

# **CPSC 661: Sampling Algorithms in ML**

---

Andre Wibisono

March 29, 2021

Yale University

## Last time

- Wasserstein distance
- Optimal transport
- Displacement interpolation as geodesic

**Today:** Otto calculus

# References

- Villani, *Topics in Optimal Transportation*, Springer, 2003
- Villani, *Optimal Transport: Old and New*, Springer, 2008
- Ambrosio, Gigli & Savaré, *Gradient Flows in Metric Spaces and in the Space of Probability Measures*, Springer, 2005

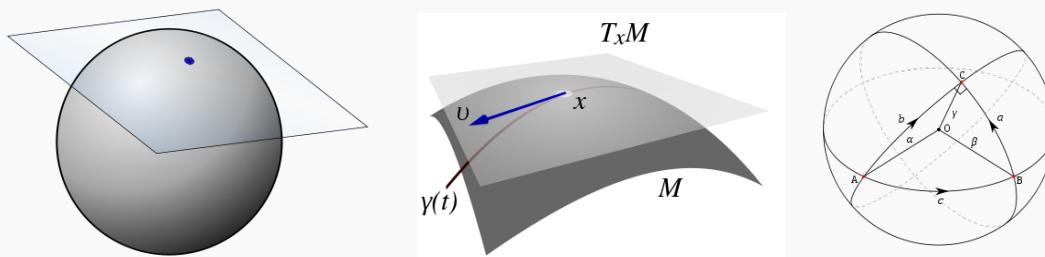
# Motivation

---

# Review: Distance on manifold

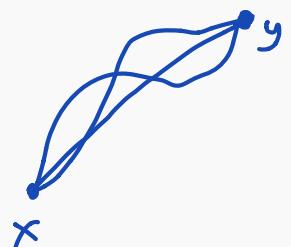
Given a manifold  $\mathcal{X}$  with metric  $g_x: T_x \mathcal{X} \times T_x \mathcal{X} \rightarrow \mathbb{R}$

$$\|v\|_x^2 = \langle v, v \rangle_x = g_x(v, v)$$



**Distance** between two points  $x, y \in \mathcal{X}$  is

$$d(x, y)^2 = \inf_{\substack{(X_t) \\ X_0=x, X_1=y}} \int_0^1 \|\dot{X}_t\|_{X_t}^2 dt$$



Shortest curve is **geodesic**:  $d(x, X_t) = t d(x, y)$

# Wasserstein distance

$$\mathcal{P}(\mathbb{R}^n) = \{ \text{probability distributions } \rho \text{ on } \mathbb{R}^n \text{ with } \mathbb{E}_\rho[\|X\|^2] < \infty \}$$

**Wasserstein distance:**

$$W_2(\rho, \nu)^2 = \inf_{\pi \in \Pi(\rho, \nu)} \mathbb{E}_\pi[\|X - Y\|^2]$$

where  $\Pi(\rho, \nu)$  is set of couplings  $(X, Y) \sim \pi$  such that  $X \sim \rho, Y \sim \nu$

# Wasserstein distance

$$\mathcal{P}(\mathbb{R}^n) = \{ \text{probability distributions } \rho \text{ on } \mathbb{R}^n \text{ with } \mathbb{E}_\rho[\|X\|^2] < \infty \}$$

**Wasserstein distance:**

$$W_2(\rho, \nu)^2 = \inf_{\pi \in \Pi(\rho, \nu)} \mathbb{E}_\pi[\|X - Y\|^2]$$

where  $\Pi(\rho, \nu)$  is set of couplings  $(X, Y) \sim \pi$  such that  $X \sim \rho$ ,  $Y \sim \nu$

Brenier: There is unique convex  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}$  s.t.  $(\nabla \phi)_\# \rho = \nu$  and

$$W_2(\rho, \nu)^2 = \mathbb{E}_\rho[\|X - \nabla \phi(X)\|^2]$$

# Wasserstein distance

$$\mathcal{P}(\mathbb{R}^n) = \{ \text{probability distributions } \rho \text{ on } \mathbb{R}^n \text{ with } \mathbb{E}_\rho[\|X\|^2] < \infty \}$$

**Wasserstein distance:**

$$W_2(\rho, \nu)^2 = \inf_{\pi \in \Pi(\rho, \nu)} \mathbb{E}_\pi[\|X - Y\|^2]$$

where  $\Pi(\rho, \nu)$  is set of couplings  $(X, Y) \sim \pi$  such that  $X \sim \rho$ ,  $Y \sim \nu$

Brenier: There is unique convex  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}$  s.t.  $(\nabla \phi)_\# \rho = \nu$  and

$$W_2(\rho, \nu)^2 = \mathbb{E}_\rho[\|X - \nabla \phi(X)\|^2]$$

McCann: Geodesic is displacement interpolation  $\rho_t = (T_t)_\# \rho$  induced by linear interpolation  $T_t(x) = (1 - t)x + t \nabla \phi(x)$

# From distance to metric

Wasserstein  $W_2$  distance:

$$W_2(\rho, \nu)^2 = \inf_{\pi \in \Pi(\rho, \nu)} \mathbb{E}_\pi[\|X - Y\|^2]$$

Question: What metric  $\|\phi\|_\rho^2$  for  $\phi \in T_\rho \mathcal{P}$  generates  $W_2$  distance?

$$W_2(\rho, \nu)^2 = \inf_{\substack{(\rho_t) \\ \rho_0=\rho, \rho_1=\nu}} \int_0^1 \|\dot{\rho}_t\|_{\rho_t}^2 dt$$

# Wasserstein metric

Tangent vector  $\phi \in T_\rho \mathcal{P}$  is  $\phi = -\nabla \cdot (\rho \nabla u)$  for some  $u: \mathbb{R}^n \rightarrow \mathbb{R}$

**Wasserstein metric:**

$$\|\phi\|_\rho^2 = \int_{\mathbb{R}^n} \|\nabla u(x)\|^2 \rho(x) dx = \mathbb{E}_\rho[\|\nabla u(X)\|^2]$$

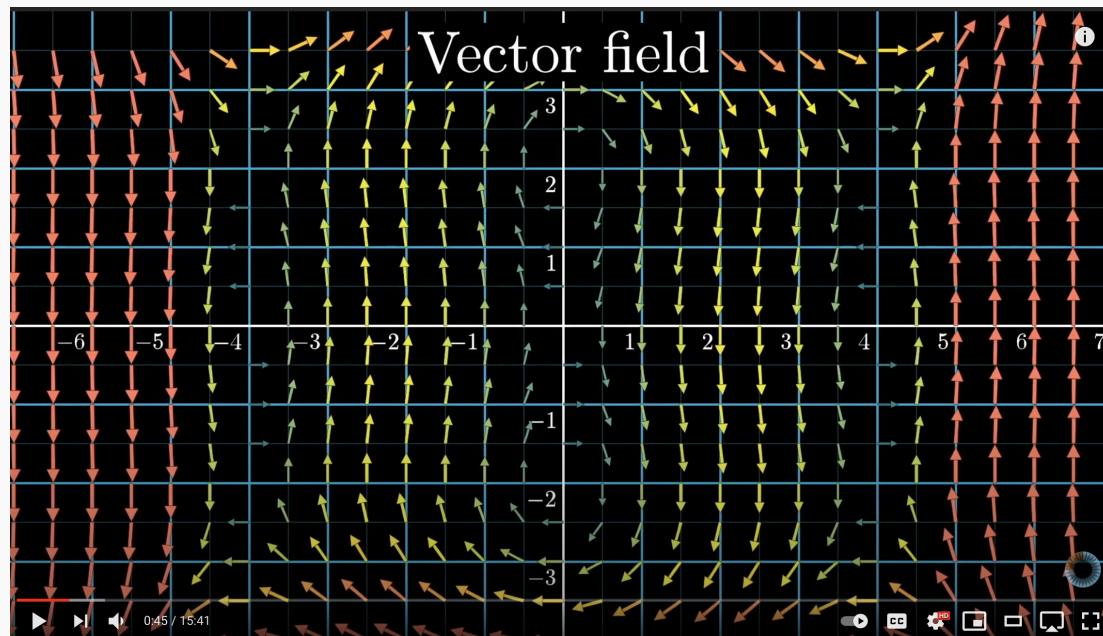
**Theorem:** Wasserstein metric generates  $W_2$  distance

$$W_2(\rho, \nu)^2 = \inf_{\substack{(\rho_t) \\ \rho_0 = \rho, \rho_1 = \nu}} \int_0^1 \|\dot{\rho}_t\|_{\rho_t}^2 dt$$

# Divergence

---

# Divergence



3Blue1Brown, *Divergence and curl: The language of Maxwell's equations, fluid flow, and more*, <https://www.youtube.com/watch?v=rB83DpBJQsE>, 2018

# Divergence

Given a vector field  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$v(x) = (v_1(x), \dots, v_n(x))$$

The **divergence**  $\nabla \cdot v: \mathbb{R}^n \rightarrow \mathbb{R}$  is

$$(\nabla \cdot v)(x) = \sum_{i=1}^n \frac{\partial v_i(x)}{\partial x_i}$$

# Divergence

Given a vector field  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$v(x) = (v_1(x), \dots, v_n(x))$$

The **divergence**  $\nabla \cdot v: \mathbb{R}^n \rightarrow \mathbb{R}$  is

$$(\nabla \cdot v)(x) = \sum_{i=1}^n \frac{\partial v_i(x)}{\partial x_i}$$

- Measures net change in flow defined by  $v$
- $\nabla \cdot v = \langle \nabla, v \rangle$  where  $\nabla = (\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_n})$
- $\nabla \cdot v = \text{Tr}(\nabla v)$  where  $\nabla v(x)$  is Jacobian matrix

$$\left( \frac{\partial v_i(x)}{\partial x_j} \right)_{i,j=1}^n$$

## Example: Laplacian

Suppose  $v$  is the **gradient** of  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$v(x) = \nabla f(x) = \left( \frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_n} \right)$$

The Jacobian is the **Hessian** of  $f$

$$\nabla(\nabla f)(x) = \nabla^2 f(x) = \left( \frac{\partial^2 f(x)}{\partial x_i \partial x_j} \right)_{i,j=1}^n$$

The divergence is the **Laplacian** of  $f$

$$(\nabla \cdot \nabla f)(x) = \Delta f(x) = \text{Tr}(\nabla^2 f(x)) = \sum_{i=1}^n \frac{\partial^2 f(x)}{\partial x_i^2}$$

# Integration by parts

For any  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $h: \mathbb{R}^n \rightarrow \mathbb{R}$  (smooth and compact support):

$$\int_{\mathbb{R}^n} h(x) (\nabla \cdot v)(x) dx = - \int_{\mathbb{R}^n} \langle \nabla h(x), v(x) \rangle dx$$


# Integration by parts

For any  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $h: \mathbb{R}^n \rightarrow \mathbb{R}$  (smooth and compact support):

$$\int_{\mathbb{R}^n} h(x) (\nabla \cdot v)(x) dx = - \int_{\mathbb{R}^n} \langle \nabla h(x), v(x) \rangle dx$$

- Divergence is *adjoint* of gradient:

$$\langle h, \nabla \cdot v \rangle_{L^2(\mathbb{R}^n \rightarrow \mathbb{R})} = \langle -\nabla h, v \rangle_{L^2(\mathbb{R}^n \rightarrow \mathbb{R}^n)}$$


# Integration by parts

For any  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $h: \mathbb{R}^n \rightarrow \mathbb{R}$  (smooth and compact support):

$$\int_{\mathbb{R}^n} h(x) (\nabla \cdot v)(x) dx = - \int_{\mathbb{R}^n} \langle \nabla h(x), v(x) \rangle dx$$

- Divergence is *adjoint* of gradient:

$$\langle h, \nabla \cdot v \rangle_{L^2(\mathbb{R}^n \rightarrow \mathbb{R})} = \langle -\nabla h, v \rangle_{L^2(\mathbb{R}^n \rightarrow \mathbb{R}^n)}$$

- Divergence has 0 integral:

$$\int_{\mathbb{R}^n} \mathbf{1} \cdot (\nabla \cdot v)(x) dx = - \int_{\mathbb{R}^n} \langle \nabla \mathbf{1}, v(x) \rangle dx = 0$$

# Divergence-free vector field

A vector field  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is **divergence-free** if

$$\nabla \cdot v = 0$$

$$(\nabla \cdot v)(x) = 0 \quad \forall x \in \mathbb{R}^n$$

- Incompressible flow
- Orthogonal to gradient: if  $\nabla \cdot v = 0$

then  $\int_{\mathbb{R}^n} \langle \nabla h(x), v(x) \rangle dx = - \int_{\mathbb{R}^n} h(x) (\underbrace{\nabla \cdot v}(x) dx = 0$

# Helmholtz decomposition

**Theorem:** (Helmholtz decomposition)

Any vector field  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$  can be written uniquely as the sum of a gradient and a divergence-free vector field:

$$v = \nabla u + w$$

for some  $u: \mathbb{R}^n \rightarrow \mathbb{R}$  and  $w: \mathbb{R}^n \rightarrow \mathbb{R}^n$  with  $\nabla \cdot w = 0$

# Helmholtz decomposition

**Theorem:** (Helmholtz decomposition)

Any vector field  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$  can be written uniquely as the sum of a gradient and a divergence-free vector field:

$$v = \nabla u + w$$

for some  $u: \mathbb{R}^n \rightarrow \mathbb{R}$  and  $w: \mathbb{R}^n \rightarrow \mathbb{R}^n$  with  $\nabla \cdot w = 0$

$$v(x) = A x$$

- Linear vector field: Any matrix  $A \in \mathbb{R}^{n \times n}$  can be written uniquely as

$$A = A_{\text{sym}} + A_{\text{ant}}$$

where  $A_{\text{sym}} = \frac{1}{2}(A + A^\top)$  is symmetric and  $A_{\text{ant}} = \frac{1}{2}(A - A^\top)$  is anti-symmetric

# Continuity equation

---

# Continuity equation

$$v(x) = (v_1(x), \dots, v_n(x))$$

Suppose  $X_t \in \mathbb{R}^n$  follows dynamics of a vector field  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$x = (x_1, \dots, x_n) \in \mathbb{R}^n$$

ODE:  $\dot{X}_t = v(X_t)$

$$\forall i: \frac{d}{dt} X_{t,i} = v_i(X_t)$$

If  $X_0 \sim \rho_0$  is random, then  $X_t \sim \rho_t$  is random

The density  $\rho_t \in \mathcal{P}(\mathbb{R}^n)$  follows the **continuity equation**

$$s_t: \mathbb{R}^n \rightarrow \mathbb{R}$$

PDE:  $\frac{\partial \rho_t}{\partial t}(x) = -\nabla \cdot (\rho_t v)(x) \Leftrightarrow \underbrace{\frac{\partial s_t}{\partial t}}_{\frac{d}{dt} s_t(x)} = -\nabla \cdot (s_t v)$

Note: Also for time-dependent vector field  $v_t(x)$

# Continuity equation

$$\frac{\partial \rho_t}{\partial t} : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$s_t v : \mathbb{R}^n \rightarrow \mathbb{R}^n$$

**Lemma:** If  $\dot{X}_t = v(X_t)$ , then  $\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v)$

$$x \mapsto s_t(x) v(x)$$

$$\nabla \cdot (s_t v) : \mathbb{R}^n \rightarrow \mathbb{R}$$

Proof: For  $h : \mathbb{R}^n \rightarrow \mathbb{R}$ , compute time derivative of

$$\mathbb{E}[h(X_t)] = \int_{\mathbb{R}^n} s_t(x) h(x) dx$$

① RHS:  $\frac{d}{dt} \int_{\mathbb{R}^n} s_t(x) h(x) dx = \int_{\mathbb{R}^n} \frac{\partial s_t(x)}{\partial t} h(x) dx$

② LHS: since  $\dot{X}_t = v(X_t)$

$$\Rightarrow \frac{d}{dt} h(X_t) = \langle \nabla h(X_t), \dot{X}_t \rangle = \langle \nabla h(X_t), v(X_t) \rangle$$

Since  $\mathbb{E}$  is linear:

$$\begin{aligned}
 \frac{d}{dt} \mathbb{E}[h(x_t)] &= \mathbb{E}[\langle \nabla h(x_t), v(x_t) \rangle] \\
 &= \int_{\mathbb{R}^n} \langle \nabla h(x), v(x) \rangle \varrho_t(x) dx \\
 \text{Integration by parts:} \quad &= - \int_{\mathbb{R}^n} h(x) \nabla \cdot (\varrho_t v)(x) dx
 \end{aligned}$$

$$\begin{aligned}
 \therefore \textcircled{1} = \textcircled{2} \Rightarrow \int_{\mathbb{R}^n} \frac{\partial \varrho_t}{\partial t}(x) h(x) &= - \int_{\mathbb{R}^n} \nabla \cdot (\varrho_t v)(x) h(x) dx \\
 \Leftrightarrow \int_{\mathbb{R}^n} \left( \frac{\partial \varrho_t}{\partial t}(x) + \nabla \cdot (\varrho_t v)(x) \right) h(x) &= 0
 \end{aligned}$$

true for arbitrary  $h: \mathbb{R}^n \rightarrow \mathbb{R}$

$$\Rightarrow \frac{\partial \varrho_t}{\partial t}(x) + \nabla \cdot (\varrho_t v)(x) = 0$$

which is continuity eq.  $\square$

# Continuity equation

**Lemma:** If  $\dot{X}_t = v(X_t)$ , then  $\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v)$

*Proof:* For test function  $h: \mathbb{R}^n \rightarrow \mathbb{R}$ , compute time derivative of

$$\mathbb{E}[h(X_t)] = \int_{\mathbb{R}^n} \rho_t(x) h(x) dx$$

1. For RHS:

$$\frac{d}{dt} \int_{\mathbb{R}^n} \rho_t(x) h(x) dx = \int_{\mathbb{R}^n} \frac{\partial \rho_t}{\partial t}(x) h(x) dx$$

2. For LHS, since  $\dot{X}_t = v(X_t)$

$$\begin{aligned}\frac{d}{dt} h(X_t) &= \langle \nabla h(X_t), \dot{X}_t \rangle \\ &= \langle \nabla h(X_t), v(X_t) \rangle\end{aligned}$$

Since  $\mathbb{E}$  is linear and using integration by parts

$$\begin{aligned}\frac{d}{dt} \mathbb{E}[h(X_t)] &= \mathbb{E}[\langle \nabla h(X_t), v(X_t) \rangle] \\ &= \int_{\mathbb{R}^n} \langle \nabla h(x), v(x) \rangle \rho_t(x) dx \\ &= - \int_{\mathbb{R}^n} h(x) (\nabla \cdot (\rho_t v))(x) dx\end{aligned}$$

Both are equal:  $\frac{d}{dt} \mathbb{E}[h(X_t)] = \frac{d}{dt} \int_{\mathbb{R}^n} \rho_t(x) h(x) dx$ . Therefore

$$\int_{\mathbb{R}^n} \frac{\partial \rho_t}{\partial t}(x) h(x) dx = - \int_{\mathbb{R}^n} h(x) (\nabla \cdot (\rho_t v))(x) dx$$

Since  $h$  is arbitrary, this implies

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

which is the continuity equation. □

## Conservation of mass

Continuity equation:  $\dot{X}_t = v(X_t)$   $\Rightarrow \frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v)$

implies **conservation of mass**:

$$\int_{\mathbb{R}^n} \rho_t(x) dx = \int_{\mathbb{R}^n} \rho_0(x) dx$$

# Conservation of mass

Continuity equation:  $\dot{X}_t = v(X_t)$   $\Rightarrow \frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v)$

implies **conservation of mass**:

$$\int_{\mathbb{R}^n} \rho_t(x) dx = \int_{\mathbb{R}^n} \rho_0(x) dx$$

Proof:

$$\frac{d}{dt} \int_{\mathbb{R}^n} \rho_t(x) dx = \int_{\mathbb{R}^n} \frac{\partial \rho_t(x)}{\partial t} dx = - \int_{\mathbb{R}^n} \nabla \cdot (\rho_t v)(x) dx = 0$$

## Example: Gradient flow

If  $X_t$  follows **gradient flow**:

$$\dot{X}_t = -\nabla f(X_t) \quad v = -\nabla f$$

then its density  $\rho_t$  follows the continuity equation

$$\frac{\partial \rho_t}{\partial t} = \nabla \cdot (\rho_t \nabla f) = -\nabla \cdot (\mathcal{S}_t v)$$

## **Time-dependent formulation of $W_2$**

---

## Kinetic energy

Suppose  $X_t \in \mathbb{R}^n$  follows dynamics generated by  $v_t: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$\dot{X}_t = v_t(X_t)$$

Let  $X_t \sim \rho_t$ , so  $\rho_t: \mathbb{R}^n \rightarrow \mathbb{R}$  follows the continuity equation

$$\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v_t)$$

## Kinetic energy

Suppose  $X_t \in \mathbb{R}^n$  follows dynamics generated by  $v_t: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$\dot{X}_t = v_t(X_t)$$

Let  $X_t \sim \rho_t$ , so  $\rho_t: \mathbb{R}^n \rightarrow \mathbb{R}$  follows the continuity equation

$$\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v_t)$$

At time  $t \geq 0$ , define **kinetic energy**:

$$E(t) = \int_{\mathbb{R}^n} \rho_t(x) \|v_t(x)\|_2^2 dx = \mathbb{E}_{g_t} [\|v_t\|^2]$$

# Action

**Action** of dynamics  $\rho = (\rho_t)_{0 \leq t \leq 1}$  generated by  $v = (v_t)_{0 \leq t \leq 1}$

$$\mathcal{A}(\rho, v) = \int_0^1 E(t) dt = \int_0^1 \int_{\mathbb{R}^n} \rho_t(x) \|v_t(x)\|_2^2 dx dt$$

$$x_0 \sim s_0$$

$$\dot{x}_t = v_t(x_t) \Leftrightarrow \frac{\partial x_t}{\partial t} = -\nabla \cdot (s_t v_t)$$

## Time-dependent formulation of $W_2$

$$W_2(s_0, s_1)^2 = \inf_{\pi \in \Pi(s_0, s_1)} \mathbb{E}[\|X - Y\|^2]$$

**Theorem:** (Benamou-Brenier formula)

$$W_2(\rho_0, \rho_1)^2 = \inf_{(\rho, v) \in V(\rho_0, \rho_1)} \mathcal{A}(\rho, v)$$

where  $V(\rho_0, \rho_1) = \{(\rho_t, v_t)_{0 \leq t \leq 1} \text{ from } \rho_0 \text{ to } \rho_1 \text{ via } \frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v_t)\}$

- See [Villani 2003, Theorem 8.1]

# Wasserstein metric

---

## Tangent space

Let  $\mathcal{P}(\mathbb{R}^n) = \left\{ \rho: \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0} \mid \int_{\mathbb{R}^n} \rho(x) dx = 1, \mathbb{E}_\rho[\|X\|^2] < \infty \right\}$

- Tangent space  $T_\rho \mathcal{P}$  contains  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\int_{\mathbb{R}^n} \phi(x) dx = 0$

## Tangent space

Let  $\mathcal{P}(\mathbb{R}^n) = \left\{ \rho: \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0} \mid \int_{\mathbb{R}^n} \rho(x) dx = 1, \mathbb{E}_\rho[\|X\|^2] < \infty \right\}$

- Tangent space  $T_\rho \mathcal{P}$  contains  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\int_{\mathbb{R}^n} \phi(x) dx = 0$
- Want  $\phi \in T_\rho \mathcal{P}$  to represent velocity  $\frac{\partial \rho_t}{\partial t}$  at  $\rho_0 = \rho$  from motion  $\dot{X}_t = v_t(X_t)$ , so  $\rho_t$  follows continuity equation  $\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v_t)$
- So  $\phi \in T_\rho \mathcal{P}$  is of the form  $\phi = -\nabla \cdot (\rho v)$  for some  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$

# Tangent space

Let  $\mathcal{P}(\mathbb{R}^n) = \left\{ \rho: \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0} \mid \int_{\mathbb{R}^n} \rho(x) dx = 1, \mathbb{E}_\rho[\|X\|^2] < \infty \right\}$

- Tangent space  $T_\rho \mathcal{P}$  contains  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\int_{\mathbb{R}^n} \phi(x) dx = 0$
- Want  $\phi \in T_\rho \mathcal{P}$  to represent velocity  $\frac{\partial \rho_t}{\partial t}$  at  $\rho_0 = \rho$  from motion  $\dot{X}_t = v_t(X_t)$ , so  $\rho_t$  follows continuity equation  $\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t v_t)$
- So  $\phi \in T_\rho \mathcal{P}$  is of the form  $\phi = -\nabla \cdot (\rho v)$  for some  $v: \mathbb{R}^n \rightarrow \mathbb{R}^n$
- Not unique: if  $v$  is a solution, then so is  $v + \frac{w}{\rho}$  for any  $\nabla \cdot w = 0$
- Which one? Choose  $v$  that minimizes *kinetic energy*  $\mathbb{E}_\rho[\|v\|^2]$

if  $\nabla \cdot (\rho v) = -\phi$  and  $\nabla \cdot w = 0$

then  $\nabla \cdot (\rho(v + \frac{w}{\rho})) = \nabla \cdot (\rho v) + \cancel{\nabla \cdot w} \stackrel{\textcircled{D}}{\approx} \nabla \cdot (\rho v) = -\phi$

# Wasserstein metric

Let  $\phi \in T_\rho \mathcal{P}$ , so  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\int_{\mathbb{R}^n} \phi(x) dx = 0$

**Wasserstein metric:**

$$\|\phi\|_\rho^2 = \inf_{\substack{v: \mathbb{R}^n \rightarrow \mathbb{R}^n \\ \nabla \cdot (\rho v) = -\phi}} \int_{\mathbb{R}^n} \rho(x) \|v(x)\|_2^2 dx = \mathbb{E}_\rho [\|v\|^2]$$

## Minimizer is gradient

**Lemma:** Minimizer is a *gradient*  $v^* = \nabla u$  for some  $u: \mathbb{R}^n \rightarrow \mathbb{R}$

# Minimizer is gradient

**Lemma:** Minimizer is a *gradient*  $v^* = \nabla u$  for some  $u: \mathbb{R}^n \rightarrow \mathbb{R}$

*Proof:* Let  $v^*$  be minimizer. For any  $w$  with  $\nabla \cdot w = 0$  and  $\eta \neq 0$ :

$$\int_{\mathbb{R}^n} \rho(x) \|v^*(x)\|_2^2 dx \leq \int_{\mathbb{R}^n} \rho(x) \left\| v^*(x) + \eta \frac{w(x)}{\rho(x)} \right\|_2^2 dx$$

Expand and simplify, divide by  $\eta$ , and let  $\eta \rightarrow 0$ :

$$\int_{\mathbb{R}^n} \langle v^*(x), w(x) \rangle dx = 0$$

Since  $w$  is arbitrary with  $\nabla \cdot w = 0$ , this implies  $v^* = \nabla u$  for some  $u: \mathbb{R}^n \rightarrow \mathbb{R}$ . *by Helmholtz decomposition.* □

## Wasserstein metric

$$\phi: \mathbb{R}^n \rightarrow \mathbb{R} \text{ with } \int_{\mathbb{R}^n} \phi(x) dx = 0$$

Tangent function  $\phi \in T_\rho \mathcal{P}$  is  $\phi = -\nabla \cdot (\rho \nabla u)$  for some  $u: \mathbb{R}^n \rightarrow \mathbb{R}$

### Wasserstein metric:

$$\|\phi\|_\rho^2 = \mathbb{E}_\rho[\|\nabla u\|^2] = \int_{\mathbb{R}^n} \rho(x) \|\nabla u(x)\|_2^2 dx$$

# Wasserstein metric

Tangent function  $\phi \in T_\rho \mathcal{P}$  is  $\phi = -\nabla \cdot (\rho \nabla u)$  for some  $u: \mathbb{R}^n \rightarrow \mathbb{R}$

**Wasserstein metric:**

$$\|\phi\|_\rho^2 = \mathbb{E}_\rho[\|\nabla u\|^2] = \int_{\mathbb{R}^n} \rho(x) \|\nabla u(x)\|_2^2 dx$$

Benamou-Brenier formula  $\Rightarrow$  metric generates  $W_2$  distance

$$W_2(\rho, \nu)^2 = \inf_{\substack{(\rho_t) \\ \rho_0=\rho, \rho_1=\nu}} \int_0^1 \|\dot{\rho}_t\|_{\rho_t}^2 dt$$

## Inner product

For  $\phi_1, \phi_2 \in T_\rho \mathcal{P}$  with

$$\phi_1 = -\nabla \cdot (\rho \nabla u_1)$$

$$\phi_2 = -\nabla \cdot (\rho \nabla u_2)$$

**Wasserstein inner product:**

$$\begin{aligned} g_\rho(\phi_1, \phi_2) &= \langle \phi_1, \phi_2 \rangle_\rho = \mathbb{E}_\rho[\langle \nabla u_1, \nabla u_2 \rangle] \\ &= \int_{\mathbb{R}^n} \rho(x) \langle \nabla u_1(x), \nabla u_2(x) \rangle dx \end{aligned}$$

# Optimization on $\mathcal{P}(\mathbb{R}^n)$

---

# Optimization on $\mathcal{P}(\mathbb{R}^n)$

Given functional  $F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$ , want to optimize

$$\min_{\rho \in \mathcal{P}(\mathbb{R}^n)} F(\rho)$$

# Optimization on $\mathcal{P}(\mathbb{R}^n)$

Given functional  $F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$ , want to optimize

$$\min_{\rho \in \mathcal{P}(\mathbb{R}^n)} F(\rho)$$

- Want to run **gradient flow**:

$$\dot{\rho}_t = -\text{grad } F(\rho_t)$$

or **gradient descent**:

$$\rho_{k+1} = \text{Exp}_{\rho_k}(-\eta \text{grad } F(\rho_k))$$

or proximal method, accelerated method, ...

# Optimization on $\mathcal{P}(\mathbb{R}^n)$

Given functional  $F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$ , want to optimize

$$\min_{\rho \in \mathcal{P}(\mathbb{R}^n)} F(\rho)$$

- Want to run **gradient flow**:

$$\dot{\rho}_t = -\text{grad } F(\rho_t)$$

or **gradient descent**:

$$\rho_{k+1} = \text{Exp}_{\rho_k}(-\eta \text{grad } F(\rho_k))$$

or proximal method, accelerated method, ...

- Want to understand convexity, smoothness, gradient domination, ...

# Examples of Functionals

$$F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$$

## 1. Potential energy:

$$F(\rho) = \int_{\mathbb{R}^n} \rho(x) f(x) dx = \mathbb{E}_\rho[f] \quad \text{for some } f: \mathbb{R}^n \rightarrow \mathbb{R}$$

or  $F(g) = \int_{\mathbb{R}^n} g(x) \|\nabla f(x)\|^2 dx = \mathbb{E}_g [\|\nabla f\|^2]$

# Examples of Functionals

$$F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$$

1. **Potential energy:**

$$F(\rho) = \int_{\mathbb{R}^n} \rho(x) f(x) dx = \mathbb{E}_\rho[f]$$

2. **Internal energy:**

$$F(\rho) = \int_{\mathbb{R}^n} U(\rho(x)) dx$$

- **Entropy:**  $H(\rho) = \int_{\mathbb{R}^n} \rho(x) \log \rho(x) dx = \mathbb{E}_\rho[\log \rho]$

# Functionals

3. Interaction energy:  $= \frac{1}{2} \mathbb{E}[W(x-y)], \quad x, y \sim \mathcal{S} \text{ iid.}$

$$F(\rho) = \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} W(x-y) \rho(x) \rho(y) dx dy$$

- Variance:  $\text{Var}(\rho) = \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} \|x-y\|^2 \rho(x) \rho(y) dx dy$

# Functionals

## 3. Interaction energy:

$$F(\rho) = \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} W(x - y) \rho(x) \rho(y) dx dy$$

- Variance:  $\text{Var}(\rho) = \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} \|x - y\|^2 \rho(x) \rho(y) dx dy$

## 4. Functionals involving derivatives

$$F(\rho) = \int_{\mathbb{R}^n} V(\rho(x), \nabla \rho(x)) dx$$

- Fisher information:  $J(\rho) = \mathbb{E}_\rho [\underbrace{\|\nabla \log \rho\|^2}_{= \frac{\nabla \rho}{\rho}}] = \int_{\mathbb{R}^n} \frac{\|\nabla \rho(x)\|^2}{\rho(x)} dx$

# Otto Calculus

---

## Review: Gradient on manifold

The **gradient** of  $f: \mathcal{X} \rightarrow \mathbb{R}$  at  $x \in \mathcal{X}$  is the tangent vector

$$\text{grad } f(x) \in T_x \mathcal{X}$$

that gives directional derivative:

$$\langle \text{grad } f(x), v \rangle_x = \frac{d}{dt} f(X_t) \Big|_{t=0}$$

where  $X_t$  is geodesic from  $X_0 = x$  along direction  $\dot{X}_0 = v$



## $L^2$ variation

The  $L^2$  variation of  $F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$  is  $\frac{\delta F}{\delta \rho}: \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$s \in \mathcal{P}(\mathbb{R}^n)$$

$$\rho: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\frac{\delta F}{\delta \rho}(x) = \frac{\partial F(\rho)}{\partial \rho(x)}$$

$$s = (\ s(x): x \in \mathbb{R}^n)$$

## $L^2$ variation

The  $L^2$  variation of  $F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$  is  $\frac{\delta F}{\delta \rho}: \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$\frac{\delta F}{\delta \rho}(x) = \frac{\partial F(\rho)}{\partial \rho(x)}$$

- Potential energy:  $F(\rho) = \int_{\mathbb{R}^n} \rho(x) f(x) dx$

$$\frac{\delta F}{\delta \rho}(x) = f(x)$$

## $L^2$ variation

The  $L^2$  variation of  $F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$  is  $\frac{\delta F}{\delta \rho}: \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$\frac{\delta F}{\delta \rho}(x) = \frac{\partial F(\rho)}{\partial \rho(x)}$$

- Potential energy:  $F(\rho) = \int_{\mathbb{R}^n} \rho(x) f(x) dx$

$$\frac{\delta F}{\delta \rho}(x) = f(x)$$

- Internal energy:  $F(\rho) = \int_{\mathbb{R}^n} U(\rho(x)) dx$

$$\frac{\delta F}{\delta \rho}(x) = U'(\rho(x))$$

# Gradient

$$T_g \mathcal{P} = \{ \phi = -\nabla \cdot (\beta \nabla u) \text{ for some } u: \mathbb{R}^n \rightarrow \mathbb{R} \}$$

Otto calculus: Rule for computing derivatives in  $W_2$  metric

**Lemma:** The *gradient* of  $F: \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$  at  $\rho$  is

$$\text{grad } F(\rho) = -\nabla \cdot \left( \rho \nabla \frac{\delta F}{\delta \rho} \right)$$

# Gradient

**Lemma:**

$$\operatorname{grad} F(\rho) = -\nabla \cdot \left( \rho \nabla \frac{\delta F}{\delta \rho} \right)$$

Proof: For any geodesic  $\rho_t$  from  $\rho_0 = \rho$  with  $\dot{\rho}_0 = -\nabla \cdot (\rho \nabla u)$

$$\begin{aligned} \frac{d}{dt} F(\rho_t) \Big|_{t=0} &= \int_{\mathbb{R}^n} \frac{\delta F(\rho)}{\delta \rho(x)} \dot{\rho}_0(x) dx \\ &= - \int_{\mathbb{R}^n} \frac{\delta F(\rho)}{\delta \rho(x)} \nabla \cdot (\rho \nabla u)(x) dx \\ &= \int_{\mathbb{R}^n} \rho(x) \left\langle \nabla \frac{\delta F(\rho)}{\delta \rho(x)}, \nabla u \right\rangle dx \\ &= \langle \operatorname{grad} F(\rho), \dot{\rho}_0 \rangle_\rho \end{aligned}$$

# Potential energy

---

# Gradient of potential energy

Potential energy:

$$F(\rho) = \int_{\mathbb{R}^n} \rho(x) f(x) dx = \mathbb{E}_\rho[f]$$

$L^2$  variation:

$$\frac{\delta F}{\delta \rho}(x) = f(x)$$

Gradient of potential energy:

$$\text{grad } F(\rho) = -\nabla \cdot (\rho \nabla f)$$

# Gradient flow of potential energy

Potential energy:  $F(\rho) = \mathbb{E}_\rho[f]$

Gradient flow of potential energy:

$$\dot{\rho}_t = -\text{grad } F(\rho_t) = \nabla \cdot (\rho_t \nabla f)$$

$$\Leftrightarrow \frac{\partial s_t}{\partial t} = \nabla \cdot (s_t \nabla f)$$

# Gradient flow of potential energy

Gradient flow of potential energy  $F(\rho) = \mathbb{E}_\rho[f]$ :

$$\frac{\partial \rho_t}{\partial t} = \nabla \cdot (\rho_t \nabla f)$$

This is the continuity equation of the **gradient flow** of  $f$ :

$$\dot{X}_t = -\nabla f(X_t)$$

Gradient flow of  $F(\rho) = \mathbb{E}_\rho[f]$  is *implemented* by gradient flow of  $f$