

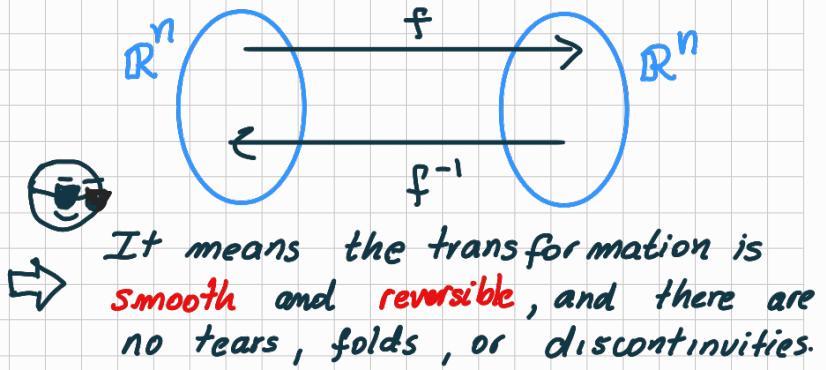
Understanding Flow Matching

by Cristian Lazo Quispe

Diffeomorphism

It is a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$
that is:

- Bijective (invertible)
- Differentiable (smooth)
- Has a differentiable inverse



Given Random Variable $x \sim p(x)$ and diffeomorphism function Φ
then

$$\tilde{x} = \Phi(x) \quad \text{and} \quad \tilde{P}(x) = \frac{p(\Phi^{-1}(x))}{\left| \det \frac{\partial \Phi}{\partial x} \right|}$$

called $\overset{\uparrow}{\text{twins}}$

Jacobian of Φ

* Note: KL-divergence (P, q) satisfy diffeomorphism variance
 $KLD(P, q) = KLD(\tilde{P}, \tilde{q})$ for any Φ

Continuity Equation



• In Physics

$$\frac{\partial \rho(x, t)}{\partial t} + \nabla \cdot (\rho(x, t) \cdot \vec{v}(x, t)) = 0$$

$\rho(x, t)$: mass density at x, t

$v(x, t)$: velocity field (how matter flows)

This means that mass is neither created nor destroyed **only** redistributed



• In Probability Distributions

$$\frac{\partial p(x, t)}{\partial t} + \nabla \cdot (p(x, t) \cdot \vec{v}(x, t)) = 0$$

$p(x, t)$: probability density at x, t

$v(x, t)$: vector field that transports samples

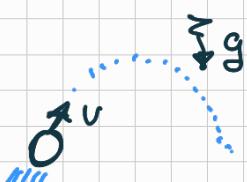
this equation ensures that the flow of samples and the evolution of the density are consistent



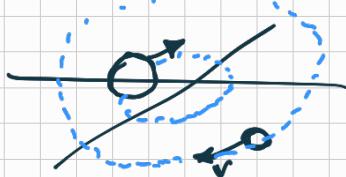
Differential Equations



$$\ddot{\theta} = -\frac{g}{l} \sin \theta - \mu \dot{\theta}(t)$$



$$\ddot{x}(t) = -g$$



$$\ddot{x}_1 = Gm_2 \left(\frac{x_2 - x_1}{\|x_2 - x_1\|^3} \right) \left(\frac{1}{\|x_2 - x_1\|^2} \right)$$

Sometimes it is hard to solve it. However, Vector field \vec{V} is often the solution

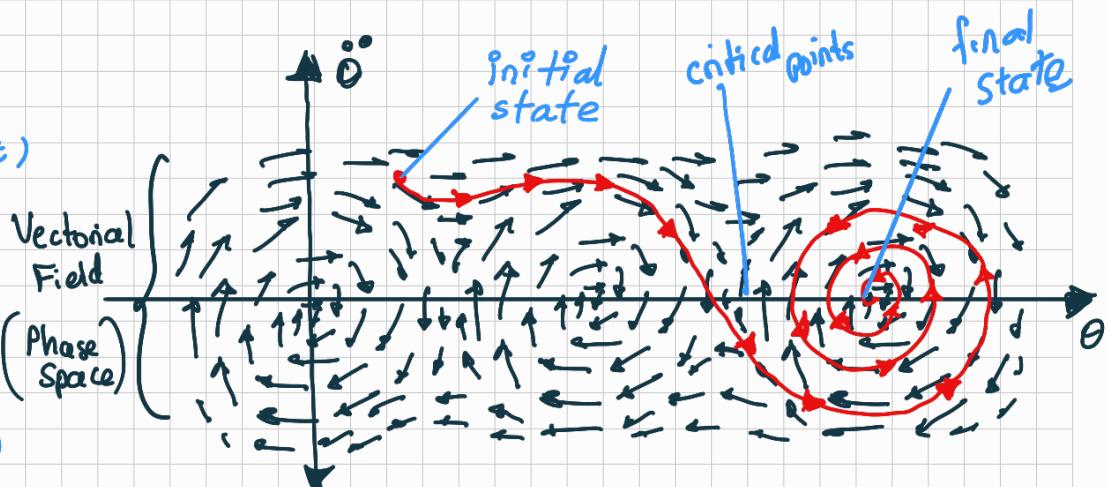
Example:



$$\ddot{\theta} = -\frac{g}{l} \sin \theta - \mu \dot{\theta}(t)$$

Second order ODE

Computational solution using Vector field



• Phase Space (Spaces encoding all kind of states of changing systems)

$$x(t) = x_0 + \int_0^t \vec{v}(x(z), z) dz$$

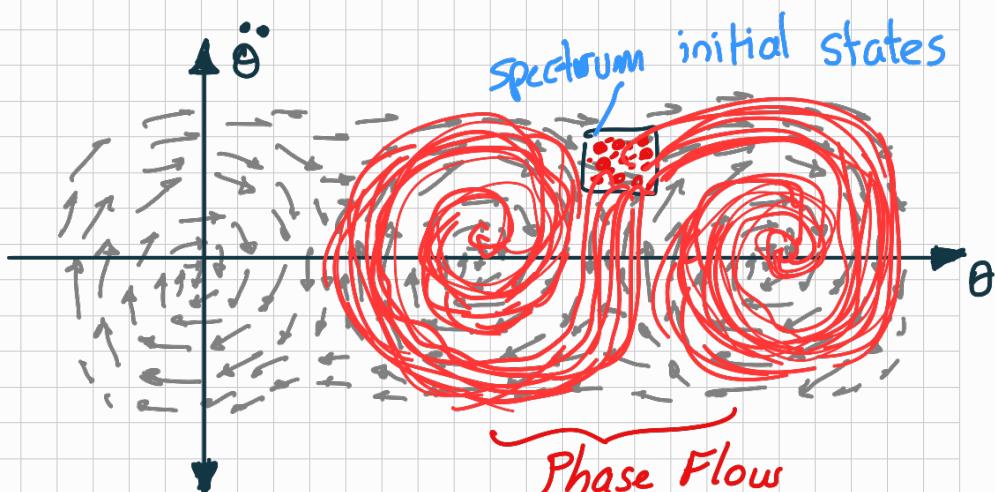
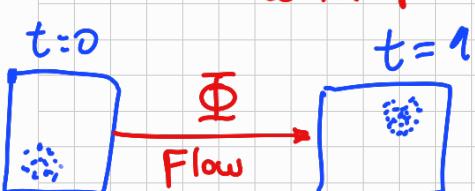
initial state

vector field

$$X(t) = \Phi^t(x_0)$$

t_0

Flow map

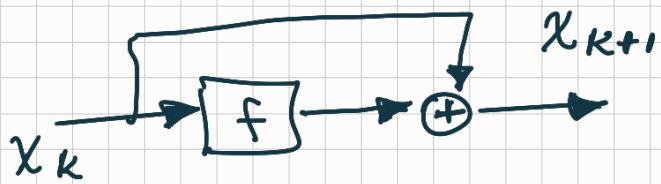


• Phase Flow (collection of all possible trajectories)

Neural ODE

Model dynamical system

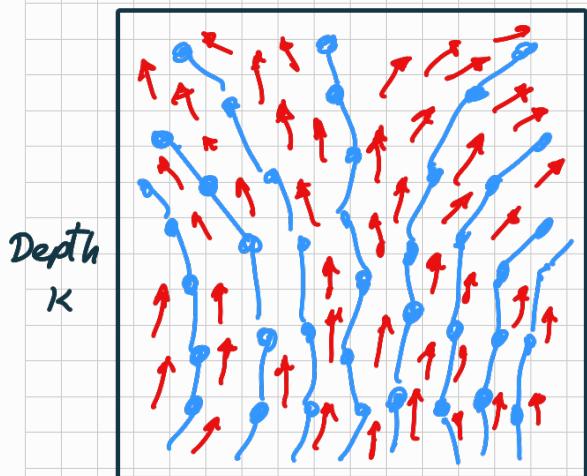
Residual Network
(Resnet) discrete



$$x_{k+1} = x_k + f(x_k)$$

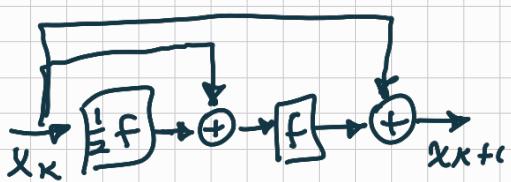
$$x_{k+1} = x_k + \Delta x_k,$$

It can be seen as Euler integrator!



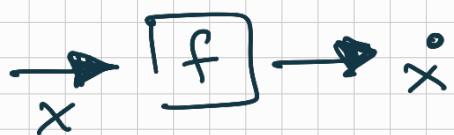
Input / hidden / output

better approximation



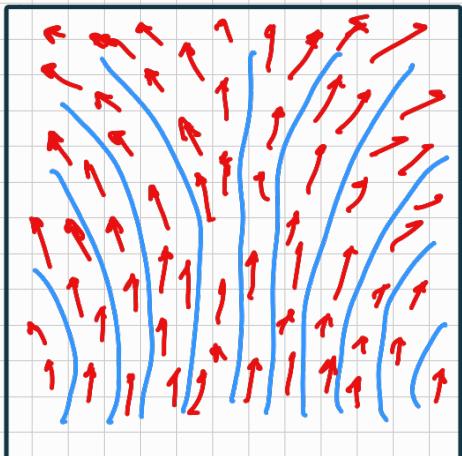
$$x_{k+1} = x_k + f(x_k + \frac{1}{2}f(x_k))$$

Neural ODE
continuous



$$\dot{x} = \frac{dx}{dt} = f(x)$$

Depth
t



Input / hidden / output

$$x_{t_0 + \Delta t} = x_{t_0} + \int_{t_0}^{t_0 + \Delta t} f(x(z)) dz$$

flow map, Φ

Free parameters: $f_\theta, t_0, \Delta t$

Hidden state: $x(z)$

$$x(z) : \Phi(x(t_0), t_0, z; \theta)$$

Optimal Transport



Problem of transforming one distribution into another in the most cost-efficient way

Particles perspective (evolution) location

$$\begin{aligned} \dot{x}_i(t) &= V(x_i(t), t) \\ x_i(0) &= x_{i,0} \quad \text{initial state} \end{aligned}$$

For V ,

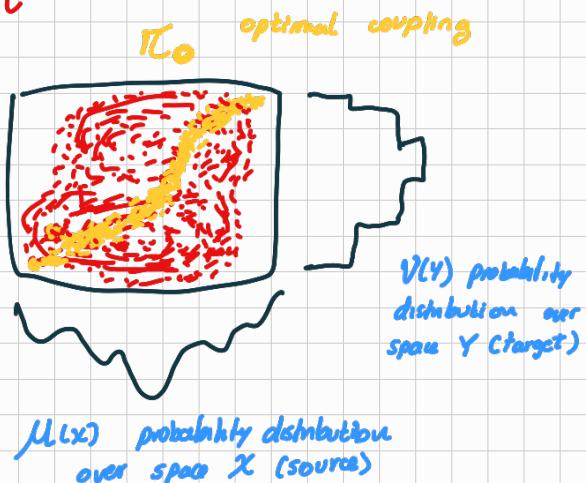
$$\text{if } p^N(0) = \frac{1}{N} \sum_{i=1}^N \delta_{x_{i,0}} \xrightarrow{N \rightarrow \infty} p_0(x), \quad \text{then } p^N(t) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i(t)} \xrightarrow{N \rightarrow \infty} p(x, t)$$

delta dirac (all probability mass at a single point x)

I have N particles moving in space. Their collective presence approximates a density

set of all joint distributions (couplings)

Π



$\Pi(\mu, \nu)$ set of all joint distributions of transport μ to ν

$C(x, y)$ cost of moving mass from x to y

the Kantorovich Formulation

$$\min_{\pi \in \Pi(\mu, \nu)} \int_{X \times Y} c(x, y) d\pi(x, y)$$

Searching over all possible joint distributions $\pi(x, y)$ that match μ and ν .

Wasserstein Distance

$$W_2(\mu, \nu) = \left(\min_{\pi \in \Pi(\mu, \nu)} \int_{X \times Y} \|x - y\|^2 d\pi(x, y) \right)^{1/2}$$

It is the optimal transport cost

• Displacement interpolation:

McCann introduces the idea of displacement interpolation a way to smoothly interpolate between 2 distributions μ and ν using the optimal transport T .

identity map

\downarrow

$$p_t = ((1-t) \cdot \text{id} + tT) \# \mu$$

$$\pi_t(x, y) = (1-t)x + ty$$

Boundaries:

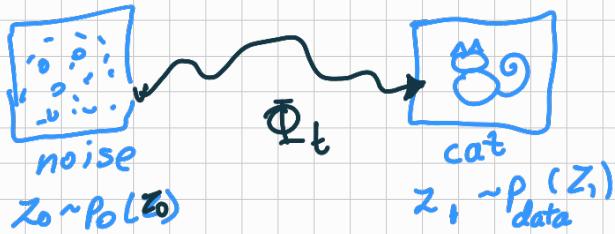
$$\pi_0(x, y) = x$$

$$\pi_1(x, y) = y$$

CNF

Continuous Normalizing Flows (CNFs)

Learn a transformation from latent $z_0 \sim p_0(z)$ to data $z(1) \sim p_{\text{data}}(z_1)$



assume
 Φ_t diffeomorphic mapping
 $\Phi = f_\theta(x, t)$
 $\Phi_0(x) = x$

Vector Field: It defines $\frac{dz(t)}{dt} = V(z(t), t) = \Phi = f_\theta(x, t)$

$$V: [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$\begin{matrix} t \\ \tilde{x} \\ \tilde{x} \end{matrix} \quad \begin{matrix} \tilde{x} \\ \tilde{x} \\ \tilde{v} \end{matrix}$$

$$z = \Phi_t(x)$$

$$\frac{d(\Phi_t(x))}{dt} = V_t(\Phi_t(x))$$

Probability density path:

$$P: [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}_{>0}$$

$$\begin{matrix} t \\ \tilde{x} \\ \tilde{x} \end{matrix} \quad \begin{matrix} \text{number} \\ \text{real} \end{matrix}$$



in instance of time
 $P_t = [\Phi_t]_* P_0(x) = P_0(\Phi_t^{-1}(x)) \det \left[\frac{\partial \Phi_t^{-1}(x)}{\partial x} \right]$
push-forward

$\underbrace{P(t, x)}_{\text{number}} \rightarrow \int_{\mathbb{R}^d} P(t, x) dx = 1 \quad \forall t \in [0, 1]$

it means all probability densities function over the time

Objective: $L_{\text{CNF}} = -\log p(z(1))$ We want to maximize the log-probability of data samples



It means learning a map from z_0 to z_1 , along with the exact densities $p(z_0)$ and $p(z_1)$, by defining a flow from z_0 to z_1 , which consequently induces a flow from $p(z_0)$ to $p(z_1)$.

how to solve:

$$\frac{\partial P(x, t)}{\partial t} + \nabla \cdot (P(x, t) \cdot \tilde{V}(x, t)) = 0$$

theorem of liouville

$$\frac{\partial P}{\partial t} + \nabla P \cdot \tilde{V} = 0$$

$$\frac{\partial P}{\partial t} + P \nabla \cdot \tilde{V} + \nabla P \cdot V = 0$$

$$\frac{\partial P}{\partial t} + \nabla P \cdot V = -P \nabla \cdot V$$

$$\frac{dp}{dt} = \frac{\partial P}{\partial t} + \frac{\partial P}{\partial z} \cdot \frac{\partial z}{\partial t} = -P \nabla \cdot V$$

$$\frac{dp}{dt} = -P \nabla \cdot V$$

$$\left\{ \frac{1}{P} \frac{dP}{dt} = -\nabla \cdot V = \frac{d \log P_t(x)}{dt} = -\nabla \cdot V_t(x) \right.$$

$$\log P(z(1)) = \log P(z(0)) - \int_0^1 \nabla \cdot V dt$$

$$L_{\text{CNF}} = -\log P(z(0)) + \int_0^1 \nabla \cdot V dt$$

$$L_{CNF} = -\log(P(z_0)) + \int_0^1 \nabla \cdot V dt$$

d : dimension
of input

$$\text{tr}(J) = \sum_{i=1}^d \underbrace{\frac{\partial V(z_i)}{\partial z_i}}_{\text{Jacobian trace}}$$

image
 $d = C \times h \cdot w$
 $RGB \rightarrow 3 \times h \cdot w$

$$J = \frac{\partial V}{\partial z} = \begin{bmatrix} \frac{\partial V}{\partial z_1} \\ \vdots \\ \frac{\partial V}{\partial z_d} \end{bmatrix}$$

expensive to compute

how to solve it \rightarrow Hutchison's trick
given $\epsilon \sim N(0, I)$

$$\nabla \cdot V = \text{tr}(J) = E_{\epsilon}[\epsilon^T J \epsilon]$$

$$\epsilon^T V = \sum \epsilon_i V_i$$

$$\nabla_z \epsilon^T V = \sum \epsilon_i \nabla_z V_i = \underbrace{\epsilon^T}_{\text{scalar}} \underbrace{V}_{\text{scalar}}$$

it could be computed with autograd over z

$$L_{CNF} = -\log P_i(z^{(1)}) = -\log(P(z_0)) + \int_0^1 \text{tr}\left(\frac{\partial V}{\partial z}\right) dt$$

$$L_{CNF} = -\log P_i(z^{(1)}) = -\log(P(z_0)) + \int_0^1 E[\epsilon^T \nabla_z (\epsilon^T V(z))] dt$$

↓

the original paper proposes this : We don't know z_0 that generates $z^{(1)}$ directly \rightarrow how is it obtained?

Doing Flow: $z^{(1)} \xrightarrow{\Phi} z^{(0)}$

$$z_1 = z_0 + \int_0^1 \vec{V}(z(z), z) dz$$

$$z_0 = z_1 - \int_0^1 \vec{V}(z(z), z) dz = z_1 + \int_1^0 \vec{V}(z(z), z) dz$$

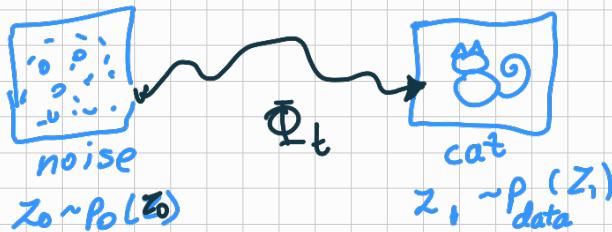
∴ Training:

$$L_{CNF} = -\log P_i(z^{(1)}) = -\log(P(z_1 + \int_1^0 \vec{V}(z(z), z) dz)) - \int_1^0 E[\epsilon^T \nabla_z (\epsilon^T V(z))] dz$$

Inference: $z_1 = z_0 + \int_0^1 \vec{V}(z(z), z) dz$

Flow Matching

Learn a transformation from latent $z_0 \sim p_0(z)$ to data $z(1) \sim p_{\text{data}}(z_1)$



assume
 Φ_t diffeomorphic
 $\Phi = f_\Theta(x, t)$
model
 $\Phi_0(x) = x$

Vector Field: It defines $\frac{dz(t)}{dt} = V(z(t), t)$ $\Phi = f_\Theta(x, t)$

$$V: \underbrace{[0, 1]}_{t} \times \underbrace{\mathbb{R}^d}_{x} \rightarrow \mathbb{R}^d$$

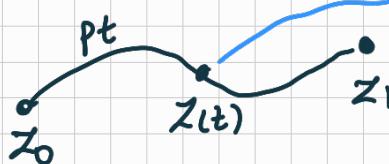
$$z = \Phi_t(x)$$

$$\frac{d(\Phi_t(x))}{dt} = V_t(\Phi_t(x))$$

Probability density path:

$$P: \underbrace{[0, 1]}_{t} \times \underbrace{\mathbb{R}^d}_{x} \rightarrow \mathbb{R}_>$$

number real



$$p_t = [\Phi_t]_* p_0(x) = p_0(\Phi_t^{-1}(x)) \det \left[\frac{\partial \Phi_t^{-1}(x)}{\partial x} \right]$$

↑ push-forward

$$\underbrace{p(t, x)}_{\text{number}} \rightarrow \int_{\mathbb{R}^d} p(t, x) dx = 1 \quad \forall t \in [0, 1]$$

it means all probability densities function over the time

Objective: $L_{FM} = E_{t, P_t(x)} \| V_t(x) - u_t(x) \|^2$

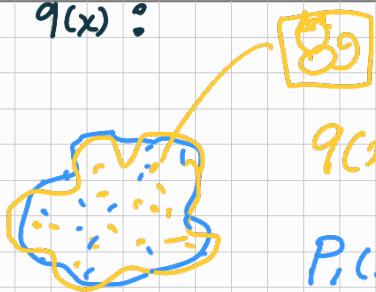
We want to reduce the mean square error of the vector field predicted over the real

It means learning a map from an intermediate point z_t to z_1 , guided by the exact flow of the vector field $V_t(x)$

How to solve it: how to get $u_t(x)$ and $P_t(x)$

We don't have access to a specific path that maps z_0 to z_1 , nor a closed-form expression for u_t that generates the desired p_t , we need to compute it somehow. In contrast, CNF don't care about the explicit path or the vector field, CNF only focus on the start and end distributions connected through a continuous flow.

Define $q(x)$:



$q(x)$: unknown distribution (world)

True world

$$P_t(x) \approx q(x)$$

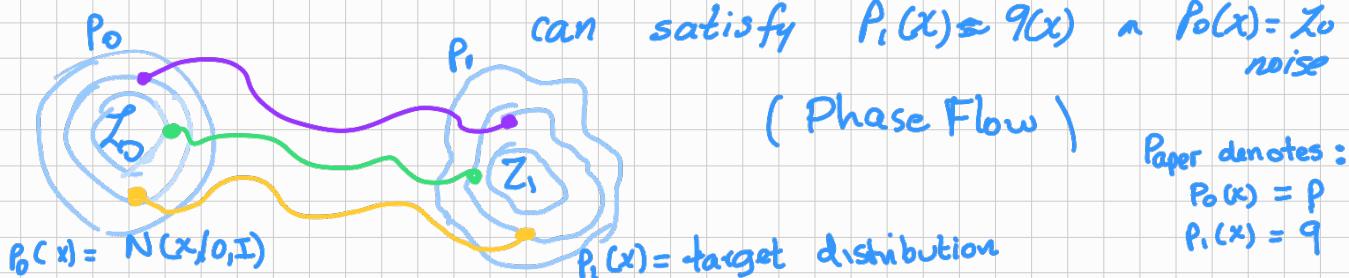
$P_t(x)$: target

we have samples

distribution (model world) parametric model

Define $P_t(x)$:

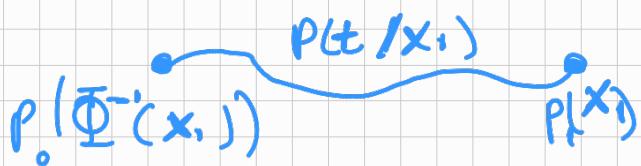
how to define $P_t(x)$? there are many probability paths that



Paper denotes:
 $P_0(x) = p$
 $P_t(x) = q$

- Define conditional probability path $P_t(x/x_i)$

We have data and x_i is a data point



Given x_i we denote

$P_t(x/x_i)$ that satisfy

$$P_0(x/x_i) = p(x), t=0$$

$P_t(x/x_i)$ = distribution concentrated around $x=x_i$

- Define marginal probability path

$$P_t(x) = \int P_t(x/x_i) dx_i \quad \rightharpoonup q(x_i)$$

$$P_t(x) = \int P_t(x/x_i) P_i(x_i) dx_i$$

$$P_t(x) = \int P_t(x/x_i) q(x_i) dx_i \quad \text{where } P_i(x) = \int P_t(x/x_i) q(x_i) dx_i \approx q(x)$$

- Define marginal vector field

$$\frac{\partial P_t(x)}{\partial t} + \nabla \cdot (P_t(x) \cdot \mu_t(x)) = 0$$

$$\frac{\partial P_t(x/x_i)}{\partial t} + \nabla \cdot (P_t(x/x_i) \cdot \mu_t(x/x_i)) = 0$$

$$q(x_i) \cdot \left[\frac{\partial P_t(x/x_i)}{\partial t} + \nabla \cdot (P_t(x/x_i) \cdot \mu_t(x/x_i)) \right] = 0$$

$$\frac{\partial P_t(x/x_i) q(x_i)}{\partial t} + \nabla \cdot (P_t(x/x_i) \cdot \mu_t(x/x_i) \cdot q(x_i)) = 0$$

$$\int \frac{\partial P_t(x/x_i) q(x_i)}{\partial t} dx_i + \nabla \cdot \left(\int P_t(x/x_i) \cdot \mu_t(x/x_i) q(x_i) dx_i \right) = 0$$

$$\int \frac{\partial P_t(x/x_1) q(x_1)}{\partial t} dx_1 + \nabla \cdot \left(\int P_t(x/x_1) \cdot u(x/x_1) q(x_1) dx_1 \right) = 0$$

$$\frac{\partial P_t(x)}{\partial t} + \nabla \cdot \left(\int P_t(x/x_1) \cdot u(x/x_1) q(x_1) dx_1 \right) = 0$$

$$\boxed{\frac{\partial P_t(x)}{\partial t} + \nabla \cdot (P_t(x) \cdot u_t(x)) = 0}$$

$$P_t(x) \cdot u_t(x) = \int P_t(x/x_1) \cdot u(x/x_1) q(x_1) dx_1$$

$$u_t(x) = \frac{1}{P_t(x)} \int P_t(x/x_1) u(x/x_1) q(x_1) dx_1$$

$$u_t(x) = \int u_\theta(x/x_1) \cdot \frac{P_t(x/x_1) q(x_1)}{P_t(x)} dx_1$$

where $u_t(x/x_1) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a conditional vector field that generates $P_t(x/x_1)$

Paper key observation:

The marginal vector field generates the marginal probability path $u_t(x)$

this generates the following theorem 1

Theorem 1:

Given vector fields $u_t(x/x_1)$ that generate conditional probability paths $P_t(x/x_1)$, for any distribution $q(x_1)$, the marginal vector field u_t generates the marginal probability path P_t , i.e., u_t and P_t satisfy the continuity equation

model

$$L_{FM} = E_{+, P_t(x)} \downarrow \| V_t(x) - u_t(x) \|^2$$

$u_t(x)$ and $P_t(x)$ are still intractable to compute

We could trying to

solve it in a simple way

$$\tilde{u}_t(x) = \int u_\theta(x/x_1) \cdot \frac{P_t(x/x_1) q(x_1) dx_1}{P_t(x)}$$



model

$$L_{FM} = E_{t, P_t(x)} \| V_t(x) - M_t(x) \|^2$$

$$E_{t, x \sim P_t(x)} (\| V_t(x) \|^2 - \langle 2 V_t(x) \cdot M_t(x) \rangle + \| M_t(x) \|^2)$$

$$\begin{aligned} E_{t, x \sim P_t(x)} \| V_t(x) \|^2 &= \int \| V_t(x) \|^2 P_t(x) dx = \int \| V_t(x) \|^2 \int P_t(x/x_1) q(x_1) dx_1 dx \\ &= \iint \| V_t(x) \|^2 P_t(x/x_1) q(x_1) dx_1 dx \end{aligned}$$

$$E_{t, x \sim P_t(x)} \| V_t(x) \|^2 = E_{t, x \sim P_t(x/x_1), x_1 \sim q(x_1)} \| V_t(x) \|^2$$

$$\begin{aligned} E_{t, x \sim P_t(x)} \| \langle 2 V_t(x) \cdot M_t(x) \rangle \| &= \int \langle 2 V_t(x) \cdot M_t(x) \rangle P_t(x) dx \\ &= 2 \int V_t(x) \cdot M_t(x) P_t(x) dx \\ &= 2 \underbrace{\int V_t(x) \cdot \frac{\int M_t(x/x_1) P_t(x/x_1) q(x_1) dx_1}{P_t(x)}}_{P_t(x)} dx \\ &= 2 \int V_t(x) \cdot \int M_t(x/x_1) P_t(x/x_1) q(x_1) dx_1 dx \\ &= 2 \iint \underbrace{V_t(x) \cdot}_{\langle 2 V_t(x) \cdot M_t(x/x_1) \rangle} M_t(x/x_1) P_t(x/x_1) q(x_1) dx_1 dx \end{aligned}$$

$$E_{t, x \sim P_t(x)} \| \langle 2 V_t(x) \cdot M_t(x) \rangle \| = E_{t, x \sim P_t(x/x_1), x_1 \sim q(x_1)} \| \langle 2 V_t(x) \cdot M_t(x/x_1) \rangle \|$$

$$\begin{aligned} E_{t, x \sim P_t(x)} \| M_t(x) \|^2 &= \int \| M_t(x) \|^2 P_t(x) dx = \int \| M_t(x) \|^2 \int P_t(x/x_1) q(x_1) dx_1 dx \\ &= \iint \| M_t(x) \|^2 P_t(x/x_1) q(x_1) dx_1 dx \end{aligned}$$

$$E_{t, x \sim P_t(x)} \| M_t(x) \|^2 = E_{t, x \sim P_t(x/x_1), x_1 \sim q(x_1)} \| M_t(x) \|^2$$

$$E_{t, P_t(x)} \| V_t(x) - M_t(x) \|^2 = E_{t, x \sim P_t(x/x_1), x_1 \sim q(x_1)} (\| V_t(x) \|^2 - \langle 2 V_t(x) \cdot M_t(x/x_1) \rangle + \| M_t(x) \|^2)$$

$$E_{t, P_t(x)} = \|V_t(x) - U_t(x)\|^2 = E_{t, x \sim P_t(x/x_i), x_i \sim q(x_i)} \left(\|V_t(x)\|^2 - \langle 2 V_t(x) \cdot U_t(x/x_i) \rangle + \|U_t(x)\|^2 \right)$$

$$E_{t, x \sim P_t(x/x_i), x_i \sim q(x_i)} \left(\|V_t(x)\|^2 - \langle 2 V_t(x) \cdot U_t(x/x_i) \rangle + \|U_t(x)\|^2 \right) + \|U_t(x/x_i)\|^2 - \|U_t(x/x_i)\|^2$$

$$E_{t, x \sim P_t(x/x_i), x_i \sim q(x_i)} \underbrace{\left(\|V_t(x)\|^2 - \langle 2 V_t(x) \cdot U_t(x/x_i) \rangle + \|U_t(x/x_i)\|^2 + \|U_t(x)\|^2 - \|U_t(x/x_i)\|^2 \right)}$$

$$E_{t, x \sim P_t(x/x_i), x_i \sim q(x_i)} \left(\|V_t(x) - U_t(x/x_i)\|^2 + \|U_t(x)\|^2 - \|U_t(x/x_i)\|^2 \right)$$

 not depend of the model
in the optimization

$$L_{CFM} = E_{t, x \sim P_t(x/x_i), x_i \sim q(x_i)} \|V_t(x) - U_t(x/x_i)\|^2 \quad t \sim U[0,1]$$

The 2nd key observation is therefore:

The L_{FM} and L_{CFM} objectives have identical gradients w.r.t θ

 this generate the following theorem 2

Theorem 2. Assuming that $p_t(x) > 0$ for all $x \in \mathbb{R}^d$ and $t \in [0,1]$, then, up to a constant independent of θ , L_{CFM} and L_{FM} are equal. Hence, $\nabla_\theta L_{FM}(\theta) = \nabla_\theta L_{CFM}(\theta)$

Now we known that we can express the objective with a conditional probability $p_t(x/x_i)$ and conditional vector field $U_t(x/x_i)$.

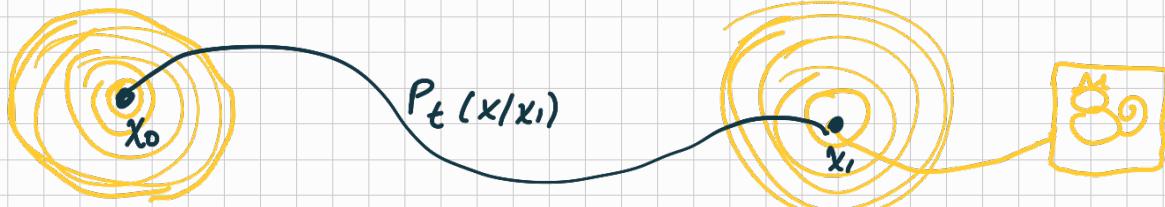
How to construct $p_t(x/x_i)$ and $U_t(x/x_i)$?



We know that LCFOT works with any choice of $P_t(x/x_i)$ and $\mu_t(x/x_i)$

Let's consider $P_t(x/x_i) = N(x | \mu_t(x_i), \sigma_t(x_i)^2 I)$
family of Gaussian conditional probability paths.

where $\mu: [0,1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the time-dependent mean of the Gaussian Distribution
 $\sigma: [0,1] \times \mathbb{R} \rightarrow \mathbb{R}_{>0}$ describes a time-dependent scalar standard deviation std



We set

$$\begin{aligned} \mu_0(x_i) &= 0 \\ \sigma_0(x_i) &= 1 \rightarrow P_0(x) = p(x) = N(x|0, I) \end{aligned}$$

noise

$$P_0(x/x_i) = P_0(x)$$

$$\begin{aligned} \mu_1(x_i) &= x_i \\ \sigma_1(x_i) &= \sigma_{\min} \quad \text{sufficiently small} \end{aligned}$$

so that

$P_t(x/x_i)$ is concentrated Gaussian distribution centered at x_i .

$$N(x|0, I) \xrightarrow{P_t(x/x_i)} N(x_i | x_i, \sigma_{\min}^2 I)$$

there are many vector fields that generate any particular probability path and many flow Φ_t because $d\frac{\Phi_t(x)}{dt} = V_t(\Phi_t(x))$

↪ let's consider a flow $\Psi_t(x)$ as a canonical transformation for Gaussian distributions (conditioned on x_i)

$$\Psi_t(x, x_i) = \Psi_t(x) = \sigma_t(x_i) \cdot x + \mu_t(x_i)$$

$$[\Psi_t] * p(x) = P_t(x/x_i)$$

$$\boxed{\frac{d}{dt} \Phi_t(x) = \mu_t(\Phi_t(x))}$$

This flow then provides a vector field that generates the conditional probability path:

$$\frac{d}{dt} \Psi_t(x) = \mu_t(\Psi_t(x)/x_i)$$

$$\begin{aligned} \Psi_0(x_0) &= x_0 & \Psi_1(x_1) &= x_1 \\ \sigma_0(x_i) \cdot x_0 + \mu_0(x_i) &= x_1 & \sigma_1(x_i) \cdot x_1 + \mu_1(x_i) &= x_1 \\ 1 \cdot x_0 + 0 &= x_1 & \sigma_{\min} \cdot x_1 + x_i &= x_1 \\ x_1 &= x_1 & \downarrow & \downarrow \\ x_1 &= x_1 & & x_1 = x_1 \end{aligned}$$

Satisfies the property of diffeomorphic transformation

$$\Psi_t(x) = \sigma_t(x_1) \cdot x + \mu_t(x_1)$$

$$x = \frac{\Psi_t(x) - \mu_t(x_1)}{\sigma_t(x_1)}$$

$$\frac{d}{dt} \Psi_t(x) = \dot{\mu}_t(\Psi_t(x)/x_1) = \frac{d}{dt} [\sigma_t(x_1) \cdot x + \mu_t(x_1)]$$

$$\frac{d\sigma_t(x_1) \cdot x}{dt} + \frac{d\mu_t(x_1)}{dt}$$

$$\frac{d}{dt} \Psi_t(x) = \dot{\mu}_t(\Psi_t(x)/x_1) = \underbrace{\frac{d\sigma_t}{dt}(x_1)}_{\text{---}} \cdot \left[\frac{\Psi_t(x) - \mu_t(x_1)}{\sigma_t(x_1)} \right] + \frac{d\mu_t(x_1)}{dt}$$

it allows to use any x value in the formulation

$$\frac{d}{dt} \Psi_t(x) = \dot{\mu}_t(\Psi_t(x)) = \dot{\mu}_t(\Psi_t(x)/x)$$

$$\dot{\mu}_t(x) = \dot{\mu}_t(x/x_1) = \frac{\sigma'_t(x_1)}{\sigma_t(x_1)} (x - \mu_t(x_1)) + \mu'_t(x_1)$$

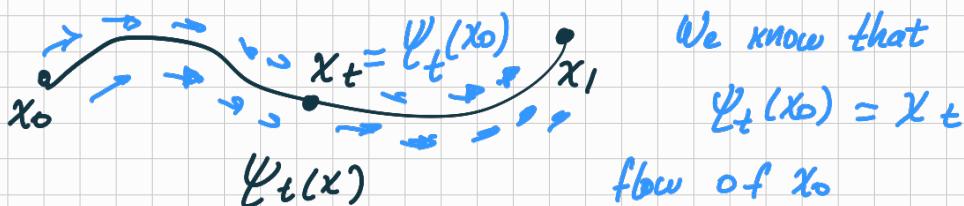
Theorem 3: Let $p_t(x/x_1)$ be a Gaussian probability path, and Ψ_t its corresponding flow map. Then, the unique vector field that defines Ψ_t has the form:

$$\dot{\mu}_t(x/x_1) = \frac{\sigma'_t(x_1)}{\sigma_t(x_1)} (x - \mu_t(x_1)) + \mu'_t(x_1)$$

Consequently $\dot{\mu}_t(x/x_1)$ generates the Gaussian path $p_t(x/x_1)$

$$L_{CFM} = E_{t_1, x \sim p_t(x/x_1), x_1 \sim q(x_1)} \| V_t(x) - \dot{\mu}_t(x/x_1) \|^2$$

Writing in terms of the intermediate point x_t



$$L_{CFM} = E_{t_1, P_t(x/x_1), q(x_1)} \| V_t(x_t) - \dot{\mu}_t(x_t/x_1) \|^2$$

$$E_{t_1, P_t(x_0/x_1), q(x_1)} \| V_t(\Psi_t(x_0)) - \dot{\mu}_t(\Psi_t(x_0)/x_1) \|^2$$

$$E_{t_1, P(x_0), q(x_1)} \| V_t(\Psi_t(x_0)) - \frac{d}{dt} \Psi_t(x_0) \|^2$$

↓ model
noise data

Remember that
 $\Psi_t(x)$ is
conditioned
on x_1

$$P_t(x/x_0) = N(x | \mu_t(x_0), \sigma_t^2 I)$$

$$\Psi_t(x/x_0) = \Psi_t(x) = \sigma_t(x_0) \cdot x + \mu_t(x_0)$$

$$W_t(x) = \mu_t(x/x_0) = \frac{\sigma_t'(x_0)}{\sigma_t(x_0)} (x - \mu_t(x_0)) + \mu_t'(x_0)$$

$$E_{t_1}, P(x_0), q(x_1) \quad ||V_t(\Psi_t(x_0)) - \frac{d}{dt} \Psi_t(x_0)||^2$$

+ ↓ ↓
noise data model

Boundaries:

$$\Psi_t(x_0) = x_0$$

$$\Psi_t(x_1) = x_1$$



Now $L_{CFN}(\theta)$ is in terms
of x_0 and $\Psi_t(x, x_1)$

This formulation is fully general for arbitrary functions $\mu_t(x_0)$ and $\sigma_t^2(x_0)$
we can set them to any differentiable function satisfying the
desired boundary conditions.

Example 1: Diffusion conditional VFs

Diffusion models start with data points and
gradually add noise until it approximates pure noise
(stochastic process)

In diffusion: the reversed (noise → data)
Variance Exploding (VE) path has the form

$$P_t(x) = N(x/x_0, \sigma_{1-t}^2 I)$$

σ_t is an increasing function
 $\sigma_0 = 0 \quad \sigma_1 \gg 1$

$$\mu_t(x_0) = x_0 \quad \sigma_t(x_0) = \sigma_{1-t}$$

$$\mu_t(x/x_0) = -\frac{\sigma_{1-t}^2}{\sigma_{1-t}} (x - x_0)$$

The reversed (noise → data) Variance Preserving
(VP) diffusion path has the form

$$P_t(x/x_0) = N(x/\alpha_{1-t} x_0, (1-\alpha_{1-t}^2) I), \text{ where}$$

$$\alpha_t = e^{-\frac{1}{2}T(t)}, T(t) = \int_0^t \beta(s) ds$$

and $\beta(t)$ is the noise scale function

$$\mu_t(x_0) = \alpha_{1-t} x_0 \quad \sigma_t(x_0) = \sqrt{1 - \alpha_{1-t}^2}$$

$$\begin{aligned} \mu_t(x/x_0) &= \frac{\alpha_{1-t}^2}{1 - \alpha_{1-t}^2} (\alpha_{1-t} x_0 - x) \\ &= -\frac{T'(1-t)}{2} \left[\frac{e^{-T(1-t)} x - e^{-\frac{1}{2}T(1-t)} x_0}{1 - e^{-T(1-t)}} \right] \end{aligned}$$

Note: These P_t were previously derived as
solutions of diffusion process, they do not
actually reach a true noise distribution
in finite time.

Example 2: Optimal Transport conditional VFs

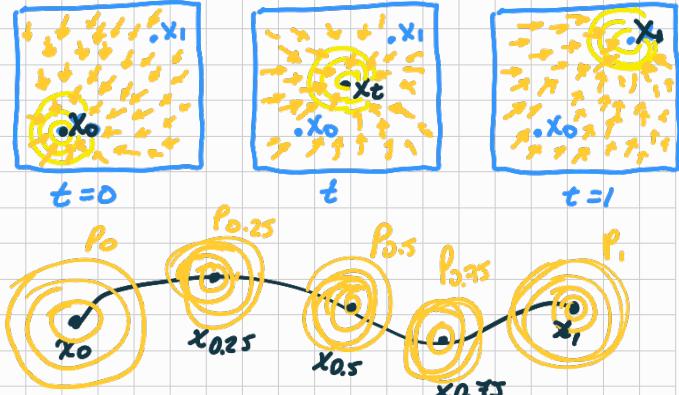
Let's define $\mu_t(x)$ and $\sigma_t^2(x)$ to
simply change linearly in time.

$$\mu_t(x) = t x_0 \quad \sigma_t^2(x) = 1 - (1 - \sigma_{\min}) t$$

Conditional Flow

$$\Psi_t(x/x_0) = \Psi_t(x) = (1 - (1 - \sigma_{\min}) t) x + t x_0$$

$$W_t(x) = \mu_t(x/x_0) = \frac{x_0 - (1 - \sigma_{\min}) x}{1 - (1 - \sigma_{\min}) t}$$



The conditional flow $\Psi_t(x)$ is in fact the
Optimal Transport (OT) displacement map
between Gaussians $P_0(x/x_0)$ and $P_1(x/x_1)$

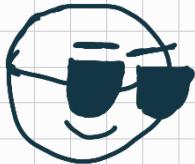
$$L_{CFN}(\theta) = E_{t_1, q(x_1), P(x_0)} \| V_t(\Psi_t(x_0)) - (x_1 - (1 - \sigma_{\min}) x_0) \|^2$$

If $\sigma_{\min} \approx 0$ then

$$\begin{aligned} \Psi_t(x) &= (1 - t)x + t x_1 \\ \Psi_t(x) &= x t \end{aligned}$$



$$L_{CFM}(\theta) = E_{t, q(x_0), p(x_0)} \left\| v_t(\Psi_t(x_0)) - (x_1 - (1-\delta_{\min})x_0) \right\|^2$$



Flow matching learns an exact vector field V_t to follow a simple, optimal transport-based path from noise to data, using only a closed-form loss over interpolated samples.
(simple and elegant generative model framework)

Personally, I consider these to be the main key ideas :

- Objective function of learning a vector field

$$V_t = \frac{d \Psi_t}{dt} \quad \text{fnn } E_{t, p_t} \|v_t(x) - M_t(x)\|^2$$

- Conditional probability path as family of Gaussian

$$P_t(x|x_0) = N(x | \mu_t(x_0), \sigma_t(x_0)^2 I)$$

- Optimal transport inspired interpolation

$$\mu_t(x) = t x_0 \quad \sigma_t(x) = 1 - (1 - \delta_{\min}) t$$

Training ($\delta_{\min}=0$) :

$$t = \text{random}(0, 1)$$

$$x_0 = \text{noise}$$

$$x_t = (1-t)x_0 + t x_1$$

$$M_t = (x_t - x_0) / (1-t)$$

$$v_t = \text{model}(x_t, t)$$

$$\text{loss} = (v_t - u_t)^2$$

Inference ($\delta_{\min}=0$): (basic solver ODE)

$$x_0 = \text{noise}$$

$$x = x_0 \quad dt = 1/N_steps$$

for i in N_steps :

$$t = i \cdot dt$$

$$v = \text{model}(x, t)$$

$$x = x + v dt$$