

Transport-Regularized Normalizing Flows¹²³⁴

“Learning by Forgetting”

J. Setpal

April 30, 2025

¹Peyré, Cuturi. [Arxiv 2020]

²Kobyzev, Prince, Brubaker. [IEEE 2021]

³Papamakarios, et. al. [Arxiv 2021]

⁴Lai, et. al. [Journal of Computational Physics 2023]

- ① Dynamic Optimal Transport
- ② Normalizing Flows
- ③ Transport-Regularized Normalizing Flows

- ① Dynamic Optimal Transport
- ② Normalizing Flows
- ③ Transport-Regularized Normalizing Flows

Dynamic Optimal Transport (1/2)

We have explored the following Optimal Transport problem:

$$L_{\mathbf{C}}(a, b) := \min_{\mathbf{P} \in \mathcal{U}(a, b)} \langle \mathbf{C}, \mathbf{P} \rangle_F = \sum_{i, j} \mathbf{C}_{i, j} \mathbf{P}_{i, j} \quad (1)$$

Dynamic Optimal Transport (1/2)

We have explored the following Optimal Transport problem:

$$L_{\mathbf{C}}(a, b) := \min_{\mathbf{P} \in \mathcal{U}(a, b)} \langle \mathbf{C}, \mathbf{P} \rangle_F = \sum_{i, j} \mathbf{C}_{i, j} \mathbf{P}_{i, j} \quad (1)$$

This notion has a couple of properties / constraints:

1. \mathcal{U} represents the set of valid couplings, which encapsulates criteria:
 - a. Mass is conserved.
 - b. Applying the coupling gets us the target measures: $\beta = \mathbf{P} \# \alpha$
2. The optimization problem is *convex*.
3. If \mathbf{C} is a distance in element space, $L_{\mathbf{C}}$ is a distance in measure space.

Dynamic Optimal Transport (1/2)

We have explored the following Optimal Transport problem:

$$L_{\mathbf{C}}(a, b) := \min_{\mathbf{P} \in \mathcal{U}(a, b)} \langle \mathbf{C}, \mathbf{P} \rangle_F = \sum_{i, j} \mathbf{C}_{i, j} \mathbf{P}_{i, j} \quad (1)$$

This notion has a couple of properties / constraints:

1. \mathcal{U} represents the set of valid couplings, which encapsulates criteria:
 - a. Mass is conserved.
 - b. Applying the coupling gets us the target measures: $\beta = \mathbf{P} \# \alpha$
2. The optimization problem is *convex*.
3. If \mathbf{C} is a distance in element space, $L_{\mathbf{C}}$ is a distance in measure space.

Observation: $\mathbf{C} = \|\cdot\|_2^2 \implies L_{\mathbf{C}}$ is **squared geodesic distance**.

Dynamic Optimal Transport (1/2)

We have explored the following Optimal Transport problem:

$$L_{\mathbf{C}}(a, b) := \min_{\mathbf{P} \in \mathcal{U}(a, b)} \langle \mathbf{C}, \mathbf{P} \rangle_F = \sum_{i, j} \mathbf{C}_{i, j} \mathbf{P}_{i, j} \quad (1)$$

This notion has a couple of properties / constraints:

1. \mathcal{U} represents the set of valid couplings, which encapsulates criteria:
 - a. Mass is conserved.
 - b. Applying the coupling gets us the target measures: $\beta = \mathbf{P} \# \alpha$
2. The optimization problem is *convex*.
3. If \mathbf{C} is a distance in element space, $L_{\mathbf{C}}$ is a distance in measure space.

Observation: $\mathbf{C} = \|\cdot\|_2^2 \implies L_{\mathbf{C}}$ is **squared geodesic distance**.

Implication: Solving for Endpoints \rightarrow Interpolatable Transport.

Dynamic Optimal Transport (1/2)

We have explored the following Optimal Transport problem:

$$L_C(a, b) := \min_{\mathbf{P} \in \mathcal{U}(a, b)} \langle \mathbf{C}, \mathbf{P} \rangle_F = \sum_{i, j} \mathbf{C}_{i, j} \mathbf{P}_{i, j} \quad (1)$$

This notion has a couple of properties / constraints:

1. \mathcal{U} represents the set of valid couplings, which encapsulates criteria:
 - a. Mass is conserved.
 - b. Applying the coupling gets us the target measures: $\beta = \mathbf{P} \# \alpha$
2. The optimization problem is *convex*.
3. If \mathbf{C} is a distance in element space, L_C is a distance in measure space.

Observation: $\mathbf{C} = \|\cdot\|_2^2 \implies L_C$ is **squared geodesic distance**.

Implication: Solving for Endpoints \rightarrow Interpolatable Transport.

How? Using fluid dynamics!

Dynamic Optimal Transport (2/2)

The Dynamic⁵ Optimal Transport enables us to borrow fluid dynamics literature, and understand how the measure evolves as *time* progresses:

$$\text{Let } \mathcal{X}, \mathcal{Y} \in \mathbb{R}^d, \quad \mathbf{C}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^2 \quad (2)$$

$$\text{With measures } \alpha_t \quad \text{s.t.} \quad T_{\#} \alpha_0 = \alpha_1 \quad \forall t \in [0, 1] \quad (3)$$

⁵Adjective AND Noun.

Dynamic Optimal Transport (2/2)

The Dynamic⁵ Optimal Transport enables us to borrow fluid dynamics literature, and understand how the measure evolves as *time* progresses:

$$\text{Let } \mathcal{X}, \mathcal{Y} \in \mathbb{R}^d, \quad \mathbf{C}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^2 \quad (2)$$

$$\text{With measures } \alpha_t \quad \text{s.t.} \quad T_{\#} \alpha_0 = \alpha_1 \quad \forall t \in [0, 1] \quad (3)$$

We can describe the path in continuous time using by moving α_t along a vector field v_t .

⁵Adjective AND Noun.

Dynamic Optimal Transport (2/2)

The Dynamic⁵ Optimal Transport enables us to borrow fluid dynamics literature, and understand how the measure evolves as *time* progresses:

$$\text{Let } \mathcal{X}, \mathcal{Y} \in \mathbb{R}^d, \quad \mathbf{C}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^2 \quad (2)$$

$$\text{With measures } \alpha_t \quad \text{s.t.} \quad T_{\#} \alpha_0 = \alpha_1 \quad \forall t \in [0, 1] \quad (3)$$

We can describe the path in continuous time using by moving α_t along a vector field v_t . Infinitesimally, our transport cost can be computed as:

$$\|v_t\|_{L^2(\alpha_t)} = \left(\int_{\mathbb{R}^d} \|v_t(\mathbf{x})\|^2 d\alpha_t(\mathbf{x}) \right)^{1/2} \quad (4)$$

⁵Adjective AND Noun.

Dynamic Optimal Transport (2/2)

The Dynamic⁵ Optimal Transport enables us to borrow fluid dynamics literature, and understand how the measure evolves as *time* progresses:

$$\text{Let } \mathcal{X}, \mathcal{Y} \in \mathbb{R}^d, \quad \mathbf{C}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^2 \quad (2)$$

$$\text{With measures } \alpha_t \quad \text{s.t.} \quad T_{\#} \alpha_0 = \alpha_1 \quad \forall t \in [0, 1] \quad (3)$$

We can describe the path in continuous time using by moving α_t along a vector field v_t . Infinitesimally, our transport cost can be computed as:

$$\|v_t\|_{L^2(\alpha_t)} = \left(\int_{\mathbb{R}^d} \|v_t(\mathbf{x})\|^2 d\alpha_t(\mathbf{x}) \right)^{1/2} \quad (4)$$

Across time t , our net transport cost is:

$$W_2^2(\alpha_0, \alpha_1) = \min_{\alpha_t, v_t} \int_0^1 \int_{\mathbb{R}^d} \|v_t(\mathbf{x})\|^2 d\alpha_t(\mathbf{x}) dt \quad (5)$$

⁵Adjective AND Noun.

Dynamic Optimal Transport (2/2)

The Dynamic⁵ Optimal Transport enables us to borrow fluid dynamics literature, and understand how the measure evolves as *time* progresses:

$$\text{Let } \mathcal{X}, \mathcal{Y} \in \mathbb{R}^d, \quad \mathbf{C}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^2 \quad (2)$$

$$\text{With measures } \alpha_t \text{ s.t. } T_{\#} \alpha_0 = \alpha_1 \quad \forall t \in [0, 1] \quad (3)$$

We can describe the path in continuous time using by moving α_t along a vector field v_t . Infinitesimally, our transport cost can be computed as:

$$\|v_t\|_{L^2(\alpha_t)} = \left(\int_{\mathbb{R}^d} \|v_t(\mathbf{x})\|^2 d\alpha_t(\mathbf{x}) \right)^{1/2} \quad (4)$$

Across time t , our net transport cost is:

$$W_2^2(\alpha_0, \alpha_1) = \min_{\alpha_t, v_t} \int_0^1 \int_{\mathbb{R}^d} \|v_t(\mathbf{x})\|^2 d\alpha_t(\mathbf{x}) dt \quad (5)$$

To satisfy mass conservation, we also enforce the following constraint:

$$\partial_t \alpha_t + \text{div}(\alpha_t v_t) = 0 \quad (6)$$

⁵Adjective AND Noun.

Mean-Field Games

Problem: The mass-conservation constraint, requiring computing $\alpha_t v_t$ makes the problem non-convex.

Mean-Field Games

Problem: The mass-conservation constraint, requiring computing $\alpha_t v_t$ makes the problem non-convex.

Solution: We can reparameterize the problem using Mean-Field Games:

1. MFGs are ∞ -agent games with each agent in spatial domain trying to minimize individual cost. We assume a density function at time t .

$$\min_{\alpha_t, v_t} T(\alpha_0, \alpha_1) + \int_0^1 \int_{\mathbb{R}^d} L(\mathbf{x}, v_t(\mathbf{x}), \alpha_t(\mathbf{x})) dx dt \quad (7)$$

Mean-Field Games

Problem: The mass-conservation constraint, requiring computing $\alpha_t v_t$ makes the problem non-convex.

Solution: We can reparameterize the problem using Mean-Field Games:

1. MFGs are ∞ -agent games with each agent in spatial domain trying to minimize individual cost. We assume a density function at time t .

$$\min_{\alpha_t, v_t} T(\alpha_0, \alpha_1) + \int_0^1 \int_{\mathbb{R}^d} L(\mathbf{x}, v_t(\mathbf{x}), \alpha_t(\mathbf{x})) \, dx \, dt \quad (7)$$

2. We define agent trajectories $F : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$ satisfying:

$$\begin{cases} \partial_t F(\mathbf{x}, t) = v_t(F(\mathbf{x}, t)) & \forall \mathbf{x} \in \mathbb{R}^d, t \in [0, 1] \\ F(\mathbf{x}, 0) = \mathbf{x} \end{cases} \quad (8)$$

Mean-Field Games

Problem: The mass-conservation constraint, requiring computing $\alpha_t v_t$ makes the problem non-convex.

Solution: We can reparameterize the problem using Mean-Field Games:

1. MFGs are ∞ -agent games with each agent in spatial domain trying to minimize individual cost. We assume a density function at time t .

$$\min_{\alpha_t, v_t} T(\alpha_0, \alpha_1) + \int_0^1 \int_{\mathbb{R}^d} L(\mathbf{x}, v_t(\mathbf{x}), \alpha_t(\mathbf{x})) dx dt \quad (7)$$

2. We define agent trajectories $F : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$ satisfying:

$$\begin{cases} \partial_t F(\mathbf{x}, t) = v_t(F(\mathbf{x}, t)) & \forall \mathbf{x} \in \mathbb{R}^d, t \in [0, 1] \\ F(\mathbf{x}, 0) = \mathbf{x} \end{cases} \quad (8)$$

3. Which under the reparameterization is *convex*:

$$\min_{\alpha_t, v_t} \int_0^1 \int_{\mathbb{R}^d} \|v_t(x)\|^2 d\alpha_t(x) dt = \min_F \int_0^1 \int_{\mathbb{R}^d} \|\partial_t F(\mathbf{x}, t)\|^2 \alpha_0(\mathbf{x}) dx dt \quad (9)$$

Mean-Field Games

Problem: The mass-conservation constraint, requiring computing $\alpha_t v_t$ makes the problem non-convex.

Solution: We can reparameterize the problem using Mean-Field Games:

1. MFGs are ∞ -agent games with each agent in spatial domain trying to minimize individual cost. We assume a density function at time t .

$$\min_{\alpha_t, v_t} T(\alpha_0, \alpha_1) + \int_0^1 \int_{\mathbb{R}^d} L(\mathbf{x}, v_t(\mathbf{x}), \alpha_t(\mathbf{x})) dx dt \quad (7)$$

2. We define agent trajectories $F : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$ satisfying:

$$\begin{cases} \partial_t F(\mathbf{x}, t) = v_t(F(\mathbf{x}, t)) & \forall \mathbf{x} \in \mathbb{R}^d, t \in [0, 1] \\ F(\mathbf{x}, 0) = \mathbf{x} \end{cases} \quad (8)$$

3. Which under the reparameterization is *convex*:

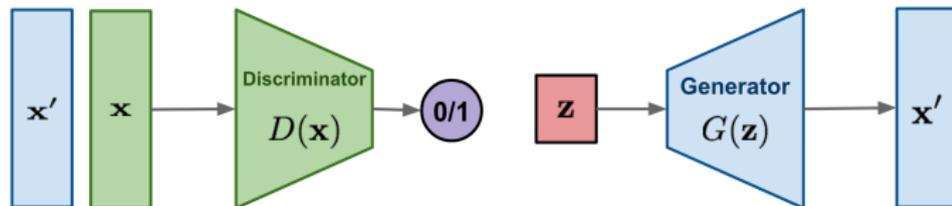
$$\min_{\alpha_t, v_t} \int_0^1 \int_{\mathbb{R}^d} \|v_t(x)\|^2 d\alpha_t(x) dt = \min_F \int_0^1 \int_{\mathbb{R}^d} \|\partial_t F(\mathbf{x}, t)\|^2 \alpha_0(\mathbf{x}) dx dt \quad (9)$$

Next, we need to find a way to learn F . Let's talk normalizing flows.

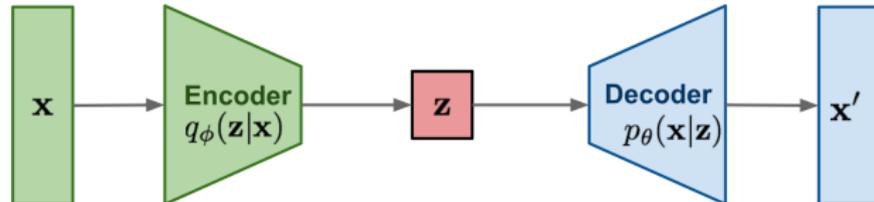
- ① Dynamic Optimal Transport
- ② Normalizing Flows
- ③ Transport-Regularized Normalizing Flows

Generative Modelling Framework

GAN: minimax the classification error loss.



VAE: maximize ELBO.



Flow-based generative models: minimize the negative log-likelihood

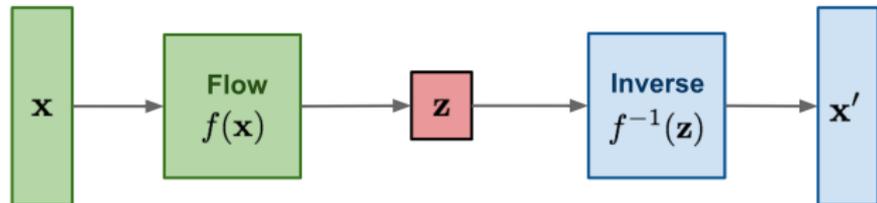


Image Credit: <https://lilianweng.github.io/posts/2018-10-13-flow-models/>

Change of Variables

Since mass vectors are sampled from Σ , α_t are *probability measures*.

Change of Variables

Since mass vectors are sampled from Σ , α_t are *probability measures*.

Setting $\mathbf{y} = T(\mathbf{x})$, we can apply the change of variables formula:

$$\alpha_1(\mathbf{y}) = \alpha_0(T_{\#}^{-1}\alpha_1(\mathbf{y})) |\det \nabla_{\mathbf{y}} T_{\#}^{-1}\alpha_1(\mathbf{y})| = \frac{\alpha_0(T_{\#}^{-1}\alpha_1(\mathbf{y}))}{|\det \nabla_{\mathbf{x}} T_{\#} \circ T_{\#}^{-1}\alpha_1(\mathbf{y})|} \quad (10)$$

To compute the PDF of α_1 .

Change of Variables

Since mass vectors are sampled from Σ , α_t are *probability measures*.

Setting $\mathbf{y} = T(\mathbf{x})$, we can apply the change of variables formula:

$$\alpha_1(\mathbf{y}) = \alpha_0(T_{\#}^{-1}\alpha_1(\mathbf{y}))|\det \nabla_{\mathbf{y}} T_{\#}^{-1}\alpha_1(\mathbf{y})| = \frac{\alpha_0(T_{\#}^{-1}\alpha_1(\mathbf{y}))}{|\det \nabla_{\mathbf{x}} T_{\#} \circ T_{\#}^{-1}\alpha_1(\mathbf{y})|} \quad (10)$$

To compute the PDF of α_1 .

We still need a known α_0 , so we use the multivariate normal $\mathcal{N}(0, \mathbf{I})$:

$$\alpha_0(\mathbf{x}) = \left(\frac{1}{2\pi}\right)^{d/2} \exp\left(-\frac{1}{2}\|\mathbf{x}\|^2\right) \quad (11)$$

Change of Variables

Since mass vectors are sampled from Σ , α_t are *probability measures*.

Setting $\mathbf{y} = T(\mathbf{x})$, we can apply the change of variables formula:

$$\alpha_1(\mathbf{y}) = \alpha_0(T_{\#}^{-1}\alpha_1(\mathbf{y})) |\det \nabla_{\mathbf{y}} T_{\#}^{-1}\alpha_1(\mathbf{y})| = \frac{\alpha_0(T_{\#}^{-1}\alpha_1(\mathbf{y}))}{|\det \nabla_{\mathbf{x}} T_{\#} \circ T_{\#}^{-1}\alpha_1(\mathbf{y})|} \quad (10)$$

To compute the PDF of α_1 .

We still need a known α_0 , so we use the multivariate normal $\mathcal{N}(0, \mathbf{I})$:

$$\alpha_0(\mathbf{x}) = \left(\frac{1}{2\pi}\right)^{d/2} \exp\left(-\frac{1}{2}\|\mathbf{x}\|^2\right) \quad (11)$$

Our goal is to learn the *inverse direction* – a function from target measure α_1 to source measure α_0 . Next, we can discuss approaches to model T .

Invertible Functions

We need three major properties from our choice of model:

1. Must be invertible, so that we can learn the normalizing direction.

Invertible Functions

We need three major properties from our choice of model:

1. Must be invertible, so that we can learn the normalizing direction.
2. Must be expressive enough to model the target distribution.

Invertible Functions

We need three major properties from our choice of model:

1. Must be invertible, so that we can learn the normalizing direction.
2. Must be expressive enough to model the target distribution.
3. Must be efficient to compute.

Invertible Functions

We need three major properties from our choice of model:

1. Must be invertible, so that we can learn the normalizing direction.
2. Must be expressive enough to model the target distribution.
3. Must be efficient to compute.

The class of functions that satisfies these properties are **diffeomorphisms**.

Invertible Functions

We need three major properties from our choice of model:

1. Must be invertible, so that we can learn the normalizing direction.
2. Must be expressive enough to model the target distribution.
3. Must be efficient to compute.

The class of functions that satisfies these properties are **diffeomorphisms**.

Diffeomorphisms are arbitrarily composable. Let $y_\ell := f_\ell$ be diffeomorphisms with inverses g_ℓ for $\ell \in \{1, \dots, L\}$:

$$F := f_L \circ f_{L-1} \circ \dots \circ f_2 \circ f_1 \quad (12)$$

$$G := g_1 \circ g_2 \circ \dots \circ g_{L-1} \circ g_L \quad (13)$$

Then we have $F = G^{-1}$.

Invertible Functions

We need three major properties from our choice of model:

1. Must be invertible, so that we can learn the normalizing direction.
2. Must be expressive enough to model the target distribution.
3. Must be efficient to compute.

The class of functions that satisfies these properties are **diffeomorphisms**.

Diffeomorphisms are arbitrarily composable. Let $y_\ell := f_\ell$ be diffeomorphisms with inverses g_ℓ for $\ell \in \{1, \dots, L\}$:

$$F := f_L \circ f_{L-1} \circ \dots \circ f_2 \circ f_1 \quad (12)$$

$$G := g_1 \circ g_2 \circ \dots \circ g_{L-1} \circ g_L \quad (13)$$

Then we have $F = G^{-1}$. Additionally, we can also compose determinants:

$$\det \nabla_{\mathbf{y}} F = \prod_{\ell=1}^L \det \nabla_{\mathbf{y}_\ell} f_\ell \quad (14)$$

Invertible Functions

We need three major properties from our choice of model:

1. Must be invertible, so that we can learn the normalizing direction.
2. Must be expressive enough to model the target distribution.
3. Must be efficient to compute.

The class of functions that satisfies these properties are **diffeomorphisms**.

Diffeomorphisms are arbitrarily composable. Let $y_\ell := f_\ell$ be diffeomorphisms with inverses g_ℓ for $\ell \in \{1, \dots, L\}$:

$$F := f_L \circ f_{L-1} \circ \dots \circ f_2 \circ f_1 \quad (12)$$

$$G := g_1 \circ g_2 \circ \dots \circ g_{L-1} \circ g_L \quad (13)$$

Then we have $F = G^{-1}$. Additionally, we can also compose determinants:

$$\det \nabla_{\mathbf{y}} F = \prod_{\ell=1}^L \det \nabla_{\mathbf{y}_\ell} f_\ell \quad (14)$$

This is huge for satisfying property 2.

Optimization by Maximum Likelihood

Learning Normalizing Flows allows us to directly maximize log-likelihood:

$$\min_{\theta} D_{\text{KL}}(\alpha_1 || T_{\#} \alpha_0) = -\mathbb{E}_{\mathbf{x} \sim \alpha_1} [\log \alpha_0(G(\mathbf{x})) + \log |\det \nabla_{\mathbf{y}} G|] + C \quad (15)$$

With some constant C , which for purpose of optimization we can ignore.

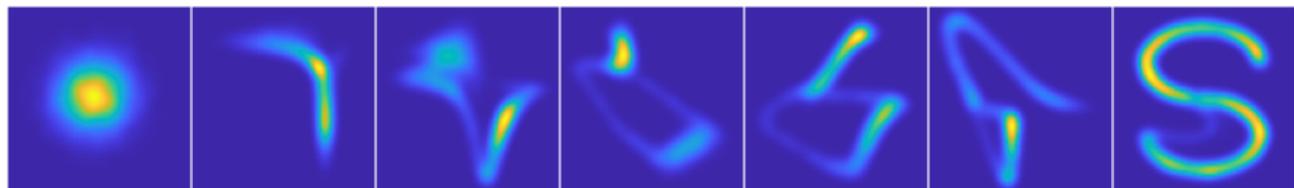
Optimization by Maximum Likelihood

Learning Normalizing Flows allows us to directly maximize log-likelihood:

$$\min_{\theta} D_{\text{KL}}(\alpha_1 || T_{\#} \alpha_0) = -\mathbb{E}_{\mathbf{x} \sim \alpha_1} [\log \alpha_0(G(\mathbf{x})) + \log |\det \nabla_{\mathbf{y}} G|] + C \quad (15)$$

With some constant C , which for purpose of optimization we can ignore.

Caveat: Models trained using D_{KL} are volatile to initialization and interpolate poorly:



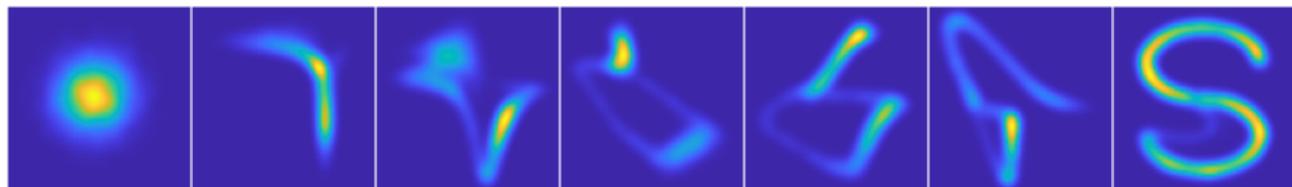
Optimization by Maximum Likelihood

Learning Normalizing Flows allows us to directly maximize log-likelihood:

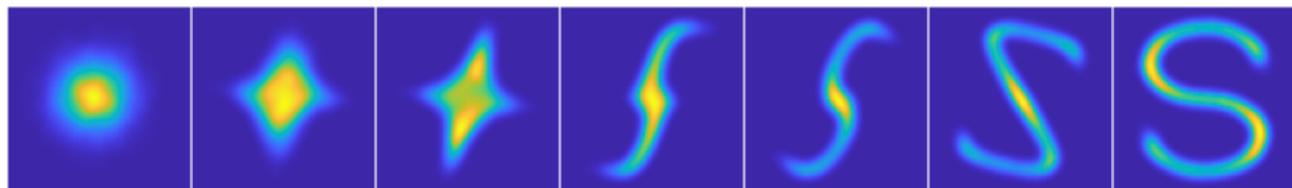
$$\min_{\theta} D_{\text{KL}}(\alpha_1 || T_{\theta} \alpha_0) = -\mathbb{E}_{\mathbf{x} \sim \alpha_1} [\log \alpha_0(G(\mathbf{x})) + \log |\det \nabla_{\mathbf{y}} G|] + C \quad (15)$$

With some constant C , which for purpose of optimization we can ignore.

Caveat: Models trained using D_{KL} are volatile to initialization and interpolate poorly:

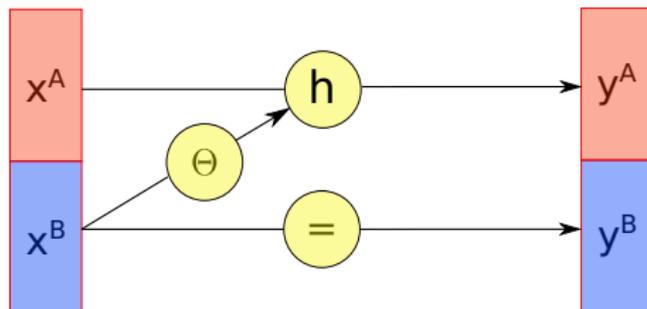


Which we can fix by regularization to transport cost:



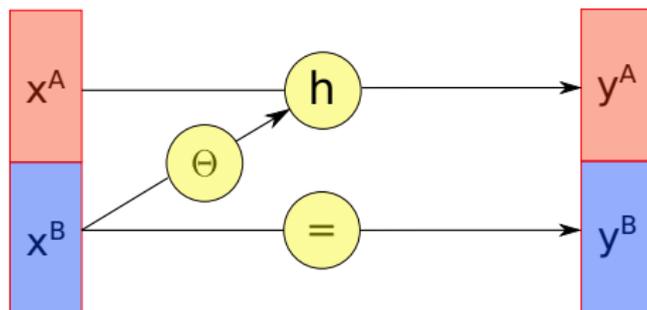
Coupling Flows

One clever diffeomorphism is a **coupling flow**:



Coupling Flows

One clever diffeomorphism is a **coupling flow**:



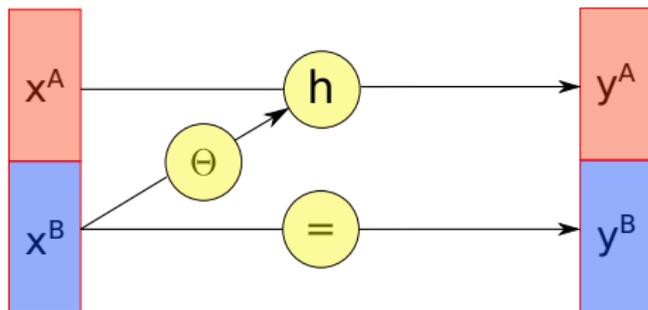
With random permutation Π , we define forward flow ($\alpha_0 \rightarrow \alpha_1$) as:

$$\mathbf{x}' := \Pi(\mathbf{x}), \quad \mathbf{w}, \mathbf{b} := h_{\Theta}(\mathbf{x}'_{D/2+1:D}) \quad (16)$$

$$f_k^{(\Theta_k)}(\mathbf{x}) := \text{Concat}([\mathbf{x}'_{1:D/2} \odot \exp(\mathbf{w}) + \mathbf{b}, \mathbf{x}'_{D/2+1:D}]) \quad (17)$$

Coupling Flows

One clever diffeomorphism is a **coupling flow**:



With random permutation Π , we define forward flow ($\alpha_0 \rightarrow \alpha_1$) as:

$$\mathbf{x}' := \Pi(\mathbf{x}), \quad \mathbf{w}, \mathbf{b} := h_{\Theta}(\mathbf{x}'_{D/2+1:D}) \quad (16)$$

$$f_k^{(\Theta_k)}(\mathbf{x}) := \text{Concat}([\mathbf{x}'_{1:D/2} \odot \exp(\mathbf{w}) + \mathbf{b}, \mathbf{x}'_{D/2+1:D}]) \quad (17)$$

Subsequently inverse flow ($\alpha_1 \rightarrow \alpha_0$) is defined as follows:

$$\mathbf{w}, \mathbf{b} := h_{\Theta}(\mathbf{y}_{D/2+1:D}) \quad (18)$$

$$f_k^{(\Theta_k)^{-1}}(\mathbf{y}) := \Pi^{-1}(\text{Concat}([\mathbf{y}_{1:D/2} - \mathbf{b}] \oslash \exp(\mathbf{w}), \mathbf{y}_{D/2+1:D}])) \quad (19)$$

Computing the log |det(·)|

Best part, the Jacobian of $f_k^{(\Theta_k)^{-1}}$ has the the following block form:

$$\nabla f_k^{(\Theta_k)^{-1}}(\mathbf{y}) = \begin{bmatrix} I_{D/2 \times D/2} & \mathbf{0}_{D/2 \times D/2} \\ \frac{\partial f_{k, D/2+1:D}^{(\Theta_k)^{-1}}}{\partial \mathbf{y}_{1:D/2}} & \text{diag}(\exp(-\mathbf{w})) \end{bmatrix}$$

Computing the log |det(·)|

Best part, the Jacobian of $f_k^{(\Theta_k)^{-1}}$ has the the following block form:

$$\nabla f_k^{(\Theta_k)^{-1}}(\mathbf{y}) = \begin{bmatrix} I_{D/2 \times D/2} & \mathbf{0}_{D/2 \times D/2} \\ \frac{\partial f_{k, D/2+1:D}^{(\Theta_k)^{-1}}}{\partial \mathbf{y}_{1:D/2}} & \text{diag}(\exp(-\mathbf{w})) \end{bmatrix} \quad (20)$$

Which enables us to compute its determinant in linear time:

$$\det \nabla f_k^{(\Theta_k)^{-1}}(\mathbf{y}) = \exp\left(\sum -\mathbf{w}_k\right) \quad (21)$$

Computing the log |det(·)|

Best part, the Jacobian of $f_k^{(\Theta_k)^{-1}}$ has the the following block form:

$$\nabla f_k^{(\Theta_k)^{-1}}(\mathbf{y}) = \begin{bmatrix} I_{D/2 \times D/2} & \mathbf{0}_{D/2 \times D/2} \\ \frac{\partial f_{k, D/2+1:D}^{(\Theta_k)^{-1}}}{\partial \mathbf{y}_{1:D/2}} & \text{diag}(\exp(-\mathbf{w})) \end{bmatrix} \quad (20)$$

Which enables us to compute its determinant in linear time:

$$\det \nabla f_k^{(\Theta_k)^{-1}}(\mathbf{y}) = \exp\left(\sum -\mathbf{w}_k\right) \quad (21)$$

Which we can further use to compute composed log det ∇F^{-1} :

$$\log \det \nabla F^{-1}(\mathbf{y}) = \sum_{\ell=0}^{L-1} \sum_{i=1}^{D/2} -\mathbf{w}_i^{(\ell)} \quad (22)$$

Computing the log |det(·)|

Best part, the Jacobian of $f_k^{(\Theta_k)^{-1}}$ has the the following block form:

$$\nabla f_k^{(\Theta_k)^{-1}}(\mathbf{y}) = \begin{bmatrix} I_{D/2 \times D/2} & \mathbf{0}_{D/2 \times D/2} \\ \frac{\partial f_{k, D/2+1:D}^{(\Theta_k)^{-1}}}{\partial \mathbf{y}_{1:D/2}} & \text{diag}(\exp(-\mathbf{w})) \end{bmatrix} \quad (20)$$

Which enables us to compute its determinant in linear time:

$$\det \nabla f_k^{(\Theta_k)^{-1}}(\mathbf{y}) = \exp\left(\sum -\mathbf{w}_k\right) \quad (21)$$

Which we can further use to compute composed log det ∇F^{-1} :

$$\log \det \nabla F^{-1}(\mathbf{y}) = \sum_{\ell=0}^{L-1} \sum_{i=1}^{D/2} -\mathbf{w}_i^{(\ell)} \quad (22)$$

Implication: h_{Θ} can be arbitrarily complex!

- ① Dynamic Optimal Transport
- ② Normalizing Flows
- ③ Transport-Regularized Normalizing Flows

Setting the Terminal Condition

The final thing left to discuss is the terminal condition. We need F s.t.:

$$F(\mathbf{x}, 1) = \mathbf{y} \sim \alpha_1 \quad \forall \mathbf{x} \in \mathbb{R}^d \quad (23)$$

Setting the Terminal Condition

The final thing left to discuss is the terminal condition. We need F s.t:

$$F(\mathbf{x}, 1) = \mathbf{y} \sim \alpha_1 \quad \forall \mathbf{x} \in \mathbb{R}^d \quad (23)$$

Instead of solving a constrained problem, we soft-constraint using D_{KL} :

$$\min_F \int_0^1 \int_{\mathbb{R}^d} \|\partial_t F(\mathbf{x}, t)\|^2 \alpha_0(\mathbf{x}) \, dx \, dt + \lambda D_{\text{KL}}(\alpha_1 \| F_{\#} \alpha_0) \quad (24)$$

D_{KL} is conventional cost; we just add transport cost, hence new formulation is called **transport-regularized normalizing flows**.

Setting the Terminal Condition

The final thing left to discuss is the terminal condition. We need F s.t:

$$F(\mathbf{x}, 1) = \mathbf{y} \sim \alpha_1 \quad \forall \mathbf{x} \in \mathbb{R}^d \quad (23)$$

Instead of solving a constrained problem, we soft-constraint using D_{KL} :

$$\min_F \int_0^1 \int_{\mathbb{R}^d} \|\partial_t F(\mathbf{x}, t)\|^2 \alpha_0(\mathbf{x}) \, dx \, dt + \lambda D_{\text{KL}}(\alpha_1 \| F_{\#} \alpha_0) \quad (24)$$

D_{KL} is conventional cost; we just add transport cost, hence new formulation is called **transport-regularized normalizing flows**.

Discretizing across a composition of L flows, we have:

$$\min_F L \cdot \mathbb{E}_{\mathbf{x} \sim \alpha_0} \left[\sum_{\ell=0}^{L-1} \|F_{\ell+1}(\mathbf{x}) - F_{\ell}(\mathbf{x})\|_2^2 \right] + \lambda D_{\text{KL}}(\alpha_1 \| F_{\#} \alpha_0) \quad (25)$$

Using this, we can train transport-regularized normalizing flows.

Implementation Details

If you can view this screen, I am making a mistake.

Thank you!

Have an awesome rest of your day!

Slides: <https://jinen.setpal.net/slides/nf.pdf>