

Flow & Diffusion Models

Theory

ODE/SDE

Goals:

- ① Flow and Diffusion models from first principles
- ② The necessary amount of mathematics.
- ③ Implement & apply these algos.

Structure:

1.
 - a. Formalize generating an image
 - b. Construct Flow & Diffusion models.
2. Define Training objective
3. Define training algorithm [score matching & flow matching]
- ④ Network architectures + Conditioning
- ⑤ Advanced Topics: Alignment, Complex Data Types, Distillation.

Lecture 1

SECTION-1 Generation to Sampling.

Data Distribution: Distribution of objects that we want to generate. $[P_{data}]$

$$P_{data}: \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0},$$

we don't know this prob. density.

$$z \mapsto P_{data}(z).$$

↓
sample.

Datasets:

A dataset consists finite number of samples from the data dist.

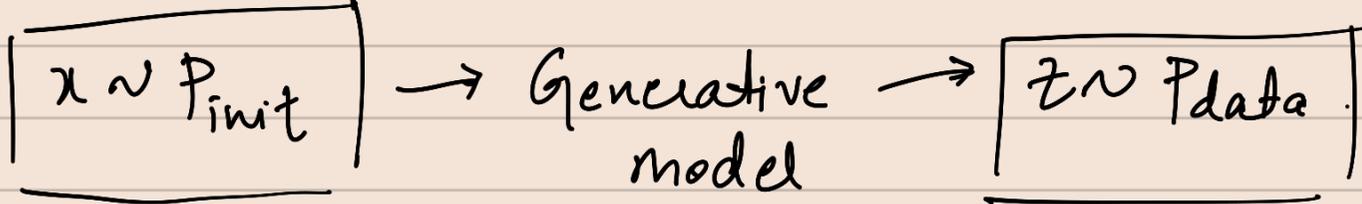
$$z_1, z_2, \dots, z_N \sim P_{data}.$$

Conditional Generation:

$$z \sim P(\cdot | y)_{data}$$

✓ First this
→ Next this

Generative Models convert samples from input dist. into samples from the data dist.

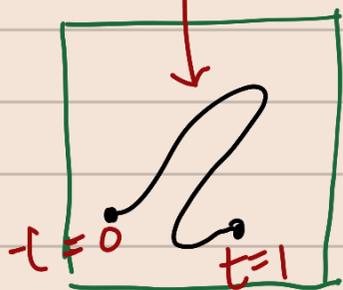


SECTION-2: Flow & Diffusion Models.

2A Flow Models

Trajectory :- $x: [0, t=1] \rightarrow \mathbb{R}^d$

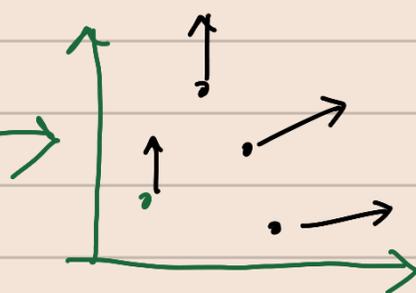
Essentially, $t \mapsto x_t$



Vector Field:

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

$$(x, t) \mapsto u_t(x)$$



Gives a dirn at every point.

Ordinary Differential Equation (ODE):

$$x_0 = x_0 \quad \left[\begin{array}{l} \text{Initial} \\ \text{Condition} \end{array} \right]$$

$$\frac{d}{dt} x_t = u_t(x_t) \quad \text{[ODE]} \rightarrow \text{This describes the direction of the } x_t, \text{ which is}$$

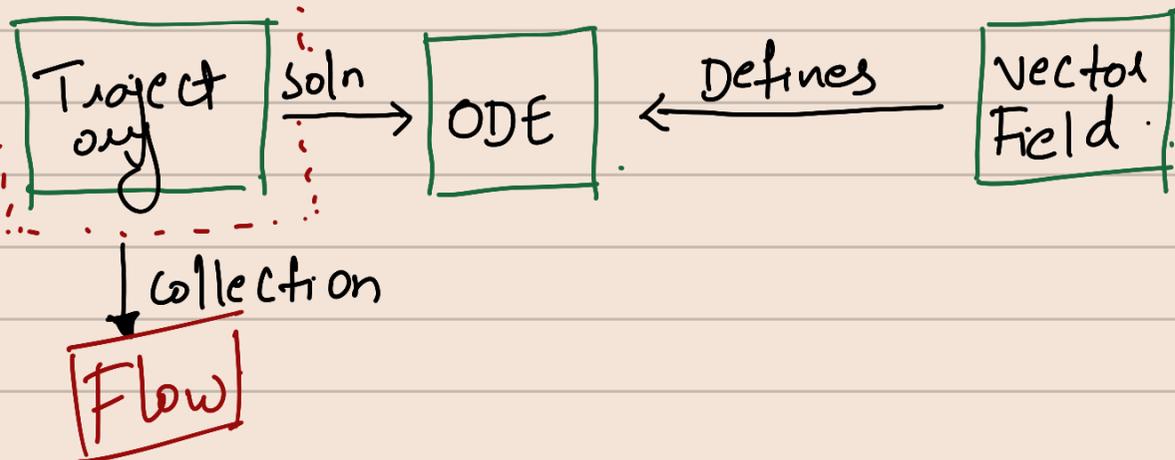
given by location of x_t , and specified by u_t

Flow: $\psi: \mathbb{R}^d \times [0, 1] \longrightarrow \mathbb{R}^d$
spatial
time

$$(x_0, t) \mapsto \psi_t(x_0)$$

$$\psi_0(x_0) = x_0 \quad \& \quad \frac{d}{dt} \psi_t(x_0) = u_t(\psi_t(x_0))$$

flow is essentially a collection of solutions to an ODE for a lot of initial conditions



Existence & Uniqueness Theorem ODE's

Picard-Lindelöf theorem: If a vector field $u_t(x)$ is continuously differentiable with bounded derivatives, then a unique solution to ODE

$$X_0 = x_0, \quad \frac{d}{dt} X_t = u_t(X_t)$$

exists.

In other words, a flow map exists. More generally this is true if vector field is Lipschitz.

Key Takeaway: In the cases of practical interest for machine learning, unique solutions to ODE/Flows exist.

Example Linear ODE

Simple Vector Field:
 $u_t(x) = -\theta x \quad (\theta > 0)$

Claim: Flow is given by.

$$\psi_t(x_0) = \exp(-\theta t) x_0$$

Proof:

Initial condition:

$$\psi_t(x_0) = \exp(0) x_0 = x_0$$

ODE:

$$\begin{aligned} \frac{d}{dt} \psi_t(x_0) &= \frac{d}{dt} \exp(-\theta t) x_0 = -\theta \exp(-\theta t) x_0 \\ &= -\theta \psi_t(x_0) \end{aligned}$$

$$= u_t(\psi_t(x_0))$$

Numerical Simulation of an ODE:

Euler Method.

For a given vector field u_t , initial condition x_0 ,

take a small step in the direction of vector field

and then return trajectory.

Flow Model:

$$P_{\text{init}} \xrightarrow{\text{ODE}} P_{\text{data}}$$

Neural Network: $u_t^\theta: \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$

Random init: $X_0 \sim P_{\text{init}}$

$$\text{ODE: } \frac{d}{dt} X_t = u_t^\theta(X_t)$$

Neural Network
Vector Field

Goal: $X_1 \sim P_{\text{data}}$

Training = Find parameters θ such that

$$X_0 \sim P_{\text{init}}, \quad dX_t = u_t^\theta(X_t) \cdot dt \quad \xrightarrow{\text{implies}} \quad X_1 \sim P_{\text{data}}$$

start with
init dist.

Follow along the
vector field

The dist.
of final point
= data dist.

Lecture 2.0

Deriving Training Target:

we need to explicitly derive it,
unlike in regression & classification.

Margin Vector Field \rightarrow Flow Matching

Marginal score F^n \rightarrow Score Matching.

Conditional Prob Path \rightarrow Conditional Vector Field \rightarrow Conditional score F^n

\rightarrow Per single Data Point

\rightarrow Across dist. of data points.

Marginal Prob path \rightarrow Marginal vector field \rightarrow Marginal score F^n

Probability Paths: The path from noise to data.

The interpolation between noise dist & data

A prob. path only specifies the marginals *dist.*
 It says nothing about the evolution of a single particle.

Dirac Dist:

$$z \in \mathbb{R}^d, \delta_z \quad X \sim \delta_z$$

Always returns the same thing

$$\Rightarrow X = z$$

Conditional probability path: $P_t(\cdot | z)$

① $P_t(\cdot | z)$ distribution over \mathbb{R}^d

② $P_0(\cdot | z) = p_{\text{init}}, P_1(\cdot | z) = \delta_z$

Example: Gaussian probability path.

$$P_t(\cdot | z) = \mathcal{N}(\alpha_t z, \beta_t^2 \mathbb{I}_d) \quad \boxed{\text{Noise scheduler}}$$

① True

② $P_0(\cdot | z) = \mathcal{N}(0, \mathbb{I}_d) \checkmark$
 $P_1(\cdot | z) = \mathcal{N}(z, 0) \checkmark$

α_t, β_t s.t

$$\alpha_0 = 0 \quad \beta_0 = 1$$

$$\alpha_1 = 1 \quad \beta_1 = 0$$

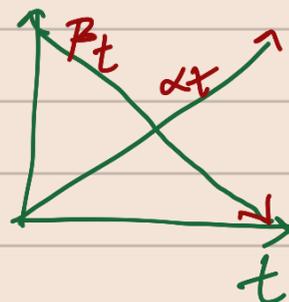
Marginal probability path: P_t

$$z \sim P_{\text{data}}, X \sim P_t(\cdot | z)$$

$$\Rightarrow \boxed{X \sim P_t}$$

forget z .

$$\textcircled{1} P_t(x) = \int P_t(x | z) P_{\text{data}}(z) \cdot dz$$



(2) $P_0 = P_{\text{init}}, P_1 = P_{\text{data}}$

Conditional & Marginal Vector Field:

Conditional Vector Field: $u_t^{\text{target}}(x|z)$
 $(0 \leq t \leq 1)$
 $x, z \in \mathbb{R}^d$

z

$P_{\text{init}} \xrightarrow[\text{ODE}]{P_t(\cdot|z)} \delta z$

such that

$x_0 \sim P_{\text{init}}, \frac{d}{dt} x_t = u_t^{\text{target}}(x_t|z)$

$P_0(\cdot|z) \Rightarrow x_t \sim P_t(\cdot|z)$
 $(0 \leq t \leq 1)$

If we follow this vec. field, we should be necessarily following the Prob. Path.

Cond. Gaussian VF:

$u_t^{\text{target}}(x|z) = \left(\alpha_t^\circ - \frac{\beta_t^\circ}{\beta_t} \alpha_t \right) z + \frac{\beta_t^\circ}{\beta_t} x$

$\alpha_t^\circ = \frac{d\alpha_t}{dt}, \beta_t^\circ$

Theorem (Marginalization Trick):

The marginal vector field by

$$u_t^{\text{target}}(\lambda) = \int u_t^{\text{target}}(\lambda|z) \frac{P_t(\lambda|z) P_{\text{data}}(z)}{P_t(\lambda)} dz$$

satisfies

$$X_0 \sim P_{\text{init}}, \frac{d}{dt} X_t = u_t^{\text{target}}(X_t) \Rightarrow X_t \sim P_t \quad (0 \leq t \leq 1).$$

Essentially, if we are following this marginal density then we are essentially following this marginal path. This means at the end we'll end up at P_{data} in the end.

Continuity Equation:

we'll use to prove the above statement.

Let us consider a flow model with vector field u_t^{target} with $X_0 \sim P_{\text{init}}$. Then $X_t \sim P_t, \forall 0 \leq t \leq 1$ i.f.f

$$\partial_t P_t \lambda = -\text{div}(P_t u_t^{\text{target}}) \lambda, \forall \lambda, 0 \leq t \leq 1$$

FLOW MATCHING: u_t^θ

Goal: $u_t^\theta \approx u_t^{\text{target}}$

Flow matching loss: $l_{fm}(\theta) = \mathbb{E} \left[\left\| u_t^\theta(x) - u_t^{\text{target}}(x) \right\|^2 \right]$

Intractable

$t \sim \text{unif}$ ← Uniform in $[0, 1]$
 $z \sim p_{\text{data}}$ ← draw data point
 $x \sim p_t(\cdot | z)$ ← draw from cond. path.

Conditional to the rescue.

Conditional Flow Matching loss:

$$l_{CFM}(\theta) = \mathbb{E} \left[\left\| u_t^\theta(x) - u_t^{\text{target}}(x|z) \right\|^2 \right]$$

$t \sim \text{unif}$
 $z \sim p_{\text{data}}$
 $x \sim p_t(\cdot | z)$

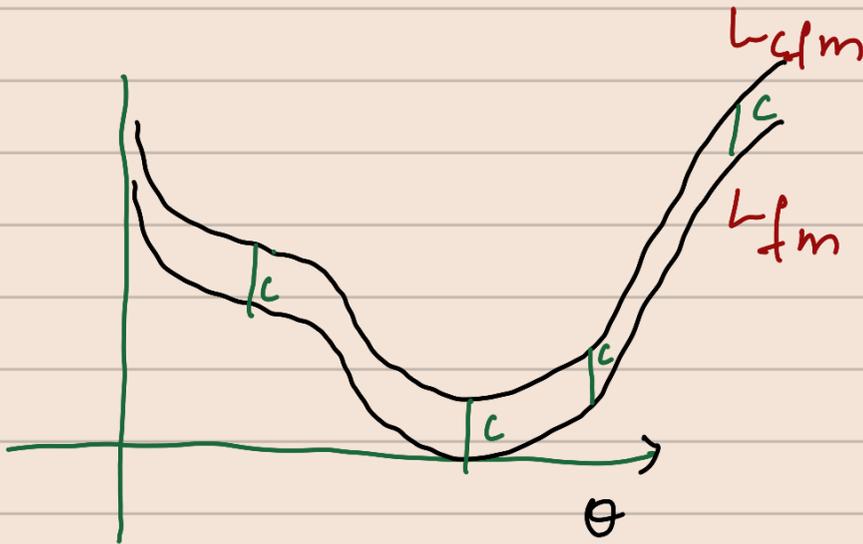
↓
 Tractable ✓
 Minimizer? $u_t^{\text{target}}(x|z)$

The beauty of FM will be seen when it will be shown that minimizing l_{CFM} is enough to minimize l_{FM} .

Theorem

$$L_{fm}(\theta) = L_{CFM}(\theta) + c.$$

for $c > 0$ & independent of θ .



⇒ (1) For minimize θ^* of L_{CFM}

$$u_t^{\theta^*} = u_t^{\text{target}}$$

(2) $\nabla_{\theta} L_{CFM} = \nabla_{\theta} L_{FM}(\theta)$

⇒ SGD the same.

Flow Matching Training Procedure

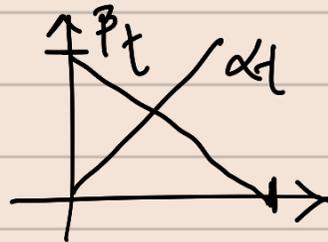
- for each mini-batch of data do
 - sample a data example z
 - sample a random time $t \sim \text{Unif}(0, 1)$
 - sample $x \sim P_t(\cdot | z)$
 - compute loss L_{CFM} .
 - update model params θ .

WCFM for Gaussian Conditional path:

$$P_t(\cdot | z) = N(\alpha_t z, \beta_t^2 \mathbb{I}_d)$$

target

$$u_t^\theta(z) = \left(\alpha_t^\bullet - \frac{\beta_t^\bullet}{\beta_t} \alpha_t \right) z + \frac{\beta_t^\bullet}{\beta_t} z$$



Sampling from $P_t(\cdot | z)$ by

$$\epsilon \sim N(0, \mathbb{I}_d) \Rightarrow \alpha_t z + \beta_t \epsilon \stackrel{\text{def}}{=} x$$

$$L_{\text{CFM}}(\theta) = \mathbb{E}_{\substack{t \sim \text{unif} \\ z \sim P_{\text{data}} \\ x \sim N(\alpha_t z, \beta_t^2 \mathbb{I}_d)}} \left[\left\| u_t^\theta(z) - \left(\alpha_t^\bullet - \frac{\beta_t^\bullet}{\beta_t} \alpha_t \right) z - \frac{\beta_t^\bullet}{\beta_t} z \right\|^2 \right] \quad x \sim P_t(\cdot | z)$$

$$x = \alpha_t z + \beta_t \epsilon$$

$$= \mathbb{E}_{\substack{t \sim \text{unif} \\ z \sim P_{\text{data}} \\ \epsilon \sim N(0, \mathbb{I}_d)}} \left[\left\| u_t^\theta(\alpha_t z + \beta_t \epsilon) - (\alpha_t^\bullet z - \beta_t^\bullet \epsilon) \right\|^2 \right]$$

$$\begin{aligned} t &\sim \text{unif} \\ z &\sim P_{\text{data}} \\ \epsilon &\sim N(0, \mathbb{I}_d) \end{aligned}$$

with specific α_t & β_t .

Cond OT Path:

$$\begin{aligned} \alpha_t &= t & \alpha_t^\bullet &= 1 \\ \beta_t &= 1-t & \beta_t^\bullet &= -1 \end{aligned}$$

$$L_{CFM} = \mathbb{E} \left[\left\| u_t^\theta(\alpha_t z + \beta_t \epsilon) - (z - \epsilon) \right\|^2 \right].$$

The OT corresponds to the optimal transport path (straight line path) from noise to data.

Simple Loss

Lecture 4

A Guided CFM Objective:

Observation 1

For fixed y , we obtain the unguided problem, and may adapt an unguided objective to obtain:

$$L_{CFM}^{\text{guided}}(\theta; y) = \mathbb{E}_{\square} \left\| u_t^\theta(z|y) - u_t^{\text{target}}(z|z) \right\|^2$$

$$\square = z \sim P_{\text{data}}(z|y), t \sim \text{Unif}(0,1), \lambda \sim P_t(\lambda|z)$$

Observation 2

By varying y , the above yields a guided objective for general y :

$$L_{CFM}^{\text{guided}}(\theta) = \mathbb{E}_{\square} \left\| u_t^{\theta}(x|y) - u_t^{\text{target}}(x|z) \right\|^2$$

$$\square = (z, y) \sim P_{\text{data}}(z, y), \quad t \sim \text{Unif}(0, 1), \\ x \sim P_t(x|z)$$

Guided sampling Procedure:

Obtain a trained guided vector field $u_t^{\theta}(x|y)$

- ① Select a prompt $y \in \mathcal{Y}$
- ② Initialize $X_0 \sim P_{\text{init}}$
- ③ Simulate $dX_t = u_t^{\theta}(X_t|y) dt$ from $t=0$ to $t=1$.

Can we do better? At least empirically, yes.

People have realized you could trade-off

diversity for perceptual quality

CFG.

Classifier Free Guidance.

We'll build this using Gaussian Paths

Recall: A Gaussian Conditional Prob path.

$$P_t(x|z) = \mathcal{N}(\alpha_t x, \beta_t^2 \text{Id})$$

where α_t, β_t are continuously diff'able, monotonic functions satisfying $\alpha_1 = \beta_0 = 1$ & $\alpha_0 = \beta_1 = 0$.

For Gaussian Probability Paths it can be shown that

$$\overset{-\text{target}}{u}_t(x|y) = \overset{-\text{target}}{u}_t(x) + b_t \nabla \log P_t(y|x),$$
$$b_t = \frac{\alpha_t \beta_t^2 - \beta_t \alpha_t}{\alpha_t}$$

For fixed w , we may define

$$\tilde{u}_t(x|y) = \overset{\text{target}}{u}_t(x) + w b_t \nabla \log P_t(y|x)$$

Rearranging yields

$$\tilde{u}_t(x|y) = (1-w) \underbrace{u_t^{\text{target}}(x)}_{\text{Unguided}} + w \underbrace{u_t^{\text{target}}(x|y)}_{\text{Guided}}.$$

CFG Training

$$u_t^{\text{-target}}(x) = u_t^{\text{target}}(x|y=\emptyset)$$

we may now train a single model $u_t^\theta(x|y)$, $y \in \{y, \emptyset\}$ by re-using $L_{\text{CFM}}^{\text{guided}}(\theta)$ & occasionally setting $y = \emptyset$.

$$L_{\text{CFM}}^{\text{CFG}}(\theta) = \mathbb{E}_{\square} \left\| u_t^\theta(x|y) - u_t^{\text{target}}(x|z) \right\|^2$$

$\square = (z, y) \sim P_{\text{data}}(z, y)$ with prob η , $y \leftarrow \emptyset$,
 $t \sim \text{Unif}[0, 1]$, $x \sim P_t(x|z)$

CFG Sampling Procedure

$$\text{Simulate } dx_t = \left[(1-w) u_t^\theta(x_t | \emptyset) + w u_t^\theta(x_t | y) \right] dt$$

from $t=0$ to $t=1$