

# CONTINUOUS NORMALIZING FLOW

## Table of Contents

Discrete Normalizing Flow (DNF)

From Discrete  $\rightarrow$  Continuous

→ Continuous Normalizing Flow (CNF)

Instantaneous Change of Variables (ICV)

Total Derivative

→ Take the Total derivative of  $\frac{d}{dt} \log p_t(x_t)$

Why the  $|_{x=x_t}$  Notation?

→ Continuity Equation & Substitute

where does the "Trace" come from?

→ Jacobian & Trace

Final Instantaneous change-of-variables

→ Loss & Training

## Discrete Normalizing Flow (DNF)

Goal: Imagine what turn a simple distribution into a more complex one

Apply a sequence of invertible transformations

$$f_1, f_2, \dots, f_k: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$z_i \rightarrow z_i'$$

$$\text{Train: } z_k = f_k \circ f_{k-1} \circ \dots \circ f_1(x)$$

$$\text{Inter: } x = f_k^{-1} \circ f_{k-1}^{-1} \circ \dots \circ f_1^{-1}(z_k)$$

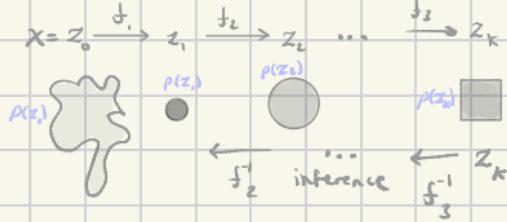
$$\text{an image } x = z_0, z_k \in \mathbb{R}^d, z_k \sim N(0, I)$$

$p(x), p(z)$

Let  $P_X(x), P_Z(z)$  be PDFs

→ Since  $f_i$  invertible,

$$f_2 : z \rightarrow z_2$$



Discrete Change-of-variables (cov)

(and log cov) formula

Gaussian

Jacobian

$$\lg p(x) = \lg p(z_k) + \sum_{i=1}^k \lg \left| \det \underbrace{\frac{\partial f_i(z_{i-1})}{\partial z_{i-1}}}_{J_{f_i}} \right| \quad (\text{A})$$

What's the problem with discrete?

- 1) Invertibility: You need specially designed invertible blocks
- 2) Expressivity: Finite blocks struggle with complex dist unless made deep/structured which gets expensive
- 3) Optimization: More layers you stack → more computationally expensive training gets

From DNF → Continuous Normalizing Flow (CNF)

Each discrete step becomes an infinitesimal transformation and the composition converges to an ODE flow

$$\frac{dx}{dt} = f_\theta(x_t, t) \quad \begin{matrix} \text{learn a vector field that} \\ \text{continuously drags a} \\ \text{single dist into the data dist} \end{matrix}$$

No more finite  
invertible layers

What is an ODE?

It tells you how something changes over time

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

The instantaneous  
RtC wrt time

tells you the slope at  
that point

# Continuous Normalizing Flow (CNF)

$x_0 \sim p_0$  is the transformed sample wwt match data  
 $x_T \sim N(0, I)$  is base distribution (say Gaussian)

Instantaneous Change-of-variables

in continuous time

$J_{t0}$

$$\frac{d}{dt} \log p_t(x_t) = -\text{Tr}\left(\frac{\partial J_{t0}(x_{t0}, t)}{\partial x_t}\right)$$

$$\log p(x_0) = \log p(x_T) - \int_0^T \text{Tr}\left(\frac{\partial J_{t0}}{\partial x_t}\right) dt$$

$p_t = \text{POF indexed}$

by time, so

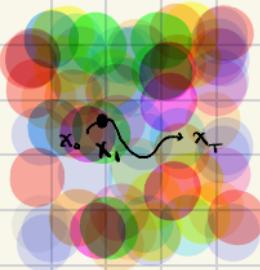
$$p_t(x_t) = p(t)(x_t) \\ = p(x_0, t)$$

whole flow is now a continuous trajectory

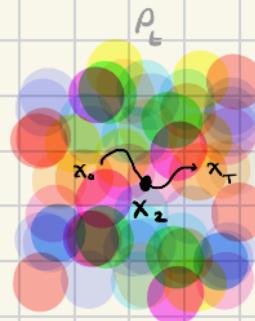
$$x_t : t \rightarrow x$$

Problem: Trace is the computational bottleneck!

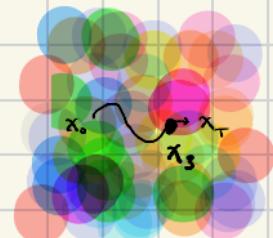
$$p_t(x) \neq p_t(x_t)$$



$$t=1 \quad p_1(x_1) = a \in \mathbb{R}$$



$$t=2$$



$$t=3$$

Think of  $p_t$   
like temperature

## Discrete Flow

Product / log-sum of Jacobian determinants

## Continuous Flows

Limit of infinitely many infinitesimal Jacobians /  
 integral of Jacobian trace

Transformation: A warp of space (invertible layer)

Continuous transformation:  
 a time-dependent velocity field that drags probability mass around

## Total Derivative

$x$  a coordinate  $\mathbb{R}^d$

$x_t$  a trajectory

and  $g(x_t, t)$

When we differentiate we have to account for 2 things

i) How the field  $g$  changes with time

at a fixed point  $\frac{\partial g}{\partial t}$

ii) How the field changes because you're moving

along the trajectory through space

$$\nabla_x g \cdot \frac{dx_t}{dt}$$

That's why, for  $x \in \mathbb{R}^d$

$$g(x_t(t), t) \quad \frac{d}{dt} g(x_t, t) = \underbrace{\frac{\partial g}{\partial x} \cdot (x_t, t)}_{\text{Lagrangian}} \cdot \frac{dx_t}{dt} + \underbrace{\frac{\partial g}{\partial t} (x_t, t)}_{\text{dot product Eulerian}}$$

$$\text{where } \frac{\partial g}{\partial x} (x_t, t) = \left( \frac{\partial g_1}{\partial x_1}, \frac{\partial g_2}{\partial x_2}, \dots, \frac{\partial g_d}{\partial x_d} \right) = \nabla_x g (x_t, t)$$

$$\text{and } \frac{dx_t}{dt} = (\dot{x}_{t,1}, \dot{x}_{t,2}, \dots, \dot{x}_{t,d}) \text{ a vector}$$

$$\Rightarrow \frac{\partial g}{\partial x} (x_t, t) \cdot \frac{dx_t}{dt} \text{ a dot product} \rightarrow \text{scalar}$$

Total Derivative

$$\underbrace{P_t(x_t)}_{\text{Lagrangian}} = P(x_t, t)$$

$$\frac{d}{dt} g(x_t, t) = \nabla_x g(x_t, t) \cdot \frac{dx_t}{dt} + \frac{\partial g}{\partial t} (x, t) \Big|_{x=x_t}$$

$$\frac{d}{dt} g_t(x) = \nabla_x g_t(x_t) \frac{dx}{dt} + \frac{\partial}{\partial t} g_t(x) \Big|_{x=x_t}$$

So we take the derivative of  $\frac{d}{dt} \log p_t(x_t)$

$$\frac{d}{dt} \log p_t(x_t) = \left. \frac{\partial}{\partial t} \log p_t(x) \right|_{x=x_t} + \nabla_x \log p_t(x_t, t) \frac{\partial x_t}{\partial t}$$

$= \left. \frac{\partial}{\partial t} \log p_t(x) \right|_{x=x_t} + \nabla_x \log p_t(x_t, t) f_\theta(x_t, t)$

How does log density of  $x$  evolve as we push it through a flow?

Recall  $\frac{\partial x}{\partial t} = f_\theta(x_t, t)$

Why the Total Derivative?

- Because 1) The density field  $p_t$  changes wrt time ( $\frac{\partial}{\partial t}$ )  
2) The sample  $x_t$  also moves ( $\nabla_x \cdot \frac{\partial x}{\partial t}$ ) through space as time evolves

Now we have the formula

$$\frac{d}{dt} \log p_t(x_t) = \left. \frac{\partial}{\partial t} \log p_t(x) \right|_{x=x_t} + \nabla_x \log p_t(x_t, t) f_\theta(x_t, t)$$

why the  $|_{x=x_t}$  notation?

when applying the chain rule, wwt be explicit that

-  $\frac{\partial}{\partial t} \log p_t(x)$  means derivative wrt  $t$  while  
treating  $x$  constant

- Then evaluate the expression at the actual  
moving location  $x_t$

$$\frac{d}{dt} \log p_t(x_t) = \left. \frac{\partial}{\partial t} \log p_t(x) \right|_{x=x_t} + \nabla_x \log p_t(x_t, t) f_\theta(x_t, t)$$

$$+ \nabla_x \log p_t(x_t, t) f_\theta(x_t, t)$$

continuity equation

$$\frac{\partial}{\partial t} \log p_t(x) = -\nabla_x f_\theta(x_t, t) - f_\theta(x_t, t) \nabla_x \log p_t(x)$$

$$= \left[ -\nabla_x f_\theta(x_t, t) - f_\theta(x_t, t) \nabla_x \log p_t(x) \right]_{x=x_t} + \nabla_x \log p_t(x_t, t) f_\theta(x_t, t)$$

$$\therefore \frac{d}{dt} \log p_t(x_t) = -\nabla_x \cdot f_\theta(x_t, t)$$

How Do We Get  $\frac{\partial}{\partial t} p_t(x)$ ?

### Conservation of Mass (Continuous)

Continuity Equation

divergence of flux

$$\frac{\partial \rho}{\partial t} = -\nabla \cdot (\rho v)$$

local  
rate of  $\rho$

dot product

(Fluid Dynamics)

$\rho(x, t)$  density at point  $x$  and time  $t$

$v(x, t)$  the velocity field moving their density at that point and time

"stuff can't vanish or appear, only move"

State mass  
diverge out  
converge in

Continuity Equation using PDFs

$$\frac{\partial \rho}{\partial t} = -\nabla(\rho v)$$

$$\frac{\partial}{\partial t} p_t(x) + \nabla_x \cdot (f_\theta(x_t, t) p_t(x)) = 0$$

$$\Rightarrow \frac{\partial}{\partial t} p_t(x) = -\nabla_x (f_\theta(x_t, t) p_t(x))$$

log density  $\rightarrow \frac{\partial}{\partial t} \log p_t(x) = \frac{\partial}{\partial t} \log p_t(x) \frac{\partial p_t(x)}{\partial t}$

$$= \frac{1}{p_t(x)} \frac{\partial}{\partial t} p_t(x)$$

by chain rule

$$= \frac{1}{p_t(x)} \left[ -\nabla_x \cdot (f_\theta(x_t, t) p_t(x)) \right]$$

$$= -\frac{1}{p_t(x)} \left( \nabla_x f_\theta(x_t, t) p_t(x) + f_\theta(x_t, t) \nabla_x p_t(x) \right)$$

product rule

$$= - \frac{\nabla_x \cdot f_\theta(x_t, t) \cancel{P_t(x)}}{P_t(x)} - \frac{f_\theta(x_t, t) \nabla_x P_t(x)}{P_t(x)}$$

$$= - \nabla_x \cdot f_\theta(x_t, t) - \frac{f_\theta(x_t, t) \nabla_x P_t(x)}{P_t(x)}$$

$\nabla_x$  is gradient operator w.r.t  
 $x = (x_1, \dots, x_d)$

$$\nabla_x g(x) = \sum \frac{\partial g_i(x)}{\partial x_i}$$

$$\frac{\nabla_x P_t(x)}{P_t(x)} = \frac{1}{P_t(x)} \nabla_x P_t(x) = \nabla_x \log P_t(x)$$

Multivariate Chain Rule

$$\nabla_x (g \circ h)(x) = g'(h(x)) \nabla_x h(x)$$

$$\frac{\partial}{\partial t} \log P_t(x) = - \nabla_x \cdot f_\theta(x_t, t) - \frac{f_\theta(x_t, t) \nabla_x P_t(x)}{P_t(x)}$$

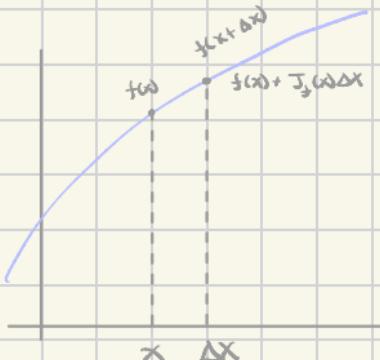
$$\frac{\partial}{\partial t} \log P_t(x) = - \nabla_x \cdot f_\theta(x_t, t) - f_\theta(x_t, t) \nabla_x \log P_t(x)$$

## Jacobian & Trace

$$x \in \mathbb{R}^d \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

Jacobian: linear approximation of  $f$  around a small neighborhood of  $x$

$$f(x + \Delta x) \approx f(x) + J_f(x) \Delta x$$



## Jacobian Matrix

$$J_f(x) = \frac{\partial f(x)}{\partial x}$$

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_d}{\partial x_1} & \dots & \frac{\partial f_d}{\partial x_d} \end{pmatrix}$$

Each diagonal term measures local stretching/compression along its own axis

Entry  $(i,j)$ : how much the  $i$ -th velocity component changes when you nudge the  $j$ -th coordinate

$\frac{\partial f_i}{\partial x_j} \rightarrow$  If you nudge  $x_j$ , how does the flow in the  $f_i$  direction speed up or slow down

## A Little More to Get Trace

Total Derivative

$$\frac{d}{dt} \log p_t(x_t) = -\nabla_x \cdot f_t(x_t, t)$$

off diagonals  
of Jacobian tell us  
rotation/shear

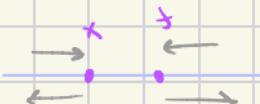
$$= -\sum_i \frac{\partial f_i(x_t, t)}{\partial x_i} = -\text{Tr}\left(\frac{\partial f_t}{\partial x_t}\right)$$

## Divergence

$$\nabla \cdot f(x) = \sum_{i=1}^d \frac{\partial f_i}{\partial x_i}$$

dot product

The diagonals of the Jacobian  
are the only entries of the  
Jacobian that directly tells us  
how volumes expand/contract  
along each axis



## Fundamental Theorem of Calculus (FTC)

$$F'(x) = f(x) \Rightarrow F(b) - F(a) = \int_a^b f(x) dx$$

equivalently

$$\frac{dy}{dt} = g(t) \Rightarrow y(T) - y(0) = \int_0^T g(t) dt$$



So now we have

$$\frac{d}{dt} \log p_t(x_t) = -\nabla_x \cdot f_0(x_t, t) = -\text{Tr}\left(\frac{\partial f_0}{\partial x_t}\right)$$

And

$$\log p(x_T) - \log p(x_0) = -\int_0^T \text{Tr}\left(\frac{\partial f_0}{\partial x_t}\right)$$

By FTC

Instantaneous Change-of-variables  
in continuous time

$$\lg p(x_t) = \lg p(x_T) - \int_0^T \text{Tr}\left(\frac{\partial f_0}{\partial x_t}\right)$$

Problem: Trace is  
the computational  
bottleneck!

Solution:

Hutchison  
Estimator  $\rightarrow$  still  
expensive

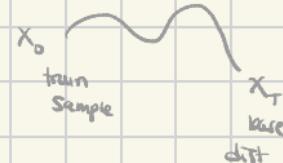
## Objective/Loss Function

Let  $x_T$  be base density (Gaussian)

$x_0$  be a training image

$$l(x; \theta) = \lg p(x_T) - \int_0^T \text{Tr}\left(\frac{\partial f_0}{\partial x_t}\right)$$

$$\text{MLE } \max_{\theta} \mathbb{E}_{x \sim \text{data}} [\log p(x)]$$



## Training -

- $x_0 \in \mathbb{R}^d$  sample

- Define trajectory via ODE

$$\frac{dx_t}{dt} = f_\theta(x_t, t) \quad x(0) = x_0$$

- Numerically integrate with ODE solver (PyTorch)  
to get  $x_T$

- choose simple base PDF  $p(x_T) = N(0, I)$

- compute

$$\log p(x_T)$$

- Compute instantaneous log-density change

$$\frac{d}{dt} \log p(x_t) = -\text{Tr}\left(\frac{\partial f_\theta}{\partial x_t}\right)$$

- Integrate  $\int_0^T \text{Tr}\left(\frac{\partial f_\theta}{\partial x_t}\right) dt$

- Compute

$$\log p(x_0) = \log p(x_T) - \int_0^T \text{Tr}\left(\frac{\partial f_\theta}{\partial x_t}\right) dt$$