

Advanced Machine Learning Normalizing Flow

Alexandre Vérine

Master Intelligence Artificielle, Systèmes, Données

18/02/2022



1 Normalizing Flow in theory

- How does it work ?
- What can we do with it ?
- How is it trained ?

2 Structure of Normalizing flow

- Theoretical Structures
- First Structures
- Coupling - Glow
- Residual Network - ResFlow

3 Current Limits of Normalizing Flow

A Normalizing Flow is usually seen as:

- a *generative model*,
- a *bijective mapping*,
- an *invertible neural network*,
- a *density estimator*.

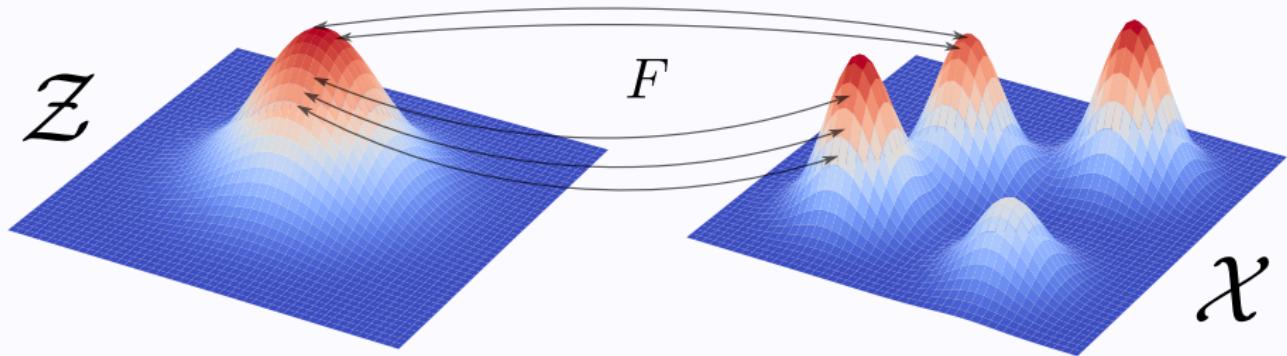


Figure: A mapping between two probability distributions
Point to point

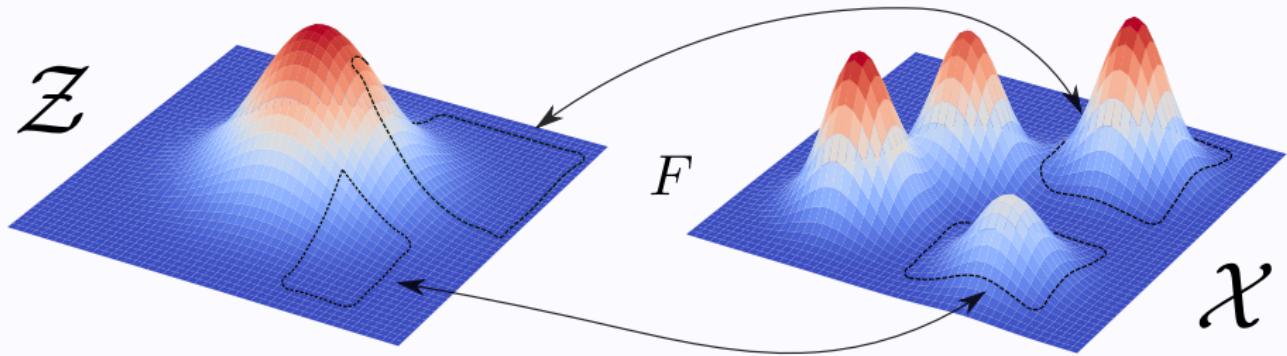


Figure: A mapping between two probability distributions
Subset to subset

Normalizing Flow

A Normalizing Flow is a bijective function between a data space \mathcal{X} and a latent space \mathcal{Z} , both subset of \mathbb{R}^d .

$$\begin{aligned} F : \quad \mathcal{X} &\longmapsto \mathcal{Z} \\ x &\longmapsto z = F(x) \end{aligned}$$

Data and Latent Distributions

In theory, a NF maps a target distribution P , ie the data distribution to a simple latent distribution Q .

Usually, Q is set to be a Normal Gaussian multivariate distribution $\mathcal{N}(0_d, I_d)$. p and q are respectively the probability densities of P and Q .

How does it work ?

In practice, the mapping is *not perfect*. P^* induces a distribution Q and similarly, the latent distribution Q induces \hat{P} , which is the learned distribution. The forward pass F is called the *Normalizing* direction while the inverse pass F^{-1} is called the *Generative* direction.

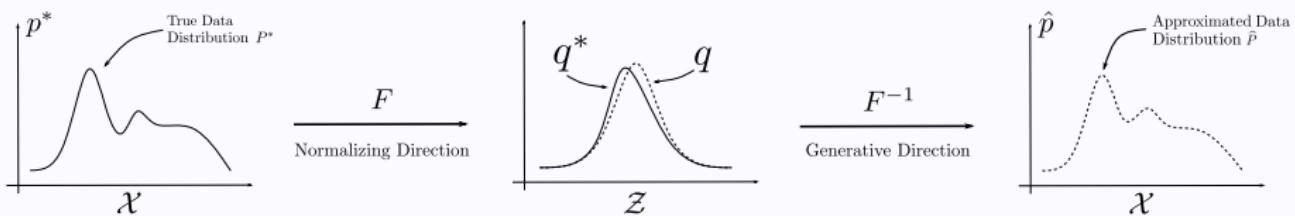


Figure: 1D Normalizing Flow process.

Induced Probabilities ?

Change of Variable Formula

For a bijective and continuous function F and a latent distribution Q , the distribution induced by Q and F is defined through the *change of variable formula*:

$$\forall \mathbf{x} \in \mathcal{X}, \quad \hat{p}(\mathbf{x}) = |\det \text{Jac}_F(\mathbf{x})| q(F(\mathbf{x})). \quad (1)$$

Induced Probabilities ?

$$\forall \mathbf{x} \in \mathcal{X}, \quad \hat{p}(\mathbf{x}) = |\det \text{Jac}_F(\mathbf{x})| q(F(\mathbf{x})).$$

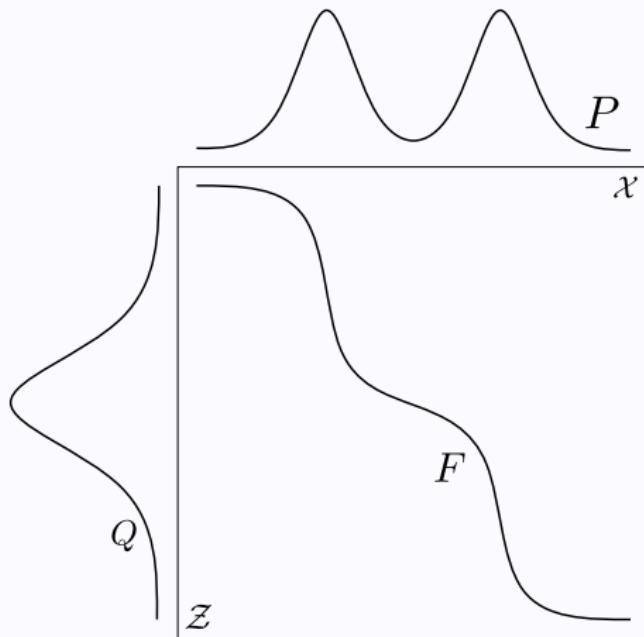


Figure: Example of 1D mapping

Induced Probabilities ?

Induced Probabilities ?

$$\forall \mathbf{x} \in \mathcal{X}, \quad \hat{p}(\mathbf{x}) = |\det \text{Jac}_F(\mathbf{x})| q(F(\mathbf{x})).$$

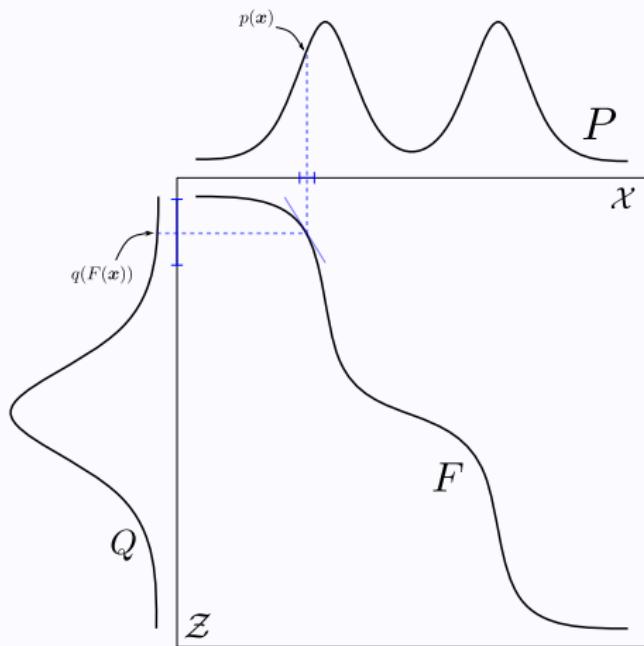


Figure: Example of 1D mapping

Induced Probabilities ?

$$\forall \mathbf{x} \in \mathcal{X}, \quad \hat{p}(\mathbf{x}) = |\det \text{Jac}_F(\mathbf{x})| q(F(\mathbf{x})).$$

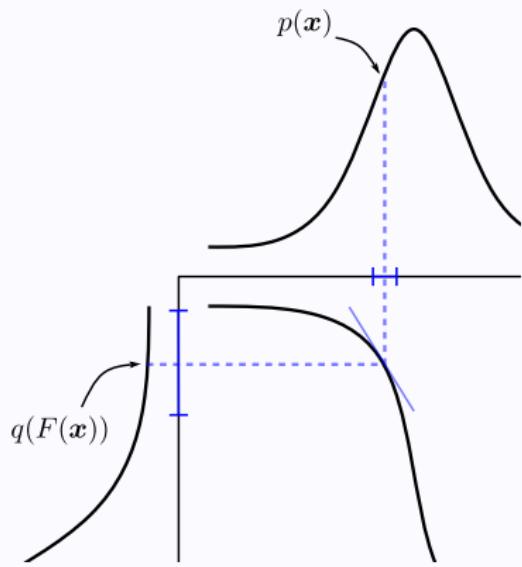


Figure: Example of 1D mapping

Induced Probabilities ?

$$\forall \mathbf{x} \in \mathcal{X}, \quad \hat{p}(\mathbf{x}) = |\det \text{Jac}_F(\mathbf{x})| q(F(\mathbf{x})).$$

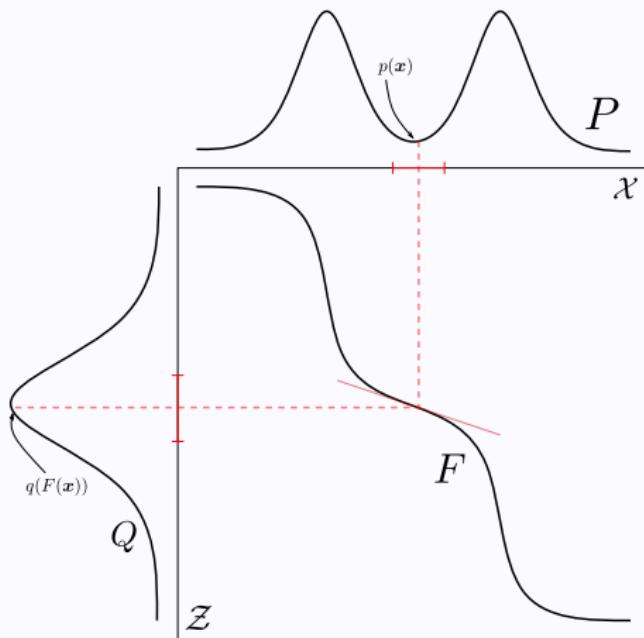


Figure: Example of 1D mapping

Induced Probabilities ?

$$\forall \mathbf{x} \in \mathcal{X}, \quad \hat{p}(\mathbf{x}) = |\det \text{Jac}_F(\mathbf{x})| q(F(\mathbf{x})).$$

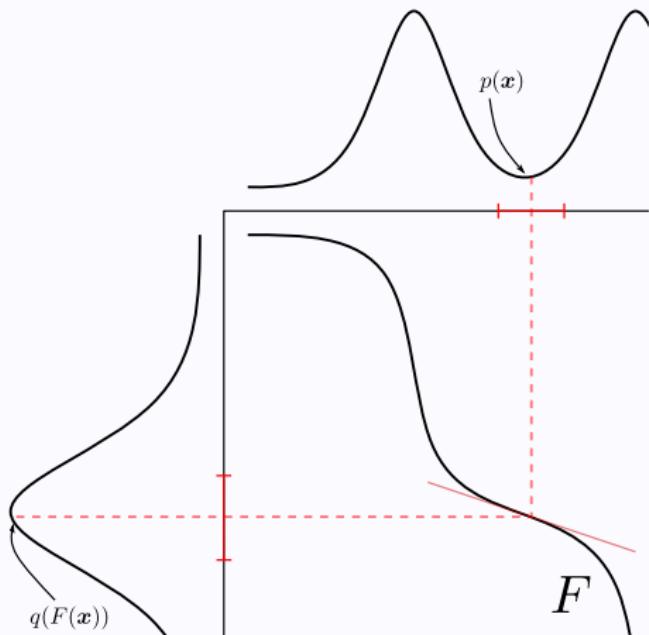


Figure: Example of 1D mapping

Density Estimation

To perform density estimation:

- ① Draw $\mathbf{x} \sim P^*$,
- ② Compute $F(\mathbf{x})$ and $|\det \text{Jac}_F(\mathbf{x})|$,
- ③ Compute $\hat{p}(\mathbf{x}) = q(F(\mathbf{x})) |\det \text{Jac}_F(\mathbf{x})|$.

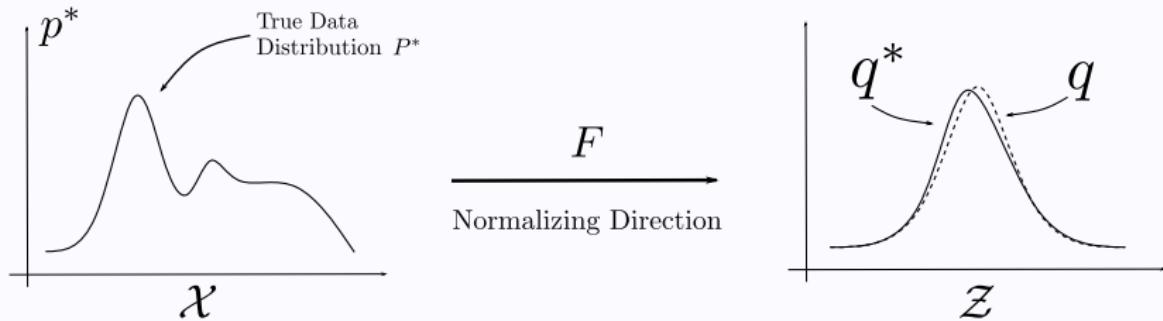


Figure: 1D Normalizing Process of Density Estimation.

Data Generation

To perform data generation:

- ① Draw $\mathbf{z} \sim Q$,
- ② Compute $\mathbf{x} = F^{-1}(\mathbf{z})$.

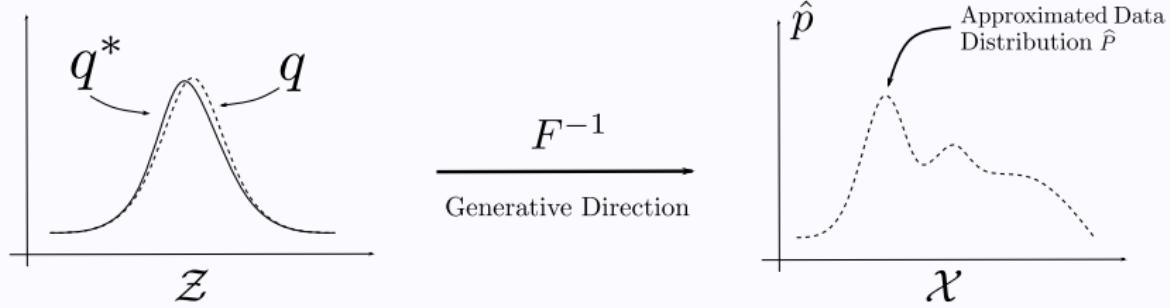


Figure: 1D Normalizing Flow process of Generation.

How is it trained ? - Generative models

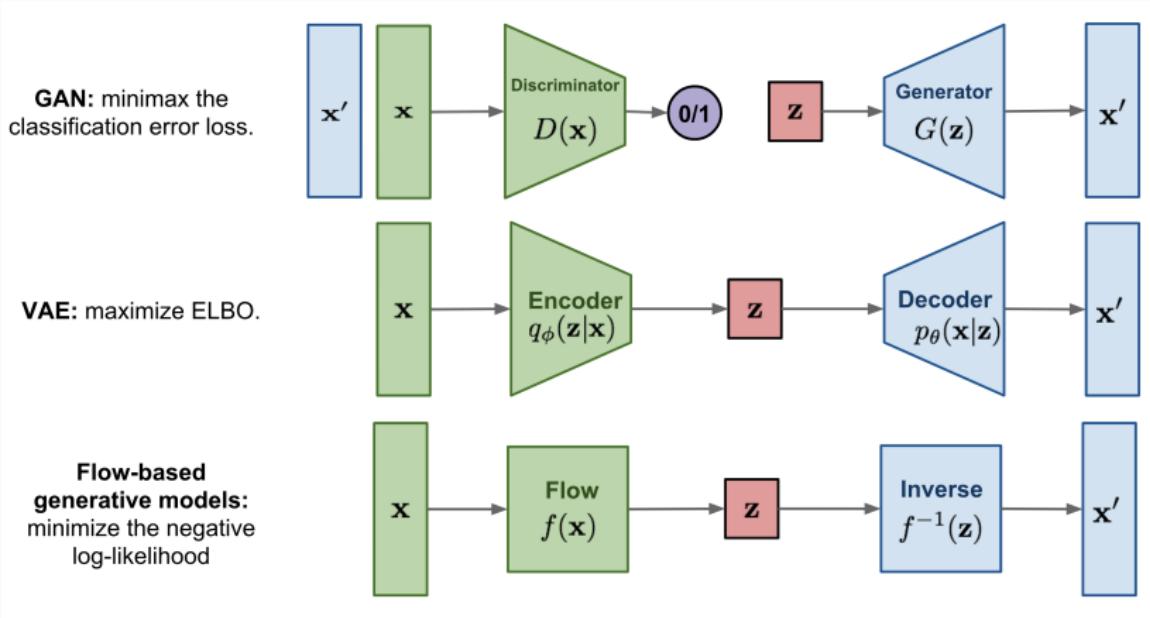


Figure: Different generative models [5]

Loss

The objective is to approximate P^* with \hat{P} . We can minimize the Kullback-Leiber Divergence :

$$\theta = \arg \min_{\theta} \mathcal{D}_{\text{KL}}(P^* \parallel \hat{P}).$$

This is equivalent to maximizing the log likelihood :

$$\theta = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{X}} [\log \hat{p}(x)].$$

$$\mathcal{D}_{\text{KL}}(P^* \parallel \hat{P}) = \int_{\mathcal{X}} p^*(\mathbf{x}) \log \left(\frac{p^*(\mathbf{x})}{\hat{p}(\mathbf{x})} \right) d\mathbf{x}$$

$$\text{nll} = -\mathbb{E}_{\mathbf{x} \sim \mathcal{X}} [\log \hat{p}(\mathbf{x})]$$

Learning steps

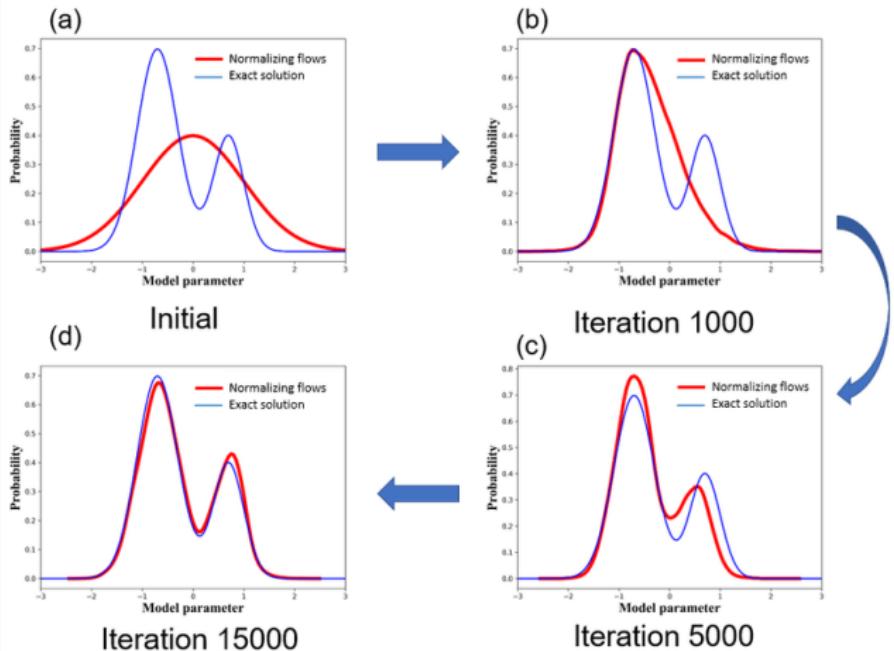


Figure: Learning Process for a 1D Normalizing Flow [6].

Figure: Learning Process for a Normalizing Flow for the MNIST dataset.

The entropy can be used as an intuitive understanding of the loss. We have :

$$H(P^*, \hat{P}) = D_{\text{KL}}(P^* \parallel \hat{P}) + H(P^*) \geq H(P^*).$$

In practice, the entropy per dimension is used to compare different models on a same dataset :

$$\text{bpd} = \frac{\sum_{\mathbf{x} \in \mathcal{X}} \log(\hat{p}(\mathbf{x})/2^K)}{d} = \frac{\text{nll}}{dK \log(2)}$$

where K is the number of different bits to encode the data. For instance, we usually have $256 = 2^8$ different values per pixel per channel, thus $K = 8$.

Normalizing Flow

A Normalizing Flow is a Neural Network for which we can efficiently compute :

- The forward pass F , ie the normalizing direction,
- the inverse pass F^{-1} , ie the generative direction,
- the determinant of the Jacobian matrix Jac_F for every \mathbf{x} .

A normalizing Flow is a Neural Network

Usually, a Neural Net is represented as the combination of atomic function $f_i : \mathbf{x} \mapsto \sigma(\mathbf{W}_i \mathbf{x} + \mathbf{b}_i)$ such that $F = f_n \circ \dots \circ f_0$.

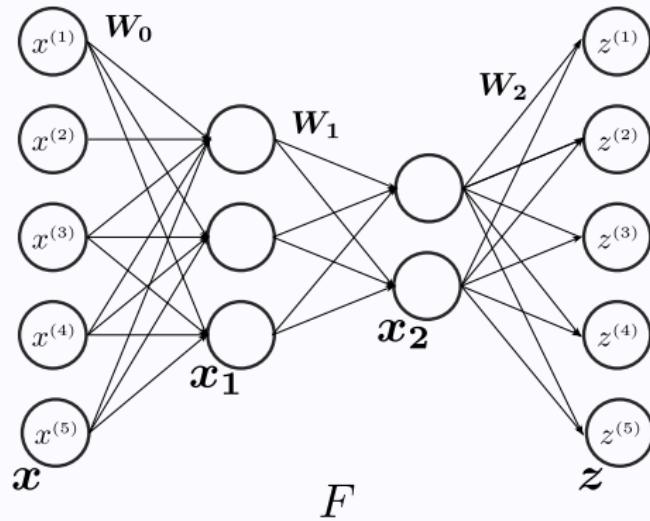


Figure: Representation of a Neural Network.

A normalizing Flow is a Neural Network

For a Normalizing Flow, every atomic block f_i satisfies the properties of F , i.e. bijective and the Jacobian must be efficiently computable.

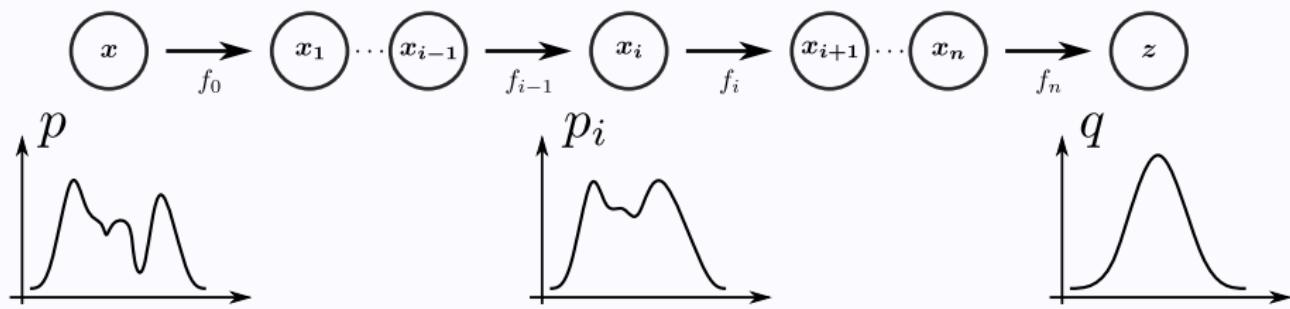


Figure: Atomic representation of the Normalizing Flow.

Therefore, we have:

$$\log \det \text{Jac}_F(\mathbf{x}) = \sum_{i=0}^n \log \det \text{Jac}_{f_i}(\mathbf{x}_i)$$

A normalizing Flow is a Neural Network

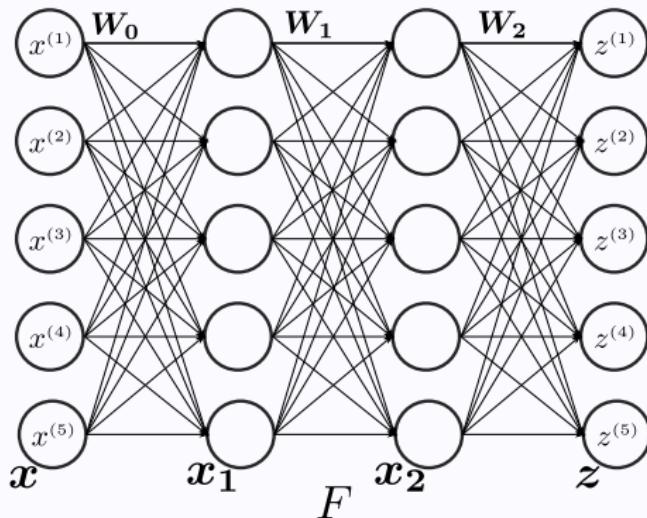


Figure: Representation of an Invertible Neural Network.

Elementwise Flow

The element wise flows are based on elementwise bijective non linear scalar function. Let $h_i : \mathbb{R} \mapsto \mathbb{R}$ be the scalar value bijection :

$$f_i(x_i^1, \dots, x_i^d) = \left(h\left(x_i^1\right), \dots, h\left(x_i^d\right) \right)^T$$

Determinant of the Jacobian matrix

Here, $\text{Jac}_{f_i}(\mathbf{x}) = \text{Diag}\left(\frac{dh}{dx}|_{x=x_1}, \dots, \frac{dh}{dx}|_{x=x_d}\right)$, so

$$|\det \text{Jac}_{f_i}(\mathbf{x})| = \prod_i^d \frac{dh}{dx}|_{x=x_i}.$$

Limits

- No mixing of variables [2].

Linear Flows

Let \mathbf{A}_i be an invertible matrix of $\mathbb{R}^{d \times d}$ and \mathbf{b}_i a vector of \mathbb{R}^d . A Linear Flow is based on the linear operation :

$$f_i(\mathbf{x}) = \mathbf{A}_i \mathbf{x} + \mathbf{b}_i$$

Determinant of the Jacobian matrix

Here, $\text{Jac}_{f_i}(\mathbf{x}) = \mathbf{A}_i$, so $|\det \text{Jac}_{f_i}(\mathbf{x})| = |\det \mathbf{A}_i| \neq 0$.

Limits

Poor Expressiveness.

Planar Flows

Let $h : \mathbb{R} \mapsto \mathbb{R}$, h' its derivative, \mathbf{u} and \mathbf{w} two vectors of \mathbb{R}^d and a scalar b . The Planar flow is based on atomic blocks f_i such that:

$$f_i(\mathbf{x}) = \mathbf{x} + \mathbf{u}h\left(\mathbf{w}^T \mathbf{x} + b\right)$$

Determinant of the Jacobian matrix

Here, $\text{Jac}_{f_i}(\mathbf{x}) = \mathbf{I}_d + \mathbf{u}h'(\mathbf{w}^T \mathbf{x} + b)\mathbf{w}^T$, so
 $|\det \text{Jac}_{f_i}(\mathbf{x})| = 1 + h'(\mathbf{w}^T \mathbf{x} + b)\mathbf{u}^T \mathbf{w}$.

Limits

No closed form of the inverse.

Radial Flows

Let \mathbf{x}_0 a vector of \mathbb{R}^d , α and β two scalars. The Radial Flow is based on atomic blocks f_i such that:

$$f_i(\mathbf{x}) = \mathbf{x} + \frac{\beta}{\alpha + \|\mathbf{x} - \mathbf{x}_0\|} (\mathbf{x} - \mathbf{x}_0)$$

Determinant of the Jacobian matrix

Here, $|\det \text{Jac}_{f_i}(\mathbf{x})| = \left[1 + \frac{\beta}{\alpha + \|\mathbf{x} - \mathbf{x}_0\|}\right]^{d-1} \left[1 + \frac{\beta}{\alpha + \|\mathbf{x} - \mathbf{x}_0\|} + \frac{\|\mathbf{x} - \mathbf{x}_0\|}{(\alpha + \|\mathbf{x} - \mathbf{x}_0\|)^2}\right]$.

Limits

No closed form of the inverse.

Radial and Planar Flows

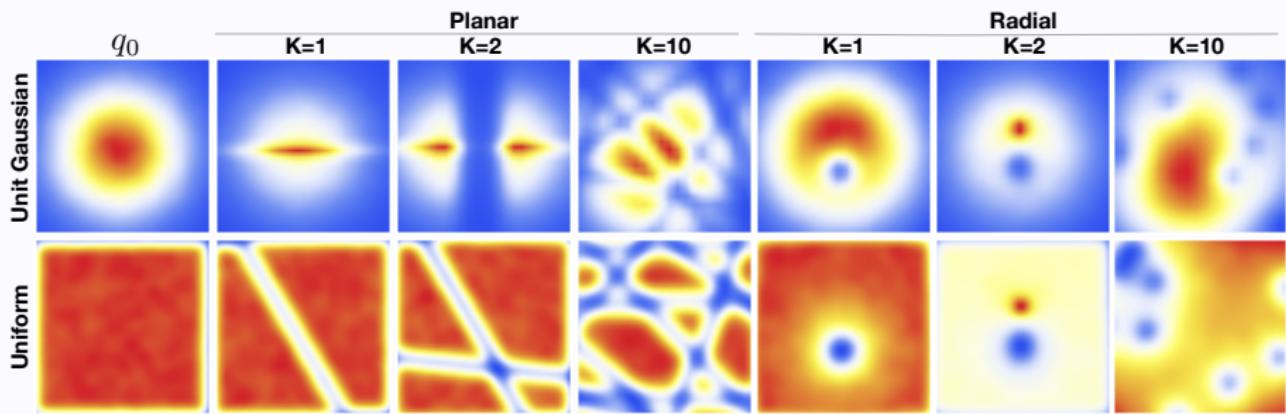


Figure: Example of Planar and Radial Flows [3].

Coupling Function

A coupling function is coupling disjoint partitions of the dataset. For an input $\mathbf{x} \in \mathbb{R}^d$, we consider a split $(\mathbf{x}^A, \mathbf{x}^B) \in \mathbb{R}^k \times \mathbb{R}^{d-k}$. Usually, $k = d/2$. We consider a bijective function, generally elementwise, $\mathbf{h} : \mathbf{x} \mapsto \mathbf{h}(\mathbf{x}, \theta) \in \mathbb{R}^k$ parameterized by θ . We can train a small neural network such that θ becomes a function of \mathbf{x}^A or \mathbf{x}^B . We can therefore define the atomic coupling function :

$$f_i\left(\begin{bmatrix} \mathbf{x}^A \\ \mathbf{x}^B \end{bmatrix}\right) = \begin{bmatrix} \mathbf{y}^A \\ \mathbf{y}^B \end{bmatrix} = \begin{bmatrix} \mathbf{h}(\mathbf{x}^A, g(\mathbf{x}^B)) \\ \mathbf{x}^B \end{bmatrix}.$$

The reverse is :

$$\mathbf{f}_i^{-1}\left(\begin{bmatrix} \mathbf{y}^A \\ \mathbf{y}^B \end{bmatrix}\right) = \begin{bmatrix} \mathbf{x}^A \\ \mathbf{x}^B \end{bmatrix} = \begin{bmatrix} \mathbf{h}^{-1}(\mathbf{y}^B, g(\mathbf{x}^B)) \\ \mathbf{x}^B \end{bmatrix}.$$

Coupling Function

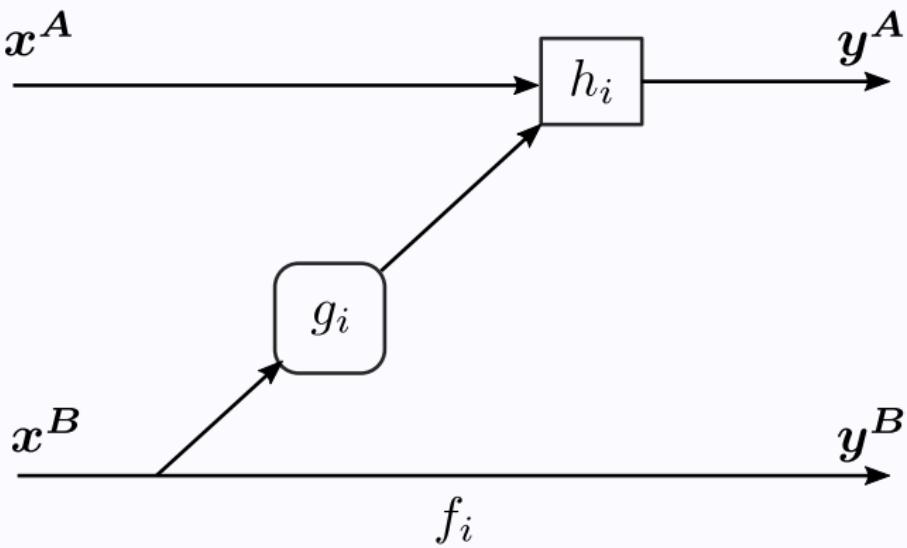


Figure: Atomic coupling block.

Coupling Function

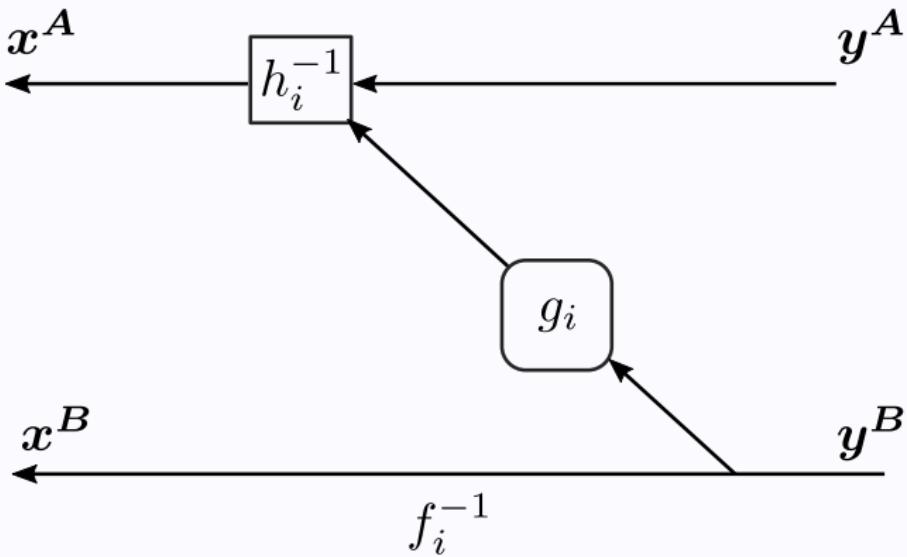


Figure: Inverse Atomic coupling block.

Coupling Function

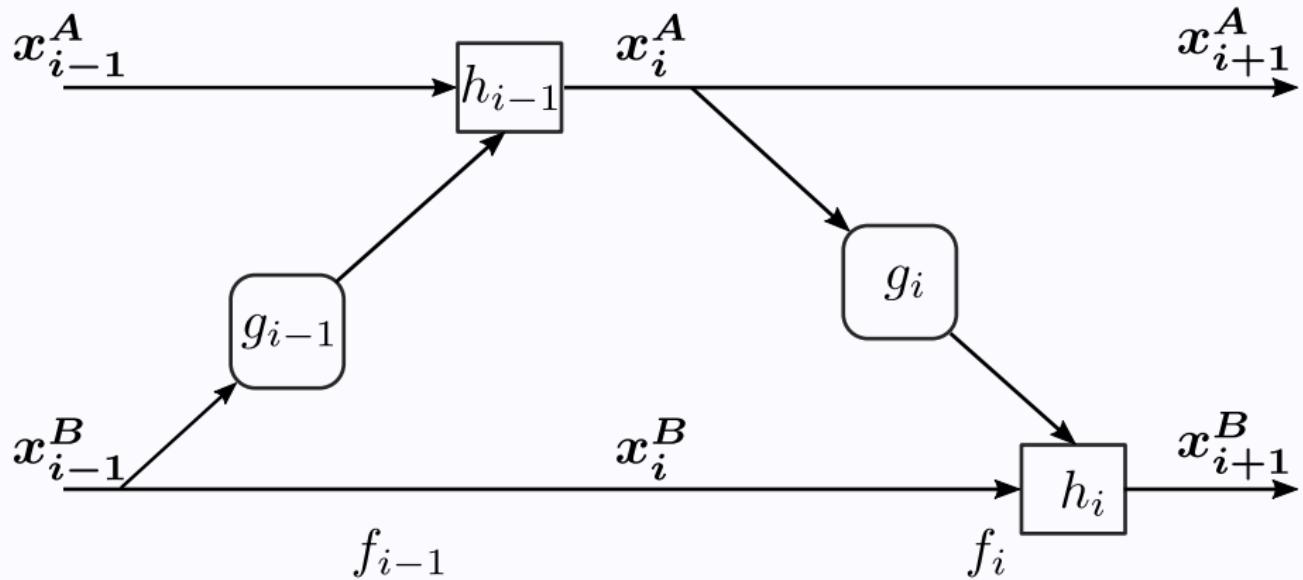


Figure: Composed Atomic coupling block.

Coupling Function

Derterminant of the Jacobian Matrix

Here, $\text{Jac}_{f_i}(\mathbf{x}) = \begin{pmatrix} \frac{\partial h}{\partial \mathbf{x}^A} |_{\mathbf{x}^A, g(\mathbf{x}^B)} & \frac{\partial h}{\partial \mathbf{x}^B} |_{\mathbf{x}^A, g(\mathbf{x}^B)} \\ \mathbf{0} & I_k \end{pmatrix}$, thus

$$|\det \text{Jac}_{f_i}(\mathbf{x})| = \left| \frac{\partial h}{\partial \mathbf{x}^A} (\mathbf{x}^A) \right|$$

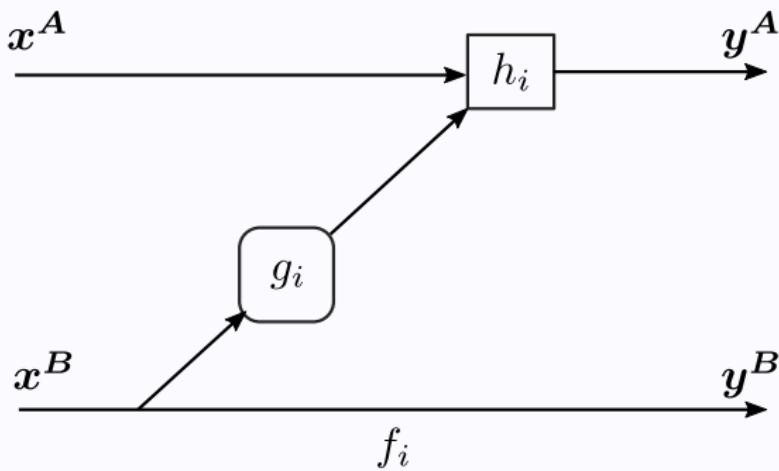


Figure: Atomic coupling block.

Additive Coupling Flow

The additive coupling flow has atomic function defined with an elementwise function $h_i : \mathbb{R} \mapsto \mathbb{R}$ and the scalar function $g_i : \mathbb{R}^d \mapsto \mathbb{R}$:

$$h_i(x_i) = x_i + g_i(\mathbf{x}^B).$$

Therefore,

$$f_i \left(\begin{bmatrix} \mathbf{x}^A \\ \mathbf{x}^B \end{bmatrix} \right) = \begin{bmatrix} \mathbf{x}_A + g(\mathbf{x}^B) \\ \mathbf{x}^B \end{bmatrix}.$$

Determinant of the Jacobian matrix

Here, $|\det \text{Jac}_{f_i}| = 1$. The flow is said to be volume preserving.

Limits

Very deep structure and poor expressiveness.

Affine Coupling Flow

The affine coupling flow has atomic function defined with an elementwise function $h_i : \mathbb{R} \mapsto \mathbb{R}$ and the scalar function $g_i : \mathbb{R}^d \mapsto \mathbb{R}$:

$$h_i(x_i) = g_i^1(\mathbf{x}^B)x_i + g_i^2(\mathbf{x}^B).$$

Therefore,

$$f_i \left(\begin{bmatrix} \mathbf{x}^A \\ \mathbf{x}^B \end{bmatrix} \right) = \begin{bmatrix} g^1(\mathbf{x}^B)\mathbf{x}_A + g^2(\mathbf{x}^B) \\ \mathbf{x}^B \end{bmatrix}.$$

Determinant of the Jacobian matrix

Here, $|\det \text{Jac}_{f_i}| = (g_i^1(\mathbf{x}^B))^k$.

Limits

Very deep structure.

Residual Flow

A residual flow is composed of residual atomic blocks f_i defined with function $\mathbf{g}_i : \mathbb{R}^d \mapsto \mathbb{R}^d$ such that:

$$f_i(\mathbf{x}) = \mathbf{x}_i + \mathbf{g}(\mathbf{x}_i).$$

In \mathbf{g} , we can stack Convolutional layers, batchnorms, activations, etc...

Determinant of the Jacobian Matrix

Here, the determinant of the log determinant of the Jacobian matrix is estimated by a Russian Roulette Estimator, and a Skilling-Hutchinson trace estimator [1]:

$$\log \det \text{Jac}_{f_i}(v\mathbf{x}) = \mathbb{E}_{n \sim \text{Geo}(p), v \sim \mathcal{N}(\mathbf{0}, I)} \left[\sum_{k=1}^n \frac{(-1)^{k+1}}{k} \frac{\mathbf{v}^T [\text{Jac}_g(\mathbf{x})^k \mathbf{v}]}{p(n \geq k)} \right]$$

Inverse

A residual Network is invertible only if $Lip(g_i) < 1$. We can compute the inverse with the Inverse-point algorithm:

Algorithm 1 Inverse via fixed-point iteration.

```
 $x_{i-1} \leftarrow x_i$ 
while NotConverged do
     $x_{i-1} \leftarrow x_i - g(x_{i-1})$ 
end while
```

Limits

Iterative inverse and approximated Jacobian.

Since Normalizing Flows are diffeomorphism, there are bi-Lipschitz :

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}, \quad \|F(\mathbf{x}_1) - F(\mathbf{x}_2)\|_2 \leq L_1 \|\mathbf{x}_1 - \mathbf{x}_2\|_2,$$

and

$$\forall \mathbf{z}_1, \mathbf{z}_2 \in \mathcal{Z}, \quad \|F^{-1}(\mathbf{z}_1) - F^{-1}(\mathbf{z}_2)\|_2 \leq L_2 \|\mathbf{z}_1 - \mathbf{z}_2\|_2.$$

Therefore, some pathological target distributions can be highlighted [4].

Expressivity of a Normalizing Flow : Dense Subset

Let F be L_1 -Lipschitz. Then:

$$\text{TV}(P^*, \hat{P}) \geq$$

$$\sup_{R, \mathbf{x}_0} \left(P^*(B_{R, \mathbf{x}_0}) - \frac{\gamma\left(\frac{d}{2}, \frac{L_1^2 R^2}{2}\right)}{\Gamma\left(\frac{d}{2}\right)} \right)$$

Therefore, if we find a ball for which the true measure satisfies

$P^*(B_{R, \mathbf{x}_0}) > \gamma\left(\frac{d}{2}, \frac{L_1^2 R^2}{2}\right)/\Gamma\left(\frac{d}{2}\right)$, then the TV is necessarily strictly positive.

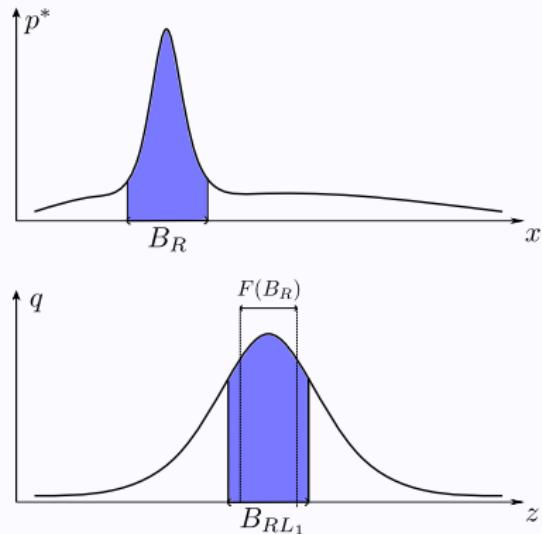


Figure: Example of a pathological target distribution: the subset B_R concentrates most of the weight in $P^*(B_R)$, but $\hat{P}(B_R) = Q(F(B_R))$ can only be as large as $Q(B_{RL_1})$.

Expressivity of a Normalizing Flow : Sparse Subset

Let F^{-1} be L_2 -Lipschitz. We consider the balls centered on $F^{-1}(0)$, we have the lower bound:

$$\text{TV}(P^*, \hat{P}) \geq \sup_R \left(\frac{\gamma\left(\frac{d}{2}, \frac{R^2}{2L_2^2}\right)}{\Gamma\left(\frac{d}{2}\right)} - P^*(B_{R,F^{-1}(0)}) \right)$$

Therefore, if we find a ball for which the true measure satisfies $P^*(B_{R,F^{-1}(0)}) < \frac{\gamma(d/2, R^2/2L_2^2)}{\Gamma(d/2)}$, then the TV is necessarily strictly positive.

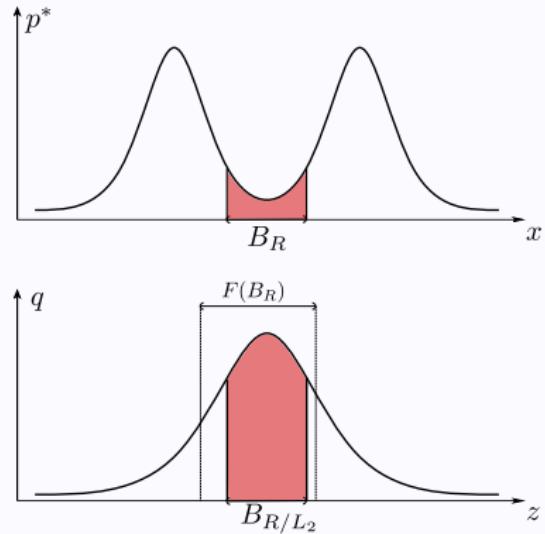


Figure: Example of a pathological target distribution: the subset B_R concentrates little weight in $P^*(B_R)$, but $\hat{P}(B_R) = Q(F(B_R))$ can only be as small as $Q(B_{R/L_2})$.

Normalizing Flow are a very powerfull tools to perform :

- Data generation,
- Density estimation.

Reversible networks have a larger spectrum of applications such as learning physics models, manifold learning...

However, they suffer from a expensive computational cost for training and inference. We will see that in the Practical session that can be found on : https://github.com/AlexVerine/AdvancedML_NF.

References I

- [1] Ricky T. Q. Chen, Jens Behrmann, David Duvenaud, and Jörn-Henrik Jacobsen. Residual Flows for Invertible Generative Modeling. In *33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*, Vancouver, Canada., July 2020. arXiv: 1906.02735.
- [2] Ivan Kobyzev, Simon J. D. Prince, and Marcus A. Brubaker. Normalizing Flows: An Introduction and Review of Current Methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020. arXiv: 1908.09257.
- [3] Danilo Jimenez Rezende and Shakir Mohamed. Variational Inference with Normalizing Flows. In *arXiv:1505.05770 [cs, stat]*, June 2016.
- [4] Alexandre Verine, Yann Chevaleyre, Fabrice Rossi, and Benjamin Negrevergne. On the expressivity of bi-Lipschitz normalizing flows. In *ICML Workshop on Invertible Neural Networks, Normalizing Flows, and Explicit Likelihood Models*, June 2021.

- [5] Lilian Weng. Flow-based deep generative models.
lilianweng.github.io/lil-log, 2018.
- [6] Xuebin Zhao, Andrew Curtis, and Xin Zhang. Bayesian Variational Seismic Tomography using Normalizing Flows. In *EGU General Assembly Conference Abstracts*, EGU General Assembly Conference Abstracts, pages EGU21-1455, April 2021.