# IV. Continuous Normalizing flows (CNF)

- For NF are based on the change-of-variables.

$$x_N = \phi(z_0), \qquad x_0 \xrightarrow{\phi_0} x_1 \xrightarrow{\phi_1} \cdots \xrightarrow{\phi_N} x_N$$

$$P_i(x_i) = P_{i-1}(x_{i-1}) \left| \text{Det} \frac{\partial \phi_i}{\partial x_{i-1}} \right|^{-1} \implies P_N(x_N) = P_0(x_0) \prod_{i=1}^{N} \left| \text{Det} \frac{\partial \phi_i}{\partial x_{i-1}} \right|^{-1}$$

$P_N$ density is the "pushforward" of $P_0$ under $f = f_N \circ \cdots \circ f_1$   "finite" flow

Drawbacks:   1) $f$ has to be invertible bijection. this may be a limitation

2) full Jacobian needs to be computed $O(d^3)$: slow

3) can be solved by restricting Jacobian matrix $\left( \begin{smallmatrix} \blacksquare & \\ 0 & \blacksquare \end{smallmatrix} \right) \sim O(d)$
this reduces expressibility...

we will see that continuous Normalizing flows (CNF) solve issues 2,3)

[R. Chen et al. "Neural Ordinary differential equations" NeurIPS 2018.]

- consider "infinitesimal flows"

$$\phi = \phi_N \circ \cdots \circ \phi_1 \longrightarrow \phi = \phi_t \quad , \quad t \in [0,1]$$

"continuous time" replaces the flow index
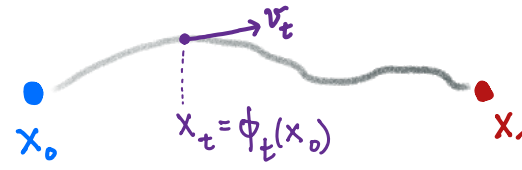
*auxiliary time*

the point of CNF is to Learn how to continuously evolve a simple distribution $q(z)$ into a complex one $P(z)$

- $\phi_t$ is a "trajectory" or "path" in time described by some <u>dynamics</u>

  we can write down an ODE

$$\begin{cases} \dfrac{d\phi_t(x)}{dt} = v_t(\phi_t(x)) \\[2mm] \phi_0(x) = x \end{cases}$$

$\phi_0 = id$

... velocity field
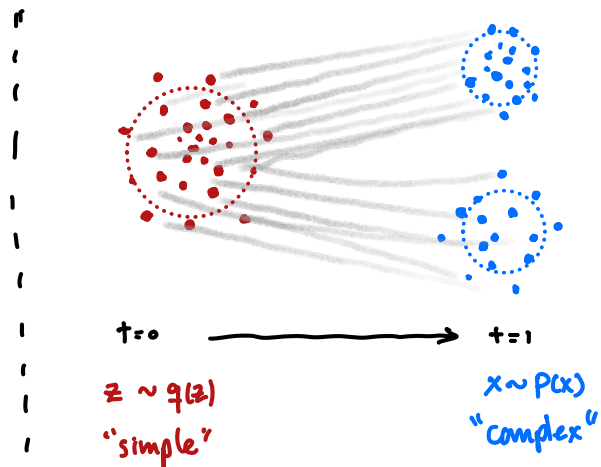


$v_t$

$x_0$    $x_t = \phi_t(x_0)$    $x_1$

the solution to this ODE is straight forward:

$$x_t = x_0 + \int_0^1 dt \; v_t(x_t)$$

$\Rightarrow$ Probability densities are now time-dependent $P_t(x)$.

we refer to $P_t(x)$ as the "probability path"

interpolating $\begin{cases} P_0 = q \\ P_1 = p \end{cases}$



$t=0$       $t=1$

$z \sim q(z)$      $x \sim p(x)$

"simple"       "complex"

- Probabilities need to be conserved over time under the vector field $v_t$

$\Rightarrow$

$$\boxed{\dfrac{\partial P_t(x)}{\partial t} + \nabla_x \cdot \dot{j}(x) = 0}$$

continuity equation or "Liouville" equation.

divergence    probability flux $\dot{j}(x) \equiv P_t(x) \, v_t(x)$

⑬

$$\Rightarrow \quad \frac{\partial P_t(x)}{\partial t} = - \sum_{i=1}^{d} \frac{\partial}{\partial x_i} \left( v_t^i(x) \, P_t(x) \right) \quad \longleftarrow$$ PDE give the evolution of the probability at a fixed point in space ⚠️

- We are interested in the change in density along trajectories $X_t$: the total derivative is:

$$d P_t(x_t) = \frac{\partial P_t(x_t)}{\partial x_t} \, dx_t + \frac{\partial P_t(x_t)}{\partial t} \, dt$$

we can use the ODE for $x_t$ and the continuity equation to get:

$$\frac{d P_t(x_t)}{dt} = \frac{\partial P_t(x_t)}{\partial x_t} v_t(x_t) - \nabla_x \cdot \left[ v_t(x_t) \, P_t(x_t) \right]$$

$$= \nabla_x \cdot \left[ v_t(x_t) \, P_t(x_t) \right] - (\nabla_x \cdot v_t(x_t)) \, P_t(x_t) - \nabla_x \cdot \left[ v_t(x_t) \, P_t(x_t) \right]$$

$$= - P_t(x_t) \, (\nabla_x \cdot v_t(x_t))$$

$$\Rightarrow \quad \frac{d}{dt} \log P_t(x_t) = - \nabla_x \cdot v_t(x_t)$$

$$\boxed{ \frac{d}{dt} \log P_t(x_t) = - \mathrm{Tr} \left[ \frac{\partial v_t}{\partial x} \right] } \quad \longleftarrow \text{ Trace of the Jacobian! matrix}$$

integrating yields the "instantaneous change-of-variables" formula:

$$\boxed{ \log P_t(x_t) = \log P_0(x_0) - \int_0^1 dt \, \mathrm{Tr} \left[ \frac{\partial v_t}{\partial x} \right] }$$

Notice that
$$\begin{cases} \boxed{NF} \longrightarrow \boxed{CNF} \\ Det\left(\frac{\partial v}{\partial x}\right) \longrightarrow Tr\left(\frac{\partial v}{\partial x}\right) \end{cases}$$

- Another way to understand why Det $\rightarrow$ Trace in the infinitesimal limit is that near the identity matrix the __determinant__ behaves like the __trace__ :
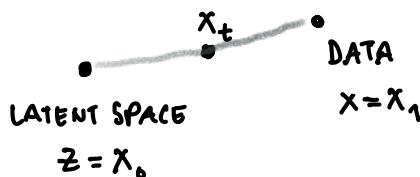
$$det(\mathbb{1} + \varepsilon A) \simeq \mathbb{1} + \varepsilon \, Tr(A) + \mathcal{O}(\varepsilon^2) \quad for \quad \varepsilon \rightarrow 0.$$

- CNF as Neural ODE (NODE) models

the CNF aims to model the continuous-time dynamical system that evolves a latent space distribution $q(z) = N(z|0,\mathbb{1})$ into the data with a "Neural ODE" :

$$\boxed{\frac{d}{dt} \phi_t(x) = v_t^\theta(\phi_t(x))}$$

$x_t$

DATA
$x = x_1$

LATENT SPACE
$z = x_0$

$$\begin{cases} P_0 = P_{base} = N(0,\mathbb{1}) \\ P_1 = P_{data} \end{cases}$$

where vector field $v_t^\theta$ is parametrized by a __Neural Network__.

- Solving the neural ODE yields the flow $\phi_1^\theta$ called a time-one map. this map transforms back and forth the data x into its latent rep. z :

$$\Rightarrow \begin{cases} x = \phi_1^\theta(z) = z + \int_0^1 dt \, v_t^\theta(x_t) \\ z = (\phi_1^\theta)^{-1}(x) = x - \int_0^1 dt \, v_t^\theta(x_t) \quad ........ \text{ inverse} \end{cases}$$
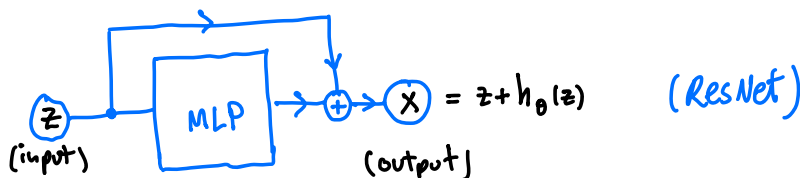
**Remark 1:**

unlike autoregressive NF, computing the inverse map for a CNF

has same complexity as forward direction !

**Remark 2:** Neural ODE can be thought as an infinitely deep

neural network. this can be seen by solving the ODE via Euler method:

$$X_t - X_{t-1} = \varepsilon \, v_t^\theta(X_n) \quad , \quad \varepsilon \ll 1$$

this has the same form as a Residual Neural Network (ResNet)
with a block $h_\theta = \varepsilon v_t^\theta$ .

[R. Chen et al. NeurIPS 2018]



$= z + h_\theta(z)$   (ResNet)

- **Training CNF's**

As with NF, we can train the CNF using Maximum likelihood
Estimation

$$\mathcal{L}(x,\theta) = \log P_\theta(x) = \log P_{base}(z) - \int_0^1 dt \, Tr\left[\frac{\partial v_t^\theta}{\partial x_t}\right]$$

**Remark 3:** computing the likelihood does not involve computing $Det(J) \sim \mathcal{O}(d^3)$

but instead $Tr[J]$ which has a better complexity of $\mathcal{O}(d)$ !

$\Rightarrow$ CNF allow for free-form Jacobian matrices

16

No restriction on architecture like for MAF or coupling flows.

$\longrightarrow$ More expressive

- In practice, we need to optimize $\mathcal{L}(\theta) \Rightarrow$ <u>Backpropagate</u> through a <u>numeric ODE solution</u> implemented via some algorithm:

$$\mathcal{L}(x,\theta) = \text{ODEsolve}(x,\theta) \qquad \nabla_\theta \mathcal{L}(x,\theta) = ? \quad \text{compute gradient}$$

$\Rightarrow$ Formulate <u>gradient computation</u> as a separate ODE in $\frac{\partial \mathcal{L}}{\partial x_t} \equiv a_t$ known as <u>adjoint</u> ODE it has solution:

$$\boxed{\nabla_\theta \mathcal{L} = -\int_1^0 dt \; a_t^T \cdot \frac{\partial v_t^\theta}{\partial \theta}}$$

solving this adjoint ODE backwards in time $t_1 \to t_0$

- <u>Fast trace computation:</u>

One can further improve the complexity of $\text{Tr}[J]$ by using Hutchinson's trace estimator:

suppose $\varepsilon$ is a noise vector such that $\begin{cases} \mathbb{E}(\varepsilon) = 0 \\ \text{Cov}(\varepsilon) = \mathbb{1}_{d\times d} \end{cases}$ e.g. $\varepsilon \sim N(0,\mathbb{1})$

$$\text{Tr}[A] = \text{Tr}\left[A \; \mathbb{E}(\varepsilon^T \varepsilon)\right] = \mathbb{E}\,\text{Tr}\left[\underbrace{\varepsilon^T A \varepsilon}_{\#}\right] = \mathbb{E}_{\varepsilon \sim p(\varepsilon)}\left[\varepsilon^T A \varepsilon\right]$$

$$\boxed{\text{Tr}[A] \approx \frac{1}{M} \sum_{i=1}^M \varepsilon_i^T A \varepsilon_i}$$ Monte Carlo

(17)

this is a stochastic estimator that and scales better!

exercise: show that $\mathbb{E}\left(\frac{1}{M}\sum_{\lambda=1}^{M} \mathcal{E}_\lambda^T A \mathcal{E}_\lambda\right) = Tr(A)$ i.e. estimation is unbiased

caveat: the jacobian matrix of the CNF flows are <u>not</u> fully <u>free-form</u>
in fact they happen to be positive definite matrices
i.e. all eigenvalues are positive. E.g. one can't write down
a function $f(x) = -x$ as a $\phi_1(x)$ time-one map

⟿ way out Augmented NODE (ANODE)

embed the data in higher dim space...
this lifts the topological obstruction!

• <u>issues with CNFs</u>

Training and Sampling the CNF model requires solving the neural
ODE, in the forward or backward direction, using Numerical ODE solvers

⟿ based on the Euler, Runge-Kutta, etc methods

→ expensive, needs many time-steps in order to give accurate
results

⟿ numerically unstable

→ In practice CNF's don't scale well to large datasets

MORAL OF STORY: WHATEVER IS GAINED IN EXPRESSIBILITY
OVER finite NF is lost IN PRACTICE!

# I. Flow-matching  [Lipman et al. ICLM 2023] {NEW}

- Flow-matching is simple training objective for CNF's that allows for scalable training ⟶ scales better than MLE objective and more stable...

$$\begin{cases} \dfrac{d\phi_t(x)}{dt} = u_t(\phi_t(x)) \quad , \quad \phi_0 = id \qquad (\phi_0(x) = x) \\[4mm] \dfrac{\partial P_t(x)}{\partial t} = -\nabla_x (u_t(x) \, P_t(x)) \end{cases}$$

- we say that $u_t$ generates the prob. paths $P_t(x)$ if the above eq. are satisfied.

the idea is to directly regress the velocity field $u_t$ with an MSE loss

$$\mathcal{L}_{FM} = \mathbb{E}_{\substack{t \sim U[0,1] \\ X \sim P_t(x)}} \| u_t^\theta(x) - u_t(x) \|^2 \quad \longleftarrow \text{ flow-matching! objective}$$

Neural Network

$t \sim U[0,1]$ is the uniform distribution.

__Huge benefit__: No need to solve ODE during training $u_t^\theta$ which usually requires going sequentially through time-steps (e.g. Euler method)

Here time can be sampled non-sequentially (Parallelized)...

109

- **problem:** How do we model the prob. path $P_t(x)$? What to take for $u_t$?

we only know that
$$\begin{cases} P_0(x) = P_{base} = N(0,1) \\ P_1(x) = P_{data} \end{cases} \quad \cdots\cdots \ \text{①}$$

- **solution?** Model simpler conditional/joint quantities that when marginalized give you back the quantities you are interested in.

  $\rightsquigarrow$ leads to "integral representations"!

Define **conditional probability** path $P_t(x|y)$ conditioned on some random variable $y$ such that:

$$\begin{cases} P_0(x|y) = P_{base}(x) = N(0,1) \quad (\text{independent of } y) \\ P_1(x|y) = \delta(x-y) \end{cases}$$

the conditional **prob. path** $P_t(x|y)$ interpolates between the std Gaussian at $t=0$ and the **delta function** centered around "$y$" at $t=1$.
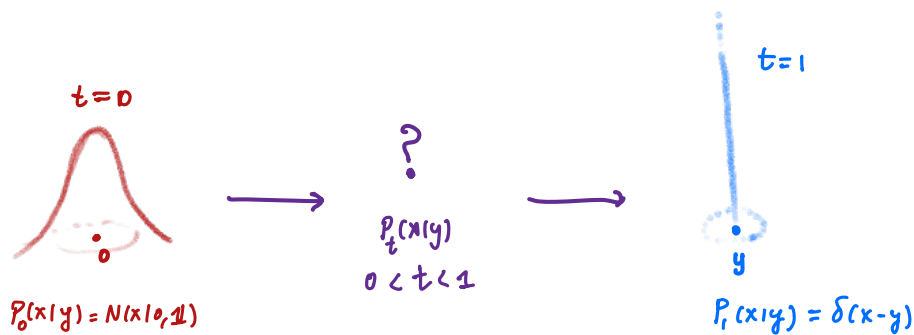
- Marginalizing over the data distribution:

$$P_t(x) = \int dx_1 \, P_t(x|x_1) \, P_{data}(x_1)$$

gives the correct boundary conditions ① above:

$$\begin{cases} \bullet \ t=0: \quad P_0(x) = \int dx_1 \, P_{base}(x) \, P_{data}(x_1) \\ \qquad\qquad = P_{base}(x) \cdot \underbrace{\int dx_1 \, P_{data}(x_1)}_{1} = P_{base}(x) \quad // \\ \\ \bullet \ t=1: \quad P_1(x) = \int dx_1 \, \delta(x-x_1) \, P_{data}(x_1) \\ \qquad\qquad = P_{data}(x) \quad // \end{cases}$$

- **Gaussian conditional probability paths:**

- Notice that the conditional paths $P_t(x|y)$ are easier to model



$t=0$

$P_0(x|y) = N(x|0,\mathbb{1})$
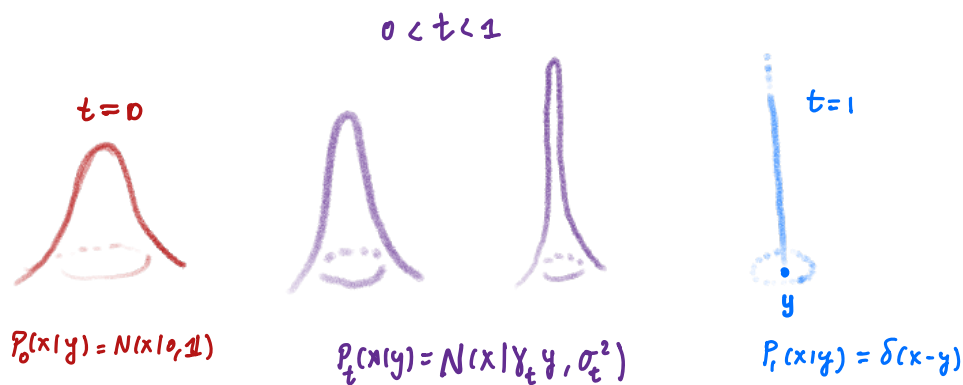
$P_t(x|y)$
$0 < t < 1$

$t=1$

$y$

$P_1(x|y) = \delta(x-y)$

- the most natural choice is to take a **Gaussian interpolation**, since the dirac delta is recoverd as the narrow width of a Gaussian.

$$\boxed{P_t(x|y) = N(x \mid \gamma_t \, y, \, \sigma_t^2 \, \mathbb{1})}$$

Gaussian conditional probabilty path!

$\gamma_t, \sigma_t$ are functions that satisfy $\begin{cases} (\gamma_0, \sigma_0) = (0, 1) \\ (\gamma_1, \sigma_1) = (1, 0) \end{cases}$

21

$$t=0 \qquad 0<t<1 \qquad t=1$$

$$P_0(x|y) = N(x|0,\mathbb{1}) \qquad P_t(x|y) = N(x|\gamma_t y, \sigma_t^2) \qquad P_1(x|y) = \delta(x-y)$$

<u>Remark:</u> Diffusion models, explained later in the course, also lead to Gaussian conditional prob. paths with particular choices of the mean $\gamma_t$ and covariance $\sigma_t$ ...

- Sampling $X_t \sim P_t(x|x_1)$ yields a cond. trajectory of the form

$$\boxed{X_t = \gamma_t X_1 + \sigma_t \mathcal{E}} \quad , \quad \mathcal{E} \sim N(0,\mathbb{1}) \qquad \text{(Reparam. trick)}$$
$$= P_{\text{base}}$$

- the conditional vector field can also be computed:

$$\boxed{u_t(x|x_1) = \dot{\gamma}_t X_1 + \dot{\sigma}_t \mathcal{E}} \quad \leftarrow \quad \begin{array}{l} \text{Recall that} \\ \dot{X}_t = u_t(x|x_1) \end{array}$$

- one can show that the vector field $v_t$ that generates $P_t(x|$ can be represented by:

$$u_t(x) = \int dx_1 \, u_t(x|x_1) \frac{P_t(x|x_1) P_1(x_1)}{P_t(x)}$$

i.e. aggregation of *conditional vector fields* $u_t(x|y)$ that generate $P_t(x|y)$ via the continuity equation.

→ Unfortunately, $u_t(x)$ can't be integrated. We still don't know the denominator $P_t(x)$ .... back to square 1?...

• Conditional Flow-matching to the rescue:

Instead of $\mathcal{L}_{FM}$, consider regressing the conditional vector field

$$\mathcal{L}_{CFM}(\theta) = \mathop{\mathbb{E}}_{\substack{t \sim U[0,1] \\ x \sim P_t(x|x_1)}} \left\| u_t^\theta(x) - u_t(x|x_1) \right\|^2$$

THEOREM: minimizing the objective $\mathcal{L}_{CFM}^\theta$ is equivalent to minimizing $\mathcal{L}_{FM}^\theta$ !

$$\nabla_\theta \hat{\mathcal{L}}_{FM}(\theta) = \nabla_\theta \mathcal{L}_{CFM}(\theta) + const.$$

- this is another comon trick in ML: if a Loss is intractable cook up a simpler loss that has the same minima!

- Notice that if we assume Gaussian conditional probability paths, where we specify $\gamma_t$ and $\sigma_t$, then $\mathcal{L}_{CFM}(\theta)$ is fully computable!

- <u>the simplest model</u>: conditional optimal transport

take: $\begin{cases} \gamma_t = t \\ \sigma_t = (1-t) \end{cases}$

$\Rightarrow \begin{cases} P_t(x|x_1) = \mathcal{N}(x| tx_1, (1-t)^2) \\ u_t(x|x_1) = \quad x_1 - \varepsilon \quad \longleftarrow \text{straight line for conditional vector field.} \\ \qquad\qquad\qquad\qquad\qquad\qquad \text{with constant speed!} \end{cases}$

$$\boxed{\mathcal{L}_{CFM}(\theta) = \mathop{\mathbb{E}}_{\substack{t \sim U[0,1] \\ \varepsilon \sim N(0,\mathbb{1}) \\ x_1 \sim P_{data}}} \left\| u_t^\theta(x) - \varepsilon - x_1 \right\|^2}$$

**Remark:** • Using the CFM objective for training is fast!

• But once we learn the conditional vector field in order to sample from the CNF we still need to solve the NODE... ← slow...

NEXT TOPICS?

• Going beyond Gaussian base distribution:

• Optimal Transport Flow-matching

• Diffusion models