

SciML 2: Generative models

E. Franck

Master CSMI, Strasbourg University

01/09/2025

1. Principe et 1eres approches

Types d'Apprentissage

- Pour le moment on a considéré l'apprentissage supervisé:
 - On a des données d'entrées et de sorties
 - On cherche à apprendre une relation déterministe qui mappe les entrées aux sorties

$$\min_{\theta} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i; \theta), y_i)$$

- Autre type d'apprentissage:
 - Apprentissage non supervisé (déterminer des structures dans les données sans sorties associées)
 - Apprentissage par renforcement (déterminer une politique d'actions pour maximiser une récompense cumulative)
- **Apprentissage non supervisé:**
 - clustering,
 - réduction de dimensionnalité,
 - dictionnaires creux,
 - estimation de densité et modèles génératifs, etc.

Modèles génératifs

Objectives (Estimation de densité): Soit des données $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ issues d'une distribution inconnue de densité $p_{\text{target}}(\mathbf{x})$. On souhaite estimer cette distribution, c'est-à-dire construire un modèle $p_{\theta}(\mathbf{x})$ qui approxime $p_{\text{target}}(\mathbf{x})$.

Objectives (Modèles génératifs): Soit des données $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ issues d'une distribution inconnue de densité $p_{\text{target}}(\mathbf{x})$. Un modèle génératif q_{θ} est processus capable de générer de nouvelles données $\tilde{\mathbf{x}}$ selon une loi qui approche $p_{\text{target}}(\mathbf{x})$.

- Un modèle d'estimation de densité est aussi un modèle génératif.
- Un modèle génératif ne donne pas forcément une estimation de densité explicite. Il s'agit d'un processus d'échantillonnage. Il caractérise donc implicitement la densité.

Objectives (modèles génératif conditionnels): On peut aussi vouloir construire des modèles génératifs conditionnels. C'est-à-dire des modèles capables de générer des données \mathbf{x} conditionnellement à une entrée \mathbf{y} . On caractérise donc implicitement la densité conditionnelle $p_{\text{target}}(\mathbf{x} \mid \mathbf{y})$.

Applications

- Génération d'images réalistes (GANs, VAEs)
- Synthèse de texte (modèles de langage, GPT)
- Modélisation de la distribution de données complexes (flux normaux, modèles variationnels)
- Synthèse de molécules en chimie et biologie computationnelle
- Simulation de systèmes physiques (modèles basés sur la physique, réseaux de neurones génératifs)
- Amélioration de la résolution d'images (super-résolution)
- Traduction automatique et génération de discours

Qu'est ce qu'on minimise ?

- **Maximum de vraisemblance:** qui consiste à maximiser la probabilité des données observées sous le modèle.

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^N p_{\theta}(z_i)$$

- **Minimisation de la divergence entre la distribution réelle et la distribution modélisée:**

$$\theta^* = \arg \min_{\theta} D(p_{\text{target}} \parallel p_{\theta}) = \arg \min_{\theta} \int p_{\text{target}}(\mathbf{x}) \log \left(\frac{p_{\text{target}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) d\mathbf{x}$$

- Les deux approches sont équivalentes lorsque la divergence choisie est la divergence de Kullback-Leibler (KL).
- Pk ?
 - En pratique on a souvent pas accès à p_{target} mais seulement à des échantillons issus de cette distribution.
 - Kl vers Maximum de vraisemblance:
 - On approxime l'intégrale par une somme sur les échantillons.

$$\theta^* = \arg \min_{\theta} \int p_{\text{target}}(\mathbf{x}) \log \left(\frac{p_{\text{target}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) d\mathbf{x} \approx \arg \min_{\theta} \sum_{i=1}^N \log \left(\frac{p_{\text{target}}(z_i)}{p_{\theta}(z_i)} \right) = \arg \max_{\theta} \sum_{i=1}^N \log(p_{\theta}(z_i))$$

1er modèle de densité: mélange de Gaussiennes

- On modélise la densité comme une combinaison linéaire de K densités gaussiennes:

$$p_{\theta}(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

où $\theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$ sont les paramètres du modèle, avec π_k les poids des composantes (somme à 1), $\boldsymbol{\mu}_k$ les moyennes et $\boldsymbol{\Sigma}_k$ les matrices de covariance.

- On peut apprendre les paramètres via l'algorithme EM (Expectation-Maximization)
- Implicitement cet algorithme maximise la vraisemblance des données sous le modèle.

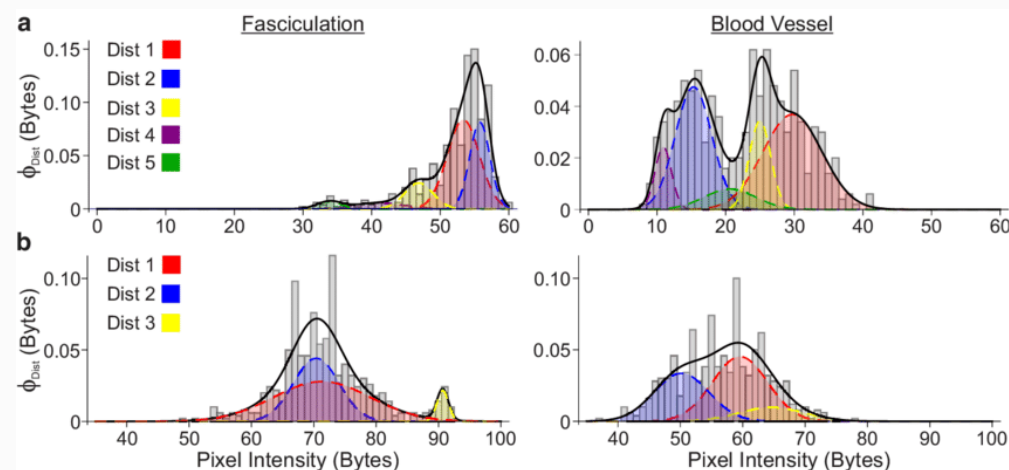


Figure 1: “Exemple de mélange de Gaussiennes pour modéliser une distribution complexe. Google Image”

1er modèle génératif: auto-encodeur variationnel (VAE)

- Un VAE est un modèle génératif probabiliste qui apprend une représentation latente des données.
- Il se compose de deux parties:
 - Un encodeur qui mappe les données d'entrée \mathbf{x} vers une distribution latente $q_{\phi}(\mathbf{z} | \mathbf{x})$.
 - Un décodeur qui mappe les échantillons de la distribution latente vers la distribution des données $p_{\theta}(\mathbf{x} | \mathbf{z})$.

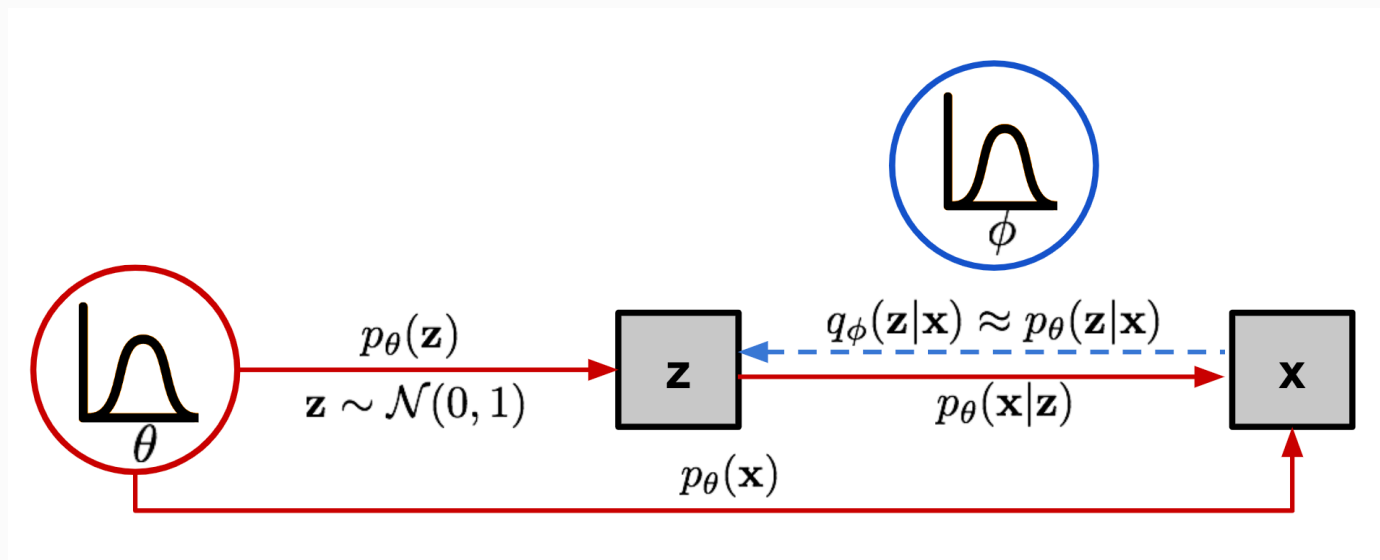


Figure 2: VAE. Lil Blog

- L'apprentissage consiste à maximiser la vraisemblance marginale des données. En pratique on utilise une estimation inférieure (ELBO).

2. Flot Normalisé

Idée principale

- Contrairement à d'autres méthodes les flots normalisés permettent de construire explicitement la densité.

Idea: Soit la densité cible $p_{\text{target}}(\mathbf{x})$. L'idée est de construire une transformation bijective $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ telle que si $\mathbf{z} \sim p_0(\mathbf{z})$ connue alors $\mathbf{x} = f(\mathbf{z})$ suit la densité cible $p_{\text{target}}(\mathbf{x})$.

- Par définition des loi de probabilité:

$$\int_{\mathbf{X}} p_{\text{target}}(\mathbf{x}) d\mathbf{x} = \int_{\mathbf{Z}} p_0(\mathbf{z}) d\mathbf{z} = \int_{f^{-1}(\mathbf{X})} p_0(\mathbf{z}) d\mathbf{z} = 1$$

- Par le changement de variable $\mathbf{z} = f^{-1}(\mathbf{x})$ on a:

$$\int_{\mathbf{X}} p_{\text{target}}(\mathbf{x}) d\mathbf{x} = \int_{\mathbf{X}} p_0(f^{-1}(\mathbf{x})) |\det J_{f^{-1}}(\mathbf{x})| d\mathbf{x} = 1$$

- On en déduit l'expression de la densité cible:

$$p_{\text{target}}(\mathbf{x}) = p_0(f^{-1}(\mathbf{x})) |\det J_{f^{-1}}(\mathbf{x})|$$

Principe 2

Remark: Si on a une transformation bijective f alors on peut donc calculer p_{target} à partir de p_0 et de f . Encore plus facilement on échantillonner de p_{target} en échantillonnant de p_0 et en appliquant f .

Idea: Une composition de transformations bijectives est aussi une transformation bijective. On peut donc construire des transformations complexes en composant des transformations simples.

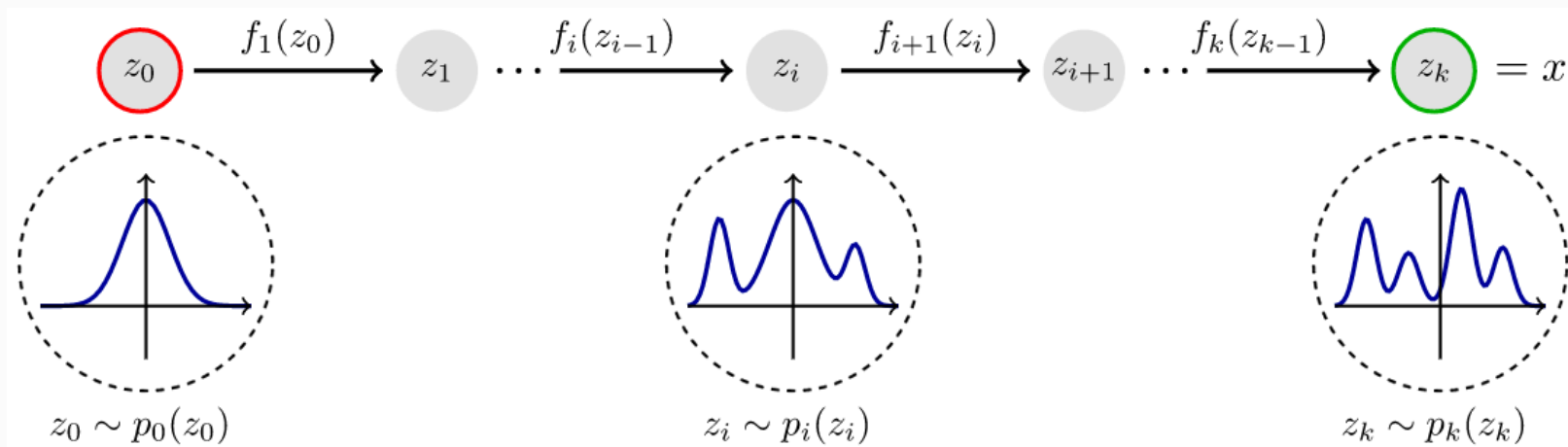


Figure 3: Flow normalisé. Lil Blog

Flot normalisé

Definition (Flots normalisés): Soit $p_0(\mathbf{z})$ une densité de probabilité simple. Un flot normalisé est une transformation bijective $f = f_L \circ \dots \circ f_1$ telle qu'un échantillon de la loi final $p_\theta(\mathbf{x})$ soit obtenue en échantillonnant $\mathbf{z} \sim p_0(\mathbf{z})$ puis en calculant $\mathbf{x} = f(\mathbf{z})$.

- On peut calculer la densité finale via la formule du changement de variable:

$$p_\theta(\mathbf{x}) = p_0(f^{-1}(\mathbf{x})) \prod_{i=1}^L |\det J_{f_i^{-1}}(\mathbf{h}_i)|$$

où $\mathbf{h}_i = f_{i-1} \circ \dots \circ f_1(\mathbf{z})$ et $\mathbf{h}_0 = \mathbf{z}$.

Proposition (log probabilité): La log probabilité de \mathbf{x} sous le modèle est donnée par:

$$\log p_\theta(\mathbf{x}) = \log p_0(f^{-1}(\mathbf{x})) + \sum_{i=1}^L \log |\det J_{f_i^{-1}}(\mathbf{h}_i)|$$

Apprentissage

- L'objectif est de minimiser la divergence entre la densité cible $p_{\text{target}}(\mathbf{x})$ et la densité modélisée $p_{\boldsymbol{\theta}}(\mathbf{x})$.

Proposition: Minimiser la divergence de KL revient à maximiser la log vraisemblance des données sous le modèle:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(p_{\text{target}} \parallel p_{\boldsymbol{\theta}}) = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{z}_i)$$

- La formule de la log probabilité permet de calculer facilement le gradient pour l'optimisation.
- On peut utiliser des algorithmes d'optimisation stochastique (SGD, Adam, etc.) pour apprendre les paramètres $\boldsymbol{\theta}$ des transformations f_i
- Comment choisir les transformations f_i ?

Remark (Besoin):

- transformations bijectives
- inversibles facilement calculables
- jacobien avec déterminant calculable efficacement

Analogie avec les EDO

- Soit $\mathbf{x}_0, \mathbf{x}_1 \in \mathbb{R}^d$. Une façon de construire une transformation bijective entre \mathbf{x}_0 et \mathbf{x}_1 est de considérer l'EDO:

$$\begin{cases} \frac{d\mathbf{x}(t)}{dt} = \mathbf{v}(\mathbf{x}(t)) \\ \mathbf{x}(0) = \mathbf{x}_0 \\ \mathbf{x}(1) = \mathbf{x}_1 \end{cases}$$

- En effet par le théorème de Cauchy-Lipschitz, si \mathbf{v} est suffisamment régulier le flot de cette EDO est réversible en temps: $\Phi_t(0, t, \mathbf{x}_0) = \Phi_{-t}(\Phi_t(\mathbf{x}_0)) = \mathbf{I}_d$

Idea:

- on paramétrise des flux d'EDO
- On discrétise l'EDO pour obtenir une couche
- Si on utilise schéma réversible ou calcul ou si peut calculer le flot exact on aura une couche **inversible**.

Exemple: Nice

- On se donne l'EDO

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}(t)),$$

- On split \mathbf{x} en deux parties $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ et obtient les deux sous équations

$$\frac{d}{dt} \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} = \begin{pmatrix} \mathbf{v}_{\boldsymbol{\theta},a}(\mathbf{x}) \\ 0 \end{pmatrix}, \quad \frac{d}{dt} \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} = \begin{pmatrix} 0 \\ \mathbf{v}_{\boldsymbol{\theta},b}(\mathbf{x}) \end{pmatrix}$$

- On choisit la paramétrisation:

$$\frac{d}{dt} \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} = \begin{pmatrix} \mathbf{v}_{\boldsymbol{\theta},a}(\mathbf{x}_b) \\ 0 \end{pmatrix}, \quad \frac{d}{dt} \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} = \begin{pmatrix} 0 \\ \mathbf{v}_{\boldsymbol{\theta},b}(\mathbf{x}_a) \end{pmatrix}$$

- On applique le schéma d'Euler explicite pour obtenir les transformations:

$$f_1 : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a + \Delta t \mathbf{v}_{\boldsymbol{\theta},a}(\mathbf{x}_b) \\ \mathbf{x}_b \end{pmatrix}, \quad f_2 : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b + \Delta t \mathbf{v}_{\boldsymbol{\theta},b}(\mathbf{x}_a) \end{pmatrix}$$

- Ses flot discrets sont exacts pour chaque sous système donc réversible en temps et donc inversible.

Exemple: Nice II

Definition (Couche Nice): Une couche Nice est une transformation bijective $f = f_2 \circ f_1$ où f_1 et f_2 sont définies par:

$$f_1 : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a + \mathbf{v}_{\theta,a}(\mathbf{x}_b) \\ \mathbf{x}_b \end{pmatrix}, \quad f_2 : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b + \mathbf{v}_{\theta,b}(\mathbf{x}_a) \end{pmatrix}$$

avec $\mathbf{v}_{\theta,a}$ et $\mathbf{v}_{\theta,b}$ des réseaux de neurones.

- Le déterminant de la jacobienne est facile à calculer: $\det J_f(\mathbf{x}) = 1$
- L'inverse est aussi facile à calculer:

$$f_1^{-1} : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a - \mathbf{v}_{\theta,a}(\mathbf{x}_b) \\ \mathbf{x}_b \end{pmatrix}, \quad f_2^{-1} : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b - \mathbf{v}_{\theta,b}(\mathbf{x}_a) \end{pmatrix}$$

- Globalement en cherchant des splitting ou des EDO nilpotentes on peut construire des transformations bijectives avec des jacobienes à déterminants calculables efficacement.

Exemple: RealNVP

$$\bullet \quad \frac{d}{dt} \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{\theta,a}(\mathbf{x}_b) \odot \mathbf{x}_a + \mathbf{s}_{\theta,a}(\mathbf{x}_b) \\ 0 \end{pmatrix}, \quad \frac{d}{dt} \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} = \begin{pmatrix} 0 \\ \mathbf{w}_{\theta,b}(\mathbf{x}_a) \odot \mathbf{x}_b + \mathbf{s}_{\theta,b}(\mathbf{x}_a) \end{pmatrix}$$

On va considère que la 1er couche mais la situation est similaire pour la seconde couche. On a

$$\frac{d}{dt} \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{\theta,a}(\mathbf{x}_b) \odot \mathbf{x}_a + \mathbf{s}_{\theta,a}(\mathbf{x}_b) \\ 0 \end{pmatrix}$$

-On résout l'équation sur \mathbf{x}_a en considérant \mathbf{x}_b comme une constante cela donne

$$\mathbf{x}_a(t) = \exp(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))t) \odot \mathbf{x}_a(0) + \int_0^{\Delta t} \exp(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))\tau) \odot \mathbf{s}_{\theta,a}(\mathbf{x}_b) d\tau$$

• On utilise le fait que si $\mathbf{w}_{\theta,a}(\mathbf{x}_b)$ est inversible on a

$$\int_0^{\Delta t} \exp(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))\tau) d\tau = \text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))^{-1} \odot (\exp(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))\Delta t) - \mathbf{1})$$

Exemple: RealNVP

- ce qui donne

$$\mathbf{x}_a(t) = \exp(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))t) \odot \mathbf{x}_a(0) + \left(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))^{-1} \odot (\exp(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))t) - \mathbf{1}) \right) \odot \mathbf{s}_{\theta,a}(\mathbf{x}_b)$$

- Vu que dans l'idée \mathbf{s} et \mathbf{w} sont des réseaux de neurones on pose $\mathbf{t}_{\theta,a}(\mathbf{x}_b) = \left(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))^{-1} \odot (\exp(\text{diag}(\mathbf{w}_{\theta,a}(\mathbf{x}_b))t) - \mathbf{1}) \right) \odot \mathbf{s}_{\theta,a}(\mathbf{x}_b)$. En faisant ça sur les deux équations on obtient les transformations:

$$f_1 : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \exp(\mathbf{w}_{\theta,a}(\mathbf{x}_b)) \odot \mathbf{x}_a + \mathbf{t}_{\theta,a}(\mathbf{x}_b) \\ \mathbf{x}_b \end{pmatrix}, \quad f_2 : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a \\ \exp(\mathbf{w}_{\theta,b}(\mathbf{x}_a)) \odot \mathbf{x}_b + \mathbf{t}_{\theta,b}(\mathbf{x}_a) \end{pmatrix}$$

qui sont obtenues en prenant le flot exact d'une EDO ce qui assure l'inversibilité.

Definition (Couche RealNVP): Une couche RealNVP est une transformation bijective $f = f_2 \circ f_1$ où f_1 et f_2 sont définies par:

$$f_1 : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a \odot \exp(\mathbf{w}_{\theta,a}(\mathbf{x}_b)) + \mathbf{t}_{\theta,a}(\mathbf{x}_b) \\ \mathbf{x}_b \end{pmatrix}, \quad f_2 : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \odot \exp(\mathbf{w}_{\theta,b}(\mathbf{x}_a)) + \mathbf{t}_{\theta,b}(\mathbf{x}_a) \end{pmatrix}$$

avec $\mathbf{w}_{\theta,a}$, $\mathbf{w}_{\theta,b}$, $\mathbf{t}_{\theta,a}$ et $\mathbf{t}_{\theta,b}$ des réseaux de neurones.

Exemple: RealNVP III

- Le déterminant de la jacobienne est facile à calculer: $\det J_{f_1}(\mathbf{x}) = \exp(\sum(s_{\theta,a}(\mathbf{x}_b)))$
- L'inverse est aussi facile à calculer:

$$f_1^{-1} : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} (\mathbf{x}_a - t_{\theta,a}(\mathbf{x}_b)) \odot \exp(-s_{\theta,a}(\mathbf{x}_b)) \\ \mathbf{x}_b \end{pmatrix}, \quad f_2^{-1} : \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x}_a \\ (\mathbf{x}_b - t_{\theta,b}(\mathbf{x}_a)) \odot \exp(-s_{\theta,b}(\mathbf{x}_a)) \end{pmatrix}$$

- Example:

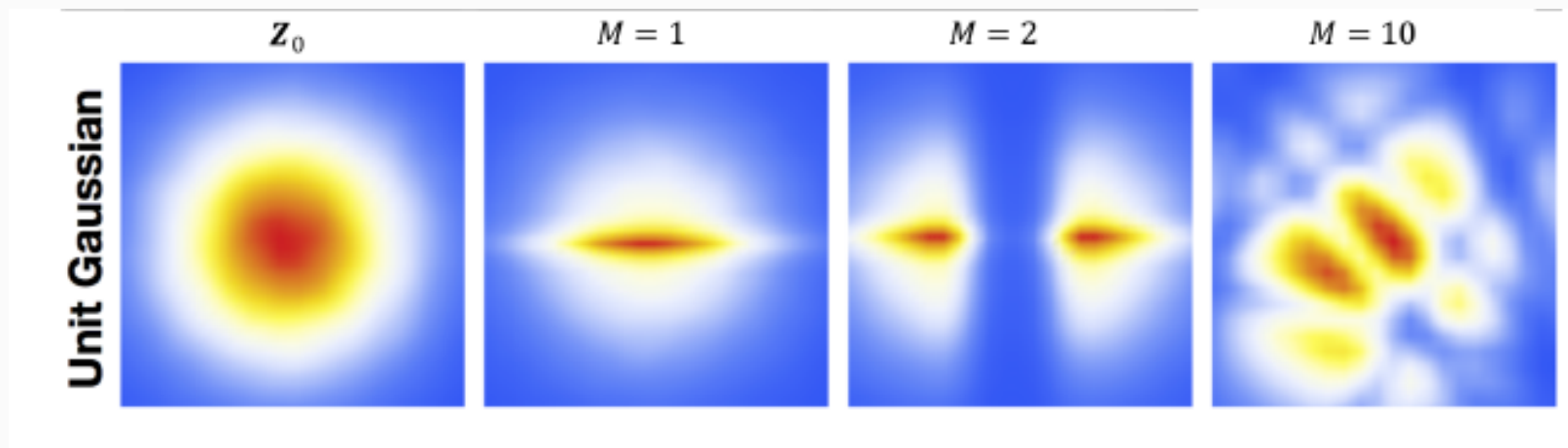


Figure 4: RealNVP. Lil Blog

Exemple: KRNet I

- idée: utiliser une transformation triangulaire:

$$\mathbf{y} = f(\mathbf{x}) = \begin{pmatrix} f_1(x_1) \\ f_2(x_1, x_2) \\ \dots \\ f_d(x_1, x_2, \dots, x_d) \end{pmatrix}$$

- Le réseau va donc enchaîner des transformations qui vont utiliser de plus en plus de variables.

$$f(\mathbf{x}) = f_N \circ f_{O,d}(\mathbf{x}) \circ f_{O,d-1}(\mathbf{x}) \circ \dots \circ f_{O,1}(\mathbf{x})$$

- L'enjeu est maintenant de construire des transformations $f_{O,i}$ bijectives avec Jacobienne à déterminant calculable efficacement.
- La transformation $f_{O,k}$ ne modifie que les composantes 1 à k

Definition (Couche de Squeezing): Une couche de Squeezing est donnée par $S_k(\mathbf{x}) = \mathbf{q} \odot \mathbf{x}$ avec $\mathbf{q} = \left(\underbrace{(1, \dots, 1)}_k, \underbrace{(0, \dots, 0)}_{d-k} \right)$ est une transformation bijective.

Exemple: KRNet II

Definition (Couche de rotation): Une couche de rotation est une transformation bijective R_{θ} définie par:

$$R_{\theta}(\mathbf{x}) = \begin{pmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_d \end{pmatrix} \begin{pmatrix} \mathbf{U} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_d \end{pmatrix} \mathbf{x}$$

avec \mathbf{L} , \mathbf{U} entraînable

Definition (Block interne): Soit $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ in $(\mathbb{R}^k, \mathbb{R}^{d-k})$ un bloc interne $f_i(\mathbf{x})$ est la compose d'une couche d'échelle donnée par:

$$f_{\text{sca}}(\mathbf{x}) = \mathbf{a} \odot \mathbf{x} + \mathbf{b}$$

avec \mathbf{a} et \mathbf{b} apprennable et d'une couche affine

$$f_{\text{aff}} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \odot (1 + \alpha \tanh(t_{\theta}(\mathbf{x}_a)) + e^{\beta} \odot \tanh(t_{\theta}(\mathbf{x}_a))) \end{pmatrix}$$

avec α un paramètre β entraînable et t_{θ} un réseau

Exemple: KRNet III

Definition (Block externe): Soit $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ in $(\mathbb{R}^k, \mathbb{R}^{d-k})$ un bloc externe

$$f_{O,k}(\mathbf{x}) = S_k \circ f_{i,L} \dots \circ f_{i,1} \circ R_{\theta}$$

- Il reste à définir la couche final L_N qui introduit une nonlinéarité assez particulière, appliqué composante par composante donné par :

$$f_{N(x)} = \begin{cases} \beta(x - a) + a, & x \in (-\infty, -a) \\ \varphi^{-1} \circ F \circ \varphi & \\ \beta(x + a) - a, & x \in (a, \infty) \end{cases}$$

avec φ affine et

$$F(x) = \int_0^x f(y) dy$$

Exemple: Residual approach I

- On se donne l'EDO

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{v}_{\theta}(\mathbf{x}(t)),$$

- On approxime le flot de cette EDO via un schéma de type Euler explicite:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \Delta t \mathbf{v}_{\theta}(\mathbf{x}_i)$$

- Si on veut inverser on résout $g(\mathbf{x}_i) = \mathbf{x}_i$ par une méthode de **Point fixe** avec

$$g(\mathbf{x}_i) = \mathbf{x}_{i+1} - \Delta t \mathbf{v}_{\theta}(\mathbf{x}_i)$$

- Pour converger on a besoin que Δt soit suffisamment petit et que \mathbf{v}_{θ} soit Lipschitzien avec constante de Lipschitz L telle que $\Delta t L < 1$.

Definition (Couche résiduelle): Une couche résiduelle est une transformation bijective f définie par:

$$f : \mathbf{x} \mapsto \mathbf{x} + \mathbf{v}_{\theta}(\mathbf{x})$$

où \mathbf{v}_{θ} est un réseau de neurones **contractant**

Exemple: Residual approach II

- Comment obtenir une paramétrisation contractante ?

Definition (Spectral normalized layer): Une couche à normalisation spectrale est une couche linéaire $w_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ de la form

$$\mathbf{x}_{l+1} = W_{s,\theta} \mathbf{x}_l$$

avec $W_{s,\theta} = \frac{W}{\Sigma(W)}$ ou $\Sigma(W)$ est la plus grande valeur singulière de W et W la matrice de poids associée à la couche.

- La valeur singulière maximale est approximée via une version différentiable de méthode de la puissance.
- Ici on contraint chaque couche a être contractante avec constante de Lipschitz $L = 1$.

Exemple: Neural Splins flows I

Idea: On applique composante par composante une fonction monotone, bijective et différentiable construite comme une spline rationnelle quadratique.

- On se donne $[-c, c]$
- On se donne des noeuds $b_1 < b_2 < \dots < b_{K-1}$ dans $[-c, c]$ avec $b_0 = -c$ et $b_K = c$
- On se donne des images $y_1 < y_2 < \dots < y_{K-1}$ dans $[-c, c]$ avec $y_0 = -c$ et $y_K = c$
- On se donne des dérivées d_0, d_1, \dots, d_K en chaque noeud
- On définit la fonction spline par morceaux:

$$S(x) = \frac{(y_{i+1} - y_i)(x - b_i)^2 + d_i(x - b_i)(b_{i+1} - x)}{(b_{i+1} - b_i)^2 + \left(d_{i+1} + d_i - 2\frac{y_{i+1} - y_i}{b_{i+1} - b_i}\right)(x - b_i)(b_{i+1} - x)}$$

- Determinant jacobienne facile à calculer:
- Inverse facile à calculer via résolution d'une équation quadratique.

Exemple: Neural Splines flows II

Soit une spline vectorielle

$$\mathbf{S} = \begin{pmatrix} S_1(\mathbf{x}_1) \\ S_2(\mathbf{x}_2) \\ \vdots \\ S_{\frac{d}{2}}(\mathbf{x}_{\frac{d}{2}}) \end{pmatrix}$$

Definition (Couche Neural Spline Flow): Une couche Neural Spline Flow est une transformation bijective f définie par une spline vectorielle:

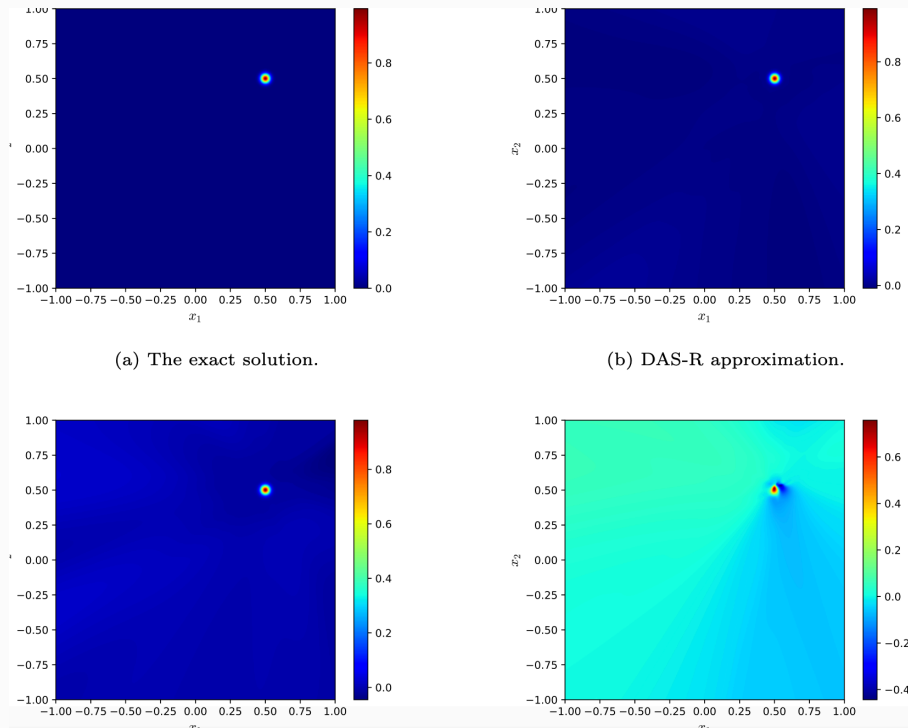
$$\mathbf{x} = \mathbf{S}(\mathbf{x}; (\mathbf{b}, \mathbf{y}, \mathbf{d})_{\mathbf{b}})$$

ou les paramètres de toutes les splines sont données par $(\mathbf{b}, \mathbf{y}, \mathbf{d})_{\mathbf{b}} = \text{MLP}_{\boldsymbol{\theta}}(\mathbf{x}_{\mathbf{a}})$ et $(\mathbf{b}, \mathbf{y}, \mathbf{d})_{\mathbf{a}}$ entraînables.

Applications aux PINNs

- Dans la version primaire des PINNs: on échantillonne uniformément les PINNs.

Remark: Pour des problèmes fortement localisé ça nécessite beaucoup de points. En pratique on aimerait échantillonner là où il faut.



- On veut échantillonner plus de point là où:
 - le résidu $|r(\mathbf{u}_x)|^2$ est grand
 - le gradient $|\nabla_x \mathbf{u}_x|^2$ est grand
 - un des critères précédent sur une autre variable
- On nomme ce critère $h(\mathbf{u}_\theta)$
- $h(\mathbf{u}_\theta)$ n'est pas une distribution de probabilité ?
comment échantillonner ?

Applications aux PINNs

Idea: apprendre à partir de la mesure de l'entraînement une densité $p_{\theta_p}(\mathbf{x})$ qui satisfait

$$\theta_p^* = \arg \min_{\theta_p} \text{KL}(h(\mathbf{u}_{\theta}) \parallel p_{\theta_p})$$

- Ici $p_{\theta_p} = \mathbb{1}_{\Omega} \circ \text{FN}_{\theta_p}$ avec FN_{θ_p} un flot normalisé et $\mathbb{1}_{\Omega}$ un algorithme qui rejette les points qui ne sont pas dans le domaine.
- On développe:

$$\arg \min_{\theta_p} \text{KL}(h(\mathbf{u}_{\theta})(\mathbf{x}) \parallel p_{\theta_p}(\mathbf{x})) = \arg \min_{\theta_p} \int_{\Omega} h(\mathbf{u}_{\theta})(\mathbf{x}) \log \left(\frac{h(\mathbf{u}_{\theta})(\mathbf{x})}{p_{\theta_p}(\mathbf{x})} \right) d\mathbf{x}$$

$$= \arg \min_{\theta_p} \int_{\Omega} h(\mathbf{u}_{\theta}) \log(h(\mathbf{u}_{\theta})) d\mathbf{x} - \arg \min_{\theta_p} \int_{\Omega} h(\mathbf{u}_{\theta}) \log(p_{\theta_p}) d\mathbf{x} = \arg \min_{\theta_p} - \int_{\Omega} h(\mathbf{u}_{\theta}) \log(p_{\theta_p}) d\mathbf{x}$$

- On utilise le fait que

$$\int_{\Omega} h(\mathbf{u}_{\theta}) \log(p_{\theta_p}) d\mathbf{x} = \mathbb{E}_{\mathcal{U}} [h(\mathbf{u}_{\theta}) \log(p_{\theta_p})] = \mathbb{E}_{p_{\theta_p}} \left[\frac{h(\mathbf{u}_{\theta})}{p_{\theta_p}} \log(p_{\theta_p}) \right]$$

Applications aux PINNs II

- On applique la méthode de Monte Carlo et donc on va minimiser

$$J_p(\theta_p) = -\frac{1}{n} \sum_{i=1}^n \frac{h(\mathbf{u}_\theta)(\mathbf{x}_i)}{p_{\theta_p}(\mathbf{x}_i)} \log(p_{\theta_p}(\mathbf{x}_i))$$

avec des points échantillonner par p_{θ_p} .

- Coté PINNs normalement on minimise

$$J(\theta) = \int_{\Omega} r(\mathbf{u}_\theta)(\mathbf{x})^2 d\mathbf{x} = \mathbb{E}_{\mathcal{U}} [r(\mathbf{u}_\theta)^2] = \mathbb{E}_{p_{\theta_p}} \left[\frac{r(\mathbf{u}_\theta)^2}{p_{\theta_p}} \right]$$

- Une fois discrétiser avec Monte Carlo on obtient le problème de Minimisation

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \frac{r(\mathbf{u}_\theta)^2(\mathbf{x}_i)}{p_{\theta_p}(\mathbf{x}_i)}$$

Remark:

- Dans l'algorithme on réentraîne régulière (toute les k itérations de gradient) le flot normalisé pour sampler de mieux en mieux
- Soit p_{θ_p} est combiné avec une loi uniforme pour assurer de garder un peu de point partout.

Applications aux PINNs III

- On peut proposer d'autre type d'optimisation.
- Par exemple:

$$\boldsymbol{\theta}^*, p_{\boldsymbol{\theta}_p}^* = \arg \min_{\boldsymbol{\theta}} \max_{p_{\boldsymbol{\theta}_p} \in V} \int_{\Omega} r(\mathbf{u}_{\boldsymbol{\theta}})^2(\mathbf{x}) p_{\boldsymbol{\theta}_p}(\mathbf{x}) d(\mathbf{x})$$

avec $V = \left\{ p_{\boldsymbol{\theta}_p}, \int p_{\boldsymbol{\theta}_p}(\mathbf{x}) d\mathbf{x} = 1, p_{\boldsymbol{\theta}_p}(\mathbf{x}) > 0 \right\}$

Remark:

- On peut montrer que résoudre ce problème fait tendre le résidu et

$$\frac{r(\mathbf{u}_{\boldsymbol{\theta}})^2(\mathbf{x})}{\int_{\Omega} r(\mathbf{u}_{\boldsymbol{\theta}})^2(\mathbf{x}) d\mathbf{x}} \rightarrow \mu_{\Omega}$$

avec μ_{Ω} la mesure de la loi uniforme. Cela pousse donc l'apprentissage à minimiser le plus uniformément possible le résidu

• Pour que le théorème marche on a besoin que la constante de Lipschitz de $p_{\boldsymbol{\theta}_p}$ soit plus petite que 1.

Applications aux PINNs IV

On résout donc le problème modifié:

$$\boldsymbol{\theta}^*, \boldsymbol{\theta}_p^* = \arg \min_{\boldsymbol{\theta}} \max_{\boldsymbol{\theta}_p} \int_{\Omega} r(\mathbf{u}_{\boldsymbol{\theta}})^2(\mathbf{x}) p_{\boldsymbol{\theta}_p}(\mathbf{x}) d(\mathbf{x}) - \beta \int |\nabla_{\mathbf{x}} p_{\boldsymbol{\theta}_p}(\mathbf{x})|^2 d\mathbf{x}$$

Remark: On va alterner l'optimisation selon $\boldsymbol{\theta}$ et $\boldsymbol{\theta}_p$

- Pour $\boldsymbol{\theta}$ on va minimiser:

$$\int_{\Omega} r(\mathbf{u}_{\boldsymbol{\theta}})^2(\mathbf{x}) p_{\boldsymbol{\theta}_p}(\mathbf{x}) \approx \sum_{i=1}^n r(\mathbf{u}_{\boldsymbol{\theta}})^2(\mathbf{x}_i)$$

- Pour $\boldsymbol{\theta}_p$ on va maximiser

$$\int_{\Omega} r(\mathbf{u}_{\boldsymbol{\theta}})^2(\mathbf{x}) p_{\boldsymbol{\theta}_p}(\mathbf{x}) d(\mathbf{x}) - \beta \int |\nabla_{\mathbf{x}} p_{\boldsymbol{\theta}_p}(\mathbf{x})|^2 d\mathbf{x} \approx \sum_{i=1}^n \frac{r(\mathbf{u}_{\boldsymbol{\theta}})^2(\mathbf{x}_i)}{p_{\boldsymbol{\theta}_p}^k(\mathbf{x}_i)} p_{\boldsymbol{\theta}_p}(\mathbf{x}_i) - \sum_{i=1}^n \frac{|\nabla_{\mathbf{x}} p_{\boldsymbol{\theta}_p}(\mathbf{x})|^2}{p_{\boldsymbol{\theta}_p}^k(\mathbf{x}_i)}$$

ou $p_{\boldsymbol{\theta}_p}^k(\mathbf{x}_i)$ est la distribution à l'itération d'avant utilisée pour échantillonner

3. Flot Normalisé continue

version continue

- On a vu que les flot normalisé se comprenaient comme des discrétisations réversibles d'EDO.
- On peut donc directement travailler avec les EDO.

Soit $\mathbf{x}_{\text{target}}$ un échantillon de la distribution cible p_{target} et \mathbf{z}_0 un échantillon de la distribution simple p_0 .

- On considère l'EDO:

$$\begin{cases} \frac{d\mathbf{x}(t)}{dt} = \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}(t), t) \\ \mathbf{x}(0) = \mathbf{z}_0 \end{cases}$$

Remark: Si on a une particule qui suit l'équation précédente alors la densité $p(\mathbf{x}, t)$ de la particule à l'instant t satisfait l'équation de continuité:

$$\frac{\partial p(\mathbf{x}, t)}{\partial t} + \nabla_{\mathbf{x}} \cdot (p(\mathbf{x}, t) \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}, t)) = 0$$

- On développe

$$\frac{\partial p(\mathbf{x}, t)}{\partial t} = -p(\mathbf{x}, t) \nabla_{\mathbf{x}} \cdot \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}, t) \cdot \nabla_{\mathbf{x}} p(\mathbf{x}, t)$$

Flot Normalisé continu

- Maintenant on regarde l'évolution de la densité pour une particule $\mathbf{x}(t)$:

$$\frac{d}{dt}p(\mathbf{x}(t), t) = \frac{\partial p(\mathbf{x}, t)}{\partial t} + \frac{d\mathbf{x}(t)}{dt} \cdot \nabla_{\mathbf{x}} p(\mathbf{x}, t)$$

- Maintenant on introduit l'équation de continuité:

$$\frac{d}{dt}p(\mathbf{x}(t), t) = -p(\mathbf{x}(t), t) \nabla_{\mathbf{x}} \cdot \mathbf{v}_{\theta}(\mathbf{x}(t), t)$$

- On multiplie par $\frac{1}{p(\mathbf{x}(t), t)}$ et on intègre en temps:

$$\frac{d}{dt} \log p(\mathbf{x}(t), t) = -\nabla_{\mathbf{x}} \cdot \mathbf{v}_{\theta}(\mathbf{x}(t), t)$$

Proposition: En intégrant on obtient la log densité de la particule à l'instant t :

$$\log p(\mathbf{x}(t), t) = \log p_0(\mathbf{x}(0)) - \int_0^t \nabla_{\mathbf{x}} \cdot \mathbf{v}_{\theta}(\mathbf{x}(\tau), \tau) d\tau = \log p_0(\mathbf{z}_0)$$

Flot Normalisé continue

- En pratique on va donc minimiser

$$\theta^* = \arg \min_{\theta} \log p_{t=1}(\mathbf{x}_{\text{target}}) = \arg \min_{\theta} -\log p_0(\mathbf{x}(0)) + \int_0^1 \nabla_{\mathbf{x}} \cdot \mathbf{v}_{\theta}(\mathbf{x}(\tau), \tau) d\tau$$

- Cela nécessite de résoudre l'EDO avant pour obtenir $\mathbf{x}(\tau)$.
- Le calcul de la divergence de \mathbf{v}_{θ} peut être coûteux en haute dimension.
- Le calcul du gradient nécessite soit de calculer l'adjoint soit de différentier par retro-propagation dans le schéma en temps.

Remark: Pad facilement utilisable en pratique

- La méthode qu'on va introduire permet de conserver l'approche ODE en simplifiant l'entraînement.

4. Flow matching

Principe

Remark: Les méthodes de flow matching utilise une edo pour transformer une loi simple en une loi cible mais propose un entraînement nettement plus simple que précédemment

Objectives: On souhaite trouver $\mathbf{v}(\mathbf{x}, t)$ tel que la densité générée p_t tel que $p_{t=0} = p_0$ la densité de base et $p_{t=1} = p_{\text{target}}$ la densité cible. Pour cela on a accès a des échantillons de p_{target} d'échantillon \mathbf{x}_1 et on échantillonne \mathbf{x}_0 de p_0 .

- On va commencer par définir l'objectif d'entraînement de la méthode flow matching.

Definition (Fonction de coût de la méthode flow matching): On souhaite minimiser la fonction de coût:

$$\mathbb{E}_{\mathbf{x}(t) \sim p_t} [\| \mathbf{v}_{t, \theta}(\mathbf{x}(t)) - \mathbf{v}_t^*(\mathbf{x}(t)) \|^2]$$

En pratique le champ de vitesse cible \mathbf{v}^* est inconnu. L'enjeu de la méthode est de remplacer cette fonctin de coût

- par une fonction de coût calculable.

Construction I

- L'idée est d'utiliser la connaissance des échantillons de p_{target}
- On va définir la trajectoire de notre distribution à partir d'une densité conditionnelle $p_{t|1}(\mathbf{x} | \mathbf{x}_1)$

$$p_t(\mathbf{x}) = \int p_{t|1}(\mathbf{x} | \mathbf{x}_1) p_{\text{target}}(\mathbf{x}_1) d\mathbf{x}_1$$

avec $p_0(\mathbf{x} | \mathbf{x}_1) = p_0(\mathbf{x})$ et $p_1(\mathbf{x} | \mathbf{x}_1) = \delta(\mathbf{x} - \mathbf{x}_1)$

Definition (Chemin de probabilité conditionnel) : Le chemin de probabilité conditionnel est défini par la densité conditionnelle $p_{t|1}(\mathbf{x} | \mathbf{x}_1)$ telle que $p_0(\mathbf{x} | \mathbf{x}_1) = p_0(\mathbf{x})$ et $p_1(\mathbf{x} | \mathbf{x}_1) = \delta(\mathbf{x} - \mathbf{x}_1)$

- Par la suite on va
 - exprimer le champ de vitesse cible \mathbf{v}^* en fonction d'un champ conditionnel et de $p_{t|1}$
 - Montrer que le champ généré à partir du champ conditionnel est le même que celui généré par le champ de vitesse cible
 - Montrer que la fonction de coût peut s'écrire en fonction du champ conditionnel uniquement.
 - proposer des chemins conditionnels simples pour lesquels on peut calculer le champ conditionnel explicitement.
 - Entraîner le modèle via la fonction de coût simplifiée.

Construction II

- On suppose que la densité conditionnelle $p_{t|1}(\mathbf{x} | \mathbf{x}_1)$ satisfait l'équation de continuité:

$$\frac{\partial p_{t|1}(\mathbf{x} | \mathbf{x}_1)}{\partial t} + \nabla_{\mathbf{x}} \cdot (p_{t|1}(\mathbf{x} | \mathbf{x}_1) \mathbf{v}_t(\mathbf{x} | \mathbf{x}_1)) = 0$$

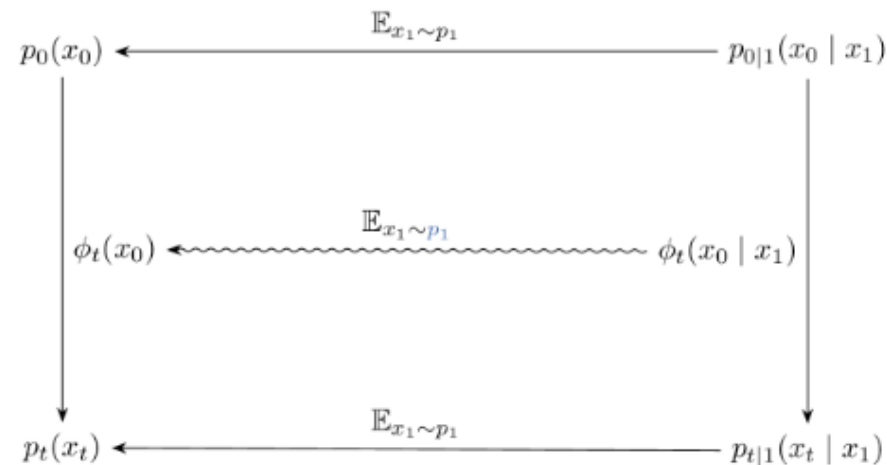


Figure 6: Schéma conditionnel. Cours flow matching Cambridge

- On peut aussi définir le champ de vitesse à partir du champ conditionnel:

$$\mathbf{v}_t(\mathbf{x}) = \mathbb{E}_{p_{1|t}(\mathbf{x}_1 | \mathbf{x})}[\mathbf{v}_t(\mathbf{x} | \mathbf{x}_1)] = \int \mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) p_{1|t}(\mathbf{x}_1 | \mathbf{x}) d\mathbf{x}_1$$

Proposition: Le champ de vitesse $\mathbf{v}_t(\mathbf{x})$ défini ci-dessus génère la même trajectoire de densité p_t que le champ de vitesse cible $\mathbf{v}_t^*(\mathbf{x})$.

Construction III

$$\begin{aligned}
\frac{d}{dt}p_t &= \int \frac{\partial p_{t|1}(\mathbf{x} | \mathbf{x}_1)}{\partial t} p_{\text{target}}(\mathbf{x}_1) d\mathbf{x}_1 \\
&= - \int \nabla_{\mathbf{x}} \cdot \left(p_{t|1}(\mathbf{x} | \mathbf{x}_1) \mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) \right) p_{\text{target}}(\mathbf{x}_1) d\mathbf{x}_1, \\
&= - \int \nabla_{\mathbf{x}} \cdot \left(p_{t|1}(\mathbf{x} | \mathbf{x}_1) \mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) p_{\text{target}}(\mathbf{x}_1) \right) d\mathbf{x}_1 \\
&= - \nabla_{\mathbf{x}} \cdot \left(\int p_{t|1}(\mathbf{x} | \mathbf{x}_1) \mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) p_{\text{target}}(\mathbf{x}_1) d\mathbf{x}_1 \right) \\
&= - \nabla_{\mathbf{x}} \cdot \left(\int \mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) \frac{p_{t|1}(\mathbf{x} | \mathbf{x}_1) p_{\text{target}}(\mathbf{x}_1)}{p_t(\mathbf{x})} p_t(\mathbf{x}) d\mathbf{x}_1 \right) \\
&= - \nabla_{\mathbf{x}} \cdot \left(\int \mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) p_{1|t}(\mathbf{x}_1 | \mathbf{x}) p_t(\mathbf{x}) d\mathbf{x}_1 \right)
\end{aligned}$$

- par le théorème de Bayes. Ensuite par la définition de u_t on a:

$$\frac{d}{dt}p_t = - \nabla_{\mathbf{x}} \cdot (\mathbf{v}_t(\mathbf{x}) p_t(\mathbf{x}))$$

Construction IV

- Maintenant qu'on a prouvé que le champ de vitesse générée par le champ conditionnel est le même que celui généré par le champ de vitesse cible on peut donc résoudre l'EDO associé au champ de vitesse conditionnel et prendre l'espérance conditionnelle pour obtenir la densité à l'instant t .

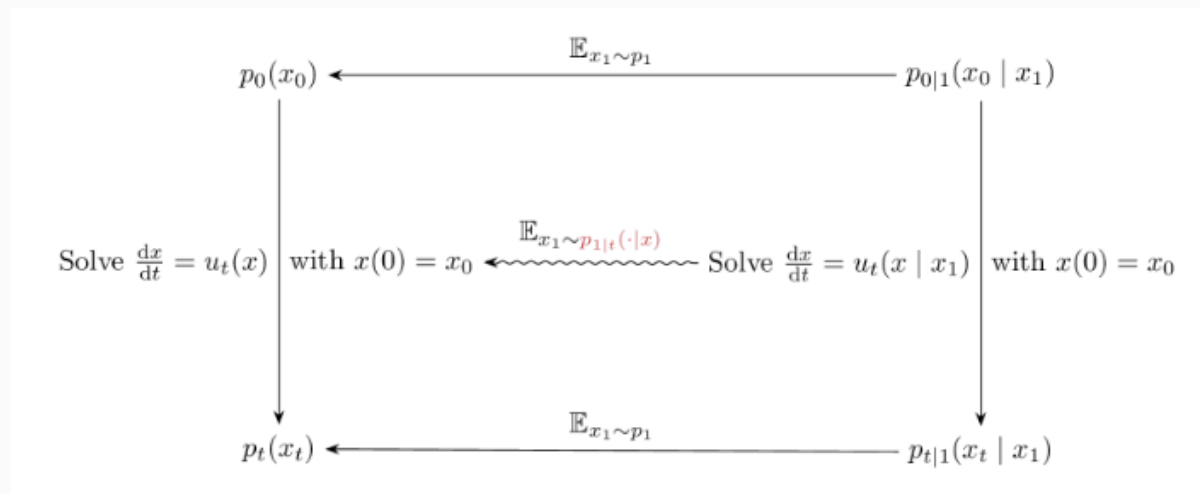


Figure 7: Schéma conditionnel. Cours flow matching Cambridge

- L'objectif de base était de minimiser

$$\mathcal{L}(\theta) = \mathbb{E}_{t \in \mathcal{U}[0,1], \mathbf{x} \sim p_t} [\| \mathbf{v}_{t,\theta}(\mathbf{x}_t) - \mathbf{v}_t^*(\mathbf{x}_t) \|^2]$$

- Puisque $\mathbf{v}(\mathbf{x}(t)) = \mathbb{E}_{p_{1|t}(\mathbf{x}_1 | \mathbf{x})}[\mathbf{v}_t(\mathbf{x} | \mathbf{x}_1)]$ génère la même trajectoire de densité que $\mathbf{v}_t^*(\mathbf{x})$ on peut minimiser

$$\mathcal{L}(\theta) = \mathbb{E}_{t \in \mathcal{U}[0,1], \mathbf{x} \sim p_t} [\| \mathbf{v}_{t,\theta}(\mathbf{x}_t) - \mathbf{v}(\mathbf{x}_t) \|^2]$$

Flow matching conditionnel

Theorem: Soit la fonction de coût

$$\mathcal{L}_{\text{FM}}(\boldsymbol{\theta}) = \mathbb{E}_{t \in \mathcal{U}[0,1], \mathbf{x} \sim p_t} [\| \mathbf{v}_{t,\boldsymbol{\theta}}(\mathbf{x}_t) - \mathbf{v}(\mathbf{x}_t) \|^2]$$

et la fonction de coût

$$\mathcal{L}_{\text{CM}}(\boldsymbol{\theta}) = \mathbb{E}_{t \in \mathcal{U}[0,1], \mathbf{x}_1 \sim p_{\text{target}}, \mathbf{x}_t \sim p_{t(\mathbf{x} | \mathbf{x}_1)}} [\| \mathbf{v}_{t,\boldsymbol{\theta}}(\mathbf{x}_t) - \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1) \|^2]$$

alors on a l'égalité

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{FM}}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{CM}}(\boldsymbol{\theta})$$

- En pratique on va donc minimiser la fonction de coût conditionnelle.

Remark: La magie de la méthode c'est en se donnant un chemin conditionnel on peut calculer explicitement le champ conditionnel $\mathbf{v}_t(\mathbf{x} | \mathbf{x}_1)$ et donc la fonction de coût

Chemin de probabilité conditionnel

Definition (Chemin de probabilité conditionnel gaussien): Le chemin de probabilité conditionnel gaussien est défini par la densité conditionnelle

$$p_{t|1}(\mathbf{x} | \mathbf{x}_1) = \mathcal{N}(\mathbf{x}; \mu_t(\mathbf{x}_1), \sigma^2(\mathbf{x}_1)\mathbf{I}_d)$$

Proposition: Le champ de vitesse conditionnel associé au chemin de probabilité conditionnel gaussien est:

$$\mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) = \frac{d\mu_t(\mathbf{x}_1)}{dt} + \frac{d\sigma_t(\mathbf{x}_1)}{dt} \sigma_t(\mathbf{x}_1)^{-1} (\mathbf{x} - \mu_t(\mathbf{x}_1))$$

• Preuve:

- Le chemin choisit implique que $\mathbf{x}_{t|1} = \mu_t(\mathbf{x}_1) + \sigma_t(\mathbf{x}_1)\mathbf{x}_0$ avec $\mathbf{x}_0 \sim p_0 = \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$
- On sait que $\frac{\partial p_{t|1}(\mathbf{x} | \mathbf{x}_1)}{\partial t} + \nabla_{\mathbf{x}} \cdot (p_{t|1}(\mathbf{x} | \mathbf{x}_1) \mathbf{v}_t(\mathbf{x} | \mathbf{x}_1)) = 0$ donc

$$\frac{d\mathbf{x}_{t|1}}{dt} = \mathbf{v}_t(\mathbf{x}_{t|1} | \mathbf{x}_1)$$

Chemin de probabilité conditionnel II

- On a donc

$$\frac{d\mathbf{x}_{t|1}}{dt} = \mathbf{v}_t(\mathbf{x}_{t|1} | \mathbf{x}_1)$$

- Il suffit de calculer la dérivée de $\mathbf{x}_{t|1}$ on obtient

$$\mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) = \frac{d\mu_t(\mathbf{x}_1)}{dt} + \frac{d\sigma_t(\mathbf{x}_1)}{dt} \mathbf{x}_0$$

- On utilise le chemin pour exprimer \mathbf{x}_0 en fonction de \mathbf{x} : $\mathbf{x}_0 = \sigma_t(\mathbf{x}_1)^{-1}(\mathbf{x} - \mu_t(\mathbf{x}_1))$.
- On obtient bien:

$$\mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) = \frac{d\mu_t(\mathbf{x}_1)}{dt} + \frac{d\sigma_t(\mathbf{x}_1)}{dt} \sigma_t(\mathbf{x}_1)^{-1}(\mathbf{x} - \mu_t(\mathbf{x}_1))$$

Remark: Un chemin classique est un chemin de probabilité conditionnel gaussien avec

$$\mu_t(\mathbf{x}_1) = t\mathbf{x}_1, \quad \sigma_t = \sigma_0 t + (1 - t)$$

avec μ_0 et σ_0 des hyperparamètres. On obtient

$$\mathbf{v}_t(\mathbf{x} | \mathbf{x}_1) = \mathbf{x}_1 - \mathbf{x} + \frac{\sigma_0 - 1}{\sigma_0 t + (1 - t)}(\mathbf{x} - t\mathbf{x}_1)$$

Algorithme

Algorithm:

- Initialiser les paramètres θ du modèle $\mathbf{v}_{t,\theta}(\mathbf{x})$
- Pour i de 1 à N_{iter} :
 - Échantillonner \mathbf{x}_1 de la distribution cible p_{target}
 - Échantillonner \mathbf{x}_0 de la distribution simple p_0
 - Échantillonner t uniformément dans $[0, 1]$
 - Calculer $\mathbf{x}_t = \mu_t(\mathbf{x}_1) + \sigma_t(\mathbf{x}_1)\mathbf{x}_0$
 - Calculer le champ conditionnel $\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1)$
 - Calculer la perte

$$\mathcal{L} = \| \mathbf{v}_{t,\theta}(\mathbf{x}_t) - \mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1) \|^2$$

- Mettre à jour les paramètres θ via une étape de descente de gradient sur \mathcal{L}

- On peut proposer plusieurs chemins conditionnels
- On peut traiter le cas conditionnel ou on cherche à générer une distribution conditionnelle $p_{\text{target}}(\mathbf{x} \mid \mathbf{c})$

Flot matching pour les PINNs

- Pour les pinns on construisait une densité qui approchait

$h(\mathbf{u}_\theta)$ (typiquement le résidu ou le gradient).

- Contrairement au flot normalisé on a pas accès à la densité qu'on construit donc on peut pas minimiser le KL.
- On a besoin échantillon et comment échantillonner selon $h(\mathbf{u}_\theta)$ (qui n'est pas une distribution).

Remark: On va utiliser une méthode de bootstrap à poids qui va nous permettre d'échantillonner selon une densité proportionnelle à $h(\mathbf{u}_\theta)$ sans connaître la densité. Ensuite on fera un entraînement de type flot matching.

- Bootstrap à poids:
 - On échantillonne n points uniformément dans le domaine \mathbf{x}_i .
 - On calcule les poids $w_i = \frac{h(\mathbf{u}_\theta)(\mathbf{x}_i)}{\sum_{j=1}^n h(\mathbf{u}_\theta)(\mathbf{x}_j)}$
 - On rééchantillonne n points avec remplacement selon les poids w_i pour obtenir des points $\mathbf{x}_i \sim p_{\theta_p}$
- Comment échantillonner avec les poids
 - on crée des variables cumulatives $c_i = \sum_{j=1}^i w_j$
 - on tire $p \in \mathcal{U}[0, 1]$ et on cherche i tel que $c_{i-1} \leq p < c_i$

Applications

- Super resolution d'images

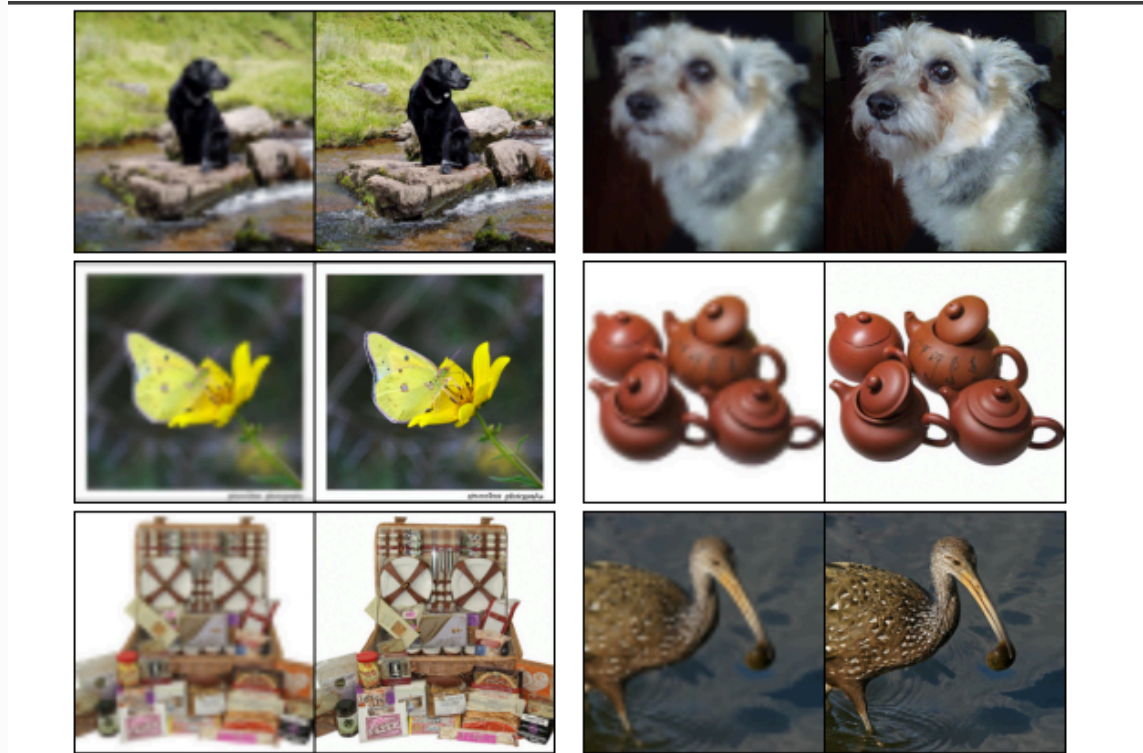


Figure 8: Super resolution d'images avec flow matching. Flow matching paper

- Génération de protéine
- Génération de turbulence etc