

Using the Gaussianized Coding Rate for Generative Modeling

Michael Psenka, Vishal Raman, Peter Tong

May 10, 2023

Abstract

Generative models all aim to model a target distribution by minimizing some “distributional distance” from the parametrized distribution to the target distribution. For the Linear Discriminative Representation (LDR) transcription framework, part of the novelty is the distributional distance used: a Gaussianized coding rate difference. Although this distributional distance hosts many benefits (efficiently computable, supports low-dimensional distributions), it is only a well-posed metric over Gaussian distributions. In this paper, we show twofold about the Gaussianized coding rate difference: (a) on a wider class of probability distributions, a class which a number of machine learning datasets fall into, the Gaussianized coding rate difference is still a well-posed metric, and (b) its simplicity allows for training stability that alternatives do not enjoy. We hope these theoretical insights guide future practitioners in how to design robust generative models using the Gaussianized coding rate distance.

1 Introduction

Distribution learning is a central task in machine learning: in its most loose definition, one aims to “learn” a probability distribution π using a finite set of samples $\{x_i\}_{i=1}^N$ from π . This is usually achieved in a parametric fashion: that is, we have a finite-dimensional parameter $\theta \in D$ and a “parametrization map” $\theta \rightarrow \pi_\theta$, and finally try to find $\theta^* \in \mathbb{R}^p$ such that $\pi_{\theta^*} \approx \pi$.

This task is usually approached through so-called *generative modeling*. Examples of generative modeling include *generative adversarial networks* (GANs) [6], *auto-encoders* (AEs) [9][7], and their various forms. A description of these methods can be found in [3] and its references.

A common concept used to both design and improve generative models is a “distributional distance”: a metric $d(\pi, \hat{\pi}) : \mathcal{D}(\mathbb{R}^d) \times \mathcal{D}(\mathbb{R}^d) \rightarrow \mathbb{R}$ between two probability distributions $\pi, \hat{\pi} \in \mathcal{D}(\mathbb{R}^d)$ over \mathbb{R}^d . This allows us to mathematically determine how close π and $\hat{\pi}$ are, adding rigor to the statement $\pi_{\theta^*} \approx \pi$: we say $\pi_{\theta^*} \approx \pi$ if $d(\pi_{\theta^*}, \pi) \leq \epsilon$ for small enough $\epsilon > 0$. Distribution learning can then be thought of minimizing a certain “distributional distance,” whose properties illuminate the properties of the respective generative model. A couple of examples:

- *KL-divergence* [11], the backbone of variational auto-encoders [8]. While this distance function is well-studied in information theory, the metric is only well-defined if both distributions $\pi, \hat{\pi} \in \mathcal{D}(\mathbb{R}^d)$ are absolutely continuous with respect to the Lebesgue measure on \mathbb{R}^d . This is broken if π or $\hat{\pi}$ have low-dimensional support, thus often in practice we need to artificially project $\pi, \hat{\pi}$ to a low-dimensional space of arbitrarily-chosen dimension.
- *Wasserstein distance* [17] (Earth-mover’s distance), the core component for the popular *Wasserstein GAN* [1]. While the Wasserstein distance is well-defined on distributions with low-dimensional support, it is computationally intractable, even on some very simple distributions such as Gaussians. In practice, the Wasserstein distance is approximated, but these approximations are not well-understood and thus lead to hidden downsides.

Finally, we turn to the main subject of this paper: LDR transcription [3], a distribution learning model that is built off of the following distributional distance:

$$\begin{aligned} \Delta R(\mathbf{Z}, \hat{\mathbf{Z}}) = & \frac{1}{2} \log \det \left(\mathbf{I} + \frac{\gamma}{2} (\mathbf{Z}\mathbf{Z}^\top + \hat{\mathbf{Z}}\hat{\mathbf{Z}}^\top) \right) \\ & - \frac{1}{4} \log \det \left(\mathbf{I} + \gamma \mathbf{Z}\mathbf{Z}^\top \right) \\ & - \frac{1}{4} \log \det \left(\mathbf{I} + \gamma \hat{\mathbf{Z}}\hat{\mathbf{Z}}^\top \right), \end{aligned} \tag{1}$$

where $\gamma := \frac{d}{n\epsilon^2}$ for some $\epsilon > 0$, and $\mathbf{Z}, \hat{\mathbf{Z}} \in \mathbb{R}^{d \times n}$ are matrices formed by collecting n samples from each of two distributions $\pi, \hat{\pi}$ supported on \mathbb{R}^d . This function comes from the coding rate of Gaussian distributions [14], where we measure the difference of coding the union of samples versus the sum of their coding rates. Eq. (1) enjoys the following nice properties:

- Eq. (1) is well defined for any distributions $\pi, \hat{\pi}$ that have finite variance, even if they are low-dimensional. While (1) is only defined on finite samples, we see by the central limit theorem that this metric does converge as we collect samples:

$$\begin{aligned} \lim_{n \rightarrow \infty} \Delta R(\mathbf{Z}, \hat{\mathbf{Z}}) &= \frac{1}{2} \log \det \left(\mathbf{I} + \frac{\alpha}{2} (\Sigma_\pi + \Sigma_{\hat{\pi}}) \right) \\ &\quad - \frac{1}{4} \log \det (\mathbf{I} + \alpha \Sigma_\pi) \\ &\quad - \frac{1}{4} \log \det (\mathbf{I} + \alpha \Sigma_{\hat{\pi}}), \end{aligned} \tag{2}$$

where $\alpha := \frac{d}{\epsilon^2}$, and $\Sigma_\pi = \mathbb{E}_\pi x x^\top := \int_{\mathbb{R}^d} x x^\top \pi(dx)$.

- Eq. (1) is interpretable, and thus practitioners can design algorithms using Eq. (1) in a principled way.
- Eq. (1), along with its gradients, are computable functions, as the function $\log \det(\cdot)$ has been studied extensively computationally.

Of course, as with any distributional distance, Eq. (1) has the following apparent (but empirically disproven) downsides:

- Eq. (1) is a “loose” distance function in the following sense: $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}}) = 0$ if and only if $\mathbf{Z}\mathbf{Z}^\top = \hat{\mathbf{Z}}\hat{\mathbf{Z}}^\top$ [18]. That is, two distributions have distance 0 only when their second moments are the same. This is well-posed between Gaussians, which are defined by their second moments, but two generic distributions may have the same second moments but wildly different behavior. This can be seen as a byproduct of the fact that $\Delta R(\cdot, \cdot)$ was derived assuming the input distributions are Gaussian.
- Difficulty to train. In the original LDR transcription paper, eq. (1) is optimized over a minimax objective, which as seen from many modern minimax distribution modelers, such as GANs, can lead to instability in training.

Of course, both of these downsides were empirically debunked in [3]: reconstruction accuracy was very high on image distributions known to be non-Gaussian, and the training process was very robust to hyperparameter tuning and initialization. However, these phenomena are not understood theoretically, thus making it difficult for outside practitioners to use the principles in [3] to make new generative models for their use case.

The purpose of this paper is to give clear, theoretical insight on why the Gaussianized coding rate difference in Eq. (1) is a good distributional distance to use for generative models, giving outside practitioners proper insight to develop their own models using Eq. (1). We show the two following ideal properties:

1. *The right amount of efficiency*: Eq. (1) is efficiently computable and interpretable because we assume the input distributions are Gaussians. However, if distributions are “close to” Gaussians in some way, then eq. (1) could be “close to” a proper metric. We show that:
 - (a) Eq. (1) is a well-defined distance over distributions with *kurtosis* close to that of Gaussians, and
 - (b) Experimental results on image datasets showing that they have kurtosis very similar to Gaussians.
Thus, Eq. (1) works properly as a metric on exactly the kinds of distributions we see in machine learning.
2. *Training stability*: Eq. (1) is stable around its critical points. Unlike in GAN training [5], where a small perturbation around a critical point can result in a massive change in loss, small changes in parameter space at a critical point of $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}})$ result in equivalently small changes in the loss, leading to more robust training.

2 The right amount of efficiency

As described in [14], Eq. (1) is derived by assuming that the distributions $\mathbf{Z}, \hat{\mathbf{Z}}$ are drawn from are Gaussian, yielding an efficient formula for computing the coding rate. However, distributions from realistic data are hardly Gaussian, so while Eq. (1) is computable on non-Gaussian distributions, it is not necessarily a well-defined metric outside of Gaussians: non-Gaussian distributions can be very different and still satisfy $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}}) = 0$.

The purpose of this section is to show that this is not a fatal problem, as many realistic data distributions are “close to Gaussian” in the right way, making $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}})$ remain a good distance of distributions while retaining is desirable computational properties. More concretely, we show this through the following:

1. We introduce the well-studied statistical quantity *kurtosis* [15] of a distribution, and describe how this measures the “closeness” of a distribution to a Gaussian distribution.
2. We show that if we consider $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}})$ where $\mathbf{Z}, \hat{\mathbf{Z}}$ are sampled from distributions with close-to-Gaussian kurtosis, then $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}}) = 0 \implies \int_{\mathbb{R}^d} |\pi - \hat{\pi}|(dx) \leq \epsilon$ for some small ϵ . In other words, $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}})$ is a proper metric on such distributions.
3. We give experimental data to show that a number of data distributions in modern machine learning have close-to-Gaussian kurtosis.
4. Finally, we show that simple neural network models preserve close-to-Gaussian kurtosis, thus training $\Delta R(f_{\theta} g_{\eta} f_{\theta}(\mathbf{X}), f_{\theta}(\mathbf{X}))$

2.1 Kurtosis

The *kurtosis* [15], denoted $k(\pi) \in \mathbb{R}$, of a distribution $\pi \in \mathcal{D}(\mathbb{R}^d)$ is a scalar measure that heuristically measures the “tailedness” of a distribution, and historically has been used to measure how close to Gaussian a distribution is: this is the property of interest for this paper.

First, we mathematically define the kurtosis. For a given dataset $X = \{x_1, x_2, x_3, \dots, x_n\}$, $X \in \mathbb{R}^{d \times n}$, we define its kurtosis to be

$$k(X) = \frac{1}{n} \sum (x_i - \bar{x})^T \Sigma^{-1} (x_i - \bar{x})^2 \quad (3)$$

where $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$, and $\Sigma := \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$.

Similarly, we define the kurtosis of a continuous distribution $\pi \in \mathcal{D}(\mathbb{R}^d)$ to be the following:

$$k(\pi) = \mathbb{E}_{\pi} \left((x - \mathbb{E}(x))^{\top} \Sigma^{-1} (x - \mathbb{E}(x))^{\top} \right)^2. \quad (4)$$

Here is a primary reason that kurtosis is used to measure how close to a Gaussian a distribution is: if π is any Gaussian over \mathbb{R}^d with invertible covariance, then $k(\pi) = d(d+2)$ (proof in Appendix A.1). Note that is only defined if Σ is invertible; this is because the kurtosis measures a scaled version of “tailedness” and is ill-defined in any direction without variance. We can fix this by projecting to the subspace with positive variance using PCA on Σ , where we then get that for *any* Gaussian $\pi = \mathcal{N}(\mu, \Sigma)$, $k(\pi) = k(k+2)$, where $k = \text{rank}(\Sigma)$.

Thus, we use the adjusted kurtosis to measure a distribution’s ratio of kurtosis compared to that of any full-rank Gaussian distribution:

$$k_{adjusted}(\pi) = \frac{k(\pi)}{d(d+2)} \quad (5)$$

The closer to 1 this number is, the closer to a Gaussian the kurtosis test indicates the distribution is [4].

2.2 ΔR on Gaussian-like distributions

When we say Gaussian-like, we mean distributions with low magnitude of *excess kurtosis*: $|k_{adjusted}(\pi) - 1| \leq \delta$. As we will show in the following section, many realistic data distributions fall into this category, even though they are non-Gaussian.

In this section of the main paper, we will give brief overviews with some intuition of the following idea: for two distributions π_1, π_2 with low excess kurtosis, equating their second moments becomes a fairly stronger condition. Thus, if ΔR is minimized on distributions with low excess kurtosis, even if they are not necessarily Gaussians, this minimization will force the distributions to converge.

The first and most obvious note is the following, immediately following from the triangle inequality:

Lemma 2.1.

Denote $B_{\delta} \subset \mathcal{D}(\mathbb{R}^d)$ the set of distributions such that $|k_{adjusted}(\pi) - 1| \leq \delta$. It follows that $|k_{adjusted}(\pi_1) - k_{adjusted}(\pi_2)| \leq 2\delta$ for all $\pi_1, \pi_2 \in B_{\delta}$.

Thus, B_{δ} does act like a δ -ball and forces member distributions to be close together. The main question to ask now is if closeness in $k_{adjusted}$ means closeness in anything stronger or more interpretable. The answer here depends on how many more assumptions we can make on the distributions π_1, π_2 :

1. Without any additional assumptions, we can assume the 4th moment of the distribution of 2-norm for each of π_1, π_2 differs by at most $2\delta C$, where C is a dimension-dependent constant. Using the theory of *cumulants* [2], we can conclude the following mouthful: the 4th order Taylor series of the log characteristic function $\log(\mathbb{E}_{\pi} e^{t \cdot X})$, where X is the 2-norm of the random variable distributed by π , differs by at most 2δ . Intuitively, the radial aggregate distribution of π_1, π_2 are similar up to a 4th order, rather than just the 2nd order similarity that equating second moments gives.

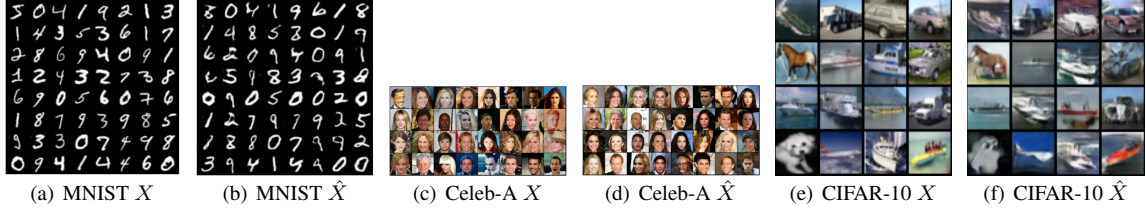


Figure 1: Visualizing the auto-encoding property of the learned ($\approx \hat{x} = g \circ f(x)$) on MNIST, Celeb-A and CIFAR-10 (zoom in for better visualization).

- For further conclusions on how close π_1, π_2 are, we need further assumptions on their regularity. For example, if we assume the distributions of π_1, π_2 don't fluctuate much, then we can establish a scaled L^1 bound $\int_{\mathbb{R}^d} \|x\|_4^4 |\pi_1 - \pi_2|(dx) \leq C\delta$, where C is a dimension-dependent constant.

The end result for this section: if our distributions for $\mathbf{Z}, \hat{\mathbf{Z}}$ both have low excess kurtosis (which we measured is true for our image distributions), and are well-behaved in another intuitive manner, then forcing $\Delta(\mathbf{Z}, \hat{\mathbf{Z}}) = 0$ imposes π_1, π_2 to be close to each other beyond their second moments.

2.3 Simulation results

For this section, we conduct experiments trying to demonstrate that the distribution of the selected dataset have kurtosis very close to Gaussian. We selected MNIST[12], Celeb-A[13] and CIFAR-10[10] datasets. LDR has performed extremely well on both the first two datasets and slightly worse on CIFAR-10. From Dai's LDR Transcription paper, we see that both MNIST and Celeb-A have great reconstruction results as shown in figure 1 only using ΔR terms. The reconstructed image looks very "real" comparing to the input data. While for CIFAR-10, although we see sample wise alignment but deviates a little in places like color and less accurate. For the simulation, we select one class from each dataset because in the ΔR term, we calculate it based on class information. All tensors in the dataset are flattened to vectors and then calculate their kurtosis and adjusted kurtosis respectively.

Dataset	Kurtosis	Gaussian Kurtosis	Adjusted Kurtosis
Celeb-A [13]	9444769	9443320	1.001
MNIST[12]	133130	123200	1.080
CIFAR-10[10]	336591	251000	1.341

Table 1: Kurtosis Result for MNIST, Celeb-A and CIFAR-10 dataset

From table 1, we see that after adjusting the kurtosis to adjusted kurtosis, both Celeb-A and MNIST are very close to 1, which indicates that both datasets are very close to Gaussian and contains most of their information in the second order. This sets a very good foundation for the analysis in the next few sections, which shows how $\Delta(R)$ can be very good objective function for these datasets. However, for the CIFAR-10 Dataset, we see that it has a much larger adjusted Kurtosis. This statistically indicates that it deviates further from the Gaussian.

2.4 Simple neural networks preserve kurtosis

In this section, we will discuss how low excess kurtosis in the dataset (i.e X) gets transferred to lower-dimensional latent space (i.e Z). In most of the simple neural networks include the one which LDR has performed the best - the Deep Convolutional Generative Adversarial Network (DCGAN) Network [16], network are consist of convolutional layers and batch norm layers, as shown in table 2.

Next we are going to prove that most parts in convolutional layers, which are fundamentally combination of linear transformation and batch norm layers preserve the kurtosis of the original data.

Theorem 2.2. Suppose for a set of data X , with A as a set of invertible linear transformation. Let $k(X)$ denote the kurtosis of for any given set of data X . Then $k(X) = k(AX)$

Proof. For convenience, we let $x = [x|1]$, $x' = [x'|1]$. In this case $x = Ax+b$ can be written as $x = Ax$.

$$S = \sum (x_i - \hat{x})(x_i - \bar{x})^T \quad (6)$$

Gray image $\in \mathbb{R}^{32 \times 32 \times 1}$
4×4 , stride=2, pad=1 conv 64 lReLU
4×4 , stride=2, pad=1 conv. BN 128 lReLU
4×4 , stride=2, pad=1 conv. BN 256 lReLU
4×4 , stride=1, pad=0 conv 128

Table 2: Encoder for MNIST.

$$k = \frac{1}{n} \sum (x_i - \bar{x})^T S^{-1} (x_i - \bar{x})^2 \quad (7)$$

since $x'_i = Ax$, $\hat{x}' = A\hat{x}$, Thus
 $S' = \sum (x'_i - \hat{x}'_i)(x'_i - \hat{x}'_i)^T = \sum A(x_i - \hat{x}_i)(x_i - \hat{x}_i)^T A^T = A \sum (x_i - \hat{x}_i)(x_i - \hat{x}_i)^T A^T = ASA^T$
Thus, we have:

$$k' = \frac{1}{n} \sum (x'_i - \bar{x}')^T S'^{-1} (x'_i - \bar{x}')^2 = \frac{1}{n} \sum (A(x_i - \bar{x}_i))^T S^{-1} (A(x_i - \bar{x}_i)) \quad (8)$$

$$= \frac{1}{n} \sum ((x_i - \bar{x}_i)^T A^T A^{-T} S^{-1} A^{-1} A (x_i - \bar{x}_i))^2 = \frac{1}{n} \sum ((x_i - \bar{x}_i)^T S^{-1} (x_i - \bar{x}_i))^2 = k \quad (9)$$

□

Theorem 2.3. Suppose for a set of data X , after a set of Batch Normalization, X' . Let $k(X)$ denote the kurtosis of for any given set of data X . Then $k(X) = k(X')$

Proof. Once again, we define $\bar{x} = \frac{1}{n} \sum x_i$, the covariance of the Data X is:

$$S = \sum (x_i - \hat{x})(x_i - \bar{x})^T \quad (10)$$

and the kurtosis, followed by definition is:

$$k = \frac{1}{n} \sum (x_i - \bar{x})^T S^{-1} (x_i - \bar{x})^2 \quad (11)$$

Since for each $x'_i = \frac{x_i - E[X]}{sd[X]}$ $S' = \sum (x'_i - \hat{x}'_i)(x'_i - \hat{x}'_i)^T = \sum x'_i x'^T_i = \sigma^2 S$ We have kurtosis being:

$$k' = \frac{1}{n} \sum (x'_i - \bar{x}')^T S'^{-1} (x'_i - \bar{x}')^2 \quad (12)$$

$$= \frac{1}{n} \sum (x_i)^T \sigma^2 S^{-1} (x_i)^2 \quad (13)$$

$$= \frac{1}{n} \sum \left(\frac{x_i - \bar{x}_i}{\sigma} \right)^T \sigma^2 S^{-1} \left(\frac{x_i - \bar{x}_i}{\sigma} \right)^2 \quad (14)$$

$$= \frac{1}{n} \sum (x_i - \bar{x}_i)^T S^{-1} (x_i - \bar{x}_i)^2 = k \quad (15)$$

□

3 Training Stability

A common problem for GAN training is stability [5], where small perturbations in the input parameters can lead to large deviations in the loss and thus lead to instability.

In this section, we show how the structure of $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}})$ naturally leads to stability at critical points. In our analysis, we consider simple linear models: $\mathbf{Z} = \mathbf{F}\mathbf{X}$, $\hat{\mathbf{Z}} = \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}$. While this is a simplification, we believe it is one that does more to highlight the true stability of ΔR . Increasing the complexity of f_θ, g_η will inevitably result in decreased stability of training ΔR : what we care about is how much additional instability training on ΔR adds on top of the instability of f_θ, g_η , which we now show is not much.

3.1 Linear Perturbation Bounds

Given $\mathbf{X} \in \mathbb{R}^D$, we define $\mathcal{T}_{\mathbf{X}}^b(\mathbf{F}, \mathbf{G}) := \Delta R(\mathbf{F}\mathbf{X}, \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X})$, the objective function we wish to train on. In this section, we give bounds on the perturbations $\mathcal{T}_{\mathbf{X}}^b(\mathbf{F} + \mathbf{M}, \mathbf{G})$ and $\mathcal{T}_{\mathbf{X}}^b(\mathbf{F}, \mathbf{G})$ where (\mathbf{F}, \mathbf{G}) are such that $\mathcal{T}_{\mathbf{X}}^b(\mathbf{F}, \mathbf{G}) = 0$ and \mathbf{M} has a bounded norm $\|\mathbf{M}\|_F \leq \delta$ for some fixed δ .

In order to do so, we will require the following facts:

Lemma 3.1. For a positive definite square matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$,

$$\text{trace}(\mathbf{I}_d - \mathbf{X}^{-1}) \leq \log \det \mathbf{X} \leq \text{trace}(\mathbf{X} - \mathbf{I}_d)$$

Lemma 3.2. For a positive semidefinite matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$, $\log \det(\mathbf{I}_d + \mathbf{X}) \geq 0$.

The proofs of the above lemmas are included in the Appendix.

Theorem 3.3. Given $\mathbf{X} \in \mathbb{R}^D$, let $\mathbf{M} \in \mathbb{R}^{D \times D}$ satisfy $\|\mathbf{M}\|_F \leq \delta < 1$ and (\mathbf{F}, \mathbf{G}) be so that $\mathbf{Z}\mathbf{Z}^* = \hat{\mathbf{Z}}\hat{\mathbf{Z}}^*$ where $\mathbf{Z} = \mathbf{F}\mathbf{X}$, $\hat{\mathbf{Z}} = \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}$. Then,

$$\mathcal{T}_{\mathbf{X}}^b(\mathbf{F}, \mathbf{G} + \mathbf{M}) \leq \frac{\gamma \|\mathbf{F}\mathbf{X}\|_F^2}{2} \left(1 + \delta \|\mathbf{F}\|_F^2 \left(1 + \frac{\|\mathbf{G}\|_F}{2} \right) \right).$$

Proof. First, note that by Lemma 3.2,

$$\begin{aligned} \mathcal{T}_{\mathbf{X}}^b(\mathbf{F}, \mathbf{G} + \mathbf{M}) &= \Delta R(\mathbf{Z}, \mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z}) \\ &= \frac{1}{2} \log \det \left(\mathbf{I} + \frac{\gamma}{2} (\mathbf{Z}\mathbf{Z}^* + \mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z}(\mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z})^*) \right) \\ &\quad - \frac{1}{4} \log \det (\mathbf{I} + \gamma \mathbf{Z}\mathbf{Z}^*) \\ &\quad - \frac{1}{4} \log \det (\mathbf{I} + \gamma (\mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z}(\mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z})^*)) \\ &\leq \frac{1}{2} \log \det \left(\mathbf{I} + \frac{\gamma}{2} (\mathbf{Z}\mathbf{Z}^* + \mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z}(\mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z})^*) \right) \end{aligned}$$

Then, note that we can write

$$\begin{aligned} \mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z}(\mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z})^* &= (\mathbf{F}\mathbf{G}\mathbf{Z} + \mathbf{F}\mathbf{M}\mathbf{Z})(\mathbf{F}\mathbf{G}\mathbf{Z} + \mathbf{F}\mathbf{M}\mathbf{Z})^* \\ &= (\hat{\mathbf{Z}} + \mathbf{F}\mathbf{M}\mathbf{Z})(\hat{\mathbf{Z}} + \mathbf{F}\mathbf{M}\mathbf{Z})^* \\ &= \hat{\mathbf{Z}}\hat{\mathbf{Z}}^* + \hat{\mathbf{Z}}(\mathbf{F}\mathbf{M}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}\hat{\mathbf{Z}}^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^* \\ &= \mathbf{Z}\mathbf{Z}^* + \mathbf{F}\mathbf{G}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{G}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^*. \end{aligned}$$

It follows that

$$\begin{aligned} \mathcal{T}_{\mathbf{X}}^b(\mathbf{F}, \mathbf{G} + \mathbf{M}) &\leq \frac{1}{2} \log \det \left(\mathbf{I} + \frac{\gamma}{2} (\mathbf{Z}\mathbf{Z}^* + \mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z}(\mathbf{F}(\mathbf{G} + \mathbf{M})\mathbf{Z})^*) \right) \\ &= \frac{1}{2} \log \det \left(\mathbf{I} + \gamma \mathbf{Z}\mathbf{Z}^* + \frac{\gamma}{2} (\mathbf{F}\mathbf{G}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{G}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^*) \right). \end{aligned}$$

By Lemma 3.1, the right-hand side is bounded by

$$\begin{aligned} &\frac{1}{2} \text{trace} \left(\gamma \mathbf{Z}\mathbf{Z}^* + \frac{\gamma}{2} (\mathbf{F}\mathbf{G}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{G}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^*) \right) \\ &= \frac{\gamma}{2} \text{trace}(\mathbf{Z}\mathbf{Z}^*) + \frac{\gamma}{4} \text{trace}(\mathbf{F}\mathbf{G}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{G}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^*) \\ &= \frac{\gamma}{2} \|\mathbf{Z}\|_F^2 + \frac{\gamma}{4} \text{trace}(\mathbf{F}\mathbf{G}\mathbf{Z}(\mathbf{F}\mathbf{M}\mathbf{Z})^* + \mathbf{F}\mathbf{M}\mathbf{Z}(\mathbf{F}\mathbf{G}\mathbf{Z})^*) + \frac{\gamma}{4} \|\mathbf{F}\mathbf{M}\mathbf{Z}\|_F^2 \\ &= \frac{\gamma}{2} \|\mathbf{Z}\|_F^2 + \frac{\gamma}{4} (\langle \mathbf{F}\mathbf{M}\mathbf{Z}, \mathbf{F}\mathbf{G}\mathbf{Z} \rangle_F + \langle \mathbf{F}\mathbf{G}\mathbf{Z}, \mathbf{F}\mathbf{M}\mathbf{Z} \rangle_F) + \frac{\gamma}{4} \|\mathbf{F}\mathbf{M}\mathbf{Z}\|_F^2 \\ &\leq \frac{\gamma}{2} \|\mathbf{Z}\|_F^2 + \frac{\gamma}{2} \|\mathbf{F}\mathbf{M}\mathbf{Z}\|_F \|\mathbf{F}\mathbf{G}\mathbf{Z}\|_F + \frac{\gamma}{4} \|\mathbf{F}\mathbf{M}\mathbf{Z}\|_F^2 \\ &\leq \frac{\gamma}{2} \|\mathbf{Z}\|_F^2 + \frac{\gamma}{2} \|\mathbf{F}\|_F^2 \|\mathbf{G}\|_F \|\mathbf{M}\|_F \|\mathbf{Z}\|_F^2 + \frac{\gamma}{4} \|\mathbf{F}\|_F^2 \|\mathbf{M}\|^2 \|\mathbf{Z}\|_F^2 \end{aligned}$$

where the last two lines follow from repeated applications of the Cauchy-Schwarz inequality.

Finally, since $\|\mathbf{M}\|_F < 1$, we have $\|\mathbf{M}\|_F^2 \leq \|\mathbf{M}\|_F \leq \delta$, it follows that

$$\begin{aligned} \frac{\gamma}{2}\|\mathbf{Z}\|_F^2 + \frac{\gamma}{2}\|\mathbf{F}\|_F^2\|\mathbf{G}\|_F\|\mathbf{M}\|_F\|\mathbf{Z}\|_F^2 + \frac{\gamma}{4}\|\mathbf{F}\|_F^2\|\mathbf{M}\|^2\|\mathbf{Z}\|_F^2 &= \frac{\gamma}{2}\|\mathbf{Z}\|_F^2 \left(1 + \|\mathbf{F}\|_F^2\|\mathbf{G}\|_F\|\mathbf{M}\|_F + \frac{1}{2}\|\mathbf{F}\|_F^2\|\mathbf{M}\|^2\right) \\ &\leq \frac{\gamma}{2}\|\mathbf{Z}\|_F^2 \left(1 + \delta\|\mathbf{F}\|_F^2\|\mathbf{G}\|_F + \frac{\delta}{2}\|\mathbf{F}\|_F^2\right) \\ &= \frac{\gamma\|\mathbf{F}\mathbf{X}\|_F^2}{2} \left(1 + \delta\|\mathbf{F}\|_F^2 \left(1 + \frac{\|\mathbf{G}\|_F}{2}\right)\right) \end{aligned}$$

□

Following similar analysis in the \mathbf{F} component, we obtain a corresponding theorem which is proved in the appendix:

Theorem 3.4. Given $\mathbf{X} \in \mathbb{R}^D$, let $\mathbf{M} \in \mathbb{R}^{D \times D}$ satisfy $\|\mathbf{M}\|_F \leq \delta < 1$ and (\mathbf{F}, \mathbf{G}) be so that $\mathbf{Z}\mathbf{Z}^* = \hat{\mathbf{Z}}\hat{\mathbf{Z}}^*$ where $\mathbf{Z} = \mathbf{F}\mathbf{X}$, $\hat{\mathbf{Z}} = \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}$. Then,

$$\mathcal{T}_{\mathbf{X}}^b(\mathbf{F} + \mathbf{M}, \mathbf{G}) \leq \frac{\gamma\|\mathbf{X}\|_F^2\|\mathbf{F}\|_F^2}{2} \left(1 + \frac{\delta\|\mathbf{G}\|_F^2\|\mathbf{F}\|_F}{2}\right).$$

Remark 3.5. By restricting to the case where $\|\mathbf{F}\|_F \leq 1$, $\|\mathbf{G}\|_F \leq 1$, $\|\mathbf{X}\|_F \leq 1$, we obtain the simpler bounds

$$\begin{aligned} \mathcal{T}_{\mathbf{X}}^b(\mathbf{F} + \mathbf{M}, \mathbf{G}) &\leq \frac{\gamma(2 + \delta)}{4}, \\ \mathcal{T}_{\mathbf{X}}^b(\mathbf{F}, \mathbf{G} + \mathbf{M}) &\leq \frac{\gamma(1 + 3\delta)}{4}. \end{aligned}$$

4 Conclusions & future directions

In this project, we have provided theoretical insights into two ideal properties - efficiency and stability - by using eq. (1) as a measurement for the distance between distributions. We introduce kurtosis to measure how close a data is to the Gaussian Distribution. showed that simple neural networks mostly preserve the kurtosis in the dataset. Furthermore, by studying the $\mathcal{T}_{\mathbf{X}}^b$ objective in the case of linear maps, we also provided an upper bound on the error when perturbing each input by a small amount. When normalizing our maps and data, we obtain bounds that are linear in the perturbation amount.

Of course, further researches and derivations are required for more complex neural networks such as the analysis on the Relu/Leaky-Relu layer's impact on kurtosis. While we chose kurtosis due to its historical significance with determining how Gaussian a distribution, there are likely better, less studied metrics out there to use in its place: our analysis highlights some possibilities for these. Additionally, further study is required to obtain tighter bounds on the linear perturbation errors.

References

- [1] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- [2] E. A. Cornish and R. A. Fisher. Moments and cumulants in the specification of distributions. *Revue de l'Institut international de Statistique*, pages 307–320, 1938.
- [3] X. Dai, S. Tong, M. Li, Z. Wu, K. H. R. Chan, P. Zhai, Y. Yu, M. Psenka, X. Yuan, H. Y. Shum, and Y. Ma. Closed-loop data transcription to an ldr via minimaxing rate reduction, 2021.
- [4] L. T. DeCarlo. On the meaning and use of kurtosis. *Psychological Methods*, 2:292–307, 1997.
- [5] F. Farnia and A. Ozdaglar. Gans may have no nash equilibria. *arXiv preprint arXiv:2002.09124*, 2020.
- [6] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [7] G. E. Hinton and R. S. Zemel. Autoencoders, minimum description length, and helmholtz free energy. *Advances in neural information processing systems*, 6:3–10, 1994.
- [8] D. P. Kingma and M. Welling. Auto-encoding variational bayes, 2014.
- [9] M. A. Kramer. Nonlinear principal component analysis using autoassociative neural networks. *AIChE journal*, 37(2):233–243, 1991.
- [10] A. Krizhevsky, V. Nair, and G. Hinton. Cifar-10 (canadian institute for advanced research).

- [11] S. Kullback and R. A. Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- [12] Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- [13] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [14] Y. Ma, H. Derksen, W. Hong, and J. Wright. Segmentation of multivariate mixed data via lossy data coding and compression. *IEEE transactions on pattern analysis and machine intelligence*, 29(9):1546–1562, 2007.
- [15] K. V. Mardia. Measures of multivariate skewness and kurtosis with applications. *Biometrika*, 57(3):519–530, 1970.
- [16] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks, 2016.
- [17] L. N. Vaserstein. Markov processes over denumerable products of spaces, describing large systems of automata. *Problemy Peredachi Informatsii*, 5(3):64–72, 1969.
- [18] Y. Yu, K. H. R. Chan, C. You, C. Song, and Y. Ma. Learning diverse and discriminative representations via the principle of maximal coding rate reduction. *arXiv preprint arXiv:2006.08558*, 2020.

A Appendix

The primary purpose of the appendix is to provide proofs for the mathematical statements made above.

A.1 Kurtosis proofs

We first address an elephant in the room: the requirement of kurtosis for the covariance to be full-rank, which is certainly not true in general. However, we can indeed generalize kurtosis to rank-deficient distributions:

Theorem A.1. $k(\pi)$ is well-defined on distributions with singular covariance matrix Σ

Proof. Consider the truncated SVD $\Sigma = U_r S_r U_r^\top$, where we truncate off all 0 singular value terms. Then $x \rightarrow U_r^\top x$ is an isometric transformation on the support of π , and the resulting covariance matrix $\Sigma' = U_r^\top \Sigma U_r$ is invertible, and we can measure the kurtosis $k(\pi)$ of the resulting distribution. \square

Theorem A.2. Denote $k(X)$ the estimated kurtosis as given in eq. (3) with a finite sample set $X \in \mathbb{R}^{d \times n}$. Then $k(X) = k(AX + b)$ for any invertible matrix $A \in GL(d)$ and any vector $b \in \mathbb{R}^d$.

Proof. $k(X) = k(X + b)$ is clear from the fact we subtract the empirical mean from X to generate $\bar{X} = X - \frac{1}{n} \sum_{i=1}^n x_i = X + b - \frac{1}{n} \sum_{i=1}^n (x_i + b)$. Thus, we only need to show $k(X) = k(AX)$. Note the effect of applying an invertible transformation on the covariance matrix:

$$S_{AX} = \sum (Ax_i - A\bar{x})(Ax_i - A\bar{x})^\top = ASA^\top \quad (16)$$

, where $S := \sum (x_i - \hat{x})(x_i - \hat{x})^\top$. Thus, we get the following for the kurtosis:

$$k(AX) = \frac{1}{n} \sum (Ax_i - A\bar{x})^\top (ASA^\top)^{-1} (Ax_i - A\bar{x}), \quad (17)$$

$$= \frac{1}{n} \sum ((x_i - \hat{x})A^\top A^{-\top} S^{-1} A^{-1} A(x_i - \hat{x}))^\top, \quad (18)$$

$$= \frac{1}{n} \sum ((x_i - \hat{x})S^{-1}(x_i - \hat{x}))^\top = k(X). \quad (19)$$

\square

We can prove the corresponding statement for the true kurtosis of continuous distributions:

Theorem A.3. Denote $k(\pi)$ the kurtosis as given in eq. (4) with respect to the distribution $\pi \in \mathcal{D}(\mathbb{R}^d)$ such that $\Sigma_\pi := \mathbb{E}_\pi(x - \bar{x})(x - \bar{x})^\top$ is invertible. Then $k(\pi)$ is invariant under affine invertible linear transformations.

Proof. We again reduce to the non-affine case since only the centered distribution is used, which is shift-invariant. For any invertible matrix A , the operation of multiplying a random variable sampled from π by A yields the distribution $\pi_A(dx) = \det(A)\pi(A^{-1}dx)$. Using the fact that $\int_{\mathbb{R}^d} f(x)\pi(A^{-1}x) = \frac{1}{\det(A)} \int_{\mathbb{R}^d} f(Ax)\pi(x)$ and $S_{\pi_A} = ASA^\top$, the proof follows exactly as in Theorem A.2. \square

Theorem A.4. Let $\pi = \mathcal{N}(\mu, \Sigma)$ be a Gaussian over \mathbb{R}^d , where Σ is invertible. Then the kurtosis $k(\pi) = d(d+2)$.

Proof. Since we have shown that kurtosis is invariant to affine linear transformations in Theorem A.3, we can assume without loss of generality that $\pi = \mathcal{N}(0, I)$. First, we make the following simplifications:

$$k(\pi) = \mathbb{E}_\pi(x^\top x)^2 = \int_{\mathbb{R}^d} (x^\top x)^2 \pi(dx), \quad (20)$$

$$= \int_{\mathbb{R}^d} x^\top x x^\top x \pi(dx), \quad (21)$$

$$= \int_{\mathbb{R}^d} \text{trace}(x^\top x x^\top x) \pi(dx), \quad (22)$$

$$= \int_{\mathbb{R}^d} \text{trace}(x x^\top x x^\top) \pi(dx), \quad (23)$$

$$= \text{trace} \left(\int_{\mathbb{R}^d} x x^\top x x^\top \pi(dx) \right), \quad (24)$$

$$= \text{trace} \left(\int_{\mathbb{R}^d} x x^\top \|x\|_2^2 \pi(dx) \right). \quad (25)$$

We now show that $K := \int_{\mathbb{R}^d} x x^\top \|x\|_2^2 \pi(dx) = (d+2)I$, and thus $k(\pi) = \text{trace}(K) = d(d+2)$. We start by showing the off-diagonal terms are 0:

$$K_{ij} = \int_{\mathbb{R}^d} x_i x_j \|x\|_2^2 \pi(dx), \quad (26)$$

$$= \int_{\mathbb{R}^d} x_i x_j \left(\sum_{k=1}^n x_k^2 \right) \pi(dx), \quad (27)$$

$$= \int_{\mathbb{R}^d} x_i^3 x_j \pi(dx) + \int_{\mathbb{R}^d} x_i x_j^3 \pi(dx) + \int_{\mathbb{R}^d} x_i x_j \left(\sum_{k \neq i, j} x_k^2 \right) \pi(dx), \quad (28)$$

$$= \int_{\mathbb{R}^d} x_i^3 \pi(dx) \left(\int_{\mathbb{R}^d} x_j \pi(dx) \right) + \int_{\mathbb{R}^d} x_i \pi(dx) \left(\int_{\mathbb{R}^d} x_j^3 \pi(dx) \right) + \int_{\mathbb{R}^d} x_i \pi(dx) \left(\int_{\mathbb{R}^d} x_j \pi(dx) \right) \left(\int_{\mathbb{R}^d} \left(\sum_{k \neq i, j} x_k^2 \right) \pi(dx) \right), \quad (29)$$

$$= 0 + 0 + 0 = 0, \quad (30)$$

since $\mathbb{E}_\pi f(x)g(x) = \mathbb{E}_\pi f(x) \mathbb{E}_\pi g(x)$ if $f(x), g(x)$ independent, and $\mathbb{E}_\pi x = 0$. Now for the diagonal terms:

$$K_{ii} = \int_{\mathbb{R}^d} x_i^2 \|x\|_2^2 \pi(dx), \quad (31)$$

$$= \int_{\mathbb{R}^d} x_i^2 \left(\sum_{k=1}^n x_k^2 \right) \pi(dx), \quad (32)$$

$$= \int_{\mathbb{R}^d} x_i^4 \pi(dx) + \int_{\mathbb{R}^d} x_i^2 \left(\sum_{k \neq i} \pi(dx) x_k^2 \right), \quad (33)$$

$$= 3 + \int_{\mathbb{R}^d} x_i^2 \pi(dx) \left(\int_{\mathbb{R}^d} \left(\sum_{k \neq i} x_k^2 \right) \pi(dx) \right), \quad (34)$$

$$= 3 + \int_{\mathbb{R}^d} \left(\sum_{k \neq i} x_k^2 \right) \pi(dx), \quad (35)$$

$$= 3 + d - 1 = d + 2, \quad (36)$$

where (34) follows since we know the 4th moment of a scalar normal distribution is 3, (35) follows from unit variance of a scalar normal distribution, and (36) follows from the fact that for a $d-1$ dimensional normal distribution $\pi_{d-1} = \mathcal{N}(0, I) \in \mathcal{D}(\mathbb{R}^{d-1})$:

$$\int_{\mathbb{R}^{d-1}} \|x\|_2^2 \pi(dx) = \int_{\mathbb{R}^{d-1}} x^\top x \pi(dx), \quad (37)$$

$$= \int_{\mathbb{R}^{d-1}} \text{trace}(x^\top x) \pi(dx), \quad (38)$$

$$= \int_{\mathbb{R}^{d-1}} \text{trace}(xx^\top) \pi(dx), \quad (39)$$

$$= \text{trace}\left(\int_{\mathbb{R}^{d-1}} xx^\top \pi(dx)\right), \quad (40)$$

$$= \text{trace}(I) = d - 1. \quad (41)$$

Thus, $K = (d + 2)I$ and $k(\pi) = \text{trace}(K) = d(d + 2)$. \square

A.2 Gaussian-like distributions proofs

Recall the definition of B_δ in Lemma 2.1. We show if the covariance matrices are normalized, then Gaussian-like distributions have similar 4th moments.

Without loss of generality, we assume all distributions are centered.

Theorem A.5. *Suppose $\pi_1, \pi_2 \in B_\delta$. Denote Σ_1, Σ_2 the covariance matrices of π_1, π_2 . If the operator norms $\|\Sigma_1\|, \|\Sigma_2\| \leq 1$, then $|\mathbb{E}_{\pi_1} \|x\|_2^4 - \mathbb{E}_{\pi_2} \|x\|_2^4| \leq 2\delta(d)(d + 1)$*

Proof. This follows from the following chain of inequalities:

$$|k_{adjusted}(\pi_1) - k_{adjusted}(\pi_2)| \leq 2\delta, \quad (42)$$

$$\implies \left| \int_{\mathbb{R}^d} (x^\top \Sigma_1^{-1} x)^2 \pi_1(dx) - \int_{\mathbb{R}^d} (x^\top \Sigma_2^{-1} x)^2 \pi_2(dx) \right| \leq 2\delta(d(d + 2)), \quad (43)$$

$$\implies \left| \int_{\mathbb{R}^d} (x^\top x)^2 \pi_1(dx) - \int_{\mathbb{R}^d} (x^\top x)^2 \pi_2(dx) \right| \leq 2\delta(d(d + 2)), \quad (44)$$

$$\implies \left| \int_{\mathbb{R}^d} \|x\|_2^4 \pi_1(dx) - \int_{\mathbb{R}^d} \|x\|_2^4 \pi_2(dx) \right| \leq 2\delta(d(d + 2)), \quad (45)$$

which follows since $\|\Sigma x\| \leq 1 \forall x \implies \|\Sigma^{-1} x\| \geq 1 \forall x$, and further by monotonicity of $(\cdot)^2$. \square

The above theorem gives the explicit connection of bounded excess kurtosis to similar 4th moments of the radial distribution. To see what this concludes about the characteristic functions, and namely similarity in 4th order Taylor polynomial expansions, we refer the reader to [2].

To make stronger assumptions, we need stronger assumptions on the distributions. One looser statement is the following:

Theorem A.6. *Suppose $\pi_1, \pi_2 \in B_\delta$ and both densities don't cross each other much: that is, either the Lebesgue measure of $\{x : \pi_1(x) \geq \pi_2(x)\}$ or $\{x : \pi_1(x) \leq \pi_2(x)\}$ is incredibly small. Then $\int_{\mathbb{R}^d} \|x\|_2^4 |\pi_1 - \pi_2|(dx) \lesssim 2\delta(d)(d + 1)$.*

Proof. From Theorem A.5, we know that $|\int_{\mathbb{R}^d} \|x\|_2^4 \pi_1(dx) - \int_{\mathbb{R}^d} \|x\|_2^4 \pi_2(dx)| \leq 2\delta(d(d + 2))$. We can decompose into the following, defining the measurable set $A := \{x : \pi_1(x) \leq \pi_2(x)\}$ and assuming that the measure of A^c is incredibly small:

$$\left| \int_{\mathbb{R}^d} \|x\|_2^4 \pi_1(dx) - \int_{\mathbb{R}^d} \|x\|_2^4 \pi_2(dx) \right| = \left| \int_A \|x\|_2^4 (\pi_1 - \pi_2)(dx) + \int_{A^c} \|x\|_2^4 (\pi_1 - \pi_2)(dx) \right|, \quad (46)$$

$$\approx \left| \int_A \|x\|_2^4 (\pi_1 - \pi_2)(dx) \right|, \quad (47)$$

$$= \int_A \|x\|_2^4 |(\pi_1 - \pi_2)|(dx) \approx \int_{\mathbb{R}^d} \|x\|_2^4 |(\pi_1 - \pi_2)|(dx). \quad (48)$$

Thus, $\int_{\mathbb{R}^d} \|x\|_2^4 |\pi_1 - \pi_2|(dx) \lesssim 2\delta(d)(d + 1)$. \square

Note that there was important hand-waving done here with the approximations, and we were not rigorous with the specific bounds. We aim in the future to improve these bounds; here, we simply want to show that with reasonable assumptions on the distributions π_1, π_2 , we can make better bounds for how different π_1, π_2 can be overall.

A.3 Stability Proofs

Lemma A.7 (Lemma 3.1). *For a positive definite square matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$,*

$$\text{trace}(\mathbf{I}_d - \mathbf{X}^{-1}) \leq \log \det \mathbf{X} \leq \text{trace}(\mathbf{X} - \mathbf{I}_d)$$

Proof. Let $\lambda_i(\mathbf{X})$ denote the eigenvalues of \mathbf{X} . We have

$$\begin{aligned} \log \det \mathbf{X} &= \log \left(\prod_{i=1}^d \lambda_i(\mathbf{X}) \right) \\ &= \sum_{i=1}^d \log \lambda_i(\mathbf{X}) \end{aligned}$$

Setting $0 < \lambda_i(X) = e^{\gamma_i}$, we have

$$\begin{aligned} \sum_{i=1}^d \log \lambda_i(\mathbf{X}) &= \sum_{i=1}^d \gamma_i \\ &\leq \sum_{i=1}^d e^{\gamma_i} - 1 \\ &= \sum_{i=1}^d \lambda_i - 1 \\ &= \text{trace}(\mathbf{X} - \mathbf{I}_d) \end{aligned}$$

by the convexity of e^x . Similarly, we have

$$\begin{aligned} \sum_{i=1}^d \log \lambda_i(\mathbf{X}) &= \sum_{i=1}^d \gamma_i \\ &\geq \sum_{i=1}^d 1 - e^{-\gamma_i} \\ &= \sum_{i=1}^d 1 - \frac{1}{\lambda_i} \\ &= \text{trace}(\mathbf{I}_d - \mathbf{X}^{-1}) \end{aligned}$$

□

Lemma A.8 (Lemma 3.2). *For a positive semidefinite matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$, $\log \det(\mathbf{I}_d + \mathbf{X}) \geq 0$.*

Proof. By Proposition 3.1, we have

$$\begin{aligned} \log \det(\mathbf{I}_d + \mathbf{X}) &\geq \text{trace}(\mathbf{I}_d - (\mathbf{I}_d + \mathbf{X})^{-1}) \\ &= d - \text{trace}((\mathbf{I} + \mathbf{X})^{-1}) \\ &= d - \sum_{i=1}^d \frac{1}{1 + \lambda_i(\mathbf{X})} \\ &\geq d - \sum_{i=1}^d 1 \\ &= 0. \end{aligned}$$

□

Theorem A.9 (Theorem 3.4). *Given $\mathbf{X} \in \mathbb{R}^D$, let $\mathbf{M} \in \mathbb{R}^{D \times D}$ satisfy $\|\mathbf{M}\|_F \leq \delta < 1$ and (\mathbf{F}, \mathbf{G}) be so that $\mathbf{Z}\mathbf{Z}^* = \hat{\mathbf{Z}}\hat{\mathbf{Z}}^*$ where $\mathbf{Z} = \mathbf{F}\mathbf{X}$, $\hat{\mathbf{Z}} = \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}$. Then,*

$$\mathcal{T}_{\mathbf{X}}^b(\mathbf{F} + \mathbf{M}, \mathbf{G}) \leq \frac{\gamma \|\mathbf{X}\|_F^2 \|\mathbf{F}\|_F^2}{2} \left(1 + \frac{\delta \|\mathbf{G}\|_F^2 \|\mathbf{F}\|_F}{2} \right).$$

Proof. Note that

$$\begin{aligned}
(\mathbf{F} + \mathbf{M})\mathbf{G}(\mathbf{F} + \mathbf{M})\mathbf{X}\mathbf{X}^*(\mathbf{F} + \mathbf{M})^*\mathbf{G}^*(\mathbf{F} + \mathbf{M})^* &= \hat{\mathbf{Z}}\hat{\mathbf{Z}}^* + \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^*
\end{aligned}$$

By Lemma 3.2, it suffices to consider the first term of $\mathcal{T}_{\mathbf{X}}^b(\mathbf{F} + \mathbf{M}, \cdot)$. Hence, we have

$$\begin{aligned}
\mathcal{T}_{\mathbf{X}}^b(\mathbf{F} + \mathbf{M}, \cdot) &\leq \frac{1}{2} \log \det(\mathbf{I} + \frac{\gamma}{2}(2\mathbf{Z}\mathbf{Z}^* + \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* + \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* + \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* + \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* + \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^*)) \\
&\leq \frac{\gamma}{4} \text{trace}(2\mathbf{Z}\mathbf{Z}^* + \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* + \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* + \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{F}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* + \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{F}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^* \\
&+ \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{F}^*\mathbf{G}^*\mathbf{M}^* + \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{F}^* + \mathbf{M}\mathbf{G}\mathbf{M}\mathbf{X}\mathbf{X}^*\mathbf{M}^*\mathbf{G}^*\mathbf{M}^*)) \\
&\leq \frac{\gamma}{2} \|\mathbf{Z}\|_F^2 + \frac{\gamma}{4} \|\mathbf{G}\|_F^2 \|\mathbf{X}\|_F^2 (4\|\mathbf{F}\|_F^3 \|\mathbf{M}\|_F + 6\|\mathbf{F}\|_F^2 \|\mathbf{M}\|_F^2 + 4\|\mathbf{F}\|_F \|\mathbf{M}\|_F^3 + \|\mathbf{M}\|_F^4) \\
&\leq \frac{\gamma \|\mathbf{X}\|_F^2}{2} \left(\|\mathbf{F}\|_F^2 + \frac{1}{2} (\|\mathbf{G}\|_F^2 (\|\mathbf{F}\|_F + \delta)^4 - \|\mathbf{F}\|_F^4) \right) \\
&\leq \frac{\gamma \|\mathbf{X}\|_F^2}{2} \left(\|\mathbf{F}\|_F^2 + \frac{1}{2} \delta \|\mathbf{G}\|_F^2 \|\mathbf{F}\|_F^3 \right) \\
&= \frac{\gamma \|\mathbf{X}\|_F^2 \|\mathbf{F}\|_F^2}{2} \left(1 + \frac{\delta \|\mathbf{G}\|_F^2 \|\mathbf{F}\|_F}{2} \right).
\end{aligned}$$

□