



Temporal Difference Flows [FPT⁺25]

Kellen Kanarios

July 10, 2025

Outline

- ① Motivation
- ② Flow Matching
- ③ Application to Reinforcement Learning

Variational inference with normalizing flows [RM]

- ▶ Want to maximize log-likelihood $\log p_{\theta}(\mathbf{x})$. Introduce latent $\mathbf{z} \sim p$.

$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{D}_{\text{KL}}[q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p(\mathbf{z})] + \mathbb{E}_q[\log p_{\theta}(\mathbf{x} \mid \mathbf{z})]. \quad (\text{ELBO})$$

- ▶ Simultaneously learn q_{ϕ} and p_{θ}

Variational inference with normalizing flows [RM]

- ▶ Want to maximize log-likelihood $\log p_{\theta}(\mathbf{x})$. Introduce latent $\mathbf{z} \sim p$.

$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{D}_{\text{KL}}[q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p(\mathbf{z})] + \mathbb{E}_q[\log p_{\theta}(\mathbf{x} \mid \mathbf{z})]. \quad (\text{ELBO})$$

- ▶ Simultaneously learn q_{ϕ} and p_{θ}

Problem. Need family q_{ϕ} to contain $p_{\theta}(z \mid x)$.

- ▶ Rarely the case.

Variational inference with normalizing flows [RM]

- ▶ Want to maximize log-likelihood $\log p_{\theta}(\mathbf{x})$. Introduce latent $\mathbf{z} \sim p$.

$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{D}_{\text{KL}}[q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p(\mathbf{z})] + \mathbb{E}_q[\log p_{\theta}(\mathbf{x} \mid \mathbf{z})]. \quad (\text{ELBO})$$

- ▶ Simultaneously learn q_{ϕ} and p_{θ}

Problem. Need family q_{ϕ} to contain $p_{\theta}(\mathbf{z} \mid \mathbf{x})$.

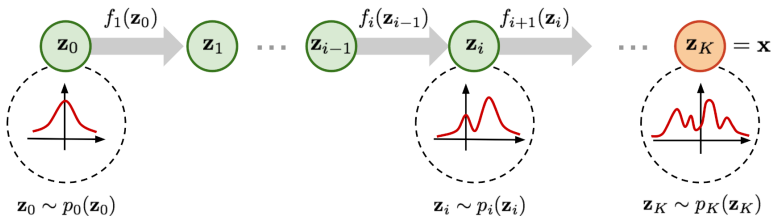
- ▶ Rarely the case.

How do we create a more expressive family that we can still optimize?

Given target distribution q . Pick some initial distribution p_0 , where we can sample $\mathbf{z}_0 \sim p_0$.

Idea: Learn sequence of functions f_1, f_2, \dots, f_K , such that

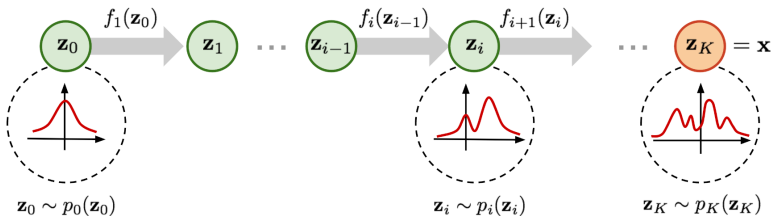
$$f_K \circ f_{K-1} \cdots \circ f_1(\mathbf{z}_0) \sim q$$



Given target distribution q . Pick some initial distribution p_0 , where we can sample $\mathbf{z}_0 \sim p_0$.

Idea: Learn sequence of functions f_1, f_2, \dots, f_K , such that

$$f_K \circ f_{K-1} \cdots \circ f_1(\mathbf{z}_0) \sim q$$

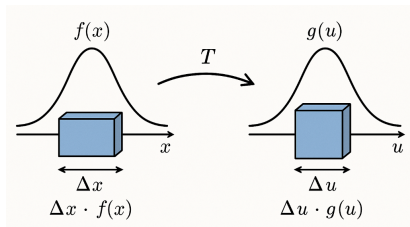


How do we learn such a sequence of functions f_i ?

Given **invertible, smooth both ways** function $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$.

Change of variables theorem

$$\int_{\mathbb{R}^d} p_0(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^d} p_0(T^{-1}(\mathbf{y})) \left| \det \frac{\partial T}{\partial \mathbf{y}} \right|^{-1} d\mathbf{y}.$$



$$\Rightarrow \text{Density of } f_1(\mathbf{z}_0) = p_0(\mathbf{z}_0) \left| \det \frac{\partial f_1}{\partial \mathbf{z}} \right|^{-1}$$

$$\Rightarrow \log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right|^{-1}$$

$$\log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right|^{-1}$$

1. Given $\mathbf{x}^{(i)}$ from dataset, compute $\mathbf{z}_0^{(i)} = (f_K \circ \dots \circ f_1)^{-1}(\mathbf{x}^{(i)})$

$$\log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \left| \det \frac{\partial \mathbf{f}_k}{\partial \mathbf{z}_{k-1}} \right|^{-1}$$

1. Given $\mathbf{x}^{(i)}$ from dataset, compute $\mathbf{x}_0^{(i)} = (\mathbf{f}_K \circ \dots \circ \mathbf{f}_1)^{-1}(\mathbf{x}^{(i)})$
2. Maximize log-likelihood

$$\max_{\theta} \left[\log p_0(\mathbf{x}_0^{(i)}(\theta)) - \sum_{k=1}^K \left| \det \frac{\partial \mathbf{f}_k(\theta)}{\partial \mathbf{x}_{k-1}^{(i)}} \right|^{-1} \right].$$

$$\log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \left| \det \frac{\partial \mathbf{f}_k}{\partial \mathbf{z}_{k-1}} \right|^{-1}$$

1. Given $\mathbf{x}^{(i)}$ from dataset, compute $\mathbf{x}_0^{(i)} = (\mathbf{f}_K \circ \dots \circ \mathbf{f}_1)^{-1}(\mathbf{x}^{(i)})$
2. Maximize log-likelihood $\cdot O(LD^3)$

$$\max_{\theta} \left[\log p_0(\mathbf{x}_0^{(i)}(\theta)) - \sum_{k=1}^K \left| \det \frac{\partial \mathbf{f}_k(\theta)}{\partial \mathbf{x}_{k-1}^{(i)}} \right|^{-1} \right].$$

$$\log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \left| \det \frac{\partial \mathbf{f}_k}{\partial \mathbf{z}_{k-1}} \right|^{-1}$$

1. Given $\mathbf{x}^{(i)}$ from dataset, compute $\mathbf{x}_0^{(i)} = (\mathbf{f}_K \circ \dots \circ \mathbf{f}_1)^{-1}(\mathbf{x}^{(i)})$
2. Maximize log-likelihood $O(LD^3)$

$$\max_{\theta} \left[\log p_0(\mathbf{x}_0^{(i)}(\theta)) - \sum_{k=1}^K \left| \det \frac{\partial \mathbf{f}_k(\theta)}{\partial \mathbf{x}_{k-1}^{(i)}} \right|^{-1} \right].$$

Need to pick **simple** transformations i.e.

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^\top \mathbf{z} + b).$$

Require **very large** K to represent complex distributions.

$$\log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \left| \det \frac{\partial \mathbf{f}_k}{\partial \mathbf{z}_{k-1}} \right|^{-1}$$

1. Given $\mathbf{x}^{(i)}$ from dataset, compute $\mathbf{x}_0^{(i)} = (\mathbf{f}_K \circ \dots \circ \mathbf{f}_1)^{-1}(\mathbf{x}^{(i)})$
2. Maximize log-likelihood

$$\max_{\theta} \left[\log p_0(\mathbf{x}_0^{(i)}(\theta)) - \sum_{k=1}^K \left| \det \frac{\partial \mathbf{f}_k(\theta)}{\partial \mathbf{x}_{k-1}^{(i)}} \right|^{-1} \right].$$

$\cdot O(LD^3)$

Need to pick **simple** transformations i.e.

$$\mathbf{f}(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^\top \mathbf{z} + b).$$

Require **very large** K to represent complex distributions.

Why not take $K \rightarrow \infty \dots$

Neural Ordinary Differential Equations [CRBD]

Replace $\mathbf{z}_{t+1} = f_t(\mathbf{z}_t)$ with the **ODE**

$$\frac{d\mathbf{z}}{dt} = f_t(\mathbf{z}_t).$$

Remark. *Instantaneous COV*

$$\frac{\partial \log p_t(\mathbf{z}_t)}{\partial t} = -\text{tr} \left(\frac{df}{d\mathbf{z}_t} \right)$$

Neural Ordinary Differential Equations [CRBD]

Replace $\mathbf{z}_{t+1} = f_t(\mathbf{z}_t)$ with the **ODE**

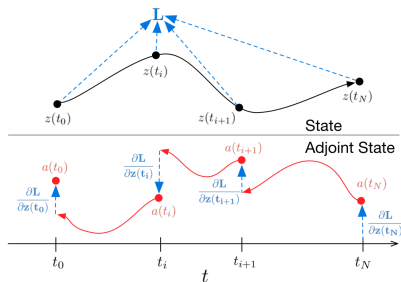
$$\frac{d\mathbf{z}}{dt} = f_t(\mathbf{z}_t).$$

Remark. *Instantaneous COV*

$$\frac{\partial \log p_t(\mathbf{z}_t)}{\partial t} = -\text{tr} \left(\frac{df}{d\mathbf{z}_t} \right)$$

Given $\mathbf{x}^{(i)}$ from dataset, now must compute $\mathbf{x}_0^{(i)} = \int_1^0 f_t^{-1}(\mathbf{x}_t^{(i)}) dt$

► Requires **exact** ODE solve for unbiased gradient.



Neural Ordinary Differential Equations [CRBD]

Replace $\mathbf{z}_{t+1} = f_t(\mathbf{z}_t)$ with the **ODE**

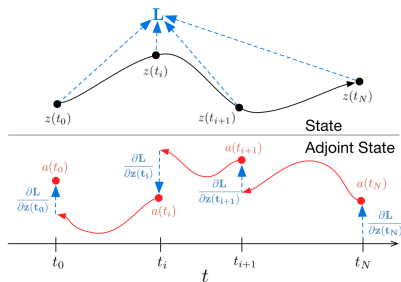
$$\frac{d\mathbf{z}}{dt} = f_t(\mathbf{z}_t).$$

Remark. *Instantaneous COV*

$$\frac{\partial \log p_t(\mathbf{z}_t)}{\partial t} = -\text{tr} \left(\frac{df}{d\mathbf{z}_t} \right)$$

Given $\mathbf{x}^{(i)}$ from dataset, now must compute $\mathbf{x}_0^{(i)} = \int_1^0 f_t^{-1}(\mathbf{x}_t^{(i)}) dt$

► Requires **exact** ODE solve for unbiased gradient.



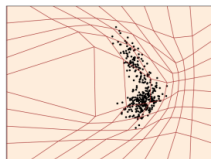
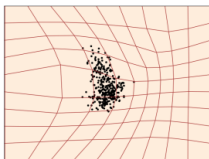
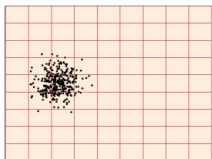
Can we learn f_t without backprop through an ODE?

Flow Matching for Generative Modeling [LCB⁺23]

Disclaimer: f_t will now often be referred to as the *flow*.

Book keeping. We say the *flow* f_t generates a probability path p_t if $X_t = f_t(X_0) \sim p_t$. Equivalently,

$$p_t(x) = [f_{t\#}p_0](x) \triangleq p_0(f_t^{-1}(y)) \left| \det \partial_y f_t^{-1}(y) \right|.$$



Key Idea: Do **not** learn the *flow* f_t , learn the **velocity** of the flow.

Def. This velocity is a *vector field* $u_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$, such that

$$\frac{df_t(\mathbf{x})}{dt} = u_t(f_t(\mathbf{x})).$$

Key Idea: Do **not** learn the *flow* f_t , learn the **velocity** of the flow.

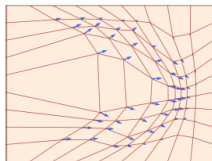
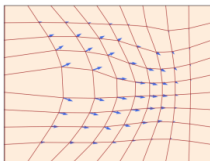
Def. This velocity is a *vector field* $u_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$, such that

$$\frac{df_t(\mathbf{x})}{dt} = u_t(f_t(\mathbf{x})).$$

To sample, solve ODE at **sample time** i.e.

$$f_{t+\Delta t}(\mathbf{x}) \approx f_t(\mathbf{x}) + \Delta t \cdot u_t(f_t(\mathbf{x})) \quad (\text{Euler method})$$

By uniqueness of ODE, (vector field u_t) \leftrightarrow (flow f_t).



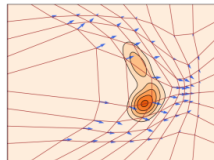
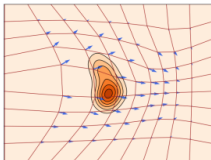
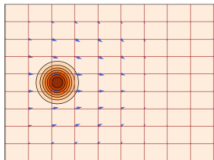
Def. A vector field u_t generates a probability path p_t if its corresponding flow f_t generates p_t .

Thm. A vector field u_t generates a probability path p_t if and only if it satisfies the continuity equation.

Continuity Equation.

$$\frac{d}{dt}p_t(x) + \operatorname{div}(p_t u_t)(x) = 0,$$

where $\operatorname{div}(v)(x) = \sum_{i=1}^d \partial_{x^i} v^i(x)$, and $v(x) = (v^1(x), \dots, v^d(x))$.



Continuity Equation.

$$\frac{d}{dt} p_t(x) + \operatorname{div}(p_t u_t)(x) = 0,$$

where $\operatorname{div}(v)(x) = \sum_{i=1}^d \partial_{x^i} v^i(x)$, and $v(x) = (v^1(x), \dots, v^d(x))$.

Proof: (If people care)

$$\begin{aligned} \frac{d}{dt} \mathbb{E} f(X_t) &= \frac{d}{dt} \int_{\mathbb{R}^d} f(x) p_t(x) dx = \int_{\mathbb{R}^d} f(x) \partial_t p_t(x) dx, \\ \mathbb{E} \frac{d}{dt} f(X_t) &= \int_{\mathbb{R}^d} (\nabla f(x_t) \cdot v_t(x_t)) p_t(x_t) dx_t \\ &= 0 - \int_{\mathbb{R}^d} f(x_t) \operatorname{div}(v_t(x_t) p_t(x_t)) dx_t \quad (\text{IBP}) \\ &= \int_{\mathbb{R}^d} -f(x) \operatorname{div}(v_t(x) p_t(x)) dx. \end{aligned}$$

The result follows from fundamental lemma of calculus of variations.

How do we actually learn the vector field u_t ?

Flow Model Loss

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, p_t(x)} \|v_t(x; \theta) - u_t(x)\|^2$$

- ▶ $p_t : \mathcal{X} \rightarrow [0, 1]$: probability density path.
- ▶ $u_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$: vector field that *generates* p_t .

Flow Model Loss

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, p_t(x)} \|v_t(x; \theta) - u_t(x)\|^2$$

- ▶ $p_t : \mathcal{X} \rightarrow [0, 1]$: probability density path.
- ▶ $u_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$: vector field that *generates* p_t .

Intractable.

1. How to choose p_t ?

Flow Model Loss

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, p_t(x)} \|v_t(x; \theta) - u_t(x)\|^2$$

- ▶ $p_t : \mathcal{X} \rightarrow [0, 1]$: probability density path.
- ▶ $u_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$: vector field that *generates* p_t .

Intractable.

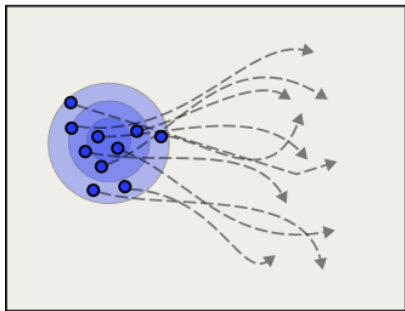
1. How to choose p_t ?
2. Given p_t need to solve continuity equation for u_t , which is a **PDE** likely without close-form.

Key Idea: Given the endpoint X_1 , we can easily construct a path between X and X_1 .

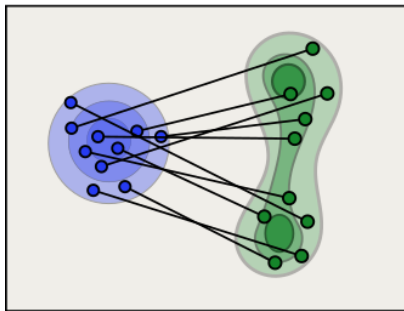
Ex: Consider

$$f_t(x \mid x_1) = tx_1 + (1 - (1 - \sigma_{\min})t)x$$

$f_t(X)$



$f_t(X \mid X_1)$



What does the conditional flow being nice have to do with our problem?

Idea 1: Given target data distribution q , approximate q by marginalizing “nice” conditionals by q i.e.

$$p_t(\mathbf{x}) = \int p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1) d\mathbf{x}_1. \quad (\text{Marginal Density})$$

For $p_1(\mathbf{x} \mid \mathbf{x}_1) \sim N(\mathbf{x}_1, \sigma_1)$, with $\sigma_1 \ll 1$, $p_1(\mathbf{x}) \approx q(\mathbf{x})$.

Idea 1: Given target data distribution q , approximate q by marginalizing “nice” conditionals by q i.e.

$$p_t(\mathbf{x}) = \int p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1) d\mathbf{x}_1. \quad (\text{Marginal Density})$$

For $p_1(\mathbf{x} \mid \mathbf{x}_1) \sim N(\mathbf{x}_1, \sigma_1)$, with $\sigma_1 \ll 1$, $p_1(\mathbf{x}) \approx q(\mathbf{x})$.

Idea 2: Use this to define an approximate vector field

$$u_t(\mathbf{x}) = \int u_t(\mathbf{x} \mid \mathbf{x}_1) \frac{p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1)}{p_t(\mathbf{x})} d\mathbf{x}_1. \quad (\text{Marginal Vector Field})$$

Here, $u_t(\mathbf{x} \mid \mathbf{x}_1)$ is the conditional vector field that generates $p_t(\mathbf{x} \mid \mathbf{x}_1)$.

Idea 1: Given target data distribution q , approximate q by marginalizing “nice” conditionals by q i.e.

$$p_t(\mathbf{x}) = \int p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1) d\mathbf{x}_1. \quad (\text{Marginal Density})$$

For $p_1(\mathbf{x} \mid \mathbf{x}_1) \sim N(\mathbf{x}_1, \sigma_1)$, with $\sigma_1 \ll 1$, $p_1(\mathbf{x}) \approx q(\mathbf{x})$.

Idea 2: Use this to define an approximate vector field

$$u_t(\mathbf{x}) = \int u_t(\mathbf{x} \mid \mathbf{x}_1) \frac{p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1)}{p_t(\mathbf{x})} d\mathbf{x}_1. \quad (\text{Marginal Vector Field})$$

Here, $u_t(\mathbf{x} \mid \mathbf{x}_1)$ is the conditional vector field that generates $p_t(\mathbf{x} \mid \mathbf{x}_1)$.

Thm. The marginal vector field u_t generates the marginal probability path p_t . (Check continuity equation)

Conditional Flow Model Loss

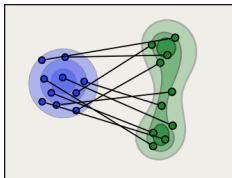
$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, q(x_1), p_t(x|x_1)} \|v_t(x; \theta) - u_t(x | x_1)\|^2$$

Conditional Flow Model Loss

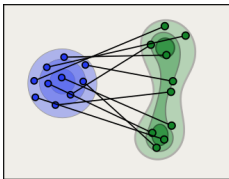
$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, q(x_1), p_t(x|x_1)} \|v_t(x; \theta) - u_t(x | x_1)\|^2$$

Thm. Up to a constant independent of θ , $\mathcal{L}_{\text{FM}}(\theta) = \mathcal{L}_{\text{CFM}}(\theta)$. In particular, $\nabla \mathcal{L}_{\text{FM}}(\theta) = \nabla \mathcal{L}_{\text{CFM}}(\theta)$.

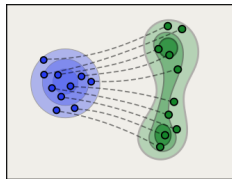
Batch i



Batch j



$\mathbb{E}[\text{Batch}]$



How do we actually use this?

Example: Optimal Transport

For $f_t(x) = \mu_t(x_1) + x\sigma_t(x_1)$, we consider

$$\mu_t = tx_1, \quad \sigma_t(x) = 1 - (1 - \sigma_{\min})t,$$

such that

Example: Optimal Transport

For $f_t(x) = \mu_t(x_1) + x\sigma_t(x_1)$, we consider

$$\mu_t = tx_1, \quad \sigma_t(x) = 1 - (1 - \sigma_{\min})t,$$

such that

$$f_t^{-1}(x) = \frac{x - tx_1}{1 - (1 - \sigma_{\min})t}, \quad \frac{d}{dt}f_t(x) = x_1 - (1 - \sigma_{\min})x.$$

Example: Optimal Transport

For $f_t(x) = \mu_t(x_1) + x\sigma_t(x_1)$, we consider

$$\mu_t = tx_1, \quad \sigma_t(x) = 1 - (1 - \sigma_{\min})t,$$

such that

$$f_t^{-1}(x) = \frac{x - tx_1}{1 - (1 - \sigma_{\min})t}, \quad \frac{d}{dt}f_t(x) = x_1 - (1 - \sigma_{\min})x.$$

Recall

$$\frac{d}{dt}f_t(x) = u(f_t(x) \mid x_1) \implies u(x \mid x_1) = \frac{d}{dt}f_t(f_t^{-1}(x)).$$

Example: Optimal Transport

Recall

$$\frac{d}{dt}f_t(x) = u(f_t(x) \mid x_1) \implies u(x \mid x_1) = \frac{d}{dt}f_t(f_t^{-1}(x)).$$

Thus,

$$u_t(x \mid x_1) = \frac{x_1 - (1 - \sigma_{\min})x}{1 - (1 - \sigma_{\min})t}.$$

Example: Optimal Transport

Recall

$$\frac{d}{dt}f_t(x) = u(f_t(x) \mid x_1) \implies u(x \mid x_1) = \frac{d}{dt}f_t(f_t^{-1}(x)).$$

Thus,

$$u_t(x \mid x_1) = \frac{x_1 - (1 - \sigma_{\min})x}{1 - (1 - \sigma_{\min})t}.$$

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, q(x_1), p_t(x|x_1)} \left\| v_t(x; \theta) - \frac{x_1 - (1 - \sigma_{\min})x}{1 - (1 - \sigma_{\min})t} \right\|^2$$

Diffusion Meets Flow Matching: Two Sides of the Same Coin **[GHH⁺24]**

Recall that instead of an ODE, in diffusion, we have a **SDE**

$$d\mathbf{x}_t = f_t(\mathbf{x}_t)dt + \sigma(\mathbf{x}_t)dB_t.$$

However, for *OU* process, we can actually absorb Brownian motion term and get vector field

$$u_t(\mathbf{x}_t) = -(\mathbf{x}_t + \nabla \ln p_t(\mathbf{x}_t)).$$

Similarly, learning the score function $\nabla \ln p_t(\mathbf{x}_t)$ can be rewritten as the flow matching objective for certain choices of p_t see [LHH⁺24].

Why Flows?

Anecdotally.

- ▶ *The OT path's conditional vector field has constant direction in time and is arguably simpler to fit with a parametric model. [LCB⁺23]*
- ▶ *The deterministic nature of ODEs equips flow-matching methods with simpler learning objectives and faster inference speed [ZPLE25]*



$t = 0.0$

$t = 1/3$

$t = 2/3$

$t = 1.0$

$t = 0.0$

$t = 1/3$

$t = 2/3$

$t = 1.0$

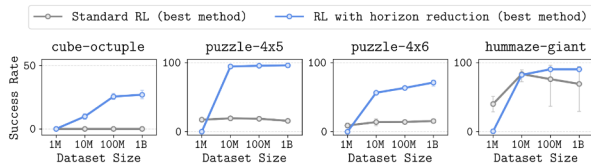
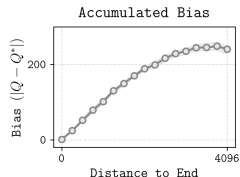
Conditional score

Conditional vector field

Q learning is not yet scalable [Par25]

Long horizon problems are hard.

$$\mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[\left(Q_{\theta}(s, a) - \underbrace{\left(r + \gamma \max_{a'} Q_{\bar{\theta}}(s', a') \right)}_{\text{Biased}} \right)^2 \right].$$



Idea: Apply flow matching to learn the *successor measure* of an MDP.

Def. We define the *successor measure* as

$$m^\pi(X \mid s, a) = (1 - \gamma) \sum_{k=0}^{\infty} \gamma^k \Pr(S_{k+1} \in X \mid S_0 = s, A_0 = a, \pi),$$

Recall. The *successor measure* is the unique fix point to the *Bellman* equation

$$\begin{aligned} m^\pi(\cdot \mid s, a) &= (\mathcal{T}^\pi m^\pi)(\cdot \mid s, a) \\ &:= (1 - \gamma)P(\cdot \mid s, a) + \gamma(P^\pi m^\pi)(\cdot \mid s, a), \end{aligned}$$

where

$$(P^\pi m)(dx \mid s, a) = \int_{s'} P(ds' \mid s, a) m(dx \mid s', \pi(s')).$$

The most straightforward idea is to just substitute m^π for q in flow matching i.e.

$$\mathcal{L}_{\text{MC-CFM}}(\theta) = \mathbb{E}_{\rho, t, Z, X_t} \left[\left\| \tilde{v}_t(X_t \mid S, A; \theta) - u_t(X_t \mid Z) \right\|^2 \right],$$

where $Z = X_1 \sim m^\pi(\cdot \mid S, A)$, $X_t \sim p_t(\cdot \mid Z)$.

Here we just use the optimal transport conditional vector field and corresponding probability density path.

The most straightforward idea is to just substitute m^π for q in flow matching i.e.

$$\mathcal{L}_{\text{MC-CFM}}(\theta) = \mathbb{E}_{\rho, t, Z, X_t} \left[\left\| \tilde{v}_t(X_t \mid S, A; \theta) - u_t(X_t \mid Z) \right\|^2 \right],$$

where $Z = X_1 \sim m^\pi(\cdot \mid S, A)$, $X_t \sim p_t(\cdot \mid Z)$.

Here we just use the optimal transport conditional vector field and corresponding probability density path.

This requires direct access to samples from m^π .

Can we learn from offline one-step transitions (S, A, S') ?

The most straightforward idea is to just substitute m^π for q in flow matching i.e.

$$\mathcal{L}_{\text{MC-CFM}}(\theta) = \mathbb{E}_{\rho, t, Z, X_t} \left[\left\| \tilde{v}_t(X_t \mid S, A; \theta) - u_t(X_t \mid Z) \right\|^2 \right],$$

where $Z = X_1 \sim m^\pi(\cdot \mid S, A), X_t \sim p_t(\cdot \mid Z)$.

Here we just use the optimal transport conditional vector field and corresponding probability density path.

This requires direct access to samples from m^π .

Can we learn from offline one-step transitions (S, A, S') ?

Leverage **recursive** structure of Bellman equation i.e.

$$\begin{aligned} X_0 &\sim p_0 \\ Z = X_1 &\sim (1 - \gamma)\delta_{S'} + \gamma\delta_{\tilde{f}_1(X_0|S', \pi(S'))}. \end{aligned} \quad (\text{TD-CFM})$$

Leverage **recursive** structure of Bellman equation i.e.

$$\begin{aligned} X_0 &\sim p_0 \\ Z = X_1 &\sim (1 - \gamma)\delta_{S'} + \gamma\delta_{\tilde{f}_1(X_0|S',\pi(S'))}. \end{aligned}$$

- ▶ With probability $(1 - \gamma)$, $X_1 = S'$
- ▶ With probability γ , sample from \tilde{m}^π by integrating \tilde{f}_t .

Can we do better?

Case 1: We sample S'

$$\vec{v}_t(x | s, a) = \int \vec{u}_t(x | x_1) \frac{\vec{p}_t(x | x_1) P(dx_1 | s, a)}{\vec{p}_t(x | s, a)},$$

$$\vec{\mathcal{L}}(\theta) = \mathbb{E}_{\rho, t, Z, \vec{X}_t} \left[\left\| \tilde{v}_t(\vec{X}_t | S, A; \theta) - \vec{u}_t(\vec{X}_t | Z) \right\|^2 \right],$$

where $Z = X_1 \sim P(\cdot | S, A)$, $\vec{X}_t \sim \vec{p}_t(\cdot | Z)$

Lemma. Assuming $v_t^{(n+1)} = \arg \min_v \mathcal{L}_{\text{TD}^2\text{-CFM}}(\text{TBD})$, then $v_t^{(n+1)}$ induces a probability path $m_t^{(n+1)}$ such that $m_0^{(n+1)} = m_0$ and $m_1^{(n+1)} = \mathcal{T}^\pi m_1^{(n)}$.

Lemma. Assuming $v_t^{(n+1)} = \arg \min_v \mathcal{L}_{\text{TD}^2-\text{CFM}}(\text{TBD})$, then $v_t^{(n+1)}$ induces a probability path $m_t^{(n+1)}$ such that $m_0^{(n+1)} = m_0$ and $m_1^{(n+1)} = \mathcal{T}^\pi m_1^{(n)}$.

Case 2: We sample future state

$$\hat{v}_t^{(n)}(x | s, a) = \int v_t^{(n)}(x | s', a') \frac{m_t^{(n)}(x | s', a') P(ds' | s, a)}{\hat{p}_t^{(n)}(x | s, a)},$$

where $\hat{p}_t^{(n)}(x | s, a) = \int m_t^{(n)}(x | s', a') P(ds' | s, a)$, and $a' = \pi(s')$.

Lemma. Assuming $\mathbf{v}_t^{(n+1)} = \arg \min_{\mathbf{v}} \mathcal{L}_{\text{TD}^2-\text{CFM}}(\text{TBD})$, then $\mathbf{v}_t^{(n+1)}$ induces a probability path $m_t^{(n+1)}$ such that $m_0^{(n+1)} = m_0$ and $m_1^{(n+1)} = \mathcal{T}^\pi m_1^{(n)}$.

Case 2: We sample future state

$$\hat{\mathbf{v}}_t^{(n)}(x | s, a) = \int \mathbf{v}_t^{(n)}(x | s', a') \frac{m_t^{(n)}(x | s', a') P(ds' | s, a)}{\hat{p}_t^{(n)}(x | s, a)},$$

where $\hat{p}_t^{(n)}(x | s, a) = \int m_t^{(n)}(x | s', a') P(ds' | s, a)$, and $a' = \pi(s')$.

$$\hat{\mathcal{L}}(\theta) = \mathbb{E}_{\rho_{t, \hat{X}_t}} \left[\left\| \tilde{\mathbf{v}}_t(\hat{X}_t | S, A; \theta) - \tilde{\mathbf{v}}_t^{(n)}(\hat{X}_t | S', \pi(S')) \right\|^2 \right],$$

where $X_0 \sim p_0$, $S' \sim P(\cdot | S, A)$, $\hat{X}_t = \tilde{f}_t^{(n)}(X_0 | S', \pi(S'))$,

Combine the **objectives**!!

$$\vec{\mathcal{L}}(\theta) = \mathbb{E}_{\rho, t, Z, \vec{X}_t} \left[\left\| \tilde{\mathbf{v}}_t(\vec{X}_t \mid S, A; \theta) - \vec{u}_t(\vec{X}_t \mid Z) \right\|^2 \right],$$

where $Z = X_1 \sim P(\cdot \mid S, A)$, $\vec{X}_t \sim \vec{p}_t(\cdot \mid Z)$

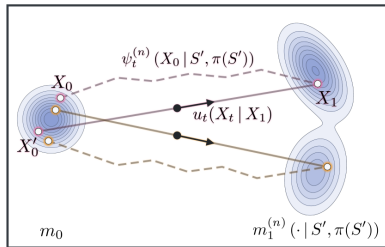
+

$$\hat{\mathcal{L}}(\theta) = \mathbb{E}_{\rho, t, \hat{X}_t} \left[\left\| \tilde{\mathbf{v}}_t(\hat{X}_t \mid S, A; \theta) - \tilde{\mathbf{v}}_t^{(n)}(\hat{X}_t \mid S', \pi(S')) \right\|^2 \right],$$

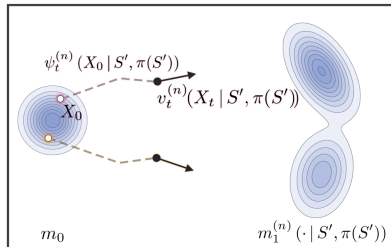
where $X_0 \sim p_0$, $S' \sim P(\cdot \mid S, A)$, $\hat{X}_t = \tilde{f}_t^{(n)}(X_0 \mid S', \pi(S'))$,

$$\mathcal{L}_{\text{TD}^2\text{-CFM}}(\theta) = (1 - \gamma)\vec{\mathcal{L}}(\theta) + \gamma\hat{\mathcal{L}}(\theta) \implies \text{Lower variance gradient.}$$

TD-CFM



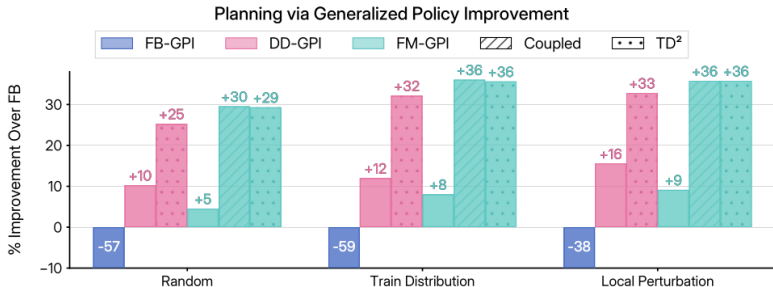
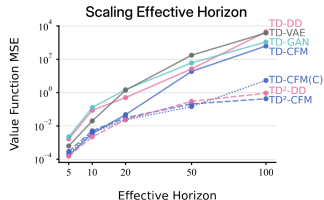
TD²-CFM



GPI:

- ▶ Train π_w with Forward backward.
- ▶ Do GPI as $w_t \in$

$$\arg \max_{w \sim D(W)} \underbrace{(1 - \gamma)^{-1} \mathbb{E}_{X \sim m^{\pi_w}(\cdot | s_t, \pi_w(s_t))} [r(X)]}_{Q^{\pi_w}(s_t, \pi_w(s_t))}.$$
- ▶ Averaged over 128 samples.



References



Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud.
Neural Ordinary Differential Equations.



Jesse Farebrother, Matteo Pirotta, Andrea Tirinzoni, Rémi Munos, Alessandro Lazaric, and Ahmed Touati.
Temporal Difference Flows, March 2025.



Ruiqi Gao, Emiel Hoogeboom, Jonathan Heek, Valentin De Bortoli, Kevin P. Murphy, and Tim Salimans.
Diffusion Meets Flow Matching: Two Sides of the Same Coin, December 2024.



Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le.
Flow Matching for Generative Modeling, February 2023.



Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky T. Q. Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat.
Flow Matching Guide and Code, December 2024.



Seohong Park.
Q-learning is not yet scalable, June 2025.



Danilo Jimenez Rezende and Shakir Mohamed.
Variational Inference with Normalizing Flows.



Chongyi Zheng, Seohong Park, Sergey Levine, and Benjamin Eysenbach.
Intention-Conditioned Flow Occupancy Models, June 2025.