 ± 0.365	0.058 ± 0.001	13.861 ± 0.152	
Mixup	0.023 ± 0.002	2.441 ± 0.286	$\overline{0.058\pm0.000}$	14.306 ± 0.048	
Mani Mixup	0.024 ± 0.004	2.475 ± 0.346	$\overline{0.058\pm0.000}$	14.556 ± 0.057	
C-Mixup	0.020 ± 0.001	2.041 ± 0.134	$\overline{\textbf{0.057} \pm \textbf{0.001}}$	13.372 ± 0.106	
ADA	0.021 ± 0.006	2.116 ± 0.689	0.059 ± 0.001	$\overline{13.464\pm0.296}$	
FOMA	$\textbf{0.013} \pm \textbf{0.000}$	$\textbf{1.262} \pm \textbf{0.037}$	$\underline{0.058\pm0.000}$	14.653 ± 0.166	
CEMS	$\underline{0.014 \pm 0.001}$	$\underline{1.269\pm0.062}$	$\underline{0.058\pm0.000}$	$\textbf{13.353} \pm \textbf{0.217}$	

	RCF (RMSE)	(RMSE) Crimes (RM		SkillCraft (RMSE)	
	Avg. ↓	Avg. ↓	Worst ↓	Avg. ↓	Worst ↓
ERM	0.164 ± 0.007	0.136 ± 0.006	0.170 ± 0.007	6.147 ± 0.407	7.906 ± 0.322
Mixup	0.159 ± 0.005	0.134 ± 0.003	0.168 ± 0.017	6.461 ± 0.426	9.834 ± 0.942
ManiMixup	0.157 ± 0.021	0.128 ± 0.003	0.155 ± 0.009	5.908 ± 0.344	9.264 ± 1.012
C-Mixup	$\overline{\textbf{0.146}\pm\textbf{0.005}}$	$\overline{\textbf{0.123}\pm\textbf{0.000}}$	$\overline{\textbf{0.146}\pm\textbf{0.002}}$	5.201 ± 0.059	7.362 ± 0.244
ADA	0.163 ± 0.014	0.130 ± 0.003	0.156 ± 0.007	$\overline{5.301\pm0.182}$	6.877 ± 1.267
FOMA	0.159 ± 0.010	$\underline{0.128\pm0.004}$	0.158 ± 0.002	-	-
CEMS	$\textbf{0.146} \pm \textbf{0.002}$	$\underline{0.128\pm0.001}$	0.159 ± 0.004	$\textbf{5.142} \pm \textbf{0.143}$	$\textbf{6.322} \pm \textbf{0.191}$

Table 12. Full results for out-of-distribution robustness. Standard deviations are derived from a 5-fold data split in PovertyMap and or calculated over 3 seeds for other datasets.

	DTI	(R)	Poverty (R)		
	Avg. ↑	Worst ↑	Avg. ↑	Worst ↑	
ERM	0.483 ± 0.008	0.439 ± 0.016	0.80 ± 0.04	0.50 ± 0.07	
Mixup	0.459 ± 0.013	0.424 ± 0.003	$\overline{\textbf{0.81}\pm\textbf{0.04}}$	0.46 ± 0.03	
ManiMixup	0.474 ± 0.004	0.431 ± 0.009	-	-	
C-Mixup	0.498 ± 0.008	0.458 ± 0.004	$\textbf{0.81} \pm \textbf{0.03}$	$\textbf{0.53} \pm \textbf{0.07}$	
ADA	0.493 ± 0.010	0.448 ± 0.009	0.79 ± 0.03	0.52 ± 0.06	
FOMA	$\underline{0.503\pm0.008}$	$\underline{0.459\pm0.010}$	0.77 ± 0.03	$\overline{0.49\pm0.05}$	
CEMS	$\textbf{5.110} \pm \textbf{0.005}$	$\textbf{0.465} \pm \textbf{0.004}$	$\textbf{0.81} \pm \textbf{0.05}$	0.50 ± 0.07	