

SPREAD DIVERGENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

For distributions p and q with different supports, the divergence $D(p||q)$ may not exist. We define a spread divergence $\tilde{D}(p||q)$ on modified p and q and describe sufficient conditions for the existence of such a divergence. We demonstrate how to maximize the discriminatory power of a given divergence by parameterizing and learning the spread. We also give examples of using a spread divergence to train and improve implicit generative models, including linear models (Independent Components Analysis) and non-linear models (Deep Generative Networks).

1 INTRODUCTION

A divergence $D(p||q)$ (see, for example Dragomir (2005)) is a measure of the difference between two distributions p and q with the property

$$D(p||q) \geq 0 \text{ and } D(p||q) = 0 \Leftrightarrow p = q \quad (1)$$

We are interested in situations in which the supports of the two distributions are different, $\text{supp}(p) \neq \text{supp}(q)$. An important class is the f -divergence, defined as

$$D_f(p||q) = \mathbb{E}_{q(x)} \left[f \left(\frac{p(x)}{q(x)} \right) \right] \quad (2)$$

where $f(x)$ is a convex function with $f(1) = 0$. A special case of an f -divergence is the well-known Kullback-Leibler divergence $\text{KL}(p||q) = \mathbb{E}_{p(x)} \left[\log \frac{p(x)}{q(x)} \right]$. By setting $p(x)$ to the empirical data distribution, maximum likelihood training of a model $q(x)$ corresponds to minimising $\text{KL}(p||q)$. However, this divergence may not be defined since the ratio $p(x)/q(x)$ can cause a division by zero.

This is a challenge since popular implicit generative models (Mohamed & Lakshminarayanan (2016)) of the form $q(x) = \int \delta(x - g_\theta(z)) p(z) dz$ only have limited support. In this case, maximum likelihood to learn the model parameters θ is not available and alternative approaches to measure the difference between distributions such as Maximum Mean Discrepancy (Gretton et al. (2012)) or Wasserstein distance (Peyré et al. (2019)) are required.

2 SPREAD DIVERGENCE

From $q(x)$ and $p(x)$ we define new distributions $\tilde{q}(y)$ and $\tilde{p}(y)$ that have the same support¹. Using the notation \int_x to denote integration $\int (\cdot) dx$ for continuous x , and $\sum_{x \in \mathcal{X}}$ for discrete x with domain \mathcal{X} , we define a random variable y with the same domain as x and distributions

$$\tilde{p}(y) = \int_x p(y|x)p(x), \quad \tilde{q}(y) = \int_x p(y|x)q(x) \quad (3)$$

where $p(y|x)$ is ‘noise’ designed to ‘spread’ the mass of p and q such that $\tilde{p}(y)$ and $\tilde{q}(y)$ have the same support. For example, if we use a Gaussian $p(y|x) = \mathcal{N}(y|x, \sigma^2)$, then \tilde{p} and \tilde{q} both have support \mathbb{R} . We also impose an additional requirement on the noise $p(y|x)$, namely that $D(\tilde{p}||\tilde{q}) = 0 \Leftrightarrow p = q$. As we show in section(2.1) this is guaranteed for certain ‘noise’ distributions. Given these requirements, we can define the Spread Divergence $\tilde{D}(p||q) \equiv D(\tilde{p}||\tilde{q})$. This satisfies the divergence requirement $\tilde{D}(p||q) \geq 0$ and $\tilde{D}(p||q) = 0 \Leftrightarrow p = q$.

¹For simplicity, we use univariate x , with the extension to the multivariate setting being straightforward.

For example, given two delta distributions $p_0(x) = \delta(x - x_0)$, $p_1(x) = \delta(x - x_1)$, the KL divergence (or f-divergence) between them is not defined. However, the spread KL divergence (or spread divergence) is defined. Assume a Gaussian noise distribution $p(y|x) = N(y|x; \sigma^2)$, the "spread" delta distributions have the form $p_0(y) = \int_{\mathcal{X}} \delta(x - x_0) N(y|x; \sigma^2) dx = N(y|x_0; \sigma^2)$, $p_1(y) = \int_{\mathcal{X}} \delta(x - x_1) N(y|x; \sigma^2) dx = N(y|x_1; \sigma^2)$. Therefore, the spread KL divergence (with Gaussian noise) between two delta distributions is equivalent to the KL divergence between two Gaussian distributions with the same variance, which has closed form (see appendix(D) for a derivation):

$$KL(p_0(x)||p_1(x)) = KL(p_0(y)||p_1(y)) = \frac{1}{2\sigma^2} \|x_0 - x_1\|_2^2 \quad (4)$$

It's worth noticing that in the case of two delta distributions, the spread KL divergence is equal to the squared 2-Wasserstein distance (see Peyré et al. (2019); Gelbrich (1990)).

2.1 NOISE REQUIREMENTS FOR A SPREAD DIVERGENCE

Our main interest is in using noise to define a new divergence in situations in which the original divergence $D(p||q)$ is itself not defined. For discrete variables $x \in \{1, \dots, n\}$, $y \in \{1, \dots, n\}$, the noise $p_{ij} = p(y = i | x = j)$ must be a distribution $\sum_j p_{ij} = 1$, $p_{ij} \geq 0$ and

$$\sum_j p_{ij} p_j = \sum_j p_{ij} q_j \quad \forall i \quad p_j = q_j \quad \forall j \quad (5)$$

which is equivalent to the requirement that the matrix is invertible. There is an additional requirement that the spread divergence exists. In the case of divergences, the spread divergence exists provided that p and q have the same support. This is guaranteed if

$$\sum_j p_{ij} p_j > 0; \quad \sum_j p_{ij} q_j > 0 \quad \forall i \quad (6)$$

which is satisfied if $p_{ij} > 0$. In general, therefore, there is a space of noise distributions $p(y|x)$ that define a valid spread divergence. For example, the 'antifreeze' method of Furnston & Barber (2009) is a special form of spread noise to define a valid KL divergence (see also Barber (2012)).

For continuous variables, in order for $D(p||q) = 0 \iff p = q$, the noise $p(y|x)$, with $\dim(Y) = \dim(X)$ must be a probability density and satisfy

$$\int_{\mathcal{Z}} p(y|x)p(x)dx = \int_{\mathcal{Z}} p(y|x)q(x)dx \quad \forall y \in \mathcal{Y} \quad p(x) = q(x) \quad \forall x \in \mathcal{X} \quad (7)$$

In the following section we discuss the special case of stationary noise for continuous systems.

3 STATIONARY SPREAD DIVERGENCES

Consider stationary noise $p(y|x) = K(y - x)$ where $K(x)$ is a probability density function with $K(x) > 0, x \in \mathbb{R}$. In this case p and q are defined as a convolution

$$p(y) = \int_{\mathcal{X}} K(y - x)p(x)dx = (K * p)(y); \quad q(y) = \int_{\mathcal{X}} K(y - x)q(x)dx = (K * q)(y) \quad (8)$$

Since $K > 0$, p and q are guaranteed to have the same support. A sufficient condition for the existence of the Fourier Transform $\mathcal{F}f$ of a function $f(x)$ for real x is that f is absolutely integrable. Since all distributions $p(x)$ are absolutely integrable, both $\mathcal{F}p$ and $\mathcal{F}q$ are guaranteed to exist. Assuming $\mathcal{F}f * K$ exists, we can then use the convolution theorem to write

$$\mathcal{F}f * p = \mathcal{F}f * K * \mathcal{F}f * p; \quad \mathcal{F}f * q = \mathcal{F}f * K * \mathcal{F}f * q \quad (9)$$

Let $\mathcal{F}f * K \neq 0$ or $\mathcal{F}f * K = 0$ on at most a countable set. Then

$$\mathcal{F}f * K * \mathcal{F}f * p = \mathcal{F}f * K * \mathcal{F}f * q \implies \mathcal{F}f * p = \mathcal{F}f * q \quad (10)$$

The proof is given in appendix(A). Using this we can write

$$D(p||q) = 0, \quad p = q \quad (11)$$

$$\mathcal{F}f * K * \mathcal{F}f * p = \mathcal{F}f * K * \mathcal{F}f * q \quad (12)$$

$$\mathcal{F}f * p = \mathcal{F}f * q, \quad p = q; \quad (13)$$

where we used the invertibility of the Fourier transform. Hence, for stationary noise $p(y|x) = K(y - x)$, we can define a valid spread divergence provided $K(x)$ is a probability density function and (ii) $\int f(x)K(x)dx > 0$ or $\int f(x)K(x)dx = 0$ on at most a countable set. Interestingly, the sufficient conditions for defining a valid spread divergence such that $D_f(p||q) = 0$, $p = q$ are analogous to the characteristic condition on kernels such that the Maximum Mean Discrepancy $MMD(p; q) = 0$, $p = q$, see Sriperumbudur et al. (2011; 2012); Gretton et al. (2012). As an example of such a noise process, consider Gaussian noise,

$$K(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}x^2}; \quad \int f(x)K(x)dx = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{if(x)} e^{-\frac{1}{2\sigma^2}x^2} dx = e^{-\frac{\sigma^2}{2}f^2} > 0 \quad (14)$$

Similarly, for Laplace noise

$$K(x) = \frac{1}{2b} e^{-\frac{|x|}{b}}; \quad \int f(x)K(x)dx = \frac{1}{2b} \int_{-\infty}^{\infty} e^{if(x)} e^{-\frac{|x|}{b}} dx > 0 \quad (15)$$

Since in both cases $\int f(x)K(x)dx > 0$ and $\int f(x)K(x)dx = 0$, Gaussian and Laplace noise can be used to define a valid spread divergence.

4 MAXIMISING DISCRIMINATORY POWER

From the data processing inequality (see appendix(B)), adding spread noise will always decrease the f -divergence $D_f(p(y)||q(y)) \leq D_f(p(x)||q(x))$. Intuitively, spreading out distributions makes them more similar. If we are to use a spread divergence to train a model using maximum likelihood (see section(5)), there is the danger that adding too much noise may make the spreaded empirical distribution and spreaded model distribution so similar that it becomes difficult to numerically distinguish them, impeding training. It is useful therefore to define spread noise that maximally discerns the difference between the two distributions $\max_x D(p(y)||q(x))$ for spread noise $p(y|x)$ parameterised by x . In general we need to constrain the spread noise to ensure that the divergence remains bounded.

We discuss below two complementary approaches to $\max_x D(p(y)||q(x))$ during training. The first approach adjusts the dimension-wise correlations (this corresponds to adjusting the covariance structure for Gaussian $p(y|x)$) and the second forms a mean transformation. In principle, both approaches can be combined and easily generalized to other noise distributions, such as Laplace noise.

4.1 LEARNING COVARIANCE STRUCTURE

Learning the covariance adjusts the shape of noise centered around the original model manifold. When we maximize the divergence between two spreaded distributions $D(p(y)||q(x))$, the learned noise will discourage overlap between the two distributions. Hence, if the data and model lie on the same manifold, the noise will be orthogonal to the manifold.

In learning the Gaussian spread distribution $p(y|x) = N(y|x; \Sigma)$, the number of parameters in the covariance matrix scales quadratically with the data dimension D . We thus define $\Sigma = \sigma^2 I + LL^T$ where $\sigma^2 > 0$ is fixed (to ensure bounded spread divergence) and L is a learnable $D \times R$ matrix with $R \leq D$. Calculating the log likelihood and sampling can then be performed efficiently using standard Woodberry identities, see appendix(J).

4.2 LEARNING THE MEAN TRANSFORM

Consider $p(y|x) = K(y - f(x))$ for injective² f and stationary K . Then, we define

$$\tilde{p}(y) = \int K(y - f(x))p_x(x)dx \quad (16)$$

Note that this is a valid spread divergence since, using change of variables,

$$\tilde{p}(y) = \int K(y - z)p_z(z)dz; \quad p_z(z) = p_x(f^{-1}(z)) = \int p_x(x)\delta(x - f^{-1}(z))dx \quad (17)$$

²Since the codomain of f is determined by its range, injective indicates invertible in this case.

where J is the absolute Jacobian of f . Hence $D(p(y|x)) = 0$, $p_z = q_z$, $p_x = q_x$. Each injective f gives a different noise $p(y|x)$, we can thus search for the best noise implicitly by learning

In our experiments we use the invertible residual network Behrmann et al. (2018) $\mathbb{R}^D \rightarrow \mathbb{R}^D$ with $f = (f^1 \dots f^T)$ denotes a ResNet with blocks $f_t = I + g_t(\cdot)$. Then f is invertible if the Lipschitz-constant $\text{Lip}(g_t) < 1$ for all $t = 1, \dots, T$. Note that when using the spread divergence for training (see section(5.2.2)) we only need samples $p(y)$ which can be obtained from equation 16 by first sampling from $p_x(x)$ and then y from $p(y|x) = K(y - f(x))$; this does not require computing the Jacobian or inverse ¹.

5 SPREAD MAXIMUM LIKELIHOOD ESTIMATION

Minimising the forward KL divergence between the empirical data distribution $p(x)$ and a model $p(x)$ is equivalent to Maximum Likelihood Estimation (MLE) of the parameters of the model. Minimising instead the forward spread KL divergence $D_{\text{spread}}(p(x)||p(x)) = \sum_{n=1}^N \log p(y_n) + \text{const.}$, where y_n are sampled i.i.d from $p(y) = \int p(y|x)p(x)$, results in a new type of estimation, namely spread MLE. In what follows, we will discuss the statistical properties of spread MLE and demonstrate how it enables the training of models where maximum likelihood is not suited.

5.1 STATISTICAL PROPERTIES

Maximum likelihood is a cherished criterion because it exhibits many favourable statistical properties, mainly consistency (convergence to the true parameters in the large data limit) and asymptotic efficiency (achieves the Cramér-Rao Lower Bound, which is a lower bound on the variance of any unbiased estimators) - see Casella & Berger (2002) for an introduction. A key desideratum for spread MLE is to analyse how these properties are affected. In appendix(E) we demonstrate that spread MLE (for a certain family of spread noise) needs weaker sufficient conditions than MLE for both consistency and asymptotic efficiency. Furthermore, a sufficient condition for the existence of MLE is that the likelihood function is continuous over a compact parameter space. We provide an example in appendix(E.1) where this compactness requirement is violated, but spread MLE is still well defined.

5.2 APPLICATIONS

As an application to show the effectiveness of spread MLE, we use it to train implicit models

$$p(x) = \int_Z p(x - g(z)) p(z) dz \tag{18}$$

where g are the parameters of the encoder. We show that, despite the likelihood not being defined (see also section(K) for a simple linear model example), we can nevertheless successfully train such models using modified EM/variational algorithms (Barber (2012)).

5.2.1 TRAINING IMPLICIT LINEAR MODELS: DETERMINISTIC ICA

ICA (Independent Components Analysis) corresponds to the model $p(x) = \int p(x|z) \prod_{i=1}^Q p(z_i)$, where the independent components follow a non-Gaussian distribution. For Gaussian noise ICA an observation x is assumed to be generated by the process $p(x|z) = \prod_{j=1}^Q N(x_j | g_j(z))$; ² where $g_j(z)$ mixes the independent latent processes in linear ICA, $g_j(z) = a_j^T z$ where a_j is the j^{th} column on the mixing matrix A . For small observation noise ², it is well known that the maximum likelihood EM algorithm to learn A from observed data is ineffective (Bermond & Cardoso, 1999; Winther & Petersen, 2007). To see this, consider $x \in \mathbb{R}^D$ (where D and Z are the dimension of the data and latents respectively) and invertible $x = Az$. At iteration k the EM algorithm has an estimate A_k of the mixing matrix. The M-step updates A_k to

$$A_{k+1} = E [xz^T] E [zz^T]^{-1} \tag{19}$$

(a) Error versus observation noise (b) Error versus number of training points.

Figure 1: Relative error $\|A_{ij}^{est} - A_{ij}^{true}\|$ versus observation noise (a) and number of training points (b). (a) For $K=20$ observations and $Z=10$ latent variables, we generate $N=20000$ datapoints from the model $x=Az$, for independent zero mean unit variance Laplace components. The elements of A used to generate the data are uniform random. We use $S_y=1$, $S_z=1000$ samples and 2000 EM iterations to estimate A . The error is averaged over all j and 10 experiments. We also plot standard errors around the mean relative error. In blue we show the error in learning the underlying parameter using the standard EM algorithm. As expected, as the error blows up as the EM algorithm 'freezes'. In orange we plot the error for EM using spread noise; no slowing down appears as the observation noise decreases. In (b), apart from very small noise, the error for the spread EM algorithm is lower than for the standard EM algorithm. Here $Z=5$, $X=10$, $S_y=1$, $S_z=1000$, $\sigma=0.2$, with 500 EM updates used. Results are averaged over 50 runs of randomly drawn

where, for zero observation noise ($\sigma \rightarrow 0$),

$$E_{x,z} xz^T = \frac{1}{N} \sum_n x_n A_k^{-1} x_n^T = \hat{S} A_k^{-1} S A_k^{-1}; \quad E_{z,z} zz^T = A_k^{-1} \hat{S} A_k^{-1} \quad (20)$$

and $\hat{S} = \frac{1}{N} \sum_n x_n x_n^T$ is the moment matrix of the data. Thus $A_{k+1} = \hat{S} A_k^{-1} A_k^{-1} \hat{S} A_k^{-1} = A_k$ and the algorithm 'freezes'. Similarly, for low noise $\sigma \rightarrow 0$ progress critically slows down.

To deal with small noise and the limiting case of a deterministic model ($\sigma \rightarrow 0$), we consider Gaussian spread noise $p(y|x) = \int_Z N(y|x; z) p(z)$ to give

$$p(y; z) = \int_X N(y|x; z) p(x) dx = \int_Y N(y_j | g_j(z); \sigma^2 + \sigma^2 I_X) \prod_i p(z_i) \quad (21)$$

Using spread noise, the empirical distribution is replaced by the spreaded empirical distribution $p(y) = \int_X N(y|x; z) p(x) dx$. The M-step has the same form as equation 19 but with modified statistics

$$E_{y,z} yz^T = \frac{1}{N} \sum_n \int_Z N(y|x^n; z) p(z) yz^T dz dy;$$

$$E_{z,z} zz^T = \frac{1}{N} \sum_n \int_Z N(y|x^n; z) p(z) zz^T dz dy: \quad (22)$$

The E-step optimally sets

$$p(z|y) = \frac{1}{Z_q(y)} N(z | (y); \sigma^2 + \sigma^2 I_X) \prod_i p(z_i); \quad Z_q(y) = \int_Y N(z | (y); \sigma^2 + \sigma^2 I_X) \prod_i p(z_i) dz \quad (23)$$

where $Z_q(y)$ is a normaliser and

$$= (\sigma^2 + \sigma^2 I_X)^{-1} A^T A^{-1}; \quad (y) = A^T A^{-1} A y: \quad (24)$$

Since the posterior $p(z|y)$ peaks around (y) , we rewrite equation 22 as

$$E_{y,z} yz^T = \frac{1}{N} \sum_n \int_Z N(y|x^n; z) N(z | (y); \sigma^2 + \sigma^2 I_X) \frac{\prod_i p(z_i)}{Z_q(y)} yz^T dz dy$$

(a) Fixed Laplace

(b) Fixed Gaussian

(c) Learned Gaussian

(d) Covariance

Figure 2: Samples from a generative model (deterministic output) trained using VAE with (a) fixed Laplace covariance, (b) fixed Gaussian covariance and (c) learned Gaussian covariance. We first train with one epoch a standard VAE as initialization to all models, and keep latent code $(z|0; I_Z)$ fixed when sampling from these models, so we can more easily compare the sample quality. Figure (d) visualizes the absolute mean of the leading 20 eigenvectors of the learned covariance.

and similarly for $E_{z \sim p(z|y)} zz^T$. Writing the expectations with respect to $p(z|y)$ allows for a simple but effective importance sampling approximation focused on regions of high probability. We implement this update by drawing S_z samples from $N(y|x_n; \Sigma_X)$ and, for each sample, we draw S_z samples from $N(z|y; \Sigma_Z)$. This scheme has the advantage over more standard variational approaches, see for example Winther & Petersen (2007), in that we obtain a consistent estimator of the M-step update for Σ_Z . We show results for a toy experiment in figure(1), learning the underlying mixing matrix in a deterministic non-square setting. Note that standard algorithms such as FastICA (Hyvärinen, 1999) fail in this setting. The spread noise is set to $\max(0.001; 2.5 \sqrt{\text{mean}(AA^T)})$. This modified EM algorithm thus learns a good approximation of the underlying A , with no critical slowing down.

5.2.2 TRAINING IMPLICIT NON-LINEAR MODELS: -VAE

A standard way to train a deep generative model $p(x) = \int p(x|z)p(z)dz$ is to use maximum likelihood (minimizing $D(p(x)||p(x))$). The likelihood equation 18 is in general intractable and it is common to use the variational lower bound (see (Kingma & Welling, 2013)). However, for a deterministic observation model $p(x|z) = \delta(x - g(z))$ and $Z < X$, this generative model describes only a low dimensional manifold in the data space and the divergence $D(p(x)||p(x))$ is not well defined. Additionally the above bound is not well defined (due to a delta function) and the variational EM approach fails, as in the deterministic ICA setting. To address this, we instead minimize the spread divergence $D(p(y)||p(y))$. For Gaussian noise with fixed diagonal noise $p(y|x) = N(y|x; \Sigma_X)$, we can write $p(y) = \int N(y|x; \Sigma_X) p(x) dx$ and

$$p(y) = \int p(y|x)p(x)dx = \int N(y|g(z); \Sigma_X) p(z)dz = \int p(yz)p(z)dz: \quad (25)$$

We then minimize the divergence

$$KL(p(y)||p(y)) = \int p(y) \log p(y) dy + \text{const}: \quad (26)$$

Typically, the integral over y is intractable, in which case we resort to a sampling estimation. Neglecting constants, the divergence estimator is $\frac{1}{N} \sum_{n=1}^N \sum_{s=1}^S \log p(y_s^n)$, where y_s^n is a spread noise sample from $p(y_n|x_n)$; for example $y_s^n \sim N(y_n^n|x_n; \Sigma_X)$. For non-linear g , the distribution $p(y)$ is usually intractable and we therefore use the variational lower bound

$$\log p(y) \approx \int q(z|y) (\log q(z|y) + \log (p(yz)p(z))) dz: \quad (27)$$

The approach is a straightforward extension of the standard variational autoencoder and in appendix(G) we also derive a lower variance objective and detail how to learn the spread noise (also see appendix(F)). We dub this model and associated spread divergence training 'VAE'.

(a) Fixed spread noise (b) Learned spread noise

Figure 3: Samples from a generative model with deterministic output trained using VAE with (a) fixed and (b) learned spread with injective function. We use a similar sampling strategy as in the MNIST experiment to facilitate sample comparison between the different models – see section (I).

Encoder-Decoder Models	FID	GAN Models	FID
VAE	63.0	WGAN GP	30.0
-VAE with fixed spread	52.7	BEGAN	38.9
-VAE with learned spread	46.5	WGAN	41.3
		DRAGAN	42.3
WAE-MMD	55.0	LSGAN	53.9
WAE-GAN	42.0	NS GAN	55.0
		MM GAN	65.6

Table 1 CelebA FID Scores. The VAE results are the average over 5 independent measurements. The scores of GAN-based models are based on a large-scale hyperparameter search and take the best FID obtained Lucic et al. (2018). The results of VAE and WAE-based model are from Tolstikhin et al. (2017).

MNIST Experiment: We trained a -VAE on MNIST (LeCun et al. (2010)) with (i) fixed Laplace spread noise, equation 15, (ii) fixed Gaussian spread noise, equation 14 and (iii) Gaussian noise with learned covariance, section (4.1) with rank = 20; see appendix (H) for details. Figures 2(a,b,c) show samples from $p(x)$ for these models; MNIST is sufficiently easy that it is hard to distinguish between the quality of the fixed and learned noise samples. However, qualitatively, the sharpness of the Laplace spread noise trained model is higher than for the Gaussian noise and motivates that the spread noise can affect the quality of the learned model. We speculate that Laplace noise improves image sharpness since the noise focuses attention on discriminating between points close to the data manifold (since the Laplace distribution is leptokurtic and has a higher probability of generating points close to the data manifold than the Gaussian distribution). Figure 2(d) visualizes the Gaussian learned covariance and shows that the learned noise is largely orthogonal to the data manifold.

CelebA Experiment: We trained a -VAE on the CelebA dataset (Liu et al., 2015) with (i) fixed and (ii) learned spread with injective function, see appendix (I). We compared to results from a standard VAE with fixed Gaussian noise $p(x|z) = N(x|g(z); 0.5I_x)$ Tolstikhin et al. (2017). For (i) the fixed spread divergence uses Gaussian noise $p(x|z) = N(x|g(z); 0.25I_x)$. For (ii) we use Gaussian noise with learned injective function ResNet $p(x|z) = I(z) + g(z)$; see appendix (I) for more details. Figure 3 shows samples from VAE trained using Gaussian spread divergence with both fixed and learned spread noise (with $g(z)$ initialised to the fixed-noise setting). It is notable how the ‘sharpness’ of the image samples substantially increases when learning the spread noise. Table 1 shows FID (Heusel et al. (2017)) score comparisons between different algorithms. The -VAE significantly improves on the standard VAE result. VAE with injective function learning also improves on the fixed-noise -VAE. Indeed the injective-VAE results are comparable to popular GAN and WAE models (Gulrajani et al. (2017); Berthelot et al. (2017); Arjovsky et al. (2017); Kodali et al. (2017); Mao et al. (2017); Fedus et al. (2017); Tolstikhin et al. (2017)). Whilst the VAE results are not fully state-of-the-art, we believe it is the first time that implicit models have been trained using a principled maximum likelihood based approach. Our expectation is that by increasing the complexity of the generative model and injective function, or using different noise such as Laplace distribution, the results will become competitive with state-of-the-art GAN models.

³FID scores were computed using github.com/bioinf-jku/TTUR

based on 10000 samples.

6 RELATED WORK

MMD versus spread f -divergence: In spite of the conditions required for defining the spread divergence being closely related to the kernel requirement of MMD (Gretton et al., 2012), we also show that MMD and spread Total Variation distance can be written as different norms (L_2 , L_1 respectively) of a common objective (see appendix(C)).

Instance noise: The instance noise trick to stabilize GAN training Roth et al. (2017); Sønderby et al. (2016) is a special case of spread divergence using fixed Gaussian noise. Whilst other similar tricks (for example Furnston & Barber (2009)) have been proposed previously, we believe that it is important to state the general utility of the spread noise approach.

-VAE versus WAE: The Wasserstein auto-encoder Tolstikhin et al. (2017) is another implicit generative model that uses an encoder-decoder architecture. The difference -VAE is based on KL divergence which is corresponding to MLE but WAE uses the Wasserstein distance.

-VAE versus denoising VAE The Denoising VAE Im et al. (2017) uses a VAE with noise added to the data only. In contrast, the VAE adds noise to both the data and model. Since the denoising VAE model only adds noise to the model, it cannot recover the true data distribution.

MMD GAN with kernel learning : The idea of learning a kernel to increase discrimination is also used in MMD GAN (Li et al. (2017)). Similar to ours, the kernel in MMD GAN is constructed by $k = k_f \circ f$ where k is a fixed kernel and f is a neural network. To ensure $k_f(p; q) = 0$, $p = q$, this requires f to be injective (Gretton et al. (2012)). However, in the MMD GAN framework, $f(x)$ usually maps x to a lower dimension. This is crucial for MMD because the amount of data required to produce a reliable estimator grows with the data dimension (Ramdas et al. (2015)) and the computation cost of MMD scales quadratically with the amount of data. Whilst using a lower-dimensional mapping makes MMD more practical it also makes it difficult to construct an injective function f . For this reason, heuristics such as the auto-encoder regularizer (Li et al. (2017)) are considered. In contrast, for the VAE, the computational cost of estimating the divergence is linear in the number of datapoints. For this reason there is no need for f to be a lower-dimensional mapping; guaranteeing that f is injective is therefore relatively straightforward for the VAE.

Flow-based generative models Invertible flow-based functions (Rezende & Mohamed (2015)) have been used to boost the representation power of generative models. Note our use of injective functions is quite distinct from the use of flow-based functions to boost generative model capacity. In our case, the injective function does not change the model – it only changes the divergence. For this reason, the spread divergence doesn't require the log determinant of the Jacobian (which is required in Rezende & Mohamed (2015); Behrmann et al. (2018)) meaning that more general invertible functions can be used to boost the discriminatory power of a spread divergence.

7 SUMMARY

We described how to define a divergence even when two distributions do not have the same support. Previous approaches (Furnston & Barber, 2009; Sønderby et al., 2016) can be seen as special cases. We showed that defining divergences this way enables us to train deterministic generative models using standard likelihood based approaches. In principle, we can learn the underlying true data generating process by the use of any valid spread divergence. In practice, however, the quality of the learned model can depend strongly on the choice of spread noise. We therefore investigated learning spread noise to maximally discriminate two distributions. We found the resulting training approach stable and that it can significantly improve the image generation results. Whilst state-of-the-art image generation is not the focus of this work, we obtained promising results. We also discussed the conditions under which spread MLE is consistent and asymptotically efficient, some of which are weaker than the equivalent MLE conditions. Perhaps the most appealing aspect of the spread noise is that it enables one to re-use standard machine learning approaches in statistics such as maximum likelihood to train models that would be otherwise unsuited to standard statistical training approaches.

⁴Total Variation distance between $p(x)$ and $q(x)$ is defined as $TV(p(x)||q(x)) = \int_{\mathcal{R}} |p(x) - q(x)| dx$, it belongs to f -divergence family (see Liese & Vajda (2006)).

REFERENCES

- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein GAN. *arXiv preprint arXiv:1701.07875*, 2017.
- D. Barber. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, New York, NY, USA, 2012. ISBN 0521518148, 9780521518147.
- J. Behrmann, D. Duvenaud, and J. Jacobsen. Invertible residual networks. *arXiv preprint arXiv:1811.00995*, 2018.
- O. Bermond and J. Cardoso. Approximate likelihood for noisy mixtures. *Proc. ICA '99*, pp. 325–330, 1999.
- D. Berthelot, T. Schumm, and L. Metz. Began: Boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717*, 2017.
- G. Casella and R. Berger. *Statistical inference*, volume 2. Duxbury Pacific Grove, CA, 2002.
- H. Cramér. *Mathematical methods of statistics*. Princeton U. Press, Princeton, pp. 500, 1946.
- S. Dragomir. Some general divergence measures for probability distributions. *Acta Mathematica Hungarica*, 109(4):331–345, Nov 2005. ISSN 1588-2632. doi: 10.1007/s10474-005-0251-6.
- W. Fedus, M. Rosca, B. Lakshminarayanan, A. Dai, S. Mohamed, and I. Goodfellow. Many paths to equilibrium: GANs do not need to decrease a divergence at every step. *arXiv preprint arXiv:1710.08446*, 2017.
- Thomas S Ferguson. An inconsistent maximum likelihood estimator. *Journal of the American Statistical Association*, 77(380):831–834, 1982.
- T. Furnston and D. Barber. Solving deterministic policy (PO) MDPs using Expectation-Maximisation and Antifreeze. In *First international workshop on learning and data mining for robotics (LEMIR)* pp. 56–70, 2009. In conjunction with ECML/PKDD-2009.
- M. Gelbrich. On a formula for the L2 Wasserstein metric between measures on euclidean and hilbert spaces. *Mathematische Nachrichten*, 147(1):185–203, 1990.
- S. Gerchinovitz, P. Ménard, and G. Stoltz. Fano's inequality for random variables. *arXiv*, 2018. doi: arXiv:1702.05985v2.
- A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.
- I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville. Improved training of Wasserstein GANs. In *Advances in Neural Information Processing Systems*, pp. 5767–5777, 2017.
- M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in Neural Information Processing Systems*, pp. 6626–6637, 2017.
- A. Hyvärinen. Fast and robust fixed-point algorithms for independent component analysis. *IEEE Transactions on Neural Networks*, 10(3):626–634, May 1999. ISSN 1045-9227. doi: 10.1109/72.761722.
- D. Im, S. Ahn, R. Memisevic, and Y. Bengio. Denoising criterion for variational auto-encoding framework. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- D. Kingma and M. Welling. Auto-Encoding Variational Bayes. *arXiv:1312.6114 [stat.ML]*, 2013.

- N. Kodali, J. Abernethy, J. Hays, and Z. Kira. On convergence and stability of gans. preprint arXiv:1705.07215, 2017.
- Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. AT&T Labs [Online]. Available: <http://yann.lecun.com/exdb/mnist>, 2010.
- E. Lehmann. Elements of large-sample theory. Springer Science & Business Media, 2004.
- C. Li, W. Chang, Y. Cheng, Y. Yang, and B. Póczos. Mmd gan: Towards deeper understanding of moment matching network. Advances in Neural Information Processing Systems, pp.2203–2213, 2017.
- F. Liese and I. Vajda. On divergences and informations in statistics and information theory. Transactions on Information Theory, 52(10):4394–4412, 2006.
- Z. Liu, P. Luo, X. Wang, and X. Tang. Deep Learning Face Attributes in the Wild. Proceedings of International Conference on Computer Vision (ICCV), 2015.
- M. Lucic, K. Kurach, M. Michalski, S. Gelly, and O. Bousquet. Are gans created equal? a large-scale study. In Advances in neural information processing systems, pp.700–709, 2018.
- X. Mao, Q. Li, H. Xie, R. Lau, Z. Wang, and Paul S. Least squares generative adversarial networks. In Proceedings of the IEEE International Conference on Computer Vision, pp.2794–2802, 2017.
- T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. arXiv preprint arXiv:1802.05957, 2018.
- S. Mohamed and B. Lakshminarayanan. Learning in implicit generative models. arXiv preprint, 2016. doi: arXiv:1610.03483.
- B. Pearlmutter. Fast exact multiplication by the hessian. Neural computation, 6(1):147–160, 1994.
- G. Peyré, M. Cuturi, et al. Computational optimal transport. Foundations and Trends in Machine Learning 11(5-6):355–607, 2019.
- A. Ramdas, S. Reddi, B. Póczos, A. Singh, and L. Wasserman. On the decreasing power of kernel and distance based nonparametric hypothesis tests in high dimensions. Twenty-Ninth AAAI Conference on Artificial Intelligence, 2015.
- J. Rezende and S. Mohamed. Variational inference with normalizing flows. arXiv preprint arXiv:1505.05770, 2015.
- K. Roth, A. Lucchi, S. Nowozin, and T. Hofmann. Stabilizing training of generative adversarial networks through regularization. Advances in neural information processing systems, pp.2018–2028, 2017.
- N. Schraudolph. Fast curvature matrix-vector products for second-order gradient descent. Neural computation, 14(7):1723–1738, 2002.
- C. Sønderby, J. Caballero, L. Theis, W. Shi, and F. Huszár. Amortised map inference for image super-resolution. arXiv preprint arXiv:1610.04490, 2016.
- B. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf, and G. Lanckriet. Hilbert space embeddings and metrics on probability measures. Journal of Machine Learning Research, 11(Apr):1517–1561, 2010.
- B. Sriperumbudur, K. Fukumizu, and G. Lanckriet. Universality, Characteristic Kernels and RKHS Embedding of Measures. Mach. Learn. Res, 12:2389–2410, July 2011. ISSN 1532-4435.
- B. Sriperumbudur, K. Fukumizu, A. Gretton, B. Schölkopf, and G. Lanckriet. On the Empirical Estimation of Integral Probability Metrics. Electronic Journal of Statistics, 5:1550–1599, 2012.
- M. Tipping and C. Bishop. Probabilistic principal component analysis. Journal of the Royal Statistical Society, Series B, 61/3:611–622, January 1999.

- I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf. Wasserstein auto-encoders, preprint arXiv:1711.01558, 2017.
- A. Wald. Note on the consistency of the maximum likelihood estimation. *The Annals of Mathematical Statistics* 20(4):595–601, 1949.
- O. Winther and K. Petersen. Bayesian independent component analysis: Variational methods and non-negative decomposition. *Digital Signal Processing* 17(5):858 – 872, 2007. ISSN 1051-2004. Special Issue on Bayesian Source Separation.
- M. Zhang, T. Bird, R. Habib, T. Xu, and D. Barber. Variational f-divergence minimization, preprint arXiv:1907.11891, 2018.

A PROOF OF THEOREM A

Theorem 1. Consider distributions p, q , and L^1 integrable function K and Fourier Transform $Ff K g$. Let $Ff K g \neq 0$ or $Ff K g = 0$ on at most a countable set. Then

$$Ff K g Ff pg = Ff K g Ff qg \implies Ff pg = Ff qg \quad (28)$$

Proof. When $Ff K g \neq 0$, $Ff K g Ff pg = Ff K g Ff qg \implies Ff pg = Ff qg$ is trivial. We first show that the Fourier transform of an L^1 function is continuous on \mathbb{R}^d . When $Ff K g = 0$ on at most a countable set, we then show that $Ff qg$ and $Ff pg$ cannot be different at a set of countable points.

Since any distribution q is in L^1 , we can write

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}^d} |Ff qg(w + \epsilon) - Ff qg(w)| dx &= \lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}^d} |q(x) e^{2ix(w + \epsilon)} - q(x) e^{2ixw}| dx \\ &= \lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}^d} |q(x)| e^{2ixw} |e^{2ix\epsilon} - 1| dx \\ &= \lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}^d} |q(x)| e^{2ixw} |e^{2ix\epsilon} - 1| dx \text{ (Dominated Convergence Theorem)} \\ &= 0 \end{aligned}$$

So $Ff qg$ is (uniformly) continuous. The same argument applies to show $Ff pg$ is uniformly continuous.

Since $Ff K g = 0$ on at most a countable set, we assume there is a point $w_0 \in \mathbb{C}$ where $Ff qg(w_0) \neq Ff pg(w_0)$. Without loss of generality, we assume $Ff qg(w_0) - Ff pg(w_0) = \delta > 0$. For points $w_0 + h$ that are not in \mathbb{C} , $Ff K g(w_0 + h) \neq 0$ and it follows therefore that $Ff K g Ff pg = Ff K g Ff qg$ implies $Ff qg(w_0 + h) - Ff pg(w_0 + h) = 0$. By continuity of $Ff pg$ and $Ff qg$, we have $Ff qg(w_0 + h) - Ff pg(w_0 + h) \neq 0$ when $h \neq 0$, which leads to a contradiction (cannot be zero). \square

B SPREAD NOISE MAKES DISTRIBUTIONS MORE SIMILAR

The data processing inequality for divergences (see for example Gerchinovitz et al. (2018)) states that $D_f(p(y)||q(y)) \leq D_f(p(x)||q(x))$. For completeness, we provide here an elementary proof of this result. We consider the following joint distributions

$$q(y; x) = p(y|x)q(x); \quad p(y; x) = p(y|x)p(x) \quad (29)$$

whose marginals are the spread distributions

$$p(y) = \int_x p(y|x)p(x); \quad q(y) = \int_x p(y|x)q(x) \quad (30)$$

The divergence between the two joint distributions is

$$D_f(p(y; x)||q(y; x)) = \int_{x,y} q(y; x) f \left(\frac{p(y|x)p(x)}{p(y|x)q(x)} \right) = D_f(p(x)||q(x)) \quad (31)$$

The f -divergence between two marginal distributions is no larger than the divergence between the joint (see also Zhang et al. (2018)). To see this, consider

$$\begin{aligned} D_f(p(u; v)||q(u; v)) &= \int_u q(u) \int_v q(v|u) f \left(\frac{p(u; v)}{q(u; v)} \right) dy du \\ &= \int_u q(u) f \left(\frac{p(u; v)}{q(v|u)q(u)} \right) dv du \\ &= \int_u q(u) f \left(\frac{p(u)}{q(u)} \right) du = D_f(p(u)||q(u)) \end{aligned}$$

Hence,

$$D_f(p(y)||q(y)) \leq D_f(p(y; x)||q(y; x)) = D_f(p(x)||q(x)) \quad (32)$$

Intuitively, spreading two distributions increases their overlap, reducing the divergence. When q do not have the same support, $D_f(p(x)||p(x))$ can be finite or not well-defined.

C RELATION TO MMD

Spread divergence can be generally constructed from divergence, we show how to build a connection to maximum mean discrepancy (Gretton et al. (2012)) by using the spread total variation distance:

For a translation invariant kernel $k(\cdot; x)$, the MMD is (we use k to denote the MMD distance with kernel k)

$$k(p(x)||q(x)) = \int_{\mathcal{H}} k(\cdot; x)p(x)dx - \int_{\mathcal{H}} k(\cdot; x)q(x)dx$$

Suppose both k and its square root of Fourier transform \hat{k} (represents the Fourier transform) are absolutely integratable, we can also rewrite the MMD distance as following (see Sriperumbudur et al. (2010) for a derivation):

$$k(p||q) = (2\pi)^{-d/2} \int_{\mathbb{R}^d} \hat{k}(\omega) \hat{p}(\omega) \hat{q}(\omega) d\omega$$

where $\hat{k} := (\mathcal{F}^{-1} k)$ (\mathcal{F} represents the inverse Fourier transform).

We can further define $\tilde{p} = \int_{\mathbb{R}^d} \hat{k}(\omega) \hat{p}(\omega) d\omega$, where $\tilde{p} = \int_{\mathbb{R}^d} \tilde{p}(x) dx$. Therefore, \tilde{p} is a probability density function and the MMD can be written as

$$k(p||q) = (2\pi)^{-d/2} \int_{\mathbb{R}^d} \tilde{p}(x) \tilde{q}(x) dx$$

Total variation distance between $p(x)$ and $q(x)$ is defined as

$$TV(p||q) = \int_{\mathbb{R}^d} |p(x) - q(x)| dx$$

In the spread total variation distance, we can define spread as the noise distribution thus

$$\Psi V(p||q) = \int_{\mathbb{R}^d} |p(x) - q(x)| dx$$

So MMD and spread Total Variation distance can be written as different norms of a common objective.

D SPREAD DIVERGENCE BETWEEN TWO DELTA DISTRIBUTIONS

Let $p_0(x) = \delta(x - \mu_0)$, $p_1(x) = \delta(x - \mu_1)$, assume Gaussian spread noise $p(x) = N(y|x; \sigma^2)$, so

$$\begin{aligned} KL(p_0(x)||p_1(x)) &= KL(p_0(y)||p_1(y)) \\ &= \int_{\mathcal{H}} p(y|x)p_0(x) \log \frac{p(y|x)p_0(x)}{p(y|x)p_1(x)} dx \\ &= \int_{\mathcal{H}} N(y|\mu_0; \sigma^2) \log \frac{N(y|\mu_0; \sigma^2)}{N(y|\mu_1; \sigma^2)} dx \\ &= \log \frac{2\pi\sigma^2}{2\pi\sigma^2} + \frac{(\mu_0 - \mu_1)^2}{2\sigma^2} \\ &= \frac{(\mu_0 - \mu_1)^2}{2\sigma^2} \end{aligned}$$

E STATISTICAL PROPERTIES OF MAXIMUM LIKELIHOOD ESTIMATOR

E.1 EXISTENCE OF SPREAD MLE

In some situations there may not exist a Maximum Likelihood Estimator (MLE) for μ , but there is a MLE for the spread mode $p(y) = \int p(y|x)p(x)dx$. For example, suppose that $N(\cdot; \sigma^2)$ ($\sigma > 0$). So $\mu = (\mu; \sigma^2) \in \mathbb{R}^+ \times \mathbb{R}^+$. Assume we only have one data point. Then the log-likelihood function is $L(\mu; \sigma^2) = \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x - \mu)^2)$. Maximising with respect to μ , we have $\mu = x$ and the log-likelihood becomes unbounded as $\sigma \rightarrow 0$. In this sense, the MLE for $(\mu; \sigma^2)$ does not exist.

On the other hand, we can check whether the MLE $\hat{\theta}_N$ exists. We assume Gaussian spread noise with fixed variance σ^2 . Since we only have one data point, the spread data distribution becomes $p(y|x) = \mathcal{N}(y|x; \sigma^2)$, and the model is $p(y) = \mathcal{N}(y; \mu^2 + \sigma^2)$. We can sample N points from the spread model, so the spread log likelihood function is (neglecting constants) $L(y_1, \dots, y_N; \mu) = \frac{N}{2} \log(\mu^2 + \sigma^2) - \frac{1}{2(\mu^2 + \sigma^2)} \sum_{i=1}^N (y_i - \mu)^2$. The MLE solution for μ is $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N y_i$; the MLE solution for σ^2 is $\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\mu})^2 + \sigma^2$, which has bounded spread likelihood value. Note that in the limit of a large number of spread samples $N \rightarrow \infty$, the MLE $\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\mu})^2 + \sigma^2$ tends to 0. Throughout, however, the (scaled) log likelihood remains bounded.

E.2 CONSISTENCY

Consistency of an estimator is an important property that guarantees the validity of the resulting estimate at convergence as the number of data points tends to infinity. In what follows, we refer to sufficient conditions for a consistent MLE estimator, before addressing the question of whether using spread MLE is also consistent and under what conditions.

E.2.1 CONSISTENCY FOR MLE

Sufficient conditions for the MLE being consistent and converging to the global maximum are given in Wald (1949). However, they are usually difficult to check even for some standard distributions. The sufficient conditions for MLE being consistent and converging to local maxima are given in Cramér (1946) and are more straightforward to check:

- C1. (Identifiable): $p(x; \theta_1) = p(x; \theta_2) \implies \theta_1 = \theta_2$.
- C2. The parameter space is an open interval (θ_1, θ_2) , $0 < \theta_1 < \theta_2 < \infty$.
- C3. $p(x; \theta)$ is continuous in θ and differentiable with respect to θ for all x .
- C4. The set $A = \{x : p(x; \theta) > 0\}$ is independent of θ .

Let X_1, X_2, \dots be iid with density $p(x; \theta_0)$ ($\theta_0 \in (\theta_1, \theta_2)$) satisfying conditions C1–C4, then there exists a sequence $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$ of local maxima of the likelihood function $l(\theta) = \prod_{i=1}^n p(x_i; \theta)$ which is consistent:

$$\hat{\theta}_n \xrightarrow{P} \theta_0 \text{ for all } n \geq 2$$

The proof can be found in Lehmann (2004) or Cramér (1946).

E.2.2 CONSISTENCY OF SPREAD MLE

We provide the necessary conditions for Spread MLE being consistent.

- C1. (Identifiable): $p(y; \theta)$ is identifiable. From section (3) it follows immediately that $p(y; \theta_1) = p(y; \theta_2) \implies p(x; \theta_1) = p(x; \theta_2) \implies \theta_1 = \theta_2$, where the final implication follows from the assumption that $p(x; \theta)$ is identifiable. Hence if $p(x; \theta)$ is identifiable, so is $p(y; \theta)$.
- C2. The parameter space is an open interval (θ_1, θ_2) , $0 < \theta_1 < \theta_2 < \infty$. This condition is unchanged for $p(y; \theta)$.
- C3. On $p(y; \theta)$, we require the same condition $p(x; \theta)$ as in MLE; $p(y; \theta)$ is continuous in θ and differentiable with respect to θ for all y .
- C4. For spread noise $p(y|x)$ who has full support on \mathbb{R}^d (for example Gaussian noise), $p(y; \theta)$ is greater than zero everywhere and hence the original condition C4 is automatically guaranteed.

The conditions that guarantee consistency for spread MLE are weaker for the spread model than for the standard model $p(x; \theta)$, since C4 is automatically satisfied. Ferguson (1982) gives an example for which MLE exists but is not consistent by violating condition C4, whereas spread MLE can be used to obtain a consistent estimator.

E.3 ASYMPTOTIC EFFICIENCY

A key desirable property of any estimator is that it is efficient. The Cramer-Rao bound places a lower bound on the variance of any unbiased estimator and an efficient estimator much reach this minimal value in the limit of a large amount of data. Under certain conditions (see below) the Maximum Likelihood Estimator attains this minimal variance value meaning that there is no better estimator possible than maximum likelihood (in the limit of a large amount of data). This is one of the reasons that the maximum likelihood is a cherished criterion.

E.3.1 ASYMPTOTIC EFFICIENCY FOR MLE

Building upon conditions C1-C4, additional conditions $p(x; \theta)$ are required to show MLE is asymptotically efficient:

- C5. For all x in its support, the density $p(x; \theta)$ is three times differentiable with respect to θ and the third derivative is continuous.
- C6. The derivatives of the integral $\int_{\mathcal{R}} p(x; \theta) dx$ respect to θ can be obtained by differentiating under the integral sign, that is: $\frac{\partial}{\partial \theta} \int_{\mathcal{R}} p(x; \theta) dx = \int_{\mathcal{R}} \frac{\partial}{\partial \theta} p(x; \theta) dx$.
- C7. There exists a positive number $c(\theta_0)$ and a function $M_{\theta_0}(x)$ such that

$$\left| \frac{\partial^j}{\partial \theta^j} \log p(x; \theta) \right| \leq M_{\theta_0}(x) \quad \text{for all } x \in \mathcal{A}; \quad |j| < c(\theta_0)$$

where \mathcal{A} is the support set of x and $E_{\theta_0}[M_{\theta_0}(x)] < 1$.

Let X_1, \dots, X_n be i.i.d with density $p(x; \theta)$ and satisfy conditions C1-C7, then any consistent sequence $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$ of roots of the likelihood equation satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0; F^{-1}(\theta_0));$$

where $F^{-1}(\theta_0)$ is the inverse of Fisher information matrix (also called Cramér-Rao Lower Bound, which is a lower bound on variance of any unbiased estimators). The conditions and proof can be found in Lehmann (2004).

E.3.2 ASYMPTOTIC EFFICIENCY FOR MLE

As with MLE above, we require further conditions $p(y; \theta)$ for ensuring spread MLE is asymptotically efficient:

- C5. On $p(y; \theta)$, we require the same condition as applied $p(x; \theta)$ in the MLE case; for all y in its support, the density $p(y; \theta)$ is three times differentiable with respect to θ and the third derivative is continuous.
- C6. For spread noise $p(y|x)$, which has full support on \mathcal{R}^d (for example Gaussian noise), the support of y is independent of x . Leibniz's rule⁵ allows us to differentiate under the integral: $\frac{\partial}{\partial \theta} \int_{\mathcal{R}} p(y|x; \theta) dy = \int_{\mathcal{R}} \frac{\partial}{\partial \theta} p(y|x; \theta) dy$, so this condition is automatically satisfied.
- C7. On $p(y; \theta)$, we require the same condition as applied $p(x; \theta)$ in the MLE case; There exist positive number $c(\theta_0)$ and a function $M_{\theta_0}(y)$ such that

$$\left| \frac{\partial^j}{\partial \theta^j} \log p(y; \theta) \right| \leq M_{\theta_0}(y) \quad \text{for all } y \in \mathcal{A}; \quad |j| < c(\theta_0)$$

where \mathcal{A} is the support set of y and $E_{\theta_0}[M_{\theta_0}(y)] < 1$.

Thus the conditions that guarantee asymptotically efficient for the spread model are weaker than for the standard model, since C4 and C6 are automatically satisfied.

⁵Leibniz's rule tells us: $\frac{d}{d\theta} \int_a^{b(\theta)} p(x; \theta) dx = \int_a^{b(\theta)} \frac{\partial}{\partial \theta} p(x; \theta) dx + p(b(\theta); \theta) \frac{d}{d\theta} b(\theta) - p(a(\theta); \theta) \frac{d}{d\theta} a(\theta)$, so if $a(\theta)$ and $b(\theta)$ are independent of θ , then $\frac{d}{d\theta} \int_a^b p(x; \theta) dx = \int_a^b \frac{\partial}{\partial \theta} p(x; \theta) dx$.

F PERTURBATION APPROXIMATION OF GAUSSIAN SPREAD

Herein, in the fixed noise setting, we derive a perturbation based approximation to the spread noise, which in principle can lead to a lower variance estimator. We can write a function with perturbed input and integrated over noise as

$$E_{p(\cdot)} [f(x + \cdot)]; \quad (33)$$

where $p(\cdot) = N(\cdot | 0; \Sigma)$. Taylor expanding around $\cdot = 0$, we have

$$E_{p(\cdot)} [f(x + \cdot)] = E_{p(\cdot)} [f(x) + \text{Tr} \nabla f(x) + \frac{1}{2} \text{Tr} \nabla^2 f(x) + O(\cdot^3)] \quad (34)$$

$$f(x) + \frac{1}{2} E_{p(\cdot)} [\text{Tr} H] \quad (35)$$

$$= f(x) + \frac{1}{2} \text{Tr} E_{p(\cdot)} [H] \quad (36)$$

$$= f(x) + \frac{1}{2} \text{Tr}(H) \quad (37)$$

Where H is the Hessian matrix $H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$. When x 's dimension is small, we can use equation 37 to explicitly calculate the trace. When the dimension is large, we can form a Monte Carlo estimation of equation 35. To do this we first sample $p(\cdot)$ and then calculate the Hessian-vector product $H \cdot$. This can be efficiently calculated by AutoDiff backward mode (Schraudolph (2002); Pearlmutter (1994)), without storing the Hessian matrix in memory.

For example, in -VAE with fixed Gaussian noise, according to the bound equation 56 (ignoring the constant):

$$\int_{\mathcal{Z}} N(y|x; \Sigma_x) \log p(y) = E_{N(x|0; \Sigma_x)} [H(\cdot(x + \cdot)))] \quad (38)$$

$$+ E_{N(x|0; \Sigma_x)} E_{N(0;1)} [\log p(x = g(\cdot(x + \cdot)) + C(x))] + \log p(z = \cdot(x + \cdot) + C(x))] \quad (39)$$

$$\underbrace{H(\cdot(x)) + E_{N(0;1)} [\log p(x = g(\cdot(x) + C(x))) + \log p(z = \cdot(x) + C(x))]}_{f(x)} + \frac{1}{2} \text{Tr}(H) \quad (40)$$

Where H is the Hessian matrix $H_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$ and $C(x)$ is the Cholesky decomposition of Σ_x .

G SPREAD DIVERGENCE FOR DETERMINISTIC DEEP GENERATIVE MODELS

Instead of minimising the likelihood, we train an implicit generative model by minimising the spread divergence

$$\min \text{KL}(p(y) || p(y)) \quad (41)$$

For Gaussian noise with fixed diagonal noise $p(y|x) = N(y|x; \Sigma_x)$, we can write

$$p(y) = \frac{1}{N} \prod_{n=1}^N N(y|x_n; \Sigma_x) \quad (42)$$

$$\text{and } \int_{\mathcal{Z}} p(y|x) p(x) dx = \int_{\mathcal{Z}} N(y|g(z); \Sigma_x) p(z) dz = \int p(y|z) p(z) dz \quad (43)$$

For the spread divergence with learned covariance Gaussian noise which is discussed in section(4.1), we can write

$$p(y|x) = N(y|x; \Sigma_x); \quad p(y) = \frac{1}{N} \prod_{n=1}^N N(y|x_n; \Sigma_x) \quad (44)$$

and spread divergence with learned injective function as discussed in section(4.2)

$$p(y|x) = N(y|f(x); \Sigma_x); \quad p(y) = \frac{1}{N} \prod_{n=1}^N N(y|f(x); \Sigma_x) \quad (45)$$

According to our general theory,

$$\min \text{KL}(p(y)||p(y)) = 0, \quad p(x) = p(x) \quad (46)$$

Here

$$\text{KL}(p(y)||p(y)) = \frac{1}{N} \sum_{n=1}^N \int p(y) \log p(y) dy + \text{const:} \quad (47)$$

Typically, the integral over y will be intractable and we resort to an unbiased sampled estimate (though see below for Gaussian). Neglecting constants, the KL divergence estimator is

$$\frac{1}{NS} \sum_{n=1}^N \sum_{s=1}^S \log p(y_s^n) \quad (48)$$

where y_s^n is a perturbed version of y . For example $y_s^n \sim N(y_s^n|x_n; \Sigma_x)$ for fixed Gaussian noise case and other cases are similar. In most cases of interest, with non-Gaussian distributions $p(y)$ is intractable. We therefore use the variational lower bound

$$\log p(y) \geq \int q(z|y) (\log q(z|y) + \log(p(y|z)p(z))) dz \quad (49)$$

Parameterising the variational distribution as a Gaussian,

$$q(z|y) = N(z|\mu(y); \Sigma(y)) \quad (50)$$

then we can reparameterise and write

$$\log p(y) \geq H(\mu(y)) + E_{N(0;I)} [\log(p(y|z = \mu(y) + C(y))p(z = \mu(y) + C(y)))] \quad (51)$$

where $H(\mu(y))$ is the entropy of a Gaussian with covariance $\Sigma(y)$. For fixed covariance Gaussian spread noise in D dimensions, this is

$$\log p(y) \geq H(\mu(y)) + E_{N(0;I)} \left[\frac{1}{(2\pi)^{D/2}} \int (y - g(\mu(y) + C(y)))^2 + \log p(z = \mu(y) + C(y)) \right] + \text{const:} \quad (52)$$

where $C(y)$ is the Cholesky decomposition of $\Sigma(y)$.

We can integrate equation 52 over y to give the bound

$$E_{N(y|x; \Sigma_x)} \log p(y) \geq E_{N(y|x; \Sigma_x)} H(\mu(y)) + E_{N(0;I)} [\log p(z = \mu(y) + C(y))] \quad (53)$$

$$\frac{1}{(2\pi)^{D/2}} E_{N(0;I)} \int E_{N(y|x; \Sigma_x)} (y - g(\mu(y) + C(y)))^2 + \text{const:} \quad (54)$$

where

$$\begin{aligned} & E_{N(y|x; \Sigma_x)} \int (y - g(\mu(y) + C(y)))^2 \\ &= \int \frac{1}{(2\pi)^{D/2}} E_{N(x|0;I_x)} [x - g(\mu(x) + C(y))] + E_{N(x|0;I_x)} \int (x - g(\mu(x) + C(y)))^2 \end{aligned} \quad (55)$$

We notice that the second term is zero, so the final bound for the fixed Gaussian spread KL divergence is (ignoring the constant)

$$\begin{aligned} & E_{N(y|x; \Sigma_x)} \log p(y) \geq E_{N(y|x; \Sigma_x)} H(\mu(y)) + E_{N(0;I)} [\log p(z = \mu(y) + C(y))] \\ & \frac{1}{(2\pi)^{D/2}} E_{N(x|0;I_x)} \int E_{N(0;I)} (x - g(\mu(x) + C(y)))^2 \end{aligned} \quad (56)$$

By analogy, for spread KL divergence with learned variance, the bound is (ignoring the constant)

$$\mathbb{E}_{N(y|x; \Sigma)} \log p(y) - \mathbb{E}_{N(y|x; \Sigma)} H(y) + \mathbb{E}_{N(z|y; \Sigma)} [\log p(z = (y) + C(y))] \quad \text{ii}$$

$$\mathbb{E}_{N(x|0; \Sigma)} \mathbb{E}_{N(z|y; \Sigma)} (x - g(x + S_x) + C(y))^T \Sigma^{-1} (x - g(x + S_x) + C(y)) \quad (57)$$

Where S is the cholesky decomposition of Σ . For specific covariance structure introduced in section 4.1, efficient methods for sampling, matrix inverting and log determinant calculation are available, see appendix J.

For spread KL divergence with learned injective function, the bound is (ignoring the constant)

$$\mathbb{E}_{N(y|x; \Sigma)} \log p(y) - \mathbb{E}_{N(y|x; \Sigma)} H(y) + \mathbb{E}_{N(z|y; \Sigma)} [\log p(z = (y) + C(y))] \quad \text{ii}$$

$$\frac{1}{(2^2)^{D=2}} \mathbb{E}_{N(x|0; I_x)} \mathbb{E}_{N(z|y; \Sigma)} (f(x) - f(g(f(x) + x) + C(y)))^2 \quad (58)$$

The overall procedure is therefore a straightforward modification of the standard VAE method Kingma & Welling (2013) with additional learning the spread to maximize the divergence:

1. Choose a noise distribution $p(y|x)$
2. Choose a tractable family for the variational distribution, for example $p(z|y) = N(z|g(y); \Sigma(y))$ and initialise Σ .
3. We then sample y_n for each datapoint (if we're using $\beta = 1$ samples)
4. If learning the spread noise:
 - (a) Draw samples to estimate $\log p(y_n)$ according to the corresponding bound.
 - (b) Do a gradient ascent step in Σ .
5. Draw samples to estimate $\log p(y_n)$ according to the corresponding bound.
6. Do a gradient ascent step in θ .
7. Go to 3 and repeat until convergence.

H MNIST EXPERIMENT

We first scaled the MNIST data to lie in $[0, 1]$. We use Laplace spread noise $\Sigma = 0.3$ and Gaussian spread noise $\Sigma = 0.3$ for the VAE. Both encoder and decoder contains 3 feed-forward layers, each layer with 400 units and ReLU activation function. The latent dimension is 64. The variational inference network $q(z|y) = N(z|g(y); \Sigma)$ has a similar structure for the mean network

$g(y)$. For fixed spread -VAE, learning was done using the Adam Kingma & Ba (2014) optimizer with learning rate $5e^{-4}$ for 200 epochs. For VAE with learned spread (learned covariance), we additionally train the covariance for 2 epochs using Adam optimizer with learning rate $5e^{-4}$ after everytime we train the model for 10 epochs.

I CELEBA EXPERIMENT

We pre-processed CelebA images by first taking 140x140 centre crops and then resizing to 64x64. Pixel values were then rescaled to lie in $[0, 1]$. For the learned spread we use Gaussian noise with learned injective function $\text{ResNet}(\cdot) = I(\cdot) + g(\cdot)$, where $g(\cdot)$ is a one layer convolutional neural net with kernel size 3 and stride 1 . We use spectral normalization Miyato et al. (2018) to satisfy the Lipschitz constraint: we replace the weight matrix of the convolution kernel by $w_{\text{SN}}(w) := c \cdot w = (w)$ where (w) is the spectral norm of w and $c = 2 / (0; 1)$. This guarantees that f is invertible – see Behrmann et al. (2018).

The encoder and decoder are 4-layer convolutional neural net with batch norm (Ioffe & Szegedy (2015)). Both encoder and decoder used fully convolutional architectures with 5x5 convolutional

layers and used vertical and horizontal strides 2 except the last deconvolution layer we used stride 1. The injective function is a 1 layer convolutional network with 3×3 kernel and stride 1. Conv_k stands for a convolution with k filters, DeConv_k for a deconvolution with k filters, BN for the batch normalization Ioffe & Szegedy (2015), Relu for the rectified linear units, FC_k for the fully connected layer mapping \mathbb{R}^k .

```

x 2 R^{64 \times 64 \times 3} ! injective_f ( ) 2 R^{64 \times 64 \times 3}
! Conv_{128} ! BN ! Relu
! Conv_{256} ! BN ! Relu
! Conv_{512} ! BN ! Relu
! Conv_{1024} ! BN ! Relu ! FC_{100}

z 2 R^{100} ! FC_{10 \times 10 \times 1024}
! DeConv_{612} ! BN ! Relu
! DeConv_{256} ! BN ! Relu
! DeConv_{128} ! BN ! Relu ! DeConv_8 ! sigmoid ( )
! injective_f ( ) 2 R^{64 \times 64 \times 3}

```

We use batch size 600 and latent dimension $d_{\text{dim}} = 100$ in all CelebA experiments. For the VAE with fixed spread, we use the fixed Gaussian noise with 0 mean and $(0.5)^2 I$ covariance. We train the model for 500 epochs using Adam optimizer with learning rate $1e^{-4}$. The learning rate decay with scaling factor 0.9 every 100000 iterations.

For the β -VAE with fixed spread We first train a β -VAE with fixed $f(x) = x$ and fixed Gaussian noise with 0 mean and $(0.5)^2 I$ diagonal covariance for 300 epochs, the learning rate decay with scaling factor 0.9 every 100000. Then we start iterative training by doing one step inner maximisation over the spread divergence's parameter using Adam optimizer with learning rate $1e^{-5}$ and one step minimization over the model parameter's β using Adam optimizer for additional 200 epochs.

We can share the first 300 epochs between two models. When we sample from two models, we first sample from a 100 dimensional standard Gaussian distribution $N(0; I)$ and use the same latent code z to get samples from both β -VAE with fixed and learned spread, so we can easily compare the sample quality between two models.

J WOODBERRY

When evaluating the log probability of this Gaussian, we use the Woodberry identity

$$1 = (I - L^T)^{-1} (I - L^T)^{-1} L (I - L^T)^{-1} L^T (I - L^T)^{-1} L (I - L^T)^{-1} L^T (I - L^T)^{-1}$$

so we only have to invert \mathbb{R}^D matrix. A similar trick is applied to calculate the log determinant:

$$\log \det(\Sigma) = \log \det(I + L^T (I - L^T)^{-1} L) + 2D \log |1|$$

The parameter β is trained using reparameterization trick. When sampling $N(\mu; LL^T + I^{-2})$, we first sample $z \in \mathbb{R}^D$ from $N(z|0; I)$ and then sample noise $N(0; I^{-2})$, thus a sample from $N(\mu; LL^T + I^{-2})$ can be represented by $Lz + \epsilon$.

K DETERMINISTIC LINEAR LATENT MODEL

Our aim here is to show how the classical deterministic PCA algorithm can be derived through a maximum-likelihood approach, rather than the classical non-probabilistic least-squares derivation. This is remarkable since the likelihood itself is not defined for this model.

For isotropic Gaussian observation noise with variance σ^2 , the Probabilistic PCA model (Tipping & Bishop, 1999) for X -dimensional observations and d -dimensional latent is

$$\begin{aligned}
 x &= Fz + \epsilon; \quad z \sim N(0; I_z); \quad \epsilon \sim N(0; I_x); \\
 p(x) &= N(y|0; FF^T + \sigma^2 I_x)
 \end{aligned} \tag{59}$$

When $\gamma = 0$, the generative mapping from z to x is deterministic and the model $p_\theta(x)$ has support only on a subset of \mathbb{R}^X and the data likelihood is in general not defined for $Z < X$.

In the following we consider general γ , later setting γ to zero throughout the calculation. To fit the model to iid data $\{x_1, \dots, x_N\}$ using maximum likelihood, the only information required from the dataset is the data covariance $\hat{\Sigma}$. For $\gamma > 0$, the maximum likelihood solution for PPCA is $F = U_Z (\Lambda_Z - \gamma^2 I_Z)^{\frac{1}{2}} R$, where Λ_Z, U_Z are the Z largest eigenvalues, eigenvectors of $\hat{\Sigma}$; R is an arbitrary orthogonal matrix. Using spread noise $p(y|x) = \mathcal{N}(y|x, \sigma^2 I_X)$, the spreaded distribution is a Gaussian $\tilde{p}_\theta(y) = \mathcal{N}(y|0, FF^\top + (\gamma^2 + \sigma^2)I_X)$. Thus, $\tilde{p}_\theta(y)$ is of the same form as PPCA, albeit with an inflated covariance matrix. Adding Gaussian spread noise to the data also simply inflates the sample covariance to $\hat{\Sigma}' = \hat{\Sigma} + \sigma^2 I_X$.

Since the eigenvalues of $\hat{\Sigma}' \equiv \hat{\Sigma} + \sigma^2 I_X$ are simply $\Lambda' = \Lambda + \sigma^2 I_X$, with unchanged eigenvectors, the optimal deterministic ($\gamma = 0$) latent linear model has solution $F = U_Z (\Lambda'_Z - \sigma^2 I_Z)^{\frac{1}{2}} R = U_Z \Lambda_Z^{\frac{1}{2}} R$.

We have thus recovered the standard PCA solution; however, the derivation is non-standard since the likelihood of the deterministic latent linear model $\gamma = 0$ is not defined. Since classical deterministic PCA cannot normally be described in terms of a likelihood, the usual derivation of PCA is to define it as the optimal least squares reconstruction solution based on a linear projection to a lower-dimensional subspace, see for example Barber (2012). Nevertheless, using the spread divergence, we learn a sensible model and recover the true data generating process if the data were exactly generated according to the deterministic model.



(a) Laplace with fixed covariance



(b) Gaussian with fixed covariance



(c) Gaussian with learned covariance

Figure 4: Samples from a generative model (deterministic output) trained using δ -VAE with (a) Laplace noise with fixed covariance, (a) Gaussian noise with fixed covariance and (c) Gaussian noise with learned covariance.

