

Groupe de lecture :

OT-Flow :

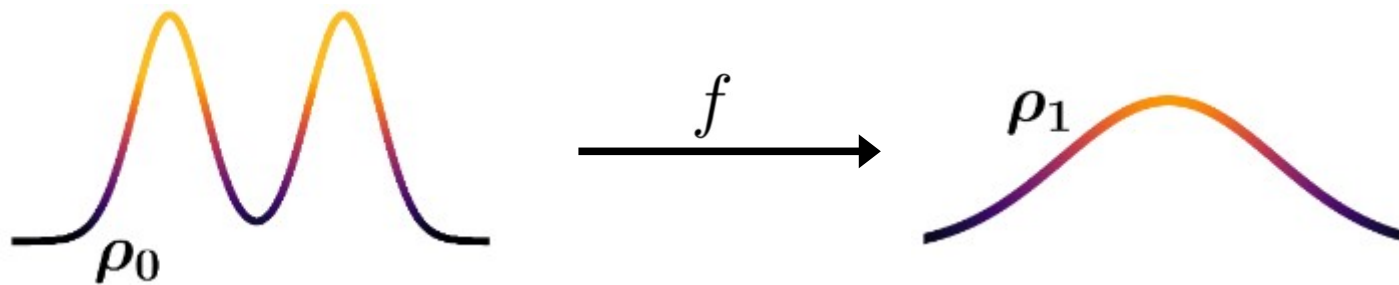
Fast and Accurate Continuous Normalizing Flows via Optimal Transport

Derek Onken, Samy Wu Fung, Xingjian Li, Lars Ruthotto 2021

Buonomo Camille,
14/05/2024,

1. Normalizing flows

Invertible mapping : $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$



Condition on f : $X \sim \rho_1, f$ s.t. $f(X) \sim \rho_0$

$$\log(\rho_0(x)) = \log(\rho_1(f(x))) + \log(\det(\mathcal{J}_f(x)))$$

1. Normalizing flows

In practice : Composition of simple invertible map

$$f = f_k \circ f_{k-1} \circ \dots \circ f_1 \circ f_0$$

$$f^{-1} = f_1^{-1} \circ \dots \circ f_k^{-1}$$

$$\det(\mathcal{J}_f)(f(x)) = \prod \det(\mathcal{J}_{f_k})(x_k)$$

1. Normalizing flows

In practice : Composition of simple invertible map

$$f = f_k \circ f_{k-1} \circ \dots \circ f_1 \circ f_0$$

Lots and lots of possibility for f_k

- Expressivity
- Inference
- Inverse
- Jacobian
- Determinant of Jacobian
- ...

[Kobyzev et Al. 2021]

Noteworthy exemple : $f_k(x) = x + hg(x; \theta)$

2. Continuous Normalizing flows

$$f_k(x) = x + hg(x; \theta)$$

Forward Euler step : $y^k = y^{k-1} + \Delta t V(y^{k-1}, t_{k-1})$

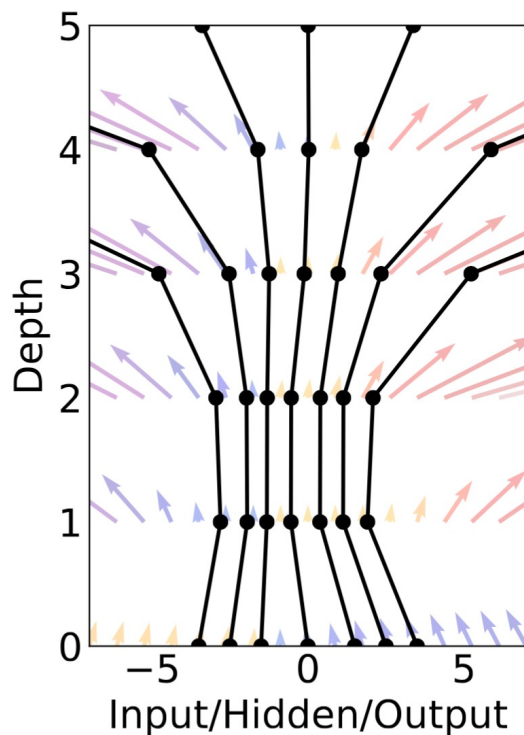
Corresponding EDO : $\partial_t z = g(z, t; \theta)$
 $z(0) = x$

→ Integration of a learned dynamic : $f(x) = z_x(T) = x + \int_0^T g(z(t), t; \theta) dt$

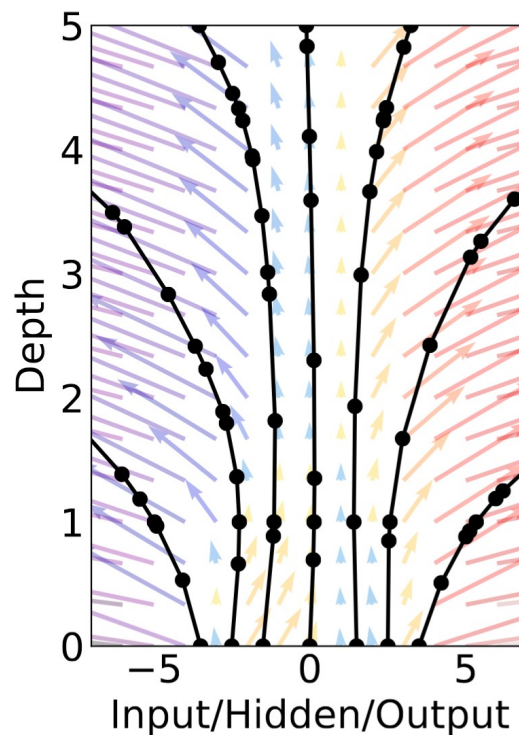
2. Continuous Normalizing flows

“Infinitely” deep neural network

Residual Network



ODE Network



[Chen et Al. 2019]

2. Continuous Normalizing flows

Instantaneous change of variable :

$$\partial_t l(x, t) = \text{Tr}(\mathcal{J}_{g(\cdot, t; \theta)}(z(x, t))) = \text{div}(g)(z(x, t), t; \theta)$$

Such that : $l(x, T) = \log(\det(\mathcal{J}_f(x)))$

Continuous backpropagation :

$$\partial_\theta L = \int_0^T \nabla_{z(t)} L \partial_\theta g(z(t), t; \theta) dt$$

[Chen et al. 2019]

2. Continuous Normalizing flows

Training : Kullback-Leibler divergence

$$\begin{aligned} \mathbb{D}_{\text{KL}} [\rho(\mathbf{z}(\mathbf{x}, T)) \parallel \rho_1(\mathbf{z}(\mathbf{x}, T))] \\ &= \int_{\mathbb{R}^d} \log \left(\frac{\rho(\mathbf{z}(\mathbf{x}, T))}{\rho_1(\mathbf{z}(\mathbf{x}, T))} \right) \rho(\mathbf{z}(\mathbf{x}, T)) \det (\nabla \mathbf{z}(\mathbf{x}, T)) \, d\mathbf{x}, \\ &= \int_{\mathbb{R}^d} \log \left(\frac{\rho_0(\mathbf{x})}{\rho_1(\mathbf{z}(\mathbf{x}, T)) \det (\nabla \mathbf{z}(\mathbf{x}, T))} \right) \rho_0(\mathbf{x}) \, d\mathbf{x}, \\ &= \int_{\mathbb{R}^d} \left[\log (\rho_0(\mathbf{x})) - \log (\rho_1(\mathbf{z}(\mathbf{x}, T))) - \log \det (\nabla \mathbf{z}(\mathbf{x}, T)) \right] \rho_0(\mathbf{x}) \, d\mathbf{x}. \end{aligned}$$

- ρ_0 is unknown
- ρ_1 is chosen as a normal distribution

2. Continuous Normalizing flows

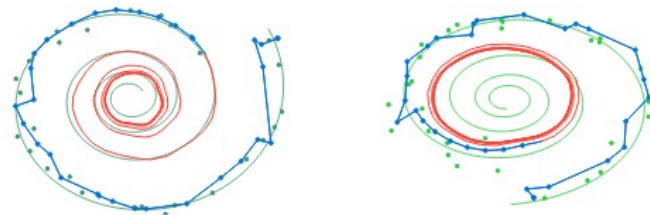
Training : log-likelihood maximisation

$$\begin{aligned}\mathbb{D}_{\text{KL}} &= \int_{\mathbb{R}^d} \left[\log(\rho_0(\mathbf{x})) - \log \det(\nabla \mathbf{z}(\mathbf{x}, T)) + \frac{1}{2} \|\mathbf{z}(\mathbf{x}, T)\|^2 + \frac{d}{2} \log(2\pi) \right] \rho_0(\mathbf{x}) \, d\mathbf{x} \\ &= \int_{\mathbb{R}^d} \left[\log(\rho_0(\mathbf{x})) + C(\mathbf{x}, T) \right] \rho_0(\mathbf{x}) \, d\mathbf{x} \quad C(x, T) = \frac{1}{2} \|z(x, T)\|^2 + l(x, T) - \frac{d}{2} \log(2\pi) \\ &= \mathbb{E}_{\rho_0(\mathbf{x})} \{ \log(\rho_0(\mathbf{x})) + C(\mathbf{x}, T) \} , \\ &\rightarrow \min_{\theta} E_{\rho_0}[C(x, T)]\end{aligned}$$

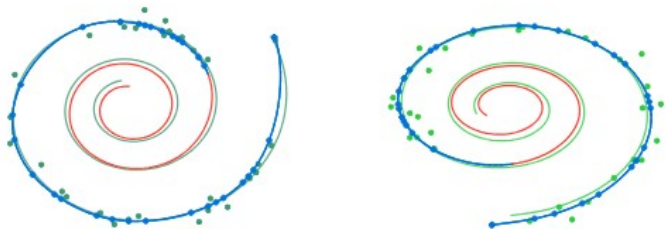
2. Continuous Normalizing flows

Drawbacks :

- Slow integration
- “High” computational cost of the score



(a) Recurrent Neural Network

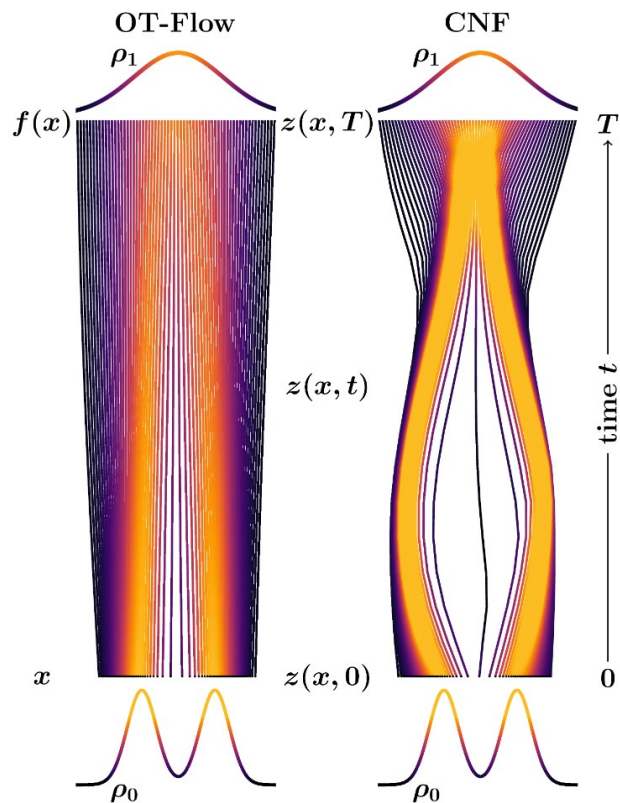


(b) Latent Neural Ordinary Differential Equation

[Chen et Al. 2019]

3. OT-Flow

Main idea : Regularization of the trajectories for a faster integration



3. OT-Flow

Main idea : Regularization of the trajectories for a faster integration

OT penalization :
$$L(x, T) = \int_0^T \frac{1}{2} \|g(z_x(t), t; \theta)\|^2 dt$$

→ Relaxed Benamou-Brenier formulation

$$\min_{\theta} E_{\rho_0} [C(x, T) + L(x, T)]$$

- “Convex optimization” problem
- Straight and non-intersecting trajectories

3. OT-Flow

Pontryagin Maximum Principle : $\exists \Phi$ s.t. $g(x, t; \theta) = -\nabla \Phi(x, t; \theta)$

→ Learning of the potential instead

Hamilton-Jacobi-Bellman equation (optimal controle theory):

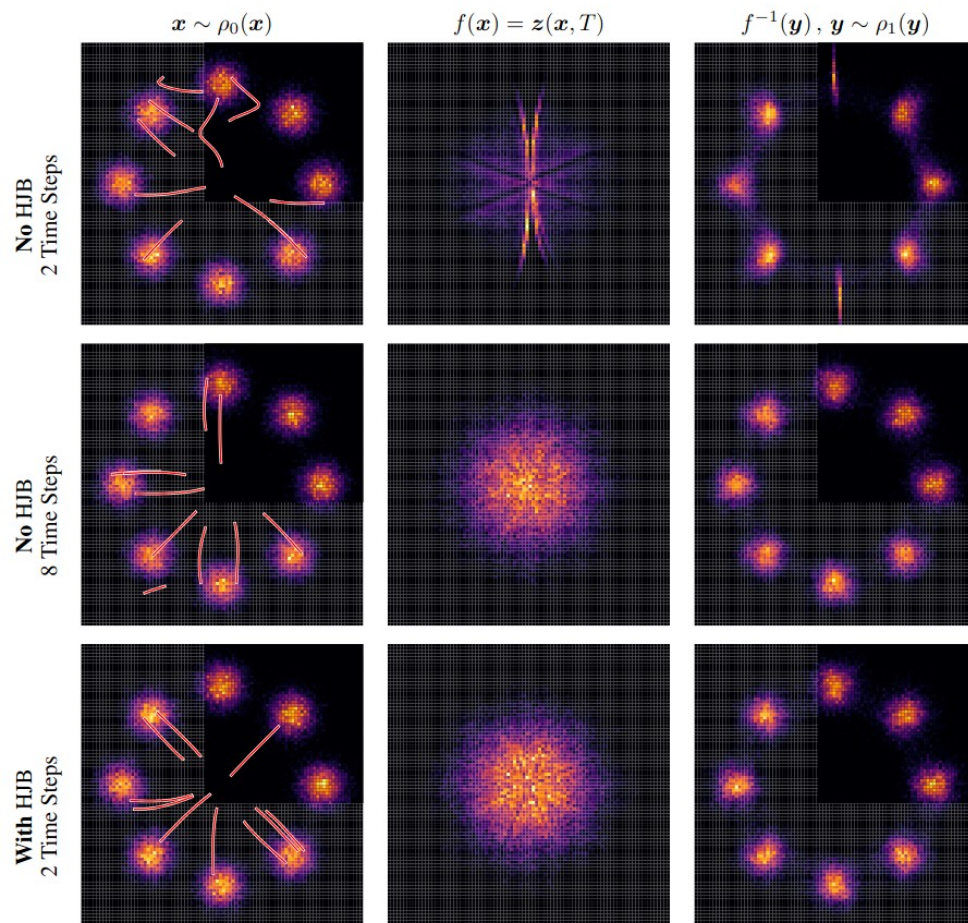
$$-\partial_t \Phi(\mathbf{x}, t) + \frac{1}{2} \|\nabla \Phi(\mathbf{z}(\mathbf{x}, t), t)\|^2 = 0$$

$$\begin{aligned} \Phi(\mathbf{x}, T) = & 1 + \log(\rho_0(\mathbf{x})) - \log(\rho_1(\mathbf{z}(\mathbf{x}, T))) \\ & - \ell(\mathbf{z}(\mathbf{x}, T), T). \end{aligned}$$

Additional constraint :

$$R(\mathbf{x}, T) = \int_0^T \left| \partial_t \Phi(\mathbf{z}(\mathbf{x}, t), t) - \frac{1}{2} \|\nabla \Phi(\mathbf{z}(\mathbf{x}, t), t)\|^2 \right| dt.$$

3. OT-Flow



3. Implementation

Network architecture :

$$\Phi(\mathbf{s}; \boldsymbol{\theta}) = \mathbf{w}^\top N(\mathbf{s}; \boldsymbol{\theta}_N) + \frac{1}{2} \mathbf{s}^\top (\mathbf{A}^\top \mathbf{A}) \mathbf{s} + \mathbf{b}^\top \mathbf{s} + c$$

ResNet :

$$\begin{aligned} \mathbf{u}_0 &= \sigma(\mathbf{K}_0 \mathbf{s} + \mathbf{b}_0) \\ N(\mathbf{s}; \boldsymbol{\theta}_N) &= \mathbf{u}_1 = \mathbf{u}_0 + h \sigma(\mathbf{K}_1 \mathbf{u}_0 + \mathbf{b}_1) \end{aligned}$$

4. Score Computation

Exact trace computation :
$$\text{tr}(\nabla^2 \Phi(\mathbf{s}; \boldsymbol{\theta})) = \text{tr} \left(\mathbf{E}^\top \nabla_{\mathbf{s}}^2 (N(\mathbf{s}; \boldsymbol{\theta}_N) \mathbf{w}) \mathbf{E} \right) + \text{tr} \left(\mathbf{E}^\top (\mathbf{A}^\top \mathbf{A}) \mathbf{E} \right),$$

$$\text{tr} \left(\mathbf{E}^\top \nabla_{\mathbf{s}}^2 (N(\mathbf{s}; \boldsymbol{\theta}_N) \mathbf{w}) \mathbf{E} \right) = t_0 + h t_1, \text{ where}$$

$$t_0 = \left(\sigma''(\mathbf{K}_0 \mathbf{s} + \mathbf{b}_0) \odot \mathbf{z}_1 \right)^\top \left((\mathbf{K}_0 \mathbf{E}) \odot (\mathbf{K}_0 \mathbf{E}) \right) \mathbf{1}.$$

$$t_1 = \left(\sigma''(\mathbf{K}_1 \mathbf{u}_0 + \mathbf{b}_1) \odot \mathbf{w} \right)^\top \left((\mathbf{K}_1 \mathbf{J}) \odot (\mathbf{K}_1 \mathbf{J}) \right) \mathbf{1}$$