

# KERNEL OF CYCLEGAN AS A PRINCIPAL HOMOGENEOUS SPACE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Unpaired image-to-image translation has attracted significant interest due to the invention of CycleGAN, a method which utilizes a combination of adversarial and cycle consistency losses to avoid the need for paired data. It is known that the CycleGAN problem might admit multiple solutions, and our goal in this paper is to analyze the space of exact solutions and to give perturbation bounds for approximate solutions. We show theoretically that the exact solution space is invariant with respect to automorphisms of the underlying probability spaces, and, furthermore, that the group of automorphisms acts freely and transitively on the space of exact solutions. We examine the case of zero ‘pure’ CycleGAN loss first in its generality, and, subsequently, expand our analysis to approximate solutions for ‘extended’ CycleGAN loss where identity loss term is included. In order to demonstrate that these results are applicable, we show that under mild conditions nontrivial smooth automorphisms exist. Furthermore, we provide empirical evidence that neural networks can learn these automorphisms with unexpected and unwanted results. We conclude that finding optimal solutions to the CycleGAN loss does not necessarily lead to the envisioned result in image-to-image translation tasks and that underlying hidden symmetries can render the result useless.

## 1 INTRODUCTION

Machine learning methods for image-to-image translation are widely studied and have applications in several fields. In medical imaging, the CycleGAN has found an important application for translating one modality to another, for instance in MR to CT translation (Han, 2017; Sjölund et al., 2015; Wolterink et al., 2017). Classically, these methods are trained in a supervised setting making their applications limited due to the a lack of good paired data. Similar issues appear in e.g. transferring the style of one artist to another (Gatys et al., 2015) or adding snow to sunny California streets (Liu et al., 2017). Unpaired image-to-image translation models such as CycleGAN (Zhu et al., 2017) promise to solve this issue by only enforcing a relationship on a distribution level, thus removing the need for paired data. However, given their widespread use, it is paramount to gain more understanding of their dynamics, to prevent unexpected things from happening, e.g., (Cohen et al., 2018). As a step in that direction, we explore the solution space of the CycleGAN in the subsequent sections of this paper.

The general task of unpaired domain translation can be informally described as follows: given two probability spaces  $X$  and  $Y$  which represent our domains, we seek to learn a mapping  $G : X \rightarrow Y$  such that a sample  $\mathbf{x} \in X$  is mapped to a sample  $G(\mathbf{x}) \in Y$  where

$$G(\mathbf{x}) \in Y \text{ is the best representative of } \mathbf{x} \text{ in } Y. \quad (1)$$

The mapping  $G$  is typically approximated by a neural network  $G_\theta$  parametrized by  $\theta$ . Without paired data, directly solving this is impossible but on a distribution level it is easily seen if  $G$  solves eq. (1) then the distribution of  $G(\mathbf{x})$  as  $\mathbf{x}$  is sampled from  $X$  is equal to that of  $Y$ . Mathematically, if  $X = (X, \mathcal{X}, \mu)$  and  $Y = (Y, \mathcal{Y}, \nu)$  are probability spaces with probability measures  $\mu$  and  $\nu$  respectively, this can be written as

$$\nu(A) = \mu(\{\mathbf{x} : G(\mathbf{x}) \in A\}) = \mu(G^{-1}(A)) \stackrel{\text{def}}{=} (G_*\mu)(A) \quad \text{for all } A \in \mathcal{Y}, \quad (2)$$

Or in words, the probability measure  $\nu$  equals the push-forward measure  $G_*\mu$ . By Jensen’s equality we can relate this to the fixed f-divergence  $D_f$ :

$$G_*\mu = \nu \text{ if and only if } D_f(G_*\mu || \nu) = 0. \quad (3)$$

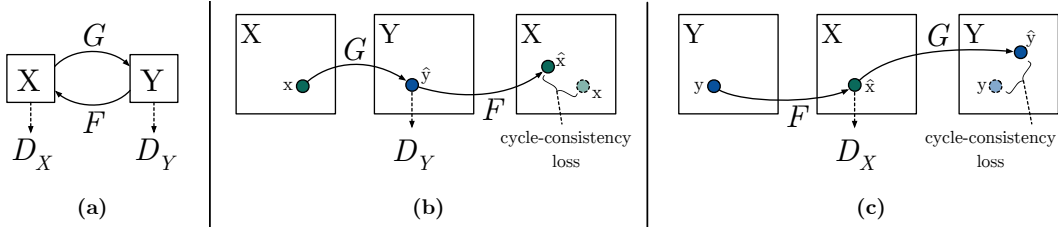


Figure 1: CycleGAN model.

While adversarial optimization techniques such as GANs can in principle solve problem eq. (3), they remain under-constrained thus not giving a reasonable solution to the original problem eq. (1).

The idea behind the *cycle consistency condition* from (Zhu et al., 2017) is to enforce additional constraints by introducing another function  $F : Y \rightarrow X$ , which is also approximated by a neural network and tries to solve the inverse task: for each  $\mathbf{y} \in Y$  find  $F(\mathbf{y}) \in X$  that would be the best translation of  $\mathbf{y}$  to  $X$ . Similar to the reasoning above, this condition would imply that

$$\mu = F_*\nu \quad \text{and} \quad D_f(F_*\nu \parallel \mu) = 0. \quad (4)$$

The goal is to enforce that  $F(G(\mathbf{x})) \approx \mathbf{x}$  for all  $\mathbf{x} \in X$  and, similarly, that  $G(F(\mathbf{y})) \approx \mathbf{y}$  for all  $\mathbf{y} \in Y$ , i.e. to minimize the following *cycle consistency loss*

$$\mathcal{L}_{\text{cyc}}(G, F) := \mathbb{E}_{\mathbf{x} \sim X} \|(F \circ G)(\mathbf{x}) - \mathbf{x}\| + \mathbb{E}_{\mathbf{y} \sim Y} \|(G \circ F)(\mathbf{y}) - \mathbf{y}\|, \quad (5)$$

where typically the  $L^1$  norm is chosen, but in principle any norm can be chosen. The authors (Zhu et al., 2017) also suggested that an adversarial loss could in principle have been used here as well, but they did not note any performance improvement.

Combining these losses, we arrive at the *CycleGAN loss* defined as

$$\mathcal{L}(G, F) := D_f(F_*\mu \parallel \nu) + D_f(G_*\nu \parallel \mu) + \alpha_{\text{cyc}} \cdot \mathcal{L}_{\text{cyc}}(G, F),$$

where the factor  $\alpha_{\text{cyc}} > 0$  determines the weight of the cycle consistency term. We illustrate the CycleGAN model in fig. 1.

Precautions with generative models have been addressed before, for example, unpaired image to image translation can hallucinate features in medical images (Cohen et al., 2018). Furthermore, it was already noted in (Zhu et al., 2017) that the CycleGAN might admit multiple solutions and that the issue of tint shift in image-to-image translation arises due to the fact that for a fixed input image  $\mathbf{x} \in X$  multiple images  $\mathbf{y}_1, \dots, \mathbf{y}_n \in Y$  with different tints might be equally plausible. Adding identity loss term was suggested in (Zhu et al., 2017) to alleviate the tint shift issue, i.e., the *extended CycleGAN loss* is defined as

$$\mathcal{L}_{\text{ext}}(G, F) := \mathcal{L}(G, F) + \alpha_{\text{id}} \cdot (\mathbb{E}_{\mathbf{y} \sim Y} \|F(\mathbf{y}) - \mathbf{y}\| + \mathbb{E}_{\mathbf{x} \sim X} \|G(\mathbf{x}) - \mathbf{x}\|),$$

where the factor  $\alpha_{\text{id}} \geq 0$  determines the weight of the identity loss term. In general, to properly define the identity loss one needs to represent both  $X$  and  $Y$  as being supported on the same manifold, which is limiting if the distributions are substantially different.

The goal of this work is to study the kernel, or null space, of the CycleGAN loss, which is the set of solutions  $(G, F)$  which have zero ‘pure’ CycleGAN loss, and to give a perturbation bounds for approximate solutions for the case of extended CycleGAN loss. We do the theoretical analysis in section 2. We show that under certain assumptions on the probability spaces  $X, Y$  the kernel has symmetries which allow for multiple possible solutions in Proposition 2.1. Furthermore, we show in Proposition 2.2 and the following remarks that the kernel admits a natural structure of a principal homogeneous space with the automorphism group  $\text{Aut}(X)$  of  $X$  acting on the set of solutions freely and transitively. Next, we expand our analysis to the case of approximate solutions for the *extended CycleGAN loss* by proving perturbation bounds in Proposition 2.3 and Corollary 2.1. We discuss the existence problem of automorphism in Proposition 2.4 and Proposition 2.6. We proceed in section 3 by showing that unexpected symmetries can be *learned* by a CycleGAN. In particular, when translating the same domain to itself CycleGAN can learn a nontrivial automorphism of the domain.

In appendix A, we briefly explain the measure-theoretic language we use heavily in the paper for those readers who are more used to working with distributions, and also remind the reader of some basic notions from differential geometry which we use as well.

## 2 THEORY

### 2.1 CYCLEGAN KERNEL AS A PRINCIPAL HOMOGENEOUS SPACE

The notions of isomorphism of probability spaces and of probability space automorphisms are central to this paper. Intuitively speaking, an isomorphism  $f : X \rightarrow Y$  of probability spaces  $X$  and  $Y$  is a bijection between  $X$  and  $Y$  such that the probability of an event  $A \subset Y$  equals the probability of event  $\{x : f(x) \in A\} \subset X$ . An isomorphism of a probability space to itself is called a probability space automorphism. For example, if our probability space consists of samples from  $n$ -dimensional spherical Gaussian distribution, then any rotation in  $SO(\mathbb{R}^n)$  is a probability space automorphism. For a precise definition we refer the reader to appendix A.

Firstly, we prove that if at least one of the probability spaces  $X, Y$  admits a nontrivial probability automorphism, then any exact solution in the kernel of CycleGAN can be altered giving a *different* solution.

**Proposition 2.1** (Invariance of the kernel). *Let  $X = (X, \mathcal{X}, \mu), Y = (Y, \mathcal{Y}, \nu)$  be probability spaces and  $\varphi : X \rightarrow X$  be a probability space automorphism. Let  $G : X \rightarrow Y$  and  $F : Y \rightarrow X$  be measurable maps satisfying*

$$\mathcal{L}(G, F) = 0. \quad (6)$$

*Then  $F, G$  are probability space isomorphisms and*

$$\mathcal{L}(G \circ \varphi, \varphi^{-1} \circ F) = 0. \quad (7)$$

*If, furthermore,  $\varphi \neq \text{id}_X$ ,<sup>1</sup> then*

$$G \circ \varphi \neq G \quad \text{and} \quad \varphi^{-1} \circ F \neq F. \quad (8)$$

*Proof.* Since  $\varphi$  is a probability space automorphism, its inverse  $\varphi^{-1}$  is an automorphism as well. In particular, it is measure-preserving since

$$\mu(\varphi(A)) = \mu(\varphi^{-1}(\varphi(A))) = \mu(A) \quad \text{for all } A \in \mathcal{X}.$$

We note that by eq. (2) and the positivity of the norms eq. (6) implies that

$$G_*\mu = \nu, \quad F_*\nu = \mu \quad (9)$$

and

$$G \circ F = \text{id}_Y \text{ a.e.}, \quad F \circ G = \text{id}_X \text{ a.e.} \quad (10)$$

Therefore both  $F$  and  $G$  are isomorphisms. By definition of  $\mathcal{L}$ ,

$$\begin{aligned} \mathcal{L}(G \circ \varphi, \varphi^{-1} \circ F) &= D_f((G \circ \varphi)_*\mu \| \nu) + D_f((\varphi^{-1} \circ F)_*\nu \| \mu) \\ &\quad + \alpha_{\text{cyc}} \cdot (\mathbb{E}_{\mathbf{x} \sim X} \|\varphi^{-1}(F(G(\varphi(\mathbf{x})))) - \mathbf{x}\| + \mathbb{E}_{\mathbf{y} \sim Y} \|G(\varphi^{-1}(F(\mathbf{y})))) - \mathbf{y}\|). \end{aligned}$$

Since  $(G \circ \varphi)_*\mu = G_*(\varphi_*\mu)$  and  $\varphi$  is measure-preserving, eq. (9) implies that  $(G \circ \varphi)_*\mu = \nu$ . Similarly,  $(\varphi^{-1} \circ F)_*\nu = \mu$  since  $\varphi^{-1}$  is measure-preserving as well. This shows that

$$D_f((G \circ \varphi)_*\mu \| \nu) = D_f((\varphi^{-1} \circ F)_*\nu \| \mu) = 0.$$

Using eq. (10) and the fact that  $\varphi^{-1} \circ \varphi = \varphi \circ \varphi^{-1} = \text{id}_X$  almost everywhere, we conclude that

$$\mathbb{E}_{\mathbf{y} \sim Y} \|G(\varphi^{-1}(F(\mathbf{y})))) - \mathbf{y}\| = \mathbb{E}_{\mathbf{y} \sim Y} \|\mathbf{y} - \mathbf{y}\| = 0.$$

and

$$\mathbb{E}_{\mathbf{x} \sim X} \|\varphi^{-1}(F(G(\varphi(\mathbf{x})))) - \mathbf{x}\| = \mathbb{E}_{\mathbf{x} \sim X} \|\mathbf{x} - \mathbf{x}\| = 0.$$

Combining these observations together, we deduce that

$$\mathcal{L}(G \circ \varphi, \varphi^{-1} \circ F) = 0$$

<sup>1</sup>Inequality should be understood in the ‘modulo null sets’ sense here, i.e., we assert that there are positive probability sets on which the maps do differ.

and the proof of eq. (7) is complete. To prove eq. (8), first note that there exists a set  $A \in \mathcal{X}$  such that  $\mu(A) > 0$  and

$$\varphi(\mathbf{x}) \neq \mathbf{x} \quad \text{for all } \mathbf{x} \in A,$$

since we assume that  $\varphi$  essentially differs from the identity mapping. If  $G \circ \varphi = G$   $\mu$ -a.e., then  $F \circ G \circ \varphi = F \circ G$   $\mu$ -a.e. as well, which implies that  $\varphi(\mathbf{x}) = \mathbf{x}$  for  $\mu$ -almost every  $\mathbf{x}$ , which is a contradiction. In a similar way one can show that  $\varphi^{-1} \circ F$  essentially differs from  $F$ .  $\square$

We provide the following converse to Proposition 2.1.

**Proposition 2.2** (Kernel as a principal homogeneous space). *Let  $X = (X, \mathcal{X}, \mu)$ ,  $Y = (Y, \mathcal{Y}, \nu)$  be probability spaces. Let  $F : X \rightarrow Y$ ,  $G : Y \rightarrow X$  and  $F' : X \rightarrow Y$ ,  $G' : Y \rightarrow X$  be measurable maps satisfying*

$$\mathcal{L}(F, G) = 0 \quad \text{and} \quad \mathcal{L}(F', G') = 0. \quad (11)$$

*Then there exists a unique probability space automorphism  $\varphi : X \rightarrow X$  such that*

$$F \circ \varphi = F' \quad \text{and} \quad \varphi^{-1} \circ G = G'.$$

For the proof it suffices to take  $\varphi := G \circ F'$ . Combined with Proposition 2.1, this allows us to say that the group  $\text{Aut}(X)$  of probability space automorphisms of  $X$  acts *freely and transitively* on the set of isomorphisms  $\text{Iso}(X, Y)$  when the latter set is nonempty. This amounts to saying that the space of solutions of CycleGAN is a principal homogeneous space. It can be helpful to view this result from the abstract category theory point of view, that is, if  $\mathcal{C}$  is a category and  $X \in \mathcal{C}$  is any fixed object, then for any object  $Y \in \mathcal{C}$  the automorphism group  $\text{Aut}(X)$  acts on the set of homomorphisms  $\text{Hom}(X, Y)$  on the right by composition, i.e. we define

$$\alpha(\phi) := \phi \circ \alpha \quad \text{for all } \phi \in \text{Hom}(X, Y), \alpha \in \text{Aut}(X).$$

This action leaves the space of isomorphisms  $\text{Iso}(X, Y) \subseteq \text{Hom}(X, Y)$  invariant, and this restricted action is transitive if  $\text{Iso}(X, Y)$  is nonempty, and, furthermore, free, i.e.  $\alpha(\phi) \neq \phi$  for all  $\alpha \neq \text{id}_X$  and all  $\phi \in \text{Iso}(X, Y)$ .

To proceed with our analysis for case of approximate solutions for extended CycleGAN loss, we first formulate a useful ‘push-forward property’ for general  $f$ -divergences between distributions<sup>2</sup> on  $\mathbb{R}^n$ .

**Lemma 2.1** (Push-forward property for  $f$ -divergences). *Let  $p, q$  be distributions on  $\mathbb{R}^n$  and  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a diffeomorphism. Then for any  $f$ -divergence  $D_f$  we have*

$$D_f(\varphi_* p \| q) = D_f(p \| (\varphi^{-1})_* q) \quad (12)$$

*Proof.* First of all, change of variables formula for the integral implies that

$$\begin{aligned} (\varphi_* p)(\mathbf{x}) &= p(\varphi^{-1}(\mathbf{x})) \left| \det \frac{\partial \varphi^{-1}}{\partial \mathbf{x}} \right|(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathbb{R}^n, \\ ((\varphi^{-1})_* q)(\mathbf{y}) &= q(\varphi(\mathbf{y})) \left| \det \frac{\partial \varphi}{\partial \mathbf{y}} \right|(\mathbf{y}) \quad \text{for all } \mathbf{y} \in \mathbb{R}^n. \end{aligned}$$

Therefore,

$$D_f(\varphi_* p \| q) = \int f\left(\frac{\varphi_* p(\mathbf{x})}{q(\mathbf{x})}\right) q(\mathbf{x}) d\mathbf{x} = \int f\left(\frac{p(\varphi^{-1}(\mathbf{x})) \left| \det \frac{\partial \varphi^{-1}}{\partial \mathbf{x}} \right|(\mathbf{x})}{q(\mathbf{x})}\right) q(\mathbf{x}) d\mathbf{x}.$$

Applying change of variables formula with  $\mathbf{x} = \varphi(\mathbf{y})$ , we get

$$\begin{aligned} & \int f\left(\frac{p(\varphi^{-1}(\mathbf{x})) \left| \det \frac{\partial \varphi^{-1}}{\partial \mathbf{x}} \right|(\mathbf{x})}{q(\mathbf{x})}\right) q(\mathbf{x}) d\mathbf{x} \\ &= \int f\left(\frac{p(\varphi^{-1}(\varphi(\mathbf{y}))) \left| \det \frac{\partial \varphi^{-1}}{\partial \mathbf{x}} \right|(\varphi(\mathbf{y}))}{q(\varphi(\mathbf{y}))}\right) q(\varphi(\mathbf{y})) \left| \det \frac{\partial \varphi}{\partial \mathbf{y}} \right|(\mathbf{y}) d\mathbf{y} \\ &\stackrel{*}{=} \int f\left(\frac{p(\mathbf{y})}{q(\varphi(\mathbf{y})) \left| \det \frac{\partial \varphi}{\partial \mathbf{y}} \right|(\mathbf{y})}\right) q(\varphi(\mathbf{y})) \left| \det \frac{\partial \varphi}{\partial \mathbf{y}} \right|(\mathbf{y}) d\mathbf{y} = D_f(p \| (\varphi^{-1})_* q), \end{aligned}$$

<sup>2</sup>While very natural to conjecture and easy to prove, we were unable to find references to it in existing ML literature, so we dubbed this property a ‘push-forward property’ and provide a proof.

where the equality in (\*) uses a general property of Jacobians of smooth invertible maps that  $\frac{\partial \varphi^{-1}}{\partial \mathbf{x}} \circ \varphi = \left(\frac{\partial \varphi}{\partial \mathbf{y}}\right)^{-1}$ . Hence  $D_f(\varphi_* p \| q) = D_f(p \| (\varphi^{-1})_* q)$ , which completes the proof.  $\square$

We are now ready to prove the perturbation bounds for approximate solutions.

**Proposition 2.3** (Perturbation bound). *Let  $X, Y$  be probability spaces with probability densities  $p_X, p_Y \in L^1(\mathbb{R}^n)$  and let  $\varphi \in \text{Aut}(X)$  be a diffeomorphic probability space automorphism. Assume that  $\varphi^{-1}$  is  $C_\varphi$ -Lipshitz, where  $C_\varphi > 0$  is some positive constant. Let  $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be measurable maps. Then the following perturbation bound holds for extended CycleGAN loss:*

$$\mathcal{L}_{\text{ext}}(G \circ \varphi, \varphi^{-1} \circ F) \leq \max(C_\varphi, 1) \cdot \mathcal{L}_{\text{ext}}(G, F) + 2 \cdot \alpha_{\text{id}} \cdot \mathbb{E}_{\mathbf{x} \sim X} \|\varphi(\mathbf{x}) - \mathbf{x}\|. \quad (13)$$

*Proof.* The proof is an adaptation of the proof of Proposition 2.1. By definition of  $\mathcal{L}_{\text{ext}}$ ,

$$\begin{aligned} \mathcal{L}_{\text{ext}}(G \circ \varphi, \varphi^{-1} \circ F) &= D_f((G \circ \varphi)_* p_X \| p_Y) + D_f((\varphi^{-1} \circ F)_* p_Y \| p_X) \\ &\quad + \alpha_{\text{cyc}} \cdot (\mathbb{E}_{\mathbf{x} \sim X} \|\varphi^{-1}(F(G(\varphi(\mathbf{x})))) - \mathbf{x}\| + \mathbb{E}_{\mathbf{y} \sim Y} \|G(\varphi^{-1}(F(\mathbf{y}))) - \mathbf{y}\|) \\ &\quad + \alpha_{\text{id}} \cdot (\mathbb{E}_{\mathbf{y} \sim Y} \|\varphi^{-1}(F(\mathbf{y})) - \mathbf{y}\| + \mathbb{E}_{\mathbf{x} \sim X} \|G(\varphi(\mathbf{x})) - \mathbf{x}\|). \end{aligned}$$

Firstly, since  $\varphi$  is measure-preserving,  $D_f((G \circ \varphi)_* p_X \| p_Y) = D_f(G_* p_X \| p_Y)$ . Using Lemma 2.1 and the fact that  $\varphi$  is measure-preserving again, we see that  $D_f((\varphi^{-1} \circ F)_* p_Y \| p_X) = D_f(F_* p_Y \| \varphi_* p_X) = D_f(F_* p_Y \| p_X)$ .

Secondly,

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim X} \|\varphi^{-1}(F(G(\varphi(\mathbf{x})))) - \mathbf{x}\| &= \mathbb{E}_{\mathbf{x} \sim X} \|\varphi^{-1}(F(G(\varphi(\mathbf{x})))) - \varphi^{-1}(\varphi(\mathbf{x}))\| \stackrel{*}{=} \\ &\stackrel{*}{=} \mathbb{E}_{\mathbf{x} \sim X} \|\varphi^{-1}(F(G(\mathbf{x}))) - \varphi^{-1}(\mathbf{x})\| \leq C_\varphi \cdot \mathbb{E}_{\mathbf{x} \sim X} \|F(G(\mathbf{x})) - \mathbf{x}\|, \end{aligned}$$

where the equality (\*) uses the fact that  $\varphi$  is measure-preserving. As in before,  $\varphi \circ \varphi^{-1} = \text{id}_X$  almost everywhere, thus  $\mathbb{E}_{\mathbf{y} \sim Y} \|G(\varphi^{-1}(F(\mathbf{y}))) - \mathbf{y}\| = \mathbb{E}_{\mathbf{y} \sim Y} \|G(F(\mathbf{y})) - \mathbf{y}\|$ .

Finally, since  $\varphi$  is a probability space automorphism and  $\varphi^{-1}$  is  $C_\varphi$ -Lipshitz, we conclude that

$$\begin{aligned} \mathbb{E}_{\mathbf{y} \sim Y} \|\varphi^{-1}(F(\mathbf{y})) - \mathbf{y}\| &\leq \mathbb{E}_{\mathbf{y} \sim Y} \|\varphi^{-1}(F(\mathbf{y})) - \varphi^{-1}(\mathbf{y})\| + \mathbb{E}_{\mathbf{y} \sim Y} \|\varphi^{-1}(\mathbf{y}) - \mathbf{y}\| \leq \\ &\leq C_\varphi \cdot \mathbb{E}_{\mathbf{y} \sim Y} \|F(\mathbf{y}) - \mathbf{y}\| + \mathbb{E}_{\mathbf{y} \sim Y} \|\varphi^{-1}(\mathbf{y}) - \mathbf{y}\| = C_\varphi \cdot \mathbb{E}_{\mathbf{y} \sim Y} \|F(\mathbf{y}) - \mathbf{y}\| + \mathbb{E}_{\mathbf{y} \sim Y} \|\varphi(\mathbf{y}) - \mathbf{y}\| \end{aligned}$$

and that

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim X} \|G(\varphi(\mathbf{x})) - \mathbf{x}\| &= \mathbb{E}_{\mathbf{x} \sim X} \|G(\varphi(\mathbf{x})) - \varphi^{-1}(\varphi(\mathbf{x}))\| = \mathbb{E}_{\mathbf{x} \sim X} \|G(\mathbf{x}) - \varphi^{-1}(\mathbf{x})\| \leq \\ &\leq \mathbb{E}_{\mathbf{x} \sim X} \|G(\mathbf{x}) - \mathbf{x}\| + \mathbb{E}_{\mathbf{x} \sim X} \|\mathbf{x} - \varphi^{-1}(\mathbf{x})\| = \mathbb{E}_{\mathbf{x} \sim X} \|G(\mathbf{x}) - \mathbf{x}\| + \mathbb{E}_{\mathbf{x} \sim X} \|\varphi(\mathbf{x}) - \mathbf{x}\|. \end{aligned}$$

Combining all these estimates together, we deduce that

$$\mathcal{L}_{\text{ext}}(G \circ \varphi, \varphi^{-1} \circ F) \leq \max(C_\varphi, 1) \cdot \mathcal{L}_{\text{ext}}(G, F) + 2 \cdot \alpha_{\text{id}} \cdot \mathbb{E}_{\mathbf{x} \sim X} \|\varphi(\mathbf{x}) - \mathbf{x}\|$$

and the proof is complete.  $\square$

**Corollary 2.1** (Asymptotic perturbation bound). *In the setting of Proposition 2.3, let  $G_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and  $F_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$  for  $i \geq 1$  be a sequence of measurable maps such that the ‘pure’ CycleGAN loss converges to zero, i.e.,*

$$\lim_{i \rightarrow \infty} \mathcal{L}(G_i, F_i) = 0$$

and let

$$\bar{\mathcal{L}}_{\text{id}} := \limsup_{i \rightarrow \infty} (\mathbb{E}_{\mathbf{y} \sim Y} \|F_i(\mathbf{y}) - \mathbf{y}\| + \mathbb{E}_{\mathbf{x} \sim X} \|G_i(\mathbf{x}) - \mathbf{x}\|).$$

Then the following asymptotic perturbation bound holds for the ‘extended’ CycleGAN loss:

$$\limsup_{i \rightarrow \infty} \mathcal{L}_{\text{ext}}(G_i \circ \varphi, \varphi^{-1} \circ F_i) \leq \max(C_\varphi, 1) \cdot \alpha_{\text{id}} \cdot \bar{\mathcal{L}}_{\text{id}} + 2 \cdot \alpha_{\text{id}} \cdot \mathbb{E}_{\mathbf{x} \sim X} \|\varphi(\mathbf{x}) - \mathbf{x}\|.$$

Corollary 2.1 has a direct practical implication. When using a CycleGAN model for translating substantially different distributions (such as different medical imaging modalities) one would be forced to pick a small value for  $\alpha_{\text{id}}$  in order for the model to produce reasonable results. Furthermore, since the distributions are substantially different, we can expect that  $\bar{\mathcal{L}}_{\text{id}} \gg 2 \cdot \mathbb{E}_{\mathbf{x} \sim X} \|\varphi(\mathbf{x}) - \mathbf{x}\|$  for many nontrivial automorphism  $\varphi$ . Therefore, the asymptotic perturbation bound automatically implies that the approximate solution space admits some symmetry, potentially leading to undesirable results.

## 2.2 EXISTENCE OF AUTOMORPHISMS

By Proposition 2.1 we see that if either space admits a nontrivial probability automorphism, then the CycleGAN problem has multiple solutions. However, for this to be a problem in practice there must actually exist such probability automorphisms, which we shall now show is the case. First of all, we state the following proposition, which says that we can transfer automorphism from an isomorphic copy of  $X$  to  $X$  itself.

**Lemma 2.2.** *Let  $f : Z \rightarrow X$  be an isomorphism of probability spaces and  $T : Z \rightarrow Z$  be an automorphism of  $Z$ . Then  $S := f \circ T \circ f^{-1}$  is an automorphism of  $X$  and the diagram*

$$\begin{array}{ccc} Z & \xleftarrow{f^{-1}} & X \\ T \downarrow & & \downarrow S \\ Z & \xrightarrow{f} & X \end{array}$$

*commutes. Furthermore, if  $Z \subset \mathbb{R}^n$ ,  $X \subset \mathbb{R}^m$  are submanifolds and  $f, T$  are diffeomorphisms, then  $S$  is a diffeomorphism as well.*

*Proof.* The first claim follows from invertibility of  $f$  and  $T$ . The second claim follows from the definition of a diffeomorphism between submanifolds, see appendix A.  $\square$

An important notion in probability theory is that of a Lebesgue probability space. Many probability spaces which emerge in practice such as  $[0, 1]^n \subset \mathbb{R}^n$  with the Lebesgue measure or  $\mathbb{R}^n$  with a Gaussian probability distribution, both defined on the respective  $\sigma$ -algebras of Lebesgue measurable sets, are instances of Lebesgue probability spaces.

**Definition 2.1.** *A probability space  $X$  is called a Lebesgue probability space if it is isomorphic as a measure space to a disjoint union  $([0, c], \lambda)$ , where  $\lambda$  is the Lebesgue measure on the  $\sigma$ -algebra of Lebesgue measurable subsets of the interval  $[0, c]$ , and at most countably many atoms of total mass  $1 - c$ .*

Informally speaking, this definition says that Lebesgue probability spaces consist of a continuous part and at most countably many Dirac deltas (=atoms). First of all, we provide an abstract result about existence of nontrivial probability space automorphisms in Lebesgue probability spaces which are either ‘not purely atomic’ or have at least two atoms with equal mass. ‘Not purely atomic’ means that the sum of the probabilities of all atoms is strictly less than 1.

**Proposition 2.4.** *Let  $X$  be a Lebesgue probability space such that at least one of the assumptions*

1.  *$X$  not purely atomic;*
2. *there exist at least two atoms  $a_j, a_k$  in  $X$  with equal mass*

*holds. Then  $X$  admits nontrivial automorphisms.*

*Proof.* If the space  $X$  is not purely atomic, we have  $X \simeq [0, c] \sqcup \bigsqcup_{i \geq 1} a_i$  for some  $c > 0$ , where  $[0, c]$  is the continuous part and  $\bigsqcup_{i \geq 1} a_i$  is the atomic part of the probability measure  $\mu$ . Interval  $[0, c]$  admits at least one nontrivial automorphism, namely the transformation  $x \mapsto c - x$  (leaving the atoms fixed), hence so does  $X$  by Lemma 2.2. In fact, there are infinitely many other automorphisms, which can be obtained by exchanging nonoverlapping subintervals  $(a, a + d), (b, b + d) \subset [0, c]$  of the same length. If there exist two atoms  $a_j, a_k$  in  $X$  with equal mass, then a transformation which transposes  $a_j$  with  $a_k$  and keeps the rest of  $X$  fixed is a nontrivial automorphism.  $\square$

Probability spaces of images which appear in real life typically have a continuous component which would correspond to continuous variations in object sizes, lighting conditions, etc. Therefore, they admit some probability space automorphisms. However, such abstract automorphisms can be highly discontinuous, which would make it questionable if neural networks can learn them. We would like to show that there are also automorphisms which are smooth, at least locally. For this, we first state the following technical claim. The proof is provided in appendix A.

**Proposition 2.5.** *Let  $\mu$  be a Borel probability measure on  $\mathbb{R}^n$  and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a continuous injective function. Then  $f : (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mu) \rightarrow (\mathbb{R}^m, \mathcal{B}(\mathbb{R}^m), f_*\mu)$  is an isomorphism of probability spaces, where  $f_*\mu$  denotes the push-forward of measure  $\mu$  to  $\mathbb{R}^m$ .*

Finally, we show the existence of smooth automorphisms under the assumption that our data manifold  $\mathcal{D} \subset \mathbb{R}^m$  can be generated by embedding  $\mathbb{R}^n$  with standard Gaussian measure into  $\mathbb{R}^m$  as a submanifold. We write  $\gamma_n$  for the standard Gaussian probability measure on the space  $\mathbb{R}^n$ .

**Proposition 2.6.** *Let  $Z := (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \gamma_n)$  be an  $n$ -dimensional standard Gaussian distribution. Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a manifold embedding. Denote by  $X$  the probability space  $(\mathbb{R}^m, \mathcal{B}(\mathbb{R}^m), f_*\gamma_n)$ . Then the following assertions hold:*

1.  $f$  is an isomorphism of probability spaces when viewed as a map  $Z \rightarrow X$ ;
2. every rotation  $T \in \text{SO}(\mathbb{R}^n)$  is a probability space automorphism and a diffeomorphism of  $Z$ .  $T$  induces a probability space automorphism of  $X$  which is, additionally, a diffeomorphism when restricted to  $\text{Im } f \subset \mathbb{R}^m$ .

*Proof.* The first claim follows directly from Proposition 2.5. For the second part, it is clear that rotations in  $\text{SO}(\mathbb{R}^n)$  preserve isotropic Gaussian distribution, and the rest follows from Lemma 2.2.  $\square$

The connection with generative models is clear if we take  $f$  to be an invertible generative model such as RealNVP Dinh et al. (2016) or Glow Kingma & Dhariwal (2018). The assumption of manifold embedding in the proposition can be seen as too limiting in general, and we explain how to ‘bypass’ it in Lemma A.1 for the interested readers. In conclusion, if we assume that the distributions we are working with could be represented by an invertible generative model, then there exists a rich space of automorphisms. Given the success of e.g. Glow, this assumption seems to be valid for natural images.

### 3 NUMERICAL RESULTS

Since we have established that the existence of automorphisms can negatively impact the results of CycleGAN, we now demonstrate how this can happen by considering a toy case with a known solution and demonstrating that CycleGAN can and does learn a nontrivial automorphism. The toy experiment which we perform is translation of MNIST dataset to itself. That is, at training time we pick two minibatches  $\text{batch}_A$  and  $\text{batch}_B$  from MNIST at random and use these as samples from  $X$  and  $Y$  respectively. The generator neural network in this case is a convolutional autoencoder with residual blocks, fully connected layer in the bottleneck and *no* skip connections from encoder to decoder. Identity loss term is not used. We also train a simple CNN for MNIST classification in order to classify CycleGAN outputs. The networks were trained using SGD. The ‘natural’ transformation in this case is, of course, the identity mapping and we expect the classification of the inputs and outputs to stay the same. But we shall see that this is not the case.

In fig. 2a–fig. 2h we show some examples for the generated fake samples and the reconstruction on test set. In fig. 3a–fig. 3b we provide the confusion matrices for the A2B and B2A generators respectively. We use these matrices to understand if e.g. the class of transformed image for A2B translation equals the source class, or if is a random variable independent of the source class, or if we can spot some deterministic permutation of classes. We have observed that in practice the identity mapping is not learned. Instead, the network leans towards producing a certain permutation of digits, rather than identity or a random assignment of classes independent of the source label. One explanation would be as follows. Suppose that we can perfectly disentangle class and style in latent digit representation Makhzani et al. (2015). Then *any* permutation in  $S_{10}$ , acting on the class part of the latent code, determines a probability space automorphism on the space of digits, which can be learned by a neural network. Further investigation of confusion matrices reveals that the networks introduce short *cycles*, e.g., mapping 2 to 6 and vice versa.

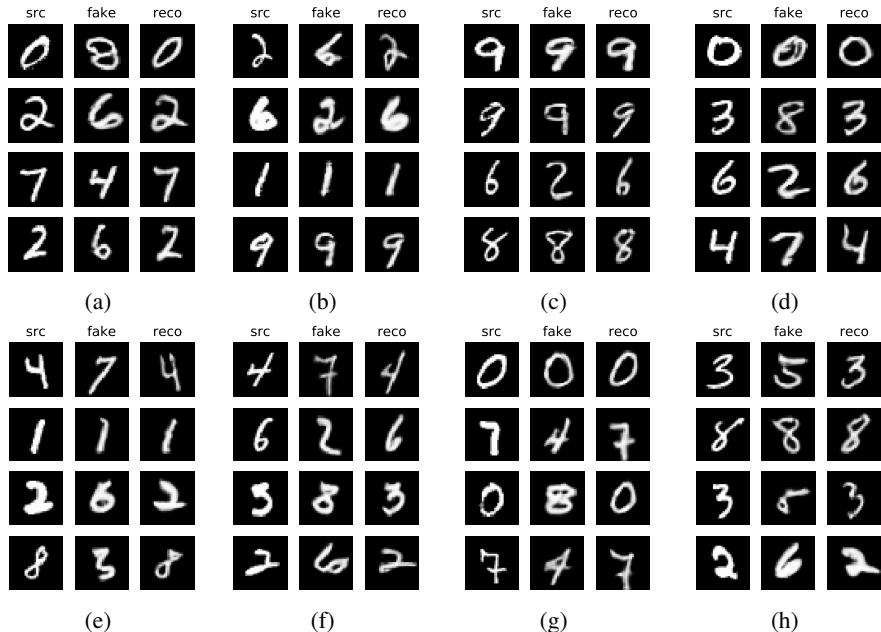


Figure 2: Examples on MNIST2MNIST task. (a)-(d) A2A translation, first column are samples from A, second column are 'fake B' and third column are reconstructions of original samples from A (e)-(h) same for B2B translation.

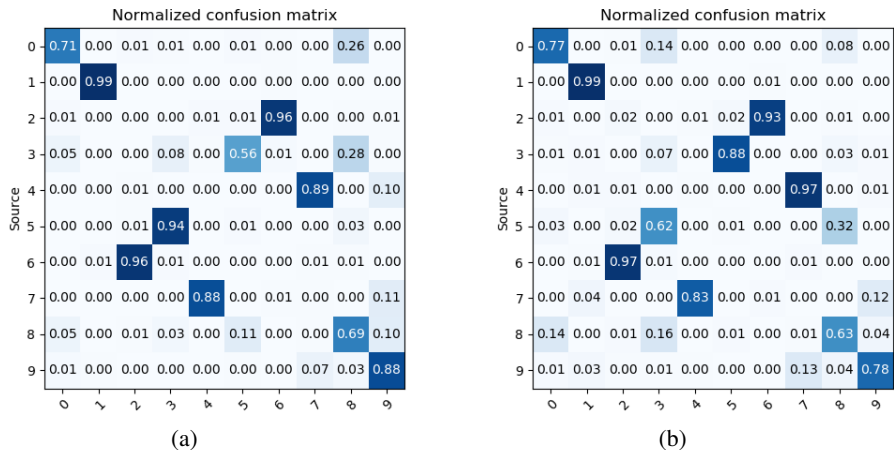


Figure 3: Normalized confusion matrices for A2B and B2A generator respectively.

#### 4 CONCLUSION

We have shown theoretically that under mild assumptions, the kernel of the CycleGAN admits nontrivial symmetries and has a natural structure of a principal homogeneous space. To show empirically that such symmetries can be learned, we have trained a CycleGAN on the task of translating a domain to itself. In particular, we show that on the MNIST2MNIST task, in contrast to the expected identity, the CycleGAN learns to permute the digits. We have therefore effectively shown that it is not the CycleGAN loss which prevents this from occurring more often, but hypothesize that skip connections in the network architecture have a major influence. Given the theoretical results in Corollary 2.1 which suggest ambiguity of solutions, even in the presence of an identity loss term, we advocate that great care should be taken when CycleGAN is used to translate between substantially different distributions in critical tasks such as medical imaging.



## REFERENCES

- V. I. Bogachev. *Measure theory. Vol. I, II*. Springer-Verlag, Berlin, 2007. ISBN 978-3-540-34513-8; 3-540-34513-2. doi: 10.1007/978-3-540-34514-5. URL <https://doi.org/10.1007/978-3-540-34514-5>.
- Joseph Paul Cohen, Margaux Luck, and Sina Honari. How to Cure Cancer (in images) with Unpaired Image Translation. In *Medical Imaging with Deep Learning (MIDL)*, volume 1, pp. 1–3, 2018.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. abs/1605.08803, 2016. URL <http://arxiv.org/abs/1605.08803>.
- Tanja Eisner, Bálint Farkas, Markus Haase, and Rainer Nagel. *Operator theoretic aspects of ergodic theory*, volume 272 of *Graduate Texts in Mathematics*. Springer, Cham, 2015. ISBN 978-3-319-16897-5; 978-3-319-16898-2. doi: 10.1007/978-3-319-16898-2. URL <https://doi.org/10.1007/978-3-319-16898-2>.
- Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. A neural algorithm of artistic style. abs/1508.06576, 2015. URL <http://arxiv.org/abs/1508.06576>.
- Xiao Han. Mr-based synthetic ct generation using a deep convolutional neural network method. *Medical physics*, 44(4):1408–1419, 2017.
- Alexander S. Kechris. *Classical descriptive set theory*, volume 156 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1995. ISBN 0-387-94374-9. doi: 10.1007/978-1-4612-4190-4. URL <https://doi.org/10.1007/978-1-4612-4190-4>.
- Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 31*, pp. 10215–10224. Curran Associates, Inc., 2018. URL <http://papers.nips.cc/paper/8224-glow-generative-flow-with-invertible-1x1-convolutions.pdf>.
- Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 700–708. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/6672-unsupervised-image-to-image-translation-networks.pdf>.
- Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, and Ian J. Goodfellow. Adversarial autoencoders. abs/1511.05644, 2015. URL <http://arxiv.org/abs/1511.05644>.
- Jens Sjölund, Daniel Forsberg, Mats Andersson, and Hans Knutsson. Generating patient specific pseudo-ct of the head from mr using atlas-based regression. *Physics in Medicine & Biology*, 60(2): 825, 2015.
- Frank W. Warner. *Foundations of differentiable manifolds and Lie groups*, volume 94 of *Graduate Texts in Mathematics*. Springer-Verlag, New York-Berlin, 1983. ISBN 0-387-90894-3. Corrected reprint of the 1971 edition.
- Jelmer M. Wolterink, Anna M. Dinkla, Mark H. F. Savenije, Peter R. Seevinck, Cornelis A. T. van den Berg, and Ivana Isgum. Deep MR to CT synthesis using unpaired data. abs/1708.01155, 2017. URL <http://arxiv.org/abs/1708.01155>.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. abs/1703.10593, 2017. URL <http://arxiv.org/abs/1703.10593>.

## A BACKGROUND

Firstly, we very briefly explain the probability theory language we use in this article, and we refer the reader to (Eisner et al., 2015; Bogachev, 2007) for more details. Formally, a *measurable space*  $(X, \mathcal{X})$  is a pair of a set  $X$  and a  $\sigma$ -algebra  $\mathcal{X}$  of subsets of  $X$ . Given a topological space  $X$  with topology  $\mathcal{U}$ , there exists the smallest  $\sigma$ -algebra  $\mathcal{B}(X)$ , which contains all open sets in  $\mathcal{U}$ . This  $\sigma$ -algebra is called *Borel  $\sigma$ -algebra* of  $X$  and its elements are called *Borel sets*. A *probability space*  $X = (X, \mathcal{X}, \mu)$  is a triple of a set  $X$ , a sigma algebra  $\mathcal{X}$  of subsets of  $X$  and a probability measure  $\mu$  defined on the sigma-algebra  $\mathcal{X}$ . Given a probability space  $(X, \mathcal{X}, \mu)$ , a measurable set  $A \in \mathcal{X}$  is called an *atom* if  $\mu(A) > 0$  and for all measurable  $B \subset A$  such that  $\mu(B) < \mu(A)$  we have  $\mu(B) = 0$ . Given measurable spaces  $(X, \mathcal{X})$  and  $(Y, \mathcal{Y})$ , we say that a mapping  $\phi : X \rightarrow Y$  is *measurable* if for any  $A \in \mathcal{Y}$  we have  $\phi^{-1}(A) \in \mathcal{X}$ . If  $X = (X, \mathcal{X}, \mu)$  and  $Y = (Y, \mathcal{Y}, \nu)$  are probability spaces and  $\phi : X \rightarrow Y$  is a measurable map, we say that  $\phi$  is *measure-preserving* if for all  $A \in \mathcal{Y}$  we have  $\mu(\phi^{-1}(A)) = \nu(A)$ . An approximation argument easily shows that a measurable transformation  $\phi : X \rightarrow X$  is measure-preserving if and only if for all nonnegative measurable functions  $f$  on  $X$  we have

$$\mathbb{E}_{\mathbf{x} \sim X} f(\mathbf{x}) d\mu = \mathbb{E}_{\mathbf{x} \sim X} (f \circ \phi)(\mathbf{x}) d\mu.$$

Given a probability space  $X$ , a measurable space  $(Y, \mathcal{Y})$  and a measurable map  $\phi : X \rightarrow Y$ , we define the *push-forward measure*  $\phi_*\mu$  on  $\mathcal{Y}$  by setting  $(\phi_*\mu)(A) := \mu(\phi^{-1}(A))$  for all  $A \in \mathcal{Y}$ .

Let  $(X, \mathcal{X}, \mu)$  and  $(Y, \mathcal{Y}, \nu)$  be probability spaces and  $f : X \rightarrow Y$  be a measure-preserving map. A measurable map  $g : Y \rightarrow X$  is called an *essential inverse* of  $f$  if  $f \circ g = \text{id}_Y$  for  $\nu$ -almost every  $\mathbf{y} \in Y$  and  $g \circ f = \text{id}_X$  for  $\mu$ -almost every  $\mathbf{x} \in X$ . One can show that essential inverse is measure preserving and uniquely defined up to equality almost everywhere. We say that  $f$  is an *isomorphism* if it admits an essential inverse. An isomorphism  $f : X \rightarrow X$  is called an *automorphism*.

We remind the reader that a *Polish space* is a separable completely metrizable topological space. A *Borel probability space* is a Polish space endowed with a probability measure  $\mu$  on its Borel  $\sigma$ -algebra, and we will also say that  $\mu$  is a *Borel probability measure*. The basic examples of Borel probability spaces would be e.g. the spaces  $[0, 1]^n \subset \mathbb{R}^n$  with its Borel  $\sigma$ -algebra  $\mathcal{B}(\mathbb{R}^n)$ , endowed with Lebesgue measure  $\lambda_n$ . A Borel  $\sigma$ -algebra of the space  $[0, 1]^n$  endowed with Lebesgue measure  $\lambda_n$  can be extended by adding all  $\lambda_n$ -measurable sets, leading to the  $\sigma$ -algebra of *Lebesgue-measurable sets*.

For the proof of Proposition 2.5 we need the following theorem, see Kechris (1995), Theorem 15.1.

**Theorem A.1** (Lusin-Souslin theorem). *Let  $X, Y$  be Polish spaces and  $f : X \rightarrow Y$  be continuous. If  $A \subset X$  is Borel and  $f|_A$  is injective, then  $f(A)$  is Borel.*

*Proof of Proposition 2.5.* Denote the image  $f(\mathbb{R}^n) \subset \mathbb{R}^m$  by  $\text{Im } f$ . Then  $\text{Im } f \subset \mathbb{R}^m$  is a Borel subset, since  $\mathbb{R}^n$  is a countable union of a compact sets and  $f$  is continuous. Furthermore, from Lusin-Souslin theorem (theorem A.1) it follows that for every Borel subset  $A \subset \mathbb{R}^n$  its image  $f(A) \subset \mathbb{R}^m$  is Borel as well. Pick a point  $x_0 \in \mathbb{R}^n$  which is not an atom of  $\mu$ . We want to define an almost everywhere inverse  $\tilde{f}$  of  $f$ . Define a function  $\tilde{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$  by

$$\tilde{f}(\mathbf{x}) = \begin{cases} f^{-1}(\mathbf{x}), & \text{if } \mathbf{x} \in \text{Im } f. \\ x_0, & \text{otherwise.} \end{cases}$$

Using the remark above it is easy to see that  $\tilde{f}$  is Borel measurable and that  $(f_*\mu)(\tilde{f}^{-1}(A)) = \mu(A)$  for every Borel  $A$ . It follows from the definition that  $\tilde{f} \circ f = \text{id}_{\mathbb{R}^n}$  and that

$$f \circ \tilde{f}(\mathbf{x}) = \begin{cases} \mathbf{x}, & \text{if } \mathbf{x} \in \text{Im } f. \\ f(x_0), & \text{otherwise.} \end{cases}$$

Since  $(f_*\mu)(\text{Im } f) = 1$ ,  $\tilde{f}$  is an almost everywhere inverse to  $f$ . We conclude that  $f$  is a probability space isomorphism.  $\square$

Secondly, we remind the reader of a couple of notions from differential geometry which we use in the text, and we refer the reader to e.g. (Warner, 1983) for more details. Given a subset  $X$  of a manifold  $M$  and a subset  $Y$  of a manifold  $N$ , a function  $f : X \rightarrow Y$  is said to be smooth if for all  $p \in X$  there is a neighborhood  $U \subset M$  of  $p$  and a smooth function  $g : U \rightarrow N$  such that  $g$  extends  $f$ , i.e.,

the restrictions agree  $g|_{U \cap X} = f|_{U \cap X}$ .  $f$  is said to be a *diffeomorphism* between  $X$  and  $Y$  if it is bijective, smooth and its inverse is smooth. Let  $M$  and  $N$  be smooth manifolds. A differentiable mapping  $f : M \rightarrow N$  is said to be an *immersion* if the tangent map  $d_p f : T_p M \rightarrow T_{f(p)} N$  is injective for all  $p \in M$ . If, in addition,  $f$  is a homeomorphism onto  $f(M) \subset N$ , where  $f(M)$  carries the subspace topology induced from  $N$ , we say that  $f$  is an *embedding*. If  $M \subset N$  and the inclusion map  $\iota : M \rightarrow N$  is an embedding, we say that  $M$  is a *submanifold* of  $N$ . Thus, the domain of an embedding is diffeomorphic to its image, and the image of an embedding is a submanifold.

We close this section with a small lemma, explaining how one can weaken the embedding assumption for generative models in Proposition 2.6.

**Lemma A.1.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be an injective manifold immersion. Let  $B_R \subset \mathbb{R}^n$  be an open ball of radius  $R > 0$  in  $\mathbb{R}^n$  and  $\overline{B}_R$  be its closure. Then  $f : B_R \rightarrow f(B_R)$  is a manifold embedding.*

*Proof.* Since  $\overline{B}_R$  is compact and  $f$  is continuous, image of every closed subset  $A \subseteq \overline{B}_R$  is compact and hence closed. This shows that  $f^{-1} : f(\overline{B}_R) \rightarrow \overline{B}_R$  is continuous and thus  $f : \overline{B}_R \rightarrow f(\overline{B}_R)$  is a homeomorphism. Restricting to the open ball  $B_R \subset \overline{B}_R$ , we conclude that  $f : B_R \rightarrow f(B_R)$  is a homeomorphism and thus a manifold embedding.  $\square$

As a consequence, for our example with spherical Gaussian latent vector one can take sufficiently large ball of radius  $R > 0$  in the latent space, truncating the latent distribution to ‘sufficiently likely’ values. This ball remains invariant under rotations, thus leading to a differentiable automorphism on the submanifold of ‘sufficiently likely’ images.