

Large Scale ML: Generative models

Gabriel V Cardoso

`gabriel.victorino_cardoso@minesparis.psl.eu`

Equipe Géostatistiques, Ecole des Mines, Univ. PSL, France



What this lesson is about

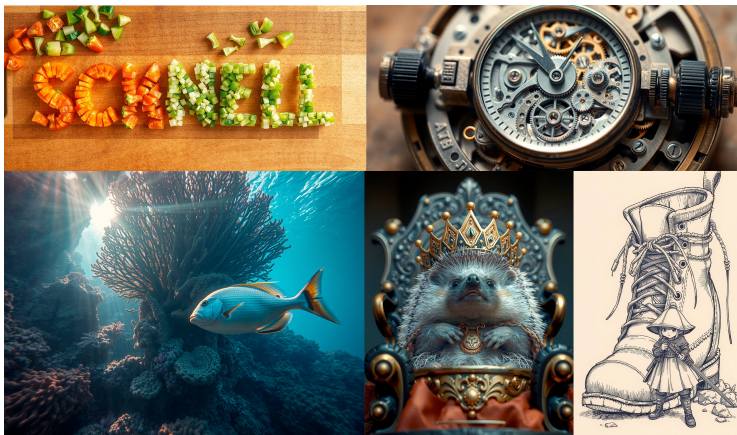


Figure: Images generated by several different prompts using the FLUX model.¹

¹<https://huggingface.co/black-forest-labs/FLUX.1-schnell>

The probabilistic point of view

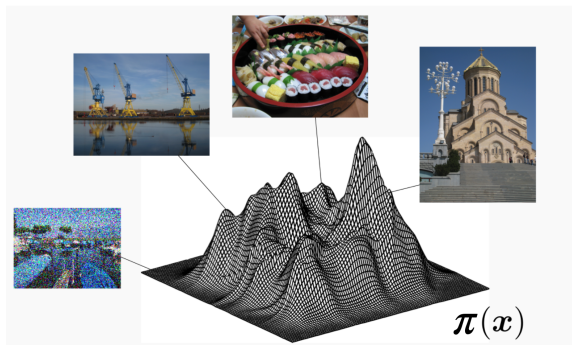


Figure: Landscape of the underlying **data** distribution ².

²Source:

https://www.charles-deledalle.fr/pages/files/ucsd_ece285_mlip/

- ▶ π lives in a **high-dimensional space** $d_x > 10^5$.
- ▶ We have access to big datasets $\mathcal{D} = \{X_1, \dots, X_n\}$, $n > 10^5$, of samples of π .
- ▶ We have no **parametric model** for our data.³

Problem: For general π , density estimation in high-dimensions requires an exponential amount of data.

³or saying it differently, we can't code a program that generates real images from a set of fixed parameters.

- ▶ π lives in a **high-dimensional space** $d_x > 10^5$.
- ▶ We have access to big datasets $\mathcal{D} = \{X_1, \dots, X_n\}$, $n > 10^5$, of samples of π .
- ▶ We have no **parametric model** for our data.³

Problem: For general π , density estimation in high-dimensions requires an exponential amount of data.

³or saying it differently, we can't code a program that generates real images from a set of fixed parameters.

The original miracle!



Figure: $X_0 \sim \pi$.



Figure: $X_\sigma = X_0 + \sigma Z$.⁴¹

- ▶ D_{θ_0} untrained CNN,
- ▶ $L_\sigma(\theta) = \|D_\theta(X_\sigma) - X_\sigma\|^2$, $L_0(\theta) = \|D_\theta(X_\sigma) - X_0\|^2$,
- ▶ $\theta_t = \theta_{t-1} - \nabla L_\sigma(\theta_{t-1})$

What happens to $\{L_\sigma(\theta_t)\}_{t=1}^N$ and $\{L_0(\theta_t)\}_{t=1}^N$?

¹ $Z \sim \mathcal{N}(0, I)$, $X_0 \perp Z$.

The original miracle!



Figure: $X_0 \sim \pi$.



Figure: $X_\sigma = X_0 + \sigma Z$.⁴¹

- ▶ D_{θ_0} untrained CNN,
- ▶ $L_\sigma(\theta) = \|D_\theta(X_\sigma) - X_\sigma\|^2$, $L_0(\theta) = \|D_\theta(X_\sigma) - X_0\|^2$,
- ▶ $\theta_t = \theta_{t-1} - \nabla L_\sigma(\theta_{t-1})$

What happens to $\{L_\sigma(\theta_t)\}_{t=1}^N$ and $\{L_0(\theta_t)\}_{t=1}^N$?

¹ $Z \sim \mathcal{N}(0, I)$, $X_0 \perp Z$.

The original miracle!



Figure: $X_0 \sim \pi$.



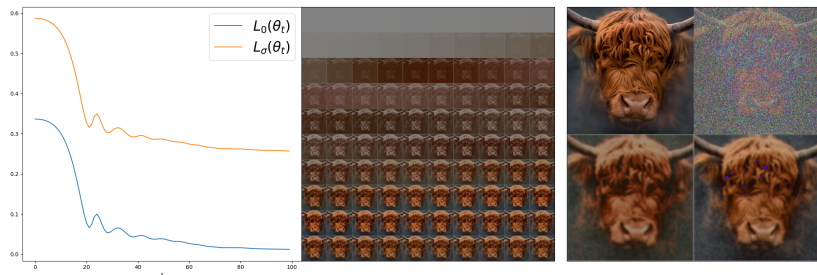
Figure: $X_\sigma = X_0 + \sigma Z$.⁴¹

- ▶ D_{θ_0} untrained CNN,
- ▶ $L_\sigma(\theta) = \|D_\theta(X_\sigma) - X_\sigma\|^2$, $L_0(\theta) = \|D_\theta(X_\sigma) - X_0\|^2$,
- ▶ $\theta_t = \theta_{t-1} - \nabla L_\sigma(\theta_{t-1})$

What happens to $\{L_\sigma(\theta_t)\}_{t=1}^N$ and $\{L_0(\theta_t)\}_{t=1}^N$?

¹ $Z \sim \mathcal{N}(0, I)$, $X_0 \perp Z$.

The original miracle ⁵.

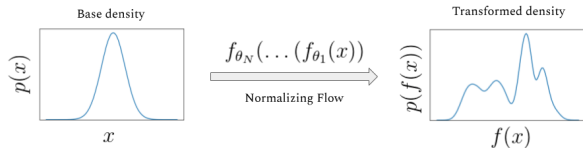


⁵This experiment is inspired by Ulyanov, Vedaldi, and Lempitsky, 2018

- 1 A bit of theory**
- 2 Normalizing flows
- 3 VAE
- 4 GANs
- 5 Score-based generative models

- 1 Define a **family of distributions** $\mathcal{F} = \{p_\theta \in \mathcal{P}(\mathbb{R}^{d_x}) | \theta \in \Theta\}$. Typically this is done by

- 1 Defining a **neural network** architecture f_θ parameterized by θ
- 2 Define a **reference distribution** λ^6
- 3 $p_\theta := \text{Law}(f_\theta(Z))$ where $Z \sim \lambda$.



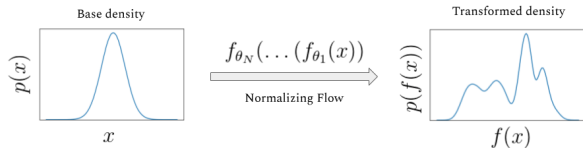
- 2 Choose a **divergence**⁷ over distributions $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_*$
- 3 Solve

$$\theta^* \in \arg \min_{\theta \in \Theta} D(\pi, p_\theta)$$

⁶Typically, a Gaussian or Uniform.

⁷Only property we need is $D(p, q) = 0 \iff p = q$.

- 1 Define a **family of distributions** $\mathcal{F} = \{p_\theta \in \mathcal{P}(\mathbb{R}^{d_x}) | \theta \in \Theta\}$. Typically this is done by
 - 1 Defining a **neural network** architecture f_θ parameterized by θ
 - 2 Define a **reference distribution** λ^6
 - 3 $p_\theta := \text{Law}(f_\theta(Z))$ where $Z \sim \lambda$.



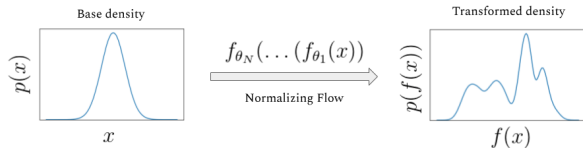
- 2 Choose a **divergence**⁷ over distributions $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_*$
- 3 Solve

$$\theta^* \in \arg \min_{\theta \in \Theta} D(\pi, p_\theta)$$

⁶Typically, a Gaussian or Uniform.

⁷Only property we need is $D(p, q) = 0 \iff p = q$.

- 1 Define a **family of distributions** $\mathcal{F} = \{p_\theta \in \mathcal{P}(\mathbb{R}^{d_x}) | \theta \in \Theta\}$. Typically this is done by
 - 1 Defining a **neural network** architecture f_θ parameterized by θ
 - 2 Define a **reference distribution** λ^6
 - 3 $p_\theta := \text{Law}(f_\theta(Z))$ where $Z \sim \lambda$.



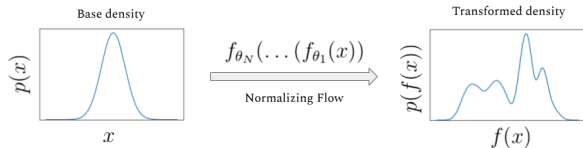
- 2 Choose a **divergence**⁷ over distributions $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_*$
- 3 Solve

$$\theta^* \in \arg \min_{\theta \in \Theta} D(\pi, p_\theta)$$

⁶Typically, a Gaussian or Uniform.

⁷Only property we need is $D(p, q) = 0 \iff p = q$.

- 1 Define a **family of distributions** $\mathcal{F} = \{p_\theta \in \mathcal{P}(\mathbb{R}^{d_x}) | \theta \in \Theta\}$. Typically this is done by
 - 1 Defining a **neural network** architecture f_θ parameterized by θ
 - 2 Define a **reference distribution** λ ⁶
 - 3 $p_\theta := \text{Law}(f_\theta(Z))$ where $Z \sim \lambda$.



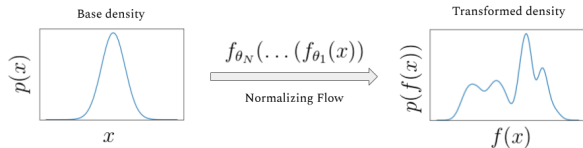
- 2 Choose a **divergence**⁷ over distributions $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_*$
- 3 Solve

$$\theta^* \in \arg \min_{\theta \in \Theta} D(\pi, p_\theta)$$

⁶Typically, a Gaussian or Uniform.

⁷Only property we need is $D(p, q) = 0 \iff p = q$.

- 1 Define a **family of distributions** $\mathcal{F} = \{p_\theta \in \mathcal{P}(\mathbb{R}^{d_x}) | \theta \in \Theta\}$. Typically this is done by
 - 1 Defining a **neural network** architecture f_θ parameterized by θ
 - 2 Define a **reference distribution** λ ⁶
 - 3 $p_\theta := \text{Law}(f_\theta(Z))$ where $Z \sim \lambda$.



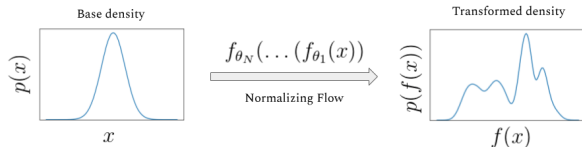
- 2 Choose a **divergence**⁷ over distributions $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_*$
- 3 Solve

$$\theta^* \in \arg \min_{\theta \in \Theta} D(\pi, p_\theta)$$

⁶Typically, a Gaussian or Uniform.

⁷Only property we need is $D(p, q) = 0 \iff p = q$.

- 1 Define a **family of distributions** $\mathcal{F} = \{p_\theta \in \mathcal{P}(\mathbb{R}^{d_x}) | \theta \in \Theta\}$. Typically this is done by
 - 1 Defining a **neural network** architecture f_θ parameterized by θ
 - 2 Define a **reference distribution** λ ⁶
 - 3 $p_\theta := \text{Law}(f_\theta(Z))$ where $Z \sim \lambda$.



- 2 Choose a **divergence**⁷ over distributions $D : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_*$
- 3 Solve

$$\theta^* \in \arg \min_{\theta \in \Theta} D(\pi, p_\theta)$$

⁶Typically, a Gaussian or Uniform.

⁷Only property we need is $D(p, q) = 0 \iff p = q$.

Kullback Leibler

Let $(p, q) \in \mathcal{P}^2$ such that $p \ll q$

$$D_{\text{KL}}(p||q) := \mathbb{E}_p \left[\log \frac{\partial p(X)}{\partial q(X)} \right]. \quad (1)$$

If they both admit a density with respect to a given measure μ then

$$D_{\text{KL}}(p||q) = \int \log \frac{p(x)}{q(x)} p(x) \mu(dx). \quad (2)$$

Some properties

1 $D_{\text{KL}}(p||q) = 0 \iff p = q$

2 $D_{\text{KL}}(p||q) \geq 0$

3 $D_{\text{KL}}(\mathcal{N}(\mu_0, \Sigma_0)||\mathcal{N}(\mu_1, \Sigma_1)) =$
 $\frac{1}{2} \left[\log \frac{|\Sigma_0|}{|\Sigma_1|} - d + \|\mu_1 - \mu_0\|_{\Sigma_1^{-1}}^2 + \text{tr}(\Sigma_1^{-1}\Sigma_0) \right].^8$

⁸ $\|v\|_A^2 = v^T A v.$

Kullback Leibler

Let $(p, q) \in \mathcal{P}^2$ such that $p \ll q$

$$D_{\text{KL}}(p||q) := \mathbb{E}_p \left[\log \frac{\partial p(X)}{\partial q(X)} \right]. \quad (1)$$

If they both admit a density with respect to a given measure μ then

$$D_{\text{KL}}(p||q) = \int \log \frac{p(x)}{q(x)} p(x) \mu(dx). \quad (2)$$

Some properties

- 1 $D_{\text{KL}}(p||q) = 0 \iff p = q$
- 2 $D_{\text{KL}}(p||q) \geq 0$
- 3 $D_{\text{KL}}(\mathcal{N}(\mu_0, \Sigma_0)||\mathcal{N}(\mu_1, \Sigma_1)) = \frac{1}{2} \left[\log \frac{|\Sigma_0|}{|\Sigma_1|} - d + \|\mu_1 - \mu_0\|_{\Sigma_1^{-1}}^2 + \text{tr}(\Sigma_1^{-1}\Sigma_0) \right].^8$

⁸ $\|v\|_A^2 = v^t A v.$

$$D_{\text{KL}}(p||q) = \int \log \frac{p(x)}{q(x)} p(x) \mu(dx) = \mathbb{E}_p [\log p(X)] - \mathbb{E}_p [\log q(X)] \quad (3)$$

Thus,

$$\theta^* \in \arg \min_{\theta \in \Theta} D_{\text{KL}}(\pi || p_{\theta}) \iff \theta^* \in \arg \max_{\theta \in \Theta} \mathbb{E}_p [\log p_{\theta}(X)] \quad (4)$$

Minimizing KL \iff doing Maximum likelihood.

So, maybe we can train our generative model by a Monte Carlo approximation of $\mathbb{E}_p [\log p_{\theta}(X)]$:

$$\arg \max_{\theta \in \Theta} n^{-1} \sum_{i=1}^n \log p_{\theta}(X_i) \quad (5)$$

$$D_{\text{KL}}(p||q) = \int \log \frac{p(x)}{q(x)} p(x) \mu(dx) = \mathbb{E}_p [\log p(X)] - \mathbb{E}_p [\log q(X)] \quad . \quad (3)$$

Thus,

$$\theta^* \in \arg \min_{\theta \in \Theta} D_{\text{KL}}(\pi||p_\theta) \iff \theta^* \in \arg \max_{\theta \in \Theta} \mathbb{E}_p [\log p_\theta(X)] \quad (4)$$

Minimizing KL \iff doing Maximum likelihood.

So, maybe we can train our generative model by a Monte Carlo approximation of $\mathbb{E}_p [\log p_\theta(X)]$:

$$\arg \max_{\theta \in \Theta} n^{-1} \sum_{i=1}^n \log p_\theta(X_i) \quad (5)$$

$$D_{\text{KL}}(p||q) = \int \log \frac{p(x)}{q(x)} p(x) \mu(dx) = \mathbb{E}_p [\log p(X)] - \mathbb{E}_p [\log q(X)] \quad (3)$$

Thus,

$$\theta^* \in \arg \min_{\theta \in \Theta} D_{\text{KL}}(\pi||p_\theta) \iff \theta^* \in \arg \max_{\theta \in \Theta} \mathbb{E}_p [\log p_\theta(X)] \quad (4)$$

Minimizing KL \iff doing Maximum likelihood.

So, maybe we can train our generative model by a Monte Carlo approximation of $\mathbb{E}_p [\log p_\theta(X)]$:

$$\arg \max_{\theta \in \Theta} n^{-1} \sum_{i=1}^n \log p_\theta(X_i) \quad (5)$$

$$D_{\text{KL}}(p||q) = \int \log \frac{p(x)}{q(x)} p(x) \mu(dx) = \mathbb{E}_p [\log p(X)] - \mathbb{E}_p [\log q(X)] \quad . \quad (3)$$

Thus,

$$\theta^* \in \arg \min_{\theta \in \Theta} D_{\text{KL}}(\pi||p_\theta) \iff \theta^* \in \arg \max_{\theta \in \Theta} \mathbb{E}_p [\log p_\theta(X)] \quad (4)$$

Minimizing KL \iff doing Maximum likelihood.

So, maybe we can train our generative model by a **Monte Carlo approximation** of $\mathbb{E}_p [\log p_\theta(X)]$:

$$\arg \max_{\theta \in \Theta} n^{-1} \sum_{i=1}^n \log p_\theta(X_i) \quad (5)$$

Can you calculate $\log p_{\theta}$?

- 1 A bit of theory
- 2 Normalizing flows**
- 3 VAE
- 4 GANs
- 5 Score-based generative models

Suppose $f_\theta : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_x}$ is invertible. Then,

$$p_\theta(x) = \lambda(f_\theta^{-1}(x)) |J_{f_\theta}(f_\theta^{-1}(x))|^{-1} . \quad (6)$$

Normalizing flow (Rezende and Mohamed, 2015)

A Normalizing flow is a Neural net

$$f_\theta(x) = \mathcal{T}_{\theta_k} \circ \cdots \circ \mathcal{T}_{\theta_0} , \quad (7)$$

where for all $\ell = 0, \dots, k$, $\mathcal{T}_{\theta_\ell}$ is invertible and has an easy to calculate $|J_{\mathcal{T}_{\theta_\ell}}|!$

Suppose $f_\theta : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_x}$ is invertible. Then,

$$p_\theta(x) = \lambda(f_\theta^{-1}(x)) |J_{f_\theta}(f_\theta^{-1}(x))|^{-1}. \quad (6)$$

Normalizing flow (Rezende and Mohamed, 2015)

A Normalizing flow is a Neural net

$$f_\theta(x) = \mathcal{T}_{\theta_k} \circ \cdots \circ \mathcal{T}_{\theta_0}, \quad (7)$$

where for all $\ell = 0, \dots, k$, $\mathcal{T}_{\theta_\ell}$ is invertible and has an easy to calculate $|J_{\mathcal{T}_{\theta_\ell}}|!$

Algorithm 1 Log likelihood with NF

```
1: procedure LOGLIK
  Input:  $x, \theta$ .
  Output:  $p_\theta(x)$ .
2:    $L = 0, z = x$ .
3:   for  $\ell = k, \dots, 0$  do
4:      $z = \mathcal{T}_{\theta_\ell}^{-1}(z)$ .
5:      $L = L - \log |\mathcal{T}_{\theta_\ell}(z)|$ .
6:   end for
7:    $L = L + \log \lambda(z)$ .
8:   Return:  $L$ 
9: end procedure
```

It is possible to

$$\arg \max_{\theta \in \Theta} n^{-1} \sum_{i=1}^n \log p_\theta(X_i) \quad (8)$$

Algorithm 2 Log likelihood with NF

```
1: procedure LOGLIK
  Input:  $x, \theta$ .
  Output:  $p_\theta(x)$ .
2:    $L = 0, z = x$ .
3:   for  $\ell = k, \dots, 0$  do
4:      $z = \mathcal{T}_{\theta_\ell}^{-1}(z)$ .
5:      $L = L - \log |\mathcal{T}_{\theta_\ell}(z)|$ .
6:   end for
7:    $L = L + \log \lambda(z)$ .
8:   Return:  $L$ 
9: end procedure
```

It is possible to

$$\arg \max_{\theta \in \Theta} n^{-1} \sum_{i=1}^n \log p_\theta(X_i) \quad (8)$$

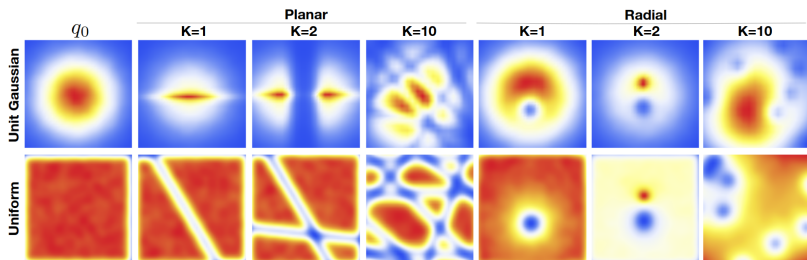


Figure: Illustration from Rezende and Mohamed, 2015 show the transformations obtained during a flow.

RealNVP (Dinh, Sohl-Dickstein, and Bengio, 2017)

$$\mathcal{T}_\theta(x[1:d], x[d:d_x]) = (x[1:d], x[d:d_x] \cdot e^{s(x[1:d];\theta)} + t(x[1:d];\theta)) , \quad (9)$$

where $s(\cdot; \theta), t(\cdot; \theta) : \mathbb{R}^d \rightarrow \mathbb{R}^{d_x-d}$ are typically neural networks.
The Jacobian of such flow is

$$J_{\mathcal{T}_\theta}(x) = \begin{bmatrix} Id & 0 \\ \dots & \text{diag}(e^{s(x[1:d];\theta)}) \end{bmatrix} , \quad (10)$$

and thus $\log |J_{\mathcal{T}_\theta}(x)| = \sum_{j=1}^{d_x-d} s(x[1:d];\theta)[j]$.

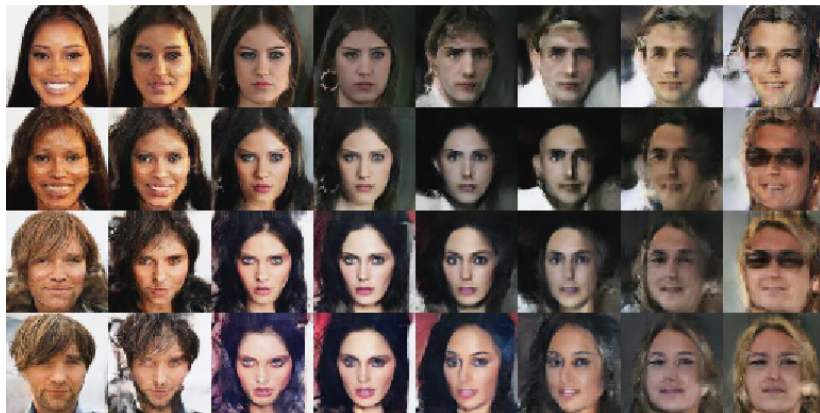


Figure: Figure taken from Dinh, Sohl-Dickstein, and Bengio, 2017 illustrating generated samples from RealNVP on the Celeb dataset.

Problems with Normalizing flow

- 1 Considerable **restriction on network architectures**.
- 2 Manifold hypothesis: p_θ admits a density with respect to the ambient Lebesgue measure!

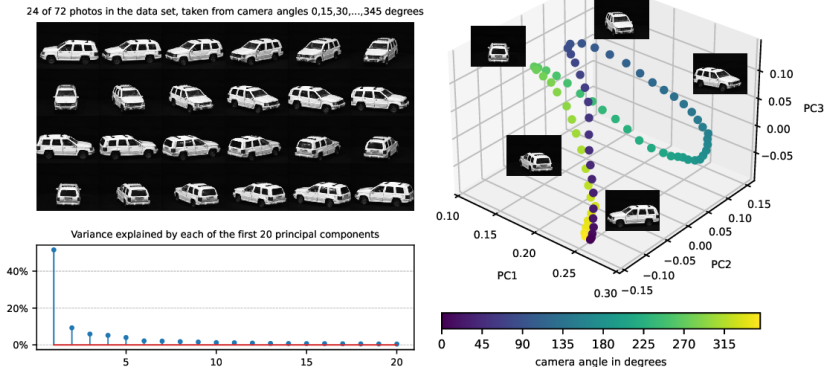


Figure: Taken from (Whiteley, Gray, and Rubin-Delanchy, 2024)

- 1 A bit of theory
- 2 Normalizing flows
- 3 VAE**
- 4 GANs
- 5 Score-based generative models

Context: We think our data is **lower dimensional**.

Context: We think our data is **lower dimensional**.

- 1 Define a latent distribution $\lambda \in \mathcal{P}(\mathbb{R}^{d_z})$,
 $d_z \ll d_x$.

Context: We think our data is **lower dimensional**.

- 1 Define a latent distribution $\lambda \in \mathcal{P}(\mathbb{R}^{d_z})$,
 $d_z \ll d_x$.
- 2 Define a **conditional distribution** over the data space:

$$z \in \mathbb{R}^{d_z} \rightarrow p_{\theta}(\cdot|z) \in \mathcal{P}(\mathbb{R}^{d_x}) .$$

1

¹Taken from Diederik P Kingma, Welling, et al., 2019.

Context: We think our data is **lower dimensional**.

- 1 Define a latent distribution $\lambda \in \mathcal{P}(\mathbb{R}^{d_z})$,
 $d_z \ll d_x$.
- 2 Define a **conditional distribution** over the data space:

$$z \in \mathbb{R}^{d_z} \rightarrow p_{\theta}(\cdot|z) \in \mathcal{P}(\mathbb{R}^{d_x}).$$

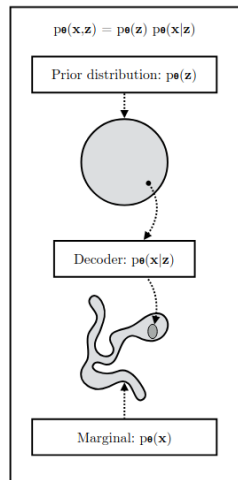


Figure: Illustration of the latent model ¹.

¹

¹Taken from Diederik P Kingma, Welling, et al., 2019.

Context: We think our data is **lower dimensional**.

- 1 Define a latent distribution $\lambda \in \mathcal{P}(\mathbb{R}^{d_z})$,
 $d_z \ll d_x$.
- 2 Define a **conditional distribution** over the data space:

$$z \in \mathbb{R}^{d_z} \rightarrow p_\theta(\cdot|z) \in \mathcal{P}(\mathbb{R}^{d_x}).$$

Gaussian Latent model

We can consider

- 1 $\lambda = \mathcal{N}(0, I)$,
- 2 $p_\theta(\cdot|z) = \mathcal{N}(\mu_\theta(z), \sigma^2 I)$.

For $\mu_\theta(z) = Az$ we recover the (Probabilistic) PCA.

In general, μ_θ will be a **Neural network**.

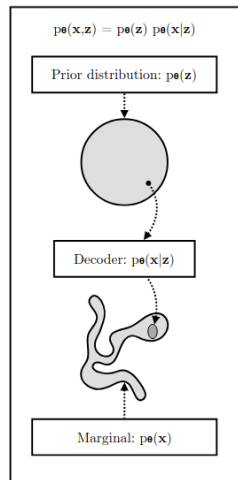


Figure: Illustration of the latent model ¹.

¹

¹Taken from Diederik P Kingma, Welling, et al., 2019.

Can you calculate $\log p_\theta$?

For the VAE,

$$\log p_\theta(x) = \log \int p_\theta(x|z) \lambda(dz) . \quad (11)$$

The answer is: Yes, but *only for linear μ_θ* !

Can you calculate $\log p_\theta$?

For the VAE,

$$\log p_\theta(x) = \log \int p_\theta(x|z) \lambda(dz) . \quad (11)$$

The answer is: Yes, but **only for linear μ_θ !**

Can you estimate
 $\log p_{\theta}(x) = \log \int p_{\theta}(x|z) \lambda(dz) ?$

Naive MC:

- 1 Sample $Z_{1:K} \sim \lambda$
- 2 Approximate $p_{\theta}(x)$ by $K^{-1} \sum_{i=1}^K p_{\theta}(x|Z_i)$.

Problem: This estimation has high variance for realistic cases!

Can you estimate

$$\log p_{\theta}(x) = \log \int p_{\theta}(x|z) \lambda(dz) ?$$

Naive MC:

- 1 Sample $Z_{1:K} \sim \lambda$
- 2 Approximate $p_{\theta}(x)$ by $K^{-1} \sum_{i=1}^K p_{\theta}(x|Z_i)$.

Problem: This estimation has high variance for realistic cases!

Can you estimate
 $\log p_{\theta}(x) = \log \int p_{\theta}(x|z) \lambda(dz) ?$

Naive MC:

- 1 Sample $Z_{1:K} \sim \lambda$
- 2 Approximate $p_{\theta}(x)$ by $K^{-1} \sum_{i=1}^K p_{\theta}(x|Z_i)$.

Problem: This estimation has high variance for realistic cases!

Introduce a **parametric family of instrumental distribution** $q_{\phi}(\cdot|x) \in \mathcal{P}(\mathbb{R}^{d_z})$ to reduce MC variance.

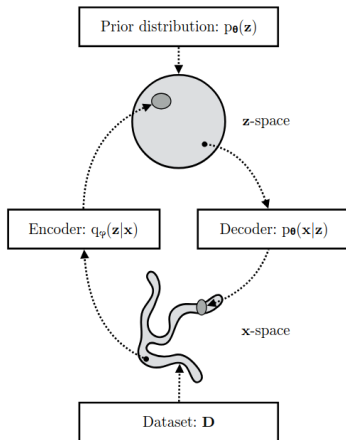


Figure: Overall view of the VAE.¹

¹Taken from Diederik P Kingma, Welling, et al., 2019.

Introduce a parametric family of instrumental distribution $q_\phi(\cdot|x) \in \mathcal{P}(\mathbb{R}^{d_z})$ to reduce MC variance.

Make it match the true (intractable) posterior:

$$p_{\theta}(z|x) = \frac{p_{\theta}(x|z) \lambda(z)}{p_{\theta}(x)}. \quad (12)$$

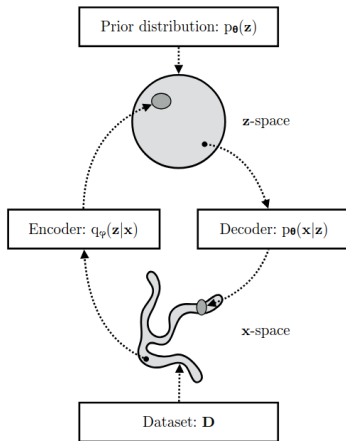


Figure: Overall view of the VAE.¹

¹Taken from Diederik P Kingma, Welling, et al., 2019.

$$D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x)) = \int \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} q_{\phi}(z|x) \, dz$$

$$\begin{aligned} D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x)) &= \int \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} q_{\phi}(z|x) dz \\ &= \int \log \frac{q_{\phi}(z|x) p_{\theta}(x)}{p_{\theta}(x|z) \lambda(z)} q_{\phi}(z|x) dz \end{aligned}$$

$$\begin{aligned} D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x)) &= \int \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} q_{\phi}(z|x) \, dz \\ &= \int \log \frac{q_{\phi}(z|x) p_{\theta}(x)}{p_{\theta}(x|z) \lambda(z)} q_{\phi}(z|x) \, dz \\ &= \log p_{\theta}(x) + \int \log \frac{q_{\phi}(z|x)}{p_{\theta}(x|z) \lambda(z)} q_{\phi}(z|x) \, dz \end{aligned}$$

$$\begin{aligned} D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x)) &= \int \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} q_{\phi}(z|x) dz \\ &= \int \log \frac{q_{\phi}(z|x) p_{\theta}(x)}{p_{\theta}(x|z) \lambda(z)} q_{\phi}(z|x) dz \\ &= \log p_{\theta}(x) + \int \log \frac{q_{\phi}(z|x)}{p_{\theta}(x|z) \lambda(z)} q_{\phi}(z|x) dz \\ &= \log p_{\theta}(x) + \int \log \frac{q_{\phi}(z|x)}{\lambda(z)} q_{\phi}(z|x) dz - \int \log p_{\theta}(x|z) q_{\phi}(z|x) dz \end{aligned}$$

$$\begin{aligned}
 D_{\text{KL}}(q_\phi(\cdot|x) || p_\theta(\cdot|x)) &= \int \log \frac{q_\phi(z|x)}{p_\theta(z|x)} q_\phi(z|x) dz \\
 &= \int \log \frac{q_\phi(z|x) p_\theta(x)}{p_\theta(x|z) \lambda(z)} q_\phi(z|x) dz \\
 &= \log p_\theta(x) + \int \log \frac{q_\phi(z|x)}{p_\theta(x|z) \lambda(z)} q_\phi(z|x) dz \\
 &= \log p_\theta(x) + \int \log \frac{q_\phi(z|x)}{\lambda(z)} q_\phi(z|x) dz - \int \log p_\theta(x|z) q_\phi(z|x) dz \\
 &= \log p_\theta(x) + D_{\text{KL}}(q_\phi(\cdot|x) || \lambda) - \mathbb{E}_{q_\phi(\cdot|x)} [\log p_\theta(x|Z)] .
 \end{aligned}$$

$$\begin{aligned} D_{\text{KL}}(q_\phi(\cdot|x) || p_\theta(\cdot|x)) &= \int \log \frac{q_\phi(z|x)}{p_\theta(z|x)} q_\phi(z|x) dz \\ &= \int \log \frac{q_\phi(z|x) p_\theta(x)}{p_\theta(x|z) \lambda(z)} q_\phi(z|x) dz \\ &= \log p_\theta(x) + \int \log \frac{q_\phi(z|x)}{p_\theta(x|z) \lambda(z)} q_\phi(z|x) dz \\ &= \log p_\theta(x) + \int \log \frac{q_\phi(z|x)}{\lambda(z)} q_\phi(z|x) dz - \int \log p_\theta(x|z) q_\phi(z|x) dz \\ &= \log p_\theta(x) + D_{\text{KL}}(q_\phi(\cdot|x) || \lambda) - \mathbb{E}_{q_\phi(\cdot|x)} [\log p_\theta(x|Z)] . \end{aligned}$$

Evidence lower bound (ELBO)

$$\begin{aligned} m(\theta, \phi; x) &:= \log p_\theta(x) - D_{\text{KL}}(q_\phi(\cdot|x) || p_\theta(\cdot|x)) \\ &= \mathbb{E}_{q_\phi(\cdot|x)} [\log p_\theta(x|Z)] - D_{\text{KL}}(q_\phi(\cdot|x) || \lambda) . \end{aligned}$$

Evidence lower bound (ELBO)

$$\begin{aligned} m(\theta, \phi; x) &:= \log p_{\theta}(x) - D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x)) \\ &= \mathbb{E}_{q_{\phi}(\cdot|x)} [\log p_{\theta}(x|Z)] - D_{\text{KL}}(q_{\phi}(\cdot|x) || \lambda) . \end{aligned}$$

Evidence lower bound (ELBO)

$$\begin{aligned} m(\theta, \phi; x) &:= \log p_{\theta}(x) - D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x)) \\ &= \mathbb{E}_{q_{\phi}(\cdot|x)} [\log p_{\theta}(x|Z)] - D_{\text{KL}}(q_{\phi}(\cdot|x) || \lambda) . \end{aligned}$$

1 Maximize $m(\theta, \phi; x)$ in $\phi \Rightarrow$ Minimize $D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x))$.

Evidence lower bound (ELBO)

$$\begin{aligned} m(\theta, \phi; x) &:= \log p_{\theta}(x) - D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x)) \\ &= \mathbb{E}_{q_{\phi}(\cdot|x)} [\log p_{\theta}(x|Z)] - D_{\text{KL}}(q_{\phi}(\cdot|x) || \lambda) . \end{aligned}$$

- 1 Maximize $m(\theta, \phi; x)$ in $\phi \Rightarrow$ Minimize $D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x))$.
- 2 Maximize $m(\theta, \phi; x)$ in $\theta \Rightarrow$ Maximize $\log p_{\theta}(x)$.

Evidence lower bound (ELBO)

$$\begin{aligned} m(\theta, \phi; x) &:= \log p_{\theta}(x) - D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x)) \\ &= \mathbb{E}_{q_{\phi}(\cdot|x)} [\log p_{\theta}(x|Z)] - D_{\text{KL}}(q_{\phi}(\cdot|x) || \lambda) . \end{aligned}$$

- 1 Maximize $m(\theta, \phi; x)$ in $\phi \Rightarrow$ Minimize $D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x))$.
- 2 Maximize $m(\theta, \phi; x)$ in $\theta \Rightarrow$ Maximize $\log p_{\theta}(x)$.

We know $\log p_{\theta}(x) - D_{\text{KL}}(q_{\phi}(\cdot|x) || p_{\theta}(\cdot|x))$ is intractable. Is $\mathbb{E}_{q_{\phi}(\cdot|x)} [\log p_{\theta}(x|Z)] - D_{\text{KL}}(q_{\phi}(\cdot|x) || \lambda)$ tractable?

$$\mathbb{E}_{q_{\phi}(\cdot|x)} [\log p_{\theta}(x|Z)] - D_{\text{KL}}(q_{\phi}(\cdot|x) \parallel \lambda)$$

$$\mathbb{E}_{q_\phi(\cdot|x)} [\log p_\theta(x|Z)] - D_{\text{KL}}(q_\phi(\cdot|x) \parallel \lambda)$$

Gaussian Gaussian Latent model

Consider

- 1 $\lambda = \mathcal{N}(0, \text{I})$,
- 2 $p_\theta(\cdot|z) = \mathcal{N}(\mu_\theta(z), \sigma^2 \text{I})$,
- 3 $q_\phi(\cdot|x) = \mathcal{N}(\mu_\phi(x), \Sigma_\phi(x))$.

Here, we consider $\Sigma_\phi(x) = \text{diag}(\sigma_\phi^2(x))$ where $\sigma_\phi^2(x) \in \mathbb{R}^{d_z}$.

$$\mathbb{E}_{q_{\phi}(\cdot|x)} [\log p_{\theta}(x|Z)] - D_{\text{KL}}(q_{\phi}(\cdot|x) \parallel \lambda)$$

Gaussian Gaussian Latent model

Consider

- 1 $\lambda = \mathcal{N}(0, \text{I}),$
- 2 $p_{\theta}(\cdot|z) = \mathcal{N}(\mu_{\theta}(z), \sigma^2 \text{I}),$
- 3 $q_{\phi}(\cdot|x) = \mathcal{N}(\mu_{\phi}(x), \Sigma_{\phi}(x)).$

Here, we consider $\Sigma_{\phi}(x) = \text{diag}(\sigma_{\phi}^2(x))$ where $\sigma_{\phi}^2(x) \in \mathbb{R}^{d_z}$.

$$D_{\text{KL}}(q_{\phi}(\cdot|x) \parallel \lambda) = \frac{1}{2} [\|\mu_{\phi}(x)\|^2 + \text{tr}(\Sigma_{\phi}(x)) + \log |\Sigma_{\phi}(x)| - d_z] \quad (13)$$

$$\mathbb{E}_{q_\phi(\cdot|x)} [\log p_\theta(x|Z)] - D_{\text{KL}}(q_\phi(\cdot|x) \parallel \lambda)$$

Gaussian Gaussian Latent model

Consider

- 1 $\lambda = \mathcal{N}(0, \text{I}),$
- 2 $p_\theta(\cdot|z) = \mathcal{N}(\mu_\theta(z), \sigma^2 \text{I}),$
- 3 $q_\phi(\cdot|x) = \mathcal{N}(\mu_\phi(x), \Sigma_\phi(x)).$

Here, we consider $\Sigma_\phi(x) = \text{diag}(\sigma_\phi^2(x))$ where $\sigma_\phi^2(x) \in \mathbb{R}^{d_z}$.

$$D_{\text{KL}}(q_\phi(\cdot|x) \parallel \lambda) = \frac{1}{2} [\|\mu_\phi(x)\|^2 + \text{tr}(\Sigma_\phi(x)) + \log |\Sigma_\phi(x)| - d_z] \quad (13)$$

$$\mathbb{E}_{q_\phi(\cdot|x)} [\log p_\theta(x|Z)] = \mathbb{E}_{q_\phi(\cdot|x)} \left[-\frac{1}{2\sigma^2} \|x - \mu_\theta(Z)\|^2 \right] + \text{cte.} \quad (14)$$

The reparametrization trick for $\mathbb{E}_{q_\phi(\cdot|x)} \left[-\frac{1}{2\sigma^2} \|x - \mu_\theta(Z)\|^2 \right]$

$$q_\phi(\cdot|x) = \mathcal{N}(\mu_\phi(x), \text{diag}(\sigma_\phi^2(x))).$$

One way of sampling:

- 1 Sample $\epsilon \sim \mathcal{N}(0, I)$
- 2 Set $Z = \mu_\phi(x) + \sigma_\phi(x) \odot \epsilon$.

Thus,

$$\mathbb{E}_{q_\phi(\cdot|x)} \left[-\frac{1}{2\sigma^2} \|x - \mu_\theta(Z)\|^2 \right] = \mathbb{E}_{\mathcal{N}(0, I)} \left[-\frac{1}{2\sigma^2} \|x - \mu_\theta(\mu_\phi(x) + \sigma_\phi(x) \odot \epsilon)\|^2 \right]$$

We can approximate with MC

$$\begin{aligned} & \mathbb{E}_{\mathcal{N}(0, I)} \left[-\frac{1}{2\sigma^2} \|x - \mu_\theta(\mu_\phi(x) + \sigma_\phi(x) \odot \epsilon)\|^2 \right] \\ & \approx K^{-1} \sum_{k=1}^K -\frac{1}{2\sigma^2} \|x - \mu_\theta(\mu_\phi(x) + \sigma_\phi(x) \odot \epsilon_k)\|^2 \end{aligned}$$

where $\epsilon_k \sim \mathcal{N}(0, I)$.

$$\mathbb{E}_{q_\phi(\cdot|x)} [\log p_\theta(x|Z)] - D_{\text{KL}}(q_\phi(\cdot|x) \parallel \lambda)$$

Gaussian Gaussian Latent model

Consider

- 1 $\lambda = \mathcal{N}(0, \text{I})$,
- 2 $p_\theta(\cdot|z) = \mathcal{N}(\mu_\theta(z), \sigma^2 \text{I})$,
- 3 $q_\phi(\cdot|x) = \mathcal{N}(\mu_\phi(x), \Sigma_\phi(x))$.

Here, we consider $\Sigma_\phi(x) = \text{diag}(\sigma_\phi^2(x))$ where $\sigma_\phi^2(x) \in \mathbb{R}^{d_z}$.

$$D_{\text{KL}}(q_\phi(\cdot|x) \parallel \lambda) = \frac{1}{2} [\|\mu_\phi(x)\|^2 + \text{tr}(\Sigma_\phi(x)) + \log |\Sigma_\phi(x)| - d_z] \quad (15)$$

$$\mathbb{E}_{q_\phi(\cdot|x)} [\log p_\theta(x|Z)] \approx K^{-1} \sum_{k=1}^K -\frac{1}{2\sigma^2} \|x - \mu_\theta(\mu_\phi(x) + \sigma_\phi(x) \odot \epsilon_k)\|^2 + \text{cte}. \quad (16)$$

We can calculate gradients!

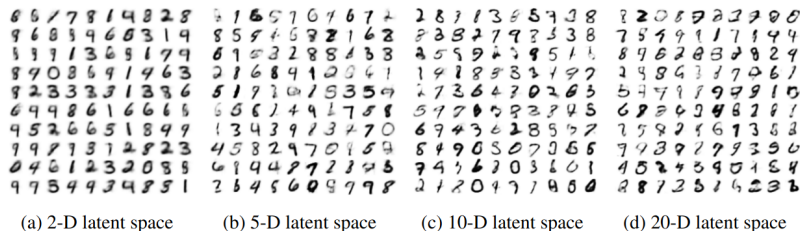


Figure: MNIST with different latent dimension.

- 1 VAE are good encoders! (Encode + Decode generalizes well!) ✓
- 2 Sampling from the prior **fails** for high dimensional latent models! ✗

$$Z \sim \lambda, \quad X \sim p_{\theta}(\cdot|Z)$$

⁸Taken from Diederik P. Kingma and Welling, 2013.

- 1 A bit of theory
- 2 Normalizing flows
- 3 VAE
- 4 GANs**
- 5 Score-based generative models

Riesz representation theorem

Every positive linear functional L in $C_c(\mathbb{R}^{d_x})$ can be represented by $L(f) = \mathbb{E}_\mu[f(X)]$ for some $\mu \in \mathcal{P}$.

Maybe, for some divergences, we can **bypass calculating log densities** by some duality theorem!

For us a duality theorem **for a given divergence D** allows to evaluate the divergence D between p_0 and p_1 **by just considering their effect on functions in a given class!**

Riesz representation theorem

Every positive linear functional L in $C_c(\mathbb{R}^{d_x})$ can be represented by $L(f) = \mathbb{E}_\mu[f(X)]$ for some $\mu \in \mathcal{P}$.

Maybe, for some divergences, we can **bypass calculating log densities** by some duality theorem!

For us a duality theorem **for a given divergence D** allows to evaluate the divergence D between p_0 and p_1 **by just considering their effect on functions in a given class!**

Riesz representation theorem

Every positive linear functional L in $C_c(\mathbb{R}^{d_x})$ can be represented by $L(f) = \mathbb{E}_\mu[f(X)]$ for some $\mu \in \mathcal{P}$.

Maybe, for some divergences, we can **bypass calculating log densities** by some duality theorem!

For us a duality theorem **for a given divergence D** allows to evaluate the divergence D between p_0 and p_1 **by just considering their effect on functions in a given class!**

Theorem (Donsker and Varadhan, 1983)

Let $(q, p) \in \mathcal{P}$, such that $p \ll q$ and $q \ll p$. Then,

$$D_{\text{KL}}(q||p) = \sup_{f \in \mathcal{G}_q} \mathbb{E}_q[f] - \log \mathbb{E}_p[\exp(f)] , \quad (17)$$

where \mathcal{G}_q is the set of measurable g such that $\mathbb{E}_q[\exp(g)] < \infty$.
Furthermore, equality hold if and only if $f = \log \frac{dq}{dp}$.

Theorem (Jensen Shannon duality Goodfellow et al., 2014)

Let $(q, p) \in \mathcal{P}$, such that $p \ll q$ and $q \ll p$. Define the *Jensen-Shannon entropy* as,

$$D_{\text{JS}}(q||p) = \frac{1}{2} (D_{\text{KL}}(q||p) + D_{\text{KL}}(p||q)) . \quad (18)$$

Then,

$$2 D_{\text{JS}}(q||p) = \sup_{D \in \mathcal{C}} \mathbb{E}_q [\log D(X)] + \mathbb{E}_p [\log(1 - D(X))] + \log(4) , \quad (19)$$

where \mathcal{C} is the set of measurable functions from $\mathbb{R}^{d_x} \rightarrow [0, 1]$.

Minimize $\sup_{D \in \mathcal{C}\ell} \mathbb{E}_q [\log D(X)] + \mathbb{E}_p [\log(1 - D(X))] \Leftrightarrow$ Minimize $D_{JS}(q||p)$.

The GAN framework

- Approximate $\mathcal{C}\ell$ by the set of classifiers defined by a Neural network architecture $\mathcal{D}_\theta := \{D_\theta : \mathbb{R}^d \rightarrow [0, 1]\}$.
- Generative model: $Z \in \mathbb{R}^d \sim \lambda$, apply a Neural net $X = f_\theta(Z)$, and p_θ the law of X .
- Train it in two steps:

Minimize $\sup_{D \in \mathcal{C}\ell} \mathbb{E}_q [\log D(X)] + \mathbb{E}_p [\log(1 - D(X))] \Leftrightarrow$ Minimize $D_{JS}(q||p)$.

The GAN framework

- 1 Approximate $\mathcal{C}\ell$ by the set of classifiers defined by a **Neural network architecture** $D_\phi = \{D_\phi : \mathbb{R}^{d_x} \rightarrow [0, 1]\}$.
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:

Minimize $\sup_{D \in \mathcal{C}\ell} \mathbb{E}_q [\log D(X)] + \mathbb{E}_p [\log(1 - D(X))] \Leftrightarrow$ Minimize $D_{JS}(q||p)$.

The GAN framework

- 1 Approximate $\mathcal{C}\ell$ by the set of classifiers defined by a **Neural network architecture** $D_\phi = \{D_\phi : \mathbb{R}^{d_x} \rightarrow [0, 1]\}$.
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:

For fixed θ , set

$$\phi = \arg \max_{\phi} \mathbb{E}_\pi [\log D_\phi(X)] + \mathbb{E}_\lambda [\log(1 - D_\phi(f_\theta(Z)))]$$

Then maximize

$$\mathbb{E}_\pi [\log D_\phi(X)] + \mathbb{E}_\lambda [\log(1 - D_\phi(f_\theta(Z)))]$$

Minimize $\sup_{D \in \mathcal{C}} \mathbb{E}_q [\log D(X)] + \mathbb{E}_p [\log(1 - D(X))] \Leftrightarrow$ Minimize $D_{\text{JS}}(q||p)$.

The GAN framework

- 1 Approximate \mathcal{C} by the set of classifiers defined by a **Neural network architecture** $D_\phi = \{D_\phi : \mathbb{R}^{d_x} \rightarrow [0, 1]\}$.
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:

1 For **fixed** θ , set

$$\phi = \arg \max_{\phi} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$$

2 For **fixed** ϕ , set

$$\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$$

Minimize $\sup_{D \in \mathcal{C}} \mathbb{E}_q [\log D(X)] + \mathbb{E}_p [\log(1 - D(X))] \Leftrightarrow$ Minimize $D_{\text{JS}}(q||p)$.

The GAN framework

- 1 Approximate \mathcal{C} by the set of classifiers defined by a **Neural network architecture** $\mathcal{D}_\phi = \{D_\phi : \mathbb{R}^{d_x} \rightarrow [0, 1]\}$.
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:
 - 1 For **fixed** θ , set

$$\phi = \arg \max_{\phi} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$$

- 2 For **fixed** ϕ , set

$$\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$$

Minimize $\sup_{D \in \mathcal{C}} \mathbb{E}_q [\log D(X)] + \mathbb{E}_p [\log(1 - D(X))] \Leftrightarrow$ Minimize $D_{JS}(q||p)$.

The GAN framework

- 1 Approximate \mathcal{C} by the set of classifiers defined by a **Neural network architecture** $D_\phi = \{D_\phi : \mathbb{R}^{d_x} \rightarrow [0, 1]\}$.
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:
 - 1 For **fixed** θ , set

$$\phi = \arg \max_{\phi} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$$

- 2 For **fixed** ϕ , set

$$\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$$

- 1** $\phi = \arg \max_{\phi} \overbrace{\mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]}^{-\text{Cross-entropy}}$
This is equivalent to training a classifier to distinguish between real data $X_r \sim \pi$ from fake data $X_f \sim p_{\theta}$.
- 2** $\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$ We are now training the generator to fool the classifier by minimizing the cross entropy.

⁸ $p_{\theta} = \text{Law}(f_{\theta}(Z))$ with $Z \sim \lambda$.

- 1 $\phi = \arg \max_{\phi} \overbrace{\mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]}^{-\text{Cross-entropy}}$
This is equivalent to training a classifier to distinguish between real data $X_r \sim \pi$ from fake data $X_f \sim p_{\theta}$.
- 2 $\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$ We are now training the generator to fool the classifier by minimizing the cross entropy.

⁸ $p_{\theta} = \text{Law}(f_{\theta}(Z))$ with $Z \sim \lambda$.

- Cross — entropy
- 1 $\phi = \arg \max_{\phi} \overbrace{\mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]}^{-\text{Cross-entropy}}$
This is equivalent to training a classifier to distinguish between real data $X_r \sim \pi$ from fake data $X_f \sim p_{\theta}$.
- 2 $\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$ We are now training the generator to "fool" the classifier by augmenting the cross entropy.

⁸ $p_{\theta} = \text{Law}(f_{\theta}(Z))$ with $Z \sim \lambda$.

- Cross — entropy
- 1 $\phi = \arg \max_{\phi} \overbrace{\mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]}^{-\text{Cross-entropy}}$
This is equivalent to training a classifier to distinguish between real data $X_r \sim \pi$ from fake data $X_f \sim p_{\theta}$.
- 2 $\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$ We are now training the generator to "fool" the classifier by augmenting the cross entropy.

⁸ $p_{\theta} = \text{Law}(f_{\theta}(Z))$ with $Z \sim \lambda$.

- Cross — entropy
- 1 $\phi = \arg \max_{\phi} \overbrace{\mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]}^{-\text{Cross-entropy}}$
This is equivalent to training a classifier to distinguish between real data $X_r \sim \pi$ from fake data $X_f \sim p_{\theta}$.
- 2 $\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]$ We are now training the generator to "fool" the classifier by augmenting the cross entropy.

⁸ $p_{\theta} = \text{Law}(f_{\theta}(Z))$ with $Z \sim \lambda$.

$$\begin{aligned} \text{1 } \phi &= \arg \max_{\phi} \overbrace{\mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))]}^{-\text{Cross-entropy}} \\ \text{2 } \theta &= \arg \min_{\theta} \mathbb{E}_{\pi} [\log D_{\phi}(X)] + \mathbb{E}_{\lambda} [\log(1 - D_{\phi}(f_{\theta}(Z)))] \end{aligned}$$

"In the proposed adversarial nets framework, the generative model is pitted against an adversary: a discriminative model that learns to determine whether a sample is from the model distribution or the data distribution. The generative model can be thought of as analogous to a team of counterfeiters, trying to produce fake currency and use it without detection, while the discriminative model is analogous to the police, trying to detect the counterfeit currency. Competition in this game drives both teams to improve their methods until the counterfeits are indistinguishable from the genuine articles."⁹

⁹from Goodfellow et al., 2014

- 1 No restriction on architecture ✓
- 2 No restriction on dimension of the latent space ✓
- 3 Training requires two networks ✗
- 4 Adversarial loss → unstable training ✗
- 5 Mode collapse ✗

- 1 No restriction on architecture ✓
- 2 No restriction on dimension of the latent space ✓
- 3 Training requires two networks ✗
- 4 Adversarial loss → unstable training ✗
- 5 Mode collapse ✗

- 1 No restriction on architecture ✓
- 2 No restriction on dimension of the latent space ✓
- 3 Training requires two networks ✗
- 4 Adversarial loss → unstable training ✗
- 5 Mode collapse ✗

- 1 No restriction on architecture ✓
- 2 No restriction on dimension of the latent space ✓
- 3 Training requires two networks ✗
- 4 Adversarial loss → unstable training ✗
- 5 Mode collapse ✗

- 1 No restriction on architecture ✓
- 2 No restriction on dimension of the latent space ✓
- 3 Training requires two networks ✗
- 4 Adversarial loss \rightarrow unstable training ✗
- 5 Mode collapse ✗

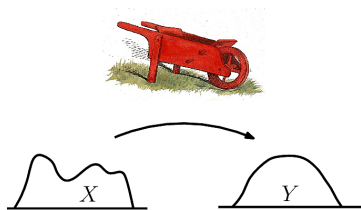


Figure: The earth mover problem

Wasserstein distance

Let $(q, p) \in (\mathcal{P}(\mathbb{R}^{d_x}))^2$, define the set of coupling distributions

$$C(q, p) = \{v \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_x}) \mid q(A) = v(A \times \mathbb{R}^{d_x}), p(A) = v(\mathbb{R}^{d_x} \times A)\}.$$

The Wasserstein distance is defined as

$$\mathcal{W}(q, p) = \inf_{v \in C(q, p)} \int \|x - y\| v(dx, dy) = \inf_{v \in C(q, p)} \mathbb{E}_{(X_0, X_1) \sim v} [\|X_0 - X_1\|]. \quad (20)$$

⁹Taken from <https://sbl.inria.fr/>.

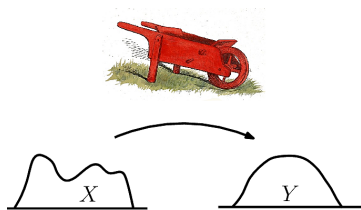


Figure: The earth mover problem

Wasserstein distance

Let $(q, p) \in (\mathcal{P}(\mathbb{R}^{d_x}))^2$, define the set of **coupling distributions**

$$C(q, p) = \{v \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_x}) \mid q(A) = v(A \times \mathbb{R}^{d_x}), p(A) = v(\mathbb{R}^{d_x} \times A)\}.$$

The **Wasserstein distance** is defined as

$$\mathcal{W}(q, p) = \inf_{v \in C(q, p)} \int \|x - y\| v(dx, dy) = \inf_{v \in C(q, p)} \mathbb{E}_{(X_0, X_1) \sim v} [\|X_0 - X_1\|]. \quad (20)$$

⁹Taken from <https://sbl.inria.fr/>.

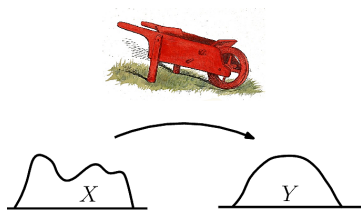


Figure: The earth mover problem

Wasserstein distance

Let $(q, p) \in (\mathcal{P}(\mathbb{R}^{d_x}))^2$, define the set of **coupling distributions**

$$C(q, p) = \{v \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_x}) \mid q(A) = v(A \times \mathbb{R}^{d_x}), p(A) = v(\mathbb{R}^{d_x} \times A)\}.$$

The **Wasserstein distance** is defined as

$$\mathcal{W}(q, p) = \inf_{v \in C(q, p)} \int \|x - y\| v(dx, dy) = \inf_{v \in C(q, p)} \mathbb{E}_{(X_0, X_1) \sim v} [\|X_0 - X_1\|]. \quad (20)$$

⁹Taken from <https://sbl.inria.fr/>.

Theorem (Kantorovich-Rubinstein theorem (see (Villani, 2021)))

$$\mathcal{W}(q, p) = \sup_{f \in \text{Lip}_1} \mathbb{E}_q[f] - \mathbb{E}_p[f] . \quad (21)$$

where $\text{Lip}_1 := \{f : \mathbb{R}^{d_x} \rightarrow \mathbb{R} \mid |f(x) - f(y)| \leq \|x - y\|\}$.

In Arjovsky, Chintala, and Bottou, 2017 proposes changing $D_{\text{JS}}(\pi \parallel p_\theta)$ by $\mathcal{W}(\pi, p_\theta)$ in the training of GANs.

Theorem (Kantorovich-Rubinstein theorem (see (Villani, 2021)))

$$\mathcal{W}(q, p) = \sup_{f \in \text{Lip}_1} \mathbb{E}_q[f] - \mathbb{E}_p[f] . \quad (21)$$

where $\text{Lip}_1 := \{f : \mathbb{R}^{d_x} \rightarrow \mathbb{R} \mid |f(x) - f(y)| \leq \|x - y\|\}$.

In Arjovsky, Chintala, and Bottou, 2017 proposes changing $D_{\text{JS}}(\pi \parallel p_\theta)$ by $\mathcal{W}(\pi, p_\theta)$ in the training of GANs.

$$\mathcal{W}(q, p) = \sup_{f \in \text{Lip}_1} \mathbb{E}_q[f] - \mathbb{E}_p[f] . \quad (22)$$

- 1 Approximate Lip_1 by a Lipschitz **Neural network architecture**
 $\mathcal{D}_\phi = \{g_\phi : \mathbb{R}^{d_x} \rightarrow \mathbb{R}\}.$
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:

$$D = \arg \max_{\phi} \left(\mathbb{E}_{x \sim p} [D_\phi(x)] - \mathbb{E}_{z \sim \lambda} [D_\phi(f_\theta(z))] \right)$$

then train θ as

$$\theta = \arg \min_{\theta} \left(\mathbb{E}_{x \sim p} [D_\phi(x)] - \mathbb{E}_{z \sim \lambda} [D_\phi(f_\theta(z))] \right)$$

$$\mathcal{W}(q, p) = \sup_{f \in \text{Lip}_1} \mathbb{E}_q[f] - \mathbb{E}_p[f] . \quad (22)$$

- 1 Approximate Lip_1 by a Lipschitz **Neural network architecture**
 $\mathcal{D}_\phi = \{g_\phi : \mathbb{R}^{d_x} \rightarrow \mathbb{R}\}.$
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:

■ For fixed θ , set

$$\phi = \arg \max_{\phi} \mathbb{E}_{p_\theta} [g_\phi(X)] - \mathbb{E}_{\lambda} [g_\phi(f_\theta(Z))]$$

■ Then, for fixed ϕ , find θ that maximizes

$$\mathbb{E}_{\lambda} [g_\phi(f_\theta(Z))] - \mathbb{E}_{p_\theta} [g_\phi(X)]$$

$$\mathcal{W}(q, p) = \sup_{f \in \text{Lip}_1} \mathbb{E}_q[f] - \mathbb{E}_p[f] . \quad (22)$$

- 1 Approximate Lip_1 by a Lipschitz **Neural network architecture**
 $\mathcal{D}_\phi = \{g_\phi : \mathbb{R}^{d_x} \rightarrow \mathbb{R}\}.$
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:

1 For **fixed** θ , set

$$\phi = \arg \max_{\phi} \mathbb{E}_{\pi} [g_{\phi}(X)] - \mathbb{E}_{\lambda} [g_{\phi}(f_{\theta}(Z))]$$

2 For **fixed** ϕ , set

$$\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [g_{\phi}(X)] - \mathbb{E}_{\lambda} [g_{\phi}(f_{\theta}(Z))]$$

$$\mathcal{W}(q, p) = \sup_{f \in \text{Lip}_1} \mathbb{E}_q[f] - \mathbb{E}_p[f] . \quad (22)$$

- 1 Approximate Lip_1 by a Lipschitz **Neural network architecture**
 $\mathcal{D}_\phi = \{g_\phi : \mathbb{R}^{d_x} \rightarrow \mathbb{R}\}.$
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:
 - 1 For **fixed** θ , set

$$\phi = \arg \max_{\phi} \mathbb{E}_{\pi} [g_{\phi}(X)] - \mathbb{E}_{\lambda} [g_{\phi}(f_{\theta}(Z))]$$

- 2 For **fixed** ϕ , set

$$\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [g_{\phi}(X)] - \mathbb{E}_{\lambda} [g_{\phi}(f_{\theta}(Z))]$$

$$\mathcal{W}(q, p) = \sup_{f \in \text{Lip}_1} \mathbb{E}_q[f] - \mathbb{E}_p[f] . \quad (22)$$

- 1 Approximate Lip_1 by a Lipschitz **Neural network architecture**
 $\mathcal{D}_\phi = \{g_\phi : \mathbb{R}^{d_x} \rightarrow \mathbb{R}\}.$
- 2 Generative model: $Z \in \mathbb{R}^{d_z} \sim \lambda$, apply a Neural net $X = f_\theta(Z)$. and p_θ the law of X .
- 3 Train it in two steps:
 - 1 For **fixed** θ , set

$$\phi = \arg \max_{\phi} \mathbb{E}_{\pi} [g_{\phi}(X)] - \mathbb{E}_{\lambda} [g_{\phi}(f_{\theta}(Z))]$$

- 2 For **fixed** ϕ , set

$$\theta = \arg \min_{\theta} \mathbb{E}_{\pi} [g_{\phi}(X)] - \mathbb{E}_{\lambda} [g_{\phi}(f_{\theta}(Z))]$$

No more mode collapse

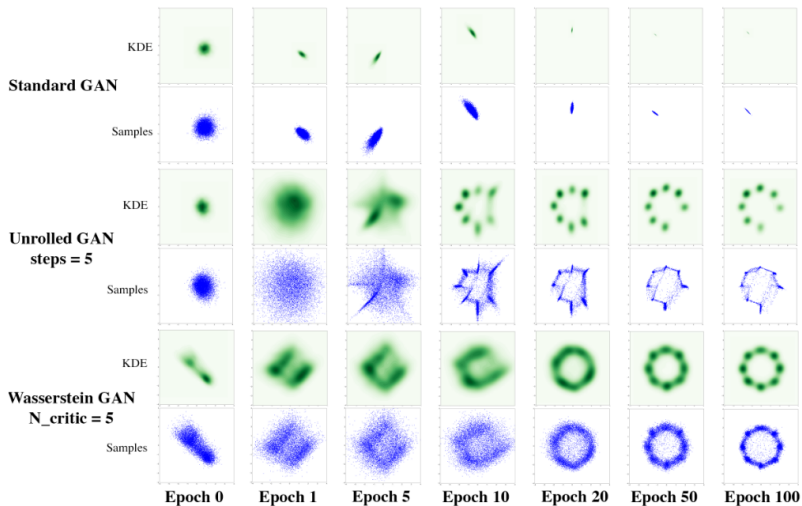


Figure: Comparison between GAN and Wasserstein GANs taken from Arjovsky, Chintala, and Bottou, 2017.

- 1 Same pros as the GANs ✓
- 2 No more mode collapse ✓
- 3 Growing literature on how to constraint Lipschitz norm of neural net ✓
- 4 still SOTA (state of the art) in some image tasks ✓
- 5 Training still requires some care !



Figure: Uncurated images from StyleGAN Karras, Laine, and Aila, 2019.

- 1 A bit of theory
- 2 Normalizing flows
- 3 VAE
- 4 GANs
- 5 Score-based generative models**



Figure: $X_0 \sim \pi$.



Figure: $X_t = X_0 + \sigma_t Z_t$.¹⁰

Suppose we have an **increasing noise schedule**

$$\sigma : [0, T] \rightarrow \sigma_t \in \mathbb{R}_* .$$

Can we **reverse the procedure**?

¹ $Z_t \perp X_0, Z_t \sim \mathcal{N}(0, \mathbf{I})$.

Suppose we have an **increasing noise schedule**

$$\sigma : [0, T] \rightarrow \sigma_t \in \mathbb{R}_* .$$

1 Noising kernel:

$$p_{t|0}(x_t|x_0) = \mathcal{N}(x_t; x_0, \sigma_t^2 \mathbf{I}) ,$$

2 Marginal:

$$p_t(x_t) = \int p_{t|0}(x_t|x_0)\pi(dx_0) .$$

Our goal is to go backwards:

$$p_T \rightarrow p_{T-\Delta t} \rightarrow \cdots \rightarrow p_0 .$$



Figure: $X_0 \sim \pi$.



Figure:
 $X_t = X_0 + \sigma_t Z_t$. ¹²1

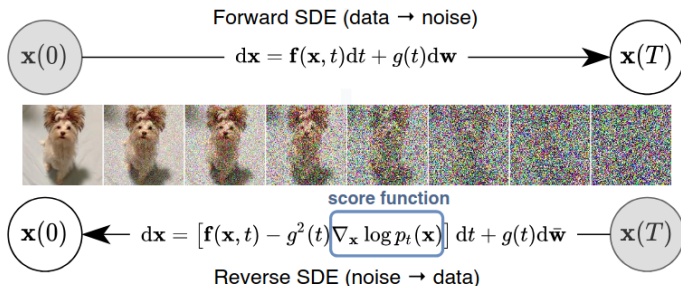


Figure: Illustration of noising and denoising with SDEs.

Our p_t is the time t marginal of the SDE

$$d\mathbf{x} = \sqrt{\frac{d\sigma_t^2}{dt}} d\mathbf{w}.$$

It admits the backward decomposition

$$d\mathbf{x} = -\frac{d\sigma_t^2}{dt} \nabla \log p_t(\mathbf{x}) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{\mathbf{w}}.$$

¹²From Y. Song et al., 2021.

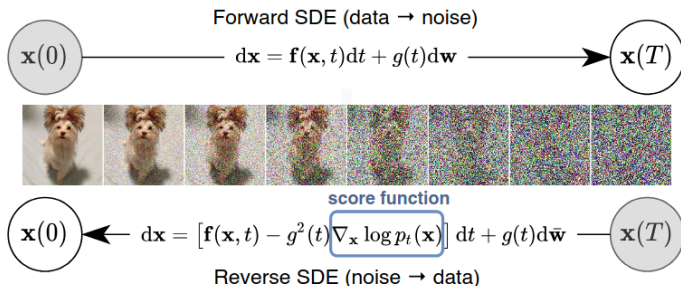


Figure: Illustration of noising and denoising with SDEs.

Our p_t is the time t marginal of the SDE

$$d\mathbf{x} = \sqrt{\frac{d\sigma_t^2}{dt}} d\mathbf{w}.$$

It admits the backward decomposition

$$d\mathbf{x} = -\frac{d\sigma_t^2}{dt} \nabla \log p_t(\mathbf{x}) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{\mathbf{w}}.$$

¹²From Y. Song et al., 2021.

$$\begin{aligned}
\nabla_{x_t} \log p_t(x_t) &= \frac{\nabla_{x_t} p_t(x_t)}{p_t(x_t)} = \frac{\int \nabla_{x_t} p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \frac{\int \nabla \log p_{t|0}(x_t|x_0) p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \int \nabla \log p_{t|0}(x_t|x_0) \frac{p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \int \nabla \log p_{t|0}(x_t|x_0) p_{0|t}(dx_0|x_t) = \mathbb{E} [\nabla \log p_{t|0}(x_t|X_0) | X_t = t] .
\end{aligned}$$

$\nabla \log p_{t|0}(x_t|x_0) = -\sigma_t^{-2}(x_t - x_0)$, thus

$$\nabla_{x_t} \log p_t(x_t) = \mathbb{E} [\nabla \log p_{t|0}(x_t|X_0) | X_t = t] = -\sigma_t^{-2}(x_t - \mathbb{E} [X_0 | X_t = x_t])$$

or equivalently:

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \nabla \log p_t(x_t) .$$

$$\begin{aligned}
\nabla_{x_t} \log p_t(x_t) &= \frac{\nabla_{x_t} p_t(x_t)}{p_t(x_t)} = \frac{\int \nabla_{x_t} p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \frac{\int \nabla \log p_{t|0}(x_t|x_0) p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \int \nabla \log p_{t|0}(x_t|x_0) \frac{p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \int \nabla \log p_{t|0}(x_t|x_0) p_{0|t}(dx_0|x_t) = \mathbb{E} [\nabla \log p_{t|0}(x_t|X_0) | X_t = t] .
\end{aligned}$$

$\nabla \log p_{t|0}(x_t|x_0) = -\sigma_t^{-2}(x_t - x_0)$, thus

$$\nabla_{x_t} \log p_t(x_t) = \mathbb{E} [\nabla \log p_{t|0}(x_t|X_0) | X_t = t] = -\sigma_t^{-2}(x_t - \mathbb{E} [X_0 | X_t = x_t])$$

or equivalently:

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \nabla \log p_t(x_t) .$$

$$\begin{aligned}
\nabla_{x_t} \log p_t(x_t) &= \frac{\nabla_{x_t} p_t(x_t)}{p_t(x_t)} = \frac{\int \nabla_{x_t} p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \frac{\int \nabla \log p_{t|0}(x_t|x_0) p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \int \nabla \log p_{t|0}(x_t|x_0) \frac{p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \int \nabla \log p_{t|0}(x_t|x_0) p_{0|t}(dx_0|x_t) = \mathbb{E} [\nabla \log p_{t|0}(x_t|X_0) | X_t = t] .
\end{aligned}$$

$\nabla \log p_{t|0}(x_t|x_0) = -\sigma_t^{-2}(x_t - x_0)$, thus

$$\nabla_{x_t} \log p_t(x_t) = \mathbb{E} [\nabla \log p_{t|0}(x_t|X_0) | X_t = t] = -\sigma_t^{-2}(x_t - \mathbb{E} [X_0 | X_t = x_t])$$

or equivalently:

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \nabla \log p_t(x_t) .$$

$$\begin{aligned}
\nabla_{x_t} \log p_t(x_t) &= \frac{\nabla_{x_t} p_t(x_t)}{p_t(x_t)} = \frac{\int \nabla_{x_t} p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \frac{\int \nabla \log p_{t|0}(x_t|x_0) p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \int \nabla \log p_{t|0}(x_t|x_0) \frac{p_{t|0}(x_t|x_0) \pi(dx_0)}{p_t(x_t)} \\
&= \int \nabla \log p_{t|0}(x_t|x_0) p_{0|t}(dx_0|x_t) = \mathbb{E} [\nabla \log p_{t|0}(x_t|X_0) | X_t = t] .
\end{aligned}$$

$\nabla \log p_{t|0}(x_t|x_0) = -\sigma_t^{-2}(x_t - x_0)$, thus

$$\nabla_{x_t} \log p_t(x_t) = \mathbb{E} [\nabla \log p_{t|0}(x_t|X_0) | X_t = t] = -\sigma_t^{-2}(x_t - \mathbb{E} [X_0 | X_t = x_t])$$

or equivalently:

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \nabla \log p_t(x_t) .$$

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \nabla \log p_t(x_t) .$$

$$\text{But } \mathbb{E} [\|X_0 - \mathbb{E} [X_0 | X_t]\|^2] = \min_{g \text{ measurable}} \mathbb{E} [\|X_0 - g(X_t)\|^2] .$$

Learning the score from data \Leftrightarrow Learn a denoiser from data.

$$L_{sm}(\theta) = \mathbb{E}_{t \sim \text{Unif}([0, T])} [\gamma_t^2 \mathbb{E}_{X_0 \sim \pi, Z \sim \mathcal{N}(0, I)} [\|X_0 - D_\theta(X_0 + \sigma_t Z, \sigma_t)\|^2]] .$$

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \nabla \log p_t(x_t) .$$

$$\text{But } \mathbb{E} [\|X_0 - \mathbb{E} [X_0 | X_t]\|^2] = \min_{g \text{ measurable}} \mathbb{E} [\|X_0 - g(X_t)\|^2] .$$

Learning the score from data \Leftrightarrow Learn a denoiser from data.

$$L_{sm}(\theta) = \mathbb{E}_{t \sim \text{Unif}([0, T])} [\gamma_t^2 \mathbb{E}_{X_0 \sim \pi, Z \sim \mathcal{N}(0, I)} [\|X_0 - D_\theta(X_0 + \sigma_t Z, \sigma_t)\|^2]] .$$

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \nabla \log p_t(x_t) .$$

$$\text{But } \mathbb{E} [\|X_0 - \mathbb{E} [X_0 | X_t]\|^2] = \min_{g \text{ measurable}} \mathbb{E} [\|X_0 - g(X_t)\|^2] .$$

Learning the score from data \Leftrightarrow Learn a denoiser from data.

$$L_{sm}(\theta) = \mathbb{E}_{t \sim \text{Unif}([0, T])} [\gamma_t^2 \mathbb{E}_{X_0 \sim \pi, Z \sim \mathcal{N}(0, I)} [\|X_0 - D_\theta(X_0 + \sigma_t Z, \sigma_t)\|^2]] .$$

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \nabla \log p_t(x_t) .$$

$$\text{But } \mathbb{E} [\|X_0 - \mathbb{E} [X_0 | X_t]\|^2] = \min_{g \text{ measurable}} \mathbb{E} [\|X_0 - g(X_t)\|^2] .$$

Learning the score from data \Leftrightarrow Learn a denoiser from data.

$$L_{sm}(\theta) = \mathbb{E}_{t \sim \text{Unif}([0, T])} \left[\gamma_t^2 \mathbb{E}_{X_0 \sim \pi, Z \sim \mathcal{N}(0, I)} [\|X_0 - D_\theta(X_0 + \sigma_t Z, \sigma_t)\|^2] \right] .$$

Suppose we have a trained denoiser $D_\theta(x, \sigma_t)$.

We now have to simulate from

$$dx = \frac{d\sigma_t^2}{dt} \sigma_t^{-2} (x - D_\theta(x, \sigma_t)) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{w}.$$

We should in theory start from p_T .

But $D_{\text{KL}}(p_T || \mathcal{N}(0, \sigma_T^2 I)) \leq \sigma_T^{-2} \mathbb{E} [\|X_0\|^2]$.

With an appropriate choice of σ_T , we can start from $\mathcal{N}(0, \sigma_T^2 I)$.

Suppose we have a trained denoiser $D_\theta(x, \sigma_t)$.

We now have to simulate from

$$dx = \frac{d\sigma_t^2}{dt} \sigma_t^{-2} (x - D_\theta(x, \sigma_t)) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{w}.$$

We should in theory start from p_T .

But $D_{\text{KL}}(p_T \| \mathcal{N}(0, \sigma_T^2 I)) \leq \sigma_T^{-2} \mathbb{E}[\|X_0\|^2]$.

With an appropriate choice of σ_T , we can start from $\mathcal{N}(0, \sigma_T^2 I)$.

Suppose we have a trained denoiser $D_\theta(x, \sigma_t)$.

We now have to simulate from

$$dx = \frac{d\sigma_t^2}{dt} \sigma_t^{-2} (x - D_\theta(x, \sigma_t)) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{w}.$$

We should in theory start from p_T .

But $D_{\text{KL}}(p_T \| \mathcal{N}(0, \sigma_T^2 I)) \leq \sigma_T^{-2} \mathbb{E} [\|X_0\|^2]$.

With an appropriate choice of σ_T , we can start from $\mathcal{N}(0, \sigma_T^2 I)$.

Suppose we have a trained denoiser $D_\theta(x, \sigma_t)$.

We now have to simulate from

$$dx = \frac{d\sigma_t^2}{dt} \sigma_t^{-2} (x - D_\theta(x, \sigma_t)) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{w}.$$

We should in theory start from p_T .

But $D_{\text{KL}}(p_T \| \mathcal{N}(0, \sigma_T^2 \mathbf{I})) \leq \sigma_T^{-2} \mathbb{E} [\|X_0\|^2]$.

With an appropriate choice of σ_T , we can start from $\mathcal{N}(0, \sigma_T^2 \mathbf{I})$.

We have to simulate from

$$dx = \frac{d\sigma_t^2}{dt} \sigma_t^{-2} (x - D_\theta(x, \sigma_t)) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{w} .$$

We consider the rather crude approximation

$$X_{t-1} = X_t + \frac{(\sigma_{t-1}^2 - \sigma_t^2)}{\sigma_t^2} (X_t - D_\theta(X_t, \sigma_t)) + \sqrt{\sigma_t^2 - \sigma_{t-1}^2} Z_t$$

We have to simulate from

$$dx = \frac{d\sigma_t^2}{dt} \sigma_t^{-2} (x - D_\theta(x, \sigma_t)) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{w} .$$

We consider the rather crude approximation

$$X_{t-1} = X_t + \frac{(\sigma_{t-1}^2 - \sigma_t^2)}{\sigma_t^2} (X_t - D_\theta(X_t, \sigma_t)) + \sqrt{\sigma_t^2 - \sigma_{t-1}^2} Z_t$$

We have to simulate from

$$dx = \frac{d\sigma_t^2}{dt} \sigma_t^{-2} (x - D_\theta(x, \sigma_t)) dt + \sqrt{\frac{d\sigma_t^2}{dt}} d\bar{w} .$$

We consider the rather crude approximation

$$X_{t-1} = X_t + \frac{(\sigma_{t-1}^2 - \sigma_t^2)}{\sigma_t^2} (X_t - D_\theta(X_t, \sigma_t)) + \sqrt{\sigma_t^2 - \sigma_{t-1}^2} Z_t$$

Algorithm 3 Backward sampling

1: **procedure** SAMPLING **Input:** N . **Output:** X_0 .2: $X_N \sim \mathcal{N}(0, I)$.3: **for** $i = N - 1, \dots, 0$ **do**4: $t_i = (i/N)T$ 5: $Z_i \sim \mathcal{N}(0, I)$.6:
$$X_i = X_{i+1} + \frac{(\sigma_{t_i}^2 - \sigma_{t_{i+1}}^2)}{\sigma_{t_{i+1}}^2} (X_{i+1} - D_\theta(X_{i+1}, \sigma_{t_{i+1}})) + \sqrt{\sigma_{t_{i+1}}^2 - \sigma_{t_i}^2} Z_i.$$
7: **end for**8: **Return:** X_0 9: **end procedure**

This is a very naive sampler. There are much better samplers, see Karras, Aittala, et al., 2022 or J. Song, Meng, and Ermon, 2021.



Figure: Illustration on the Celeb dataset from Y. Song et al., 2021.

- 1 SOTA in most image tasks ✓
- 2 No mode collapse ✓
- 3 No adversarial training ✓
- 4 Easier to condition on text (for example) ✓
- 5 Easy to conditiong à-posteriori ✓
- 6 Sampling requires more Neural function evaluations (NFE) ✗

- ▶ George Papamakarios et al. “Normalizing Flows for Probabilistic Modeling and Inference”. In: *Journal of Machine Learning Research* 22.57 (2021), pp. 1–64
- ▶ Durk P Kingma and Prafulla Dhariwal. “Glow: Generative Flow with Invertible 1x1 Convolutions”. In: *Advances in Neural Information Processing Systems*. Vol. 31. Curran Associates, Inc., 2018. (Visited on 03/03/2025)
- ▶ Aaron van den Oord et al. “Conditional Image Generation with PixelCNN Decoders”. In: *Advances in Neural Information Processing Systems*. Vol. 29. Curran Associates, Inc., 2016. (Visited on 03/03/2025)
- ▶ Louis Grenioux et al. “On Sampling with Approximate Transport Maps”. In: *Proceedings of the 40th International Conference on Machine Learning*. Ed. by Andreas Krause et al. Vol. 202. Proceedings of Machine Learning Research. PMLR, July 2023, pp. 11698–11733

- ▶ Yufeng Zhang et al. “Kullback-Leibler Divergence-Based out-of-Distribution Detection with Flow-Based Generative Models”. In: *IEEE Transactions on Knowledge and Data Engineering* 36.4 (Apr. 2024), pp. 1683–1697. ISSN: 1558-2191. DOI: 10.1109/TKDE.2023.3309853. (Visited on 02/01/2025)
- ▶ Tim Salimans et al. “PixelCNN++: Improving the PixelCNN with Discretized Logistic Mixture Likelihood and Other Modifications”. In: *International Conference on Learning Representations*. Feb. 2017. (Visited on 03/03/2025)
- ▶ Ruizhi Deng et al. “Modeling Continuous Stochastic Processes with Dynamic Normalizing Flows”. In: *Advances in Neural Information Processing Systems*. Vol. 33. Curran Associates, Inc., 2020, pp. 7805–7815. (Visited on 03/03/2025)
- ▶ Shuangfei Zhai et al. *Normalizing Flows Are Capable Generative Models*. Dec. 2024. DOI: 10.48550/arXiv.2412.06329. arXiv: 2412.06329 [cs]. (Visited on 03/03/2025)

- ▶ Diederik P Kingma, Max Welling, et al. “An Introduction to Variational Autoencoders”. In: *Foundations and Trends® in Machine Learning* 12.4 (2019), pp. 307–392
- ▶ Rong Tang and Yun Yang. “On Empirical Bayes Variational Autoencoder: An Excess Risk Bound”. In: *Proceedings of Thirty Fourth Conference on Learning Theory*. PMLR, July 2021, pp. 4068–4125. (Visited on 10/23/2024)
- ▶ Aaron Van Den Oord and Oriol Vinyals. “Neural Discrete Representation Learning”. In: *Advances in Neural Information Processing Systems* 30 (2017). (Visited on 03/03/2025)
- ▶ Irina Higgins et al. “Beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework”. In: *International Conference on Learning Representations*. Feb. 2017. (Visited on 03/03/2025)

- ▶ Jakub Tomczak and Max Welling. “VAE with a VampPrior”. In: *Proceedings of the Twenty-first International Conference on Artificial Intelligence and Statistics*. PMLR, Mar. 2018, pp. 1214–1223. (Visited on 03/03/2025)
- ▶ Arash Vahdat and Jan Kautz. “NVAE: A Deep Hierarchical Variational Autoencoder”. In: *Advances in Neural Information Processing Systems*. Vol. 33. Curran Associates, Inc., 2020, pp. 19667–19679. (Visited on 03/03/2025)
- ▶ Ali Razavi, Aaron van den Oord, and Oriol Vinyals. “Generating Diverse High-Fidelity Images with VQ-VAE-2”. In: *Advances in Neural Information Processing Systems*. Vol. 32. Curran Associates, Inc., 2019. (Visited on 03/03/2025)

If you want to know more about GANs

- ▶ Ishaan Gulrajani et al. “Improved Training of Wasserstein GANs”. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. NIPS’17. Red Hook, NY, USA: Curran Associates Inc., 2017, pp. 5769–5779. ISBN: 978-1-5108-6096-4
- ▶ Tero Karras, Samuli Laine, and Timo Aila. “A Style-Based Generator Architecture for Generative Adversarial Networks”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019, pp. 4401–4410. (Visited on 03/02/2025)
- ▶ Tero Karras, Samuli Laine, Miika Aittala, et al. *Analyzing and Improving the Image Quality of StyleGAN*. Mar. 2020. DOI: 10.48550/arXiv.1912.04958. arXiv: 1912.04958 [cs]. (Visited on 12/20/2024)
- ▶ Cem Anil, James Lucas, and Roger Grosse. *Sorting out Lipschitz Function Approximation*. June 2019. DOI: 10.48550/arXiv.1811.05381. arXiv: 1811.05381 [cs]. (Visited on 12/20/2024)

- ▶ Gérard Biau, Maxime Sangnier, and Ugo Tanielian. “Some Theoretical Insights into Wasserstein GANs”. In: *Journal of Machine Learning Research* 22.119 (2021), pp. 1–45. ISSN: 1533-7928. (Visited on 11/26/2024)
- ▶ Elen Vardanyan et al. “Statistically Optimal Generative Modeling with Maximum Deviation from the Empirical Distribution”. In: *Proceedings of the 41st International Conference on Machine Learning*. PMLR, July 2024, pp. 49203–49225. (Visited on 03/03/2025)

- ▶ Jonathan Ho, Ajay Jain, and Pieter Abbeel. “Denoising Diffusion Probabilistic Models”. In: *Advances in Neural Information Processing Systems*. Vol. 33. 2020, pp. 6840–6851
- ▶ Jiaming Song, Chenlin Meng, and Stefano Ermon. “Denoising Diffusion Implicit Models”. In: *International Conference on Learning Representations*. 2021
- ▶ Prafulla Dhariwal and Alexander Nichol. “Diffusion Models Beat Gans on Image Synthesis”. In: *Advances in neural information processing systems* 34 (2021), pp. 8780–8794
- ▶ Tero Karras, Miika Aittala, et al. “Elucidating the Design Space of Diffusion-Based Generative Models”. In: *Proc. NeurIPS*. 2022

If you want to know more about Score based generative models

- ▶ Robin Rombach et al. “High-Resolution Image Synthesis with Latent Diffusion Models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 10684–10695. (Visited on 01/28/2025)
- ▶ Zahra Kadkhodaie et al. “Generalization in Diffusion Models Arises from Geometry-Adaptive Harmonic Representation”. In: *The Twelfth International Conference on Learning Representations*. 2023
- ▶ Giovanni Conforti, Alain Durmus, and Marta Gentiloni Silveri. *KL Convergence Guarantees for Score Diffusion Models under Minimal Data Assumptions*. Sept. 2024. arXiv: 2308.12240 [math, stat]. (Visited on 09/17/2024)
- ▶ Xingchao Liu, Chengyue Gong, and Qiang Liu. “Flow Straight and Fast: Learning to Generate and Transfer Data with Rectified Flow”. In: *The Eleventh International Conference on Learning Representations*. Sept. 2022. (Visited on 03/03/2025)

- [1] Cem Anil, James Lucas, and Roger Grosse. *Sorting out Lipschitz Function Approximation*. June 2019. DOI: 10.48550/arXiv.1811.05381. arXiv: 1811.05381 [cs]. (Visited on 12/20/2024).
- [2] Martin Arjovsky, Soumith Chintala, and Léon Bottou. “Wasserstein Generative Adversarial Networks”. In: *International Conference on Machine Learning*. PMLR, 2017, pp. 214–223.
- [3] Gérard Biau, Maxime Sangnier, and Ugo Tanielian. “Some Theoretical Insights into Wasserstein GANs”. In: *Journal of Machine Learning Research* 22.119 (2021), pp. 1–45. ISSN: 1533-7928. (Visited on 11/26/2024).
- [4] Giovanni Conforti, Alain Durmus, and Marta Gentiloni Silveri. *KL Convergence Guarantees for Score Diffusion Models under Minimal Data Assumptions*. Sept. 2024. arXiv: 2308.12240 [math, stat]. (Visited on 09/17/2024).

- [5] Ruizhi Deng et al. “Modeling Continuous Stochastic Processes with Dynamic Normalizing Flows”. In: *Advances in Neural Information Processing Systems*. Vol. 33. Curran Associates, Inc., 2020, pp. 7805–7815. (Visited on 03/03/2025).
- [6] Prafulla Dhariwal and Alexander Nichol. “Diffusion Models Beat Gans on Image Synthesis”. In: *Advances in neural information processing systems* 34 (2021), pp. 8780–8794.
- [7] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. “Density Estimation Using Real NVP”. In: *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. 2017.
- [8] M. D. Donsker and S. R. S. Varadhan. “Asymptotic Evaluation of Certain Markov Process Expectations for Large Time. IV”. In: *Communications on Pure and Applied Mathematics* 36.2 (1983), pp. 183–212. ISSN: 1097-0312. DOI: 10.1002/cpa.3160360204. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/cpa.3160360204> (visited on 01/03/2025).

- [9] Ian J. Goodfellow et al. “Generative Adversarial Nets”. In: *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2. NIPS’14*. Cambridge, MA, USA: MIT Press, 2014, pp. 2672–2680.
- [10] Louis Grenioux et al. “On Sampling with Approximate Transport Maps”. In: *Proceedings of the 40th International Conference on Machine Learning*. Ed. by Andreas Krause et al. Vol. 202. Proceedings of Machine Learning Research. PMLR, July 2023, pp. 11698–11733.
- [11] Ishaan Gulrajani et al. “Improved Training of Wasserstein GANs”. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems. NIPS’17*. Red Hook, NY, USA: Curran Associates Inc., 2017, pp. 5769–5779. ISBN: 978-1-5108-6096-4.

- [12] Irina Higgins et al. “Beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework”. In: *International Conference on Learning Representations*. Feb. 2017. (Visited on 03/03/2025).
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. “Denoising Diffusion Probabilistic Models”. In: *Advances in Neural Information Processing Systems*. Vol. 33. 2020, pp. 6840–6851.
- [14] Zahra Kadkhodaie et al. “Generalization in Diffusion Models Arises from Geometry-Adaptive Harmonic Representation”. In: *The Twelfth International Conference on Learning Representations*. 2023.
- [15] Tero Karras, Miika Aittala, et al. “Elucidating the Design Space of Diffusion-Based Generative Models”. In: *Proc. NeurIPS*. 2022.

- [16] Tero Karras, Samuli Laine, and Timo Aila. “A Style-Based Generator Architecture for Generative Adversarial Networks”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019, pp. 4401–4410. (Visited on 03/02/2025).
- [17] Tero Karras, Samuli Laine, Miika Aittala, et al. *Analyzing and Improving the Image Quality of StyleGAN*. Mar. 2020. DOI: 10.48550/arXiv.1912.04958. arXiv: 1912.04958 [cs]. (Visited on 12/20/2024).
- [18] Diederik P Kingma, Max Welling, et al. “An Introduction to Variational Autoencoders”. In: *Foundations and Trends® in Machine Learning* 12.4 (2019), pp. 307–392.
- [19] Diederik P. Kingma and Max Welling. “Auto-Encoding Variational Bayes”. In: (Dec. 23, 2013). URL: <https://openreview.net/forum?id=33X9fd2-9FyZd> (visited on 01/03/2025).

- [20] Durk P Kingma and Prafulla Dhariwal. “Glow: Generative Flow with Invertible 1x1 Convolutions”. In: *Advances in Neural Information Processing Systems*. Vol. 31. Curran Associates, Inc., 2018. (Visited on 03/03/2025).
- [21] Xingchao Liu, Chengyue Gong, and Qiang Liu. “Flow Straight and Fast: Learning to Generate and Transfer Data with Rectified Flow”. In: *The Eleventh International Conference on Learning Representations*. Sept. 2022. (Visited on 03/03/2025).
- [22] George Papamakarios et al. “Normalizing Flows for Probabilistic Modeling and Inference”. In: *Journal of Machine Learning Research* 22.57 (2021), pp. 1–64.
- [23] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. “Generating Diverse High-Fidelity Images with VQ-VAE-2”. In: *Advances in Neural Information Processing Systems*. Vol. 32. Curran Associates, Inc., 2019. (Visited on 03/03/2025).

- [24] Danilo Rezende and Shakir Mohamed. “Variational Inference with Normalizing Flows”. In: *Proceedings of the 32nd International Conference on Machine Learning*. Ed. by Francis Bach and David Blei. Vol. 37. Proceedings of Machine Learning Research. Lille, France: PMLR, July 2015, pp. 1530–1538.
- [25] Robin Rombach et al. “High-Resolution Image Synthesis with Latent Diffusion Models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 10684–10695. (Visited on 01/28/2025).
- [26] Tim Salimans et al. “PixelCNN++: Improving the PixelCNN with Discretized Logistic Mixture Likelihood and Other Modifications”. In: *International Conference on Learning Representations*. Feb. 2017. (Visited on 03/03/2025).
- [27] Jiaming Song, Chenlin Meng, and Stefano Ermon. “Denoising Diffusion Implicit Models”. In: *International Conference on Learning Representations*. 2021.

- [28] Yang Song et al. “Score-Based Generative Modeling through Stochastic Differential Equations”. In: *International Conference on Learning Representations*. 2021.
- [29] Rong Tang and Yun Yang. “On Empirical Bayes Variational Autoencoder: An Excess Risk Bound”. In: *Proceedings of Thirty Fourth Conference on Learning Theory*. PMLR, July 2021, pp. 4068–4125. (Visited on 10/23/2024).
- [30] Jakub Tomczak and Max Welling. “VAE with a VampPrior”. In: *Proceedings of the Twenty-first International Conference on Artificial Intelligence and Statistics*. PMLR, Mar. 2018, pp. 1214–1223. (Visited on 03/03/2025).
- [31] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. “Deep Image Prior”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2018.

- [32] Arash Vahdat and Jan Kautz. “NVAE: A Deep Hierarchical Variational Autoencoder”. In: *Advances in Neural Information Processing Systems*. Vol. 33. Curran Associates, Inc., 2020, pp. 19667–19679. (Visited on 03/03/2025).
- [33] Aaron Van Den Oord and Oriol Vinyals. “Neural Discrete Representation Learning”. In: *Advances in Neural Information Processing Systems* 30 (2017). (Visited on 03/03/2025).
- [34] Aaron van den Oord et al. “Conditional Image Generation with PixelCNN Decoders”. In: *Advances in Neural Information Processing Systems*. Vol. 29. Curran Associates, Inc., 2016. (Visited on 03/03/2025).
- [35] Elen Vardanyan et al. “Statistically Optimal Generative Modeling with Maximum Deviation from the Empirical Distribution”. In: *Proceedings of the 41st International Conference on Machine Learning*. PMLR, July 2024, pp. 49203–49225. (Visited on 03/03/2025).

- [36] Cédric Villani. *Topics in optimal transportation*. Vol. 58. American Mathematical Soc., 2021.
- [37] Nick Whiteley, Annie Gray, and Patrick Rubin-Delanchy. *Statistical Exploration of the Manifold Hypothesis*. Feb. 2024. DOI: 10.48550/arXiv.2208.11665. arXiv: 2208.11665 [stat]. (Visited on 02/28/2025).
- [38] Shuangfei Zhai et al. *Normalizing Flows Are Capable Generative Models*. Dec. 2024. DOI: 10.48550/arXiv.2412.06329. arXiv: 2412.06329 [cs]. (Visited on 03/03/2025).
- [39] Yufeng Zhang et al. “Kullback-Leibler Divergence-Based out-of-Distribution Detection with Flow-Based Generative Models”. In: *IEEE Transactions on Knowledge and Data Engineering* 36.4 (Apr. 2024), pp. 1683–1697. ISSN: 1558-2191. DOI: 10.1109/TKDE.2023.3309853. (Visited on 02/01/2025).

Algorithm 4 Backward sampling

1: **procedure** SAMPLING

Input: N .

Output: X_0 .

2: $X_N \sim \mathcal{N}(0, I)$.

3: **for** $i = N - 1, \dots, 0$ **do**

4: $t_i = (i/N)T$

5: $Z_i \sim \mathcal{N}(0, I)$.

6: $X_i = X_{i+1} + \frac{\sigma_{t_i}^2 - \sigma_{t_{i+1}}^2}{\sigma_{t_{i+1}}^2} (X_{i+1} - D_\theta(X_{i+1}, \sigma_{t_{i+1}})) +$

$$\sqrt{\frac{\sigma_{t_i}^2}{\sigma_{t_{i+1}}^2} (\sigma_{t_{i+1}}^2 - \sigma_{t_i}^2)} Z_i.$$

7: **end for**

8: **Return:** X_0

9: **end procedure**
