

A Geometric View of Optimal Transportation and Generative Model

Na Lei ^{*} Kehua Su [†] Li Cui [‡] Shing-Tung Yau [§] David Xianfeng Gu [¶]

Abstract

In this work, we show the intrinsic relations between optimal transportation and convex geometry, especially the variational approach to solve Alexandrov problem: constructing a convex polytope with prescribed face normals and volumes. This leads to a geometric interpretation to generative models, and leads to a novel framework for generative models.

By using the optimal transportation view of GAN model, we show that the discriminator computes the Kantorovich potential, the generator calculates the transportation map. For a large class of transportation costs, the Kantorovich potential can give the optimal transportation map by a close-form formula. Therefore, it is sufficient to solely optimize the discriminator. This shows the adversarial competition can be avoided, and the computational architecture can be simplified.

Preliminary experimental results show the geometric method outperforms WGAN for approximating probability measures with multiple clusters in low dimensional space.

1 Introduction

GAN model Generative Adversarial Networks (GANs) [10] aim at learning a mapping from a simple distribution to a given distribution. A GAN model consists of a generator G and a discriminator D , both are represented as deep networks. The generator captures the data distribution and generates samples, the discriminator estimates the probability that a sample came from the training data rather than G . Both generator and the discriminator are trained simultaneously. The competition drives both of them to improve their performance until the generated samples are indistinguishable from the genuine data samples. At the Nash equilibrium [50], the distribution generated by G equals to the real data distribution. GANs have several advantages: they can automatically generate samples, and reduce the amount of real data samples; furthermore, GANs do not need the explicit expression of the distribution of given data.

Recently, GANs receive an exploding amount of attention. For example, GANs have been widely applied to numerous computer vision tasks such as image inpainting [35, 49, 28], image super resolution [26, 19], semantic segmentation [52, 31], object detection [37, 27, 47], video prediction [32, 46], image translation [20, 51, 7, 29], 3D vision [48, 34], face editing [25, 30, 36, 39, 6, 40, 18], etc. Also, in machine learning field, GANs have been applied to semi-supervised learning [33, 24, 38], clustering [41], cross domain learning [42, 22], and ensemble learning [43].

^{*}Dalian University of Technology, Dalian, China. Email: nalei@dlut.edu.cn

[†]Wuhan University, Wuhan, China. Email: skh@whu.edu.cn

[‡]Beijing Normal University, Beijing, China. Email: licui@bnu.edu.cn

[§]Harvard University, Boston, US. Email: yau@math.harvard.edu

[¶]Stony Brook University, New York, US. Email: gu@cs.stonybrook.edu.

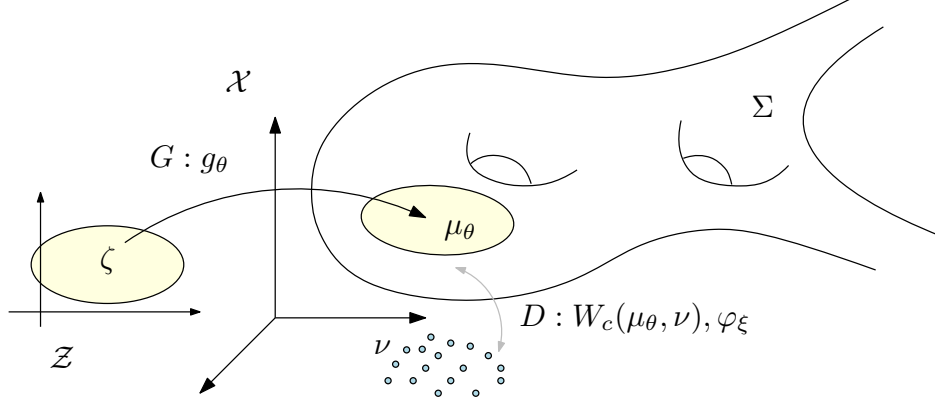


Figure 1: Wasserstein Generative Adversarial Networks (W-GAN) framework.

Optimal Transportation View Recently, optimal mass transportation theory has been applied to improve GANs. The Wasserstein distance has been adapted by GANs as the loss function as the discriminator, such as WGAN [3], WGAN-GP [13] and RWGAN [14]. When the supports of two distributions have no overlap, Wasserstein distance still provides a suitable gradient for the generator to update.

Figure 1 shows the optimal mass transportation point of view of WGAN [3]. The ambient image space is \mathcal{X} , the real data distribution is ν . The latent space is \mathcal{Z} with much lower dimension. The generator G can be treated as a mapping from the latent space to the sample space, $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$, realized by a deep network with parameter θ . Let ζ be a fixed distribution on the latent space, such as uniform distribution or Gaussian distribution. The generator G pushes forward ζ to a distribution $\mu_\theta = g_{\theta\#}\zeta$ in the ambient space \mathcal{X} . The discriminator D computes the distance between μ_θ and ν , in general using the Wasserstein distance, $W_c(\mu_\theta, \nu)$. The Wasserstein distance is equivalent to finding the so-called Kantorovich potential function φ_ξ , which is carried out by another deep network with parameter ξ . Therefore, G improves the "decoding" map g_θ to approximate ν by $g_{\theta\#}\zeta$; D improves the φ_ξ to increase the approximation accuracy to the Wasserstein distance. The generator G and the discriminator D are trained alternatively, until the competition reaches an equilibrium.

In summary, the generative model has a natural connection with the optimal mass transportation (OMT) theory:

1. In generator G , the generating map g_θ in GAN is equivalent to the optimal transportation map in OMT;
2. In discriminator D , the metric between distributions is equivalent to the Kantorovich potential φ_ξ .
3. The alternative training process of W-GAN is the min-max optimization of expectations:

$$\min_{\theta} \max_{\xi} \mathbb{E}_{z \sim \zeta} (\varphi_\xi(g_\theta(z))) + \mathbb{E}_{y \sim \nu} (\varphi_\xi^c(y)).$$

The deep nets of D and G perform the maximization and the minimization respectively.

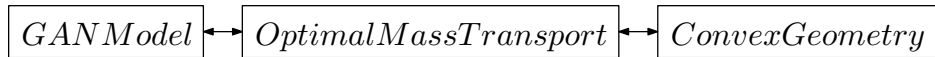


Figure 2: The GAN model, OMT theory and convex geometry has intrinsic relations.

Geometric Interpretation The optimal mass transportation theory has intrinsic connections with the convex geometry. Special OMT problem is equivalent to the Alexandrov theory in convex geometry: finding the optimal transportation map with L^2 cost is equivalent to constructing a convex polytope with user prescribed normals and face volumes. The geometric view leads to a practical algorithm, which finds the generating map g_θ by a convex optimization. Furthermore, the optimization can be carried out using Newton’s method with explicit geometric meaning. The geometric interpretation also gives the direct relation between the transportation map g_θ for G and the Kantorovich potential φ_ξ for D .

These concepts can be explained using the plain language in computational geometry [8],

1. the Kantorovich potential φ_ξ corresponds to the power distance;
2. the optimal transportation map g_θ represents the mapping from the power diagram to the power centers, each power cell is mapped to the corresponding site.

Imaginary Adversary In the current work, we use optimal mass transportation theory to show the fact that: by carefully designing the model and choosing special distance functions c , the generator map g_θ and the discriminator function (Kantorovich potential) φ_ξ are equivalent, one can be deduced from the other by a simple closed formula. Therefore, once the Kantorovich potential reaches the optimum, the generator map can be obtained directly without training. One of the deep neural net for G or D is redundant, one of the training processes is wasteful. The competition between the generator G and the discriminator D is unnecessary. In one word, the adversary is imaginary.

Contributions The major contributions of the current work are as follows:

1. Give an explicit geometric interpretation of optimal mass transportation map, and apply it for generative model;
2. Prove in theorem 3.7 that if the cost function $c(x, y) = h(x - y)$, where h is a strictly convex function, then once the optimal discriminator is obtained, the generator can be written down in an explicit formula. In this section, the competition between the discriminator and the generator is unnecessary and the computational architecture can be simplified;
3. Propose a novel framework for generative model, which uses geometric construction of the optimal mass transportation map;
4. Conduct preliminary experiments for the proof of concepts.

Organization The article is organized as follows: section 2 explains the optimal transportation view of WGAN in details; section 3 lists the main theory of OMT; section 4 gives the detailed exposition of Minkowski and Alexandrov theorems in convex geometry, and its close relation with power diagram theory in computational geometry, an explicit computational algorithm is given to solve Alexandrov’s problem; section 5 analyzes semi-discrete optimal transportation problem, and connects Alexandrov problem with the optimal transportation map; section 6 proposes a novel geometric generative model, which applies the geometric OMT map to the generative model; preliminary experiments are conducted for proof of concept, which are reported in section 7. The work concludes in the section 8.

2 Optimal Transportation View of GAN

This section, the GAN model is interpreted from the optimal transportation point of view. We show that the discriminator mainly looks for the Kantorovich potential.

Let $\mathcal{X} \subset \mathbb{R}^n$ be the (abient) image space, $\mathcal{P}(\mathcal{X})$ be the Wasserstein space of all probability measures on \mathcal{X} . Assume the data distribution is $\nu \in \mathcal{P}(\mathcal{X})$, represented as an empirical distribution

$$\nu := \frac{1}{n} \sum_{j=1}^n \delta_{y_j}, \quad (1)$$

where $y_j \in \mathcal{X}, j = 1, \dots, n$ are data samples. A *generative model* produces a parametric family of probability distributions $\mu_\theta, \theta \in \Theta$, a Minimum Kantorovitch Estimator for θ is defined as any solution to the problem

$$\min_{\theta} W_c(\mu_\theta, \nu),$$

where W_c is the Wasserstein cost on $\mathcal{P}(\mathcal{X})$ for some ground cost function $c : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$,

$$W_c(\mu, \nu) = \min_{\gamma \in \mathcal{P}(\mathcal{X} \times \mathcal{X})} \left\{ \int_{\mathcal{X} \times \mathcal{X}} c(x, y) d\gamma(x, y) \mid \pi_{x\#} \gamma = \mu, \pi_{y\#} \gamma = \nu \right\} \quad (2)$$

where π_x and π_y are projectors, $\pi_{x\#}$ and $\pi_{y\#}$ are marginalization operators. In a generative model, the image samples are encoded to a low dimensional latent space (or a feature space) $\mathcal{Z} \subset \mathbb{R}^m, m \ll n$. Let ζ be a fixed distribution supported on \mathcal{Z} . A WGAN produces a parametric mapping $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$, which is treated as a "decoding" map the latent space \mathcal{Z} to the original image space \mathcal{X} . g_θ pushes ζ forward to $\mu_\theta \in \mathcal{P}(\mathcal{X}), \mu_\theta = g_{\theta\#} \zeta$. The minimal Kantorovich estimator in WGAN is formulated as

$$\min_{\theta} E(\theta) := W_c(g_{\theta\#} \zeta, \nu).$$

According to the optimal transportation theory, the Kantorovich problem has a dual formulation

$$E(\theta) = \max_{\varphi, \psi} \left\{ \int_{\mathcal{Z}} \varphi(g_\theta(z)) d\zeta(z) + \int_{\mathcal{X}} \psi(y) d\nu(y); \varphi(x) + \psi(y) \leq c(x, y) \right\} \quad (3)$$

The gradient of the dual energy with respect to θ can be written as

$$\nabla E(\theta) = \int_{\mathcal{Z}} [\partial_\theta g_\theta(z)]^T \nabla \varphi^*(g_\theta(z)) d\zeta(z),$$

where φ^* is the optimal Kantorovich potential. In practice, ψ can be replaced by the c-transform of φ , defined as

$$\varphi^c(y) := \inf_x c(x, y) - \varphi(x).$$

The function φ is called the *Kantorovich potential*. Since ν is discrete, one can replace the continuous potential φ^c by a discrete vector $\sum_i \psi_i \delta_{y_i}$ and impose $\varphi = (\sum_i \psi_i \delta_{y_i})^c$. The optimization over $\{\psi_i\}$ can then be achieved using stochastic gradient descent, as in [9].

In WGAN [3], the dual problem Eqn. 3 is solved by approximating the Kantorovich potential φ by the so-called "adversarial" map $\varphi_\xi : \mathcal{X} \rightarrow \mathbb{R}$, where ξ is represented by a discriminative deep network. This leads to the Wasserstein-GAN problem

$$\min_{\theta} \max_{\xi} \int_{\mathcal{Z}} \varphi_\xi \circ g_\theta(z) d\zeta(z) + \frac{1}{n} \sum_{j=1}^n \varphi_\xi^c(y_j). \quad (4)$$

The generator produces g_θ , the discriminator estimates φ_ξ , by simultaneous training, the competition reaches the equilibrium. In WGAN [3], $c(x, y) = |x - y|$, then the c-transform of φ_ξ equals to $-\varphi_\xi$, subject to φ_ξ being a 1-Lipschitz function. This is used in to replace φ_ξ^c by $-\varphi_\xi$ in Eqn. 4 and use deep network made of ReLu units whose Lipschitz constant is upper-bounded by 1.

3 Optimal Mass Transport Theory

In this section, we review the classical optimal mass transportation theory. Theorem 3.7 shows the intrinsic relation between the Wasserstein distance (Kantorovich potential) and the optimal transportation map (Brenier potential), this demonstrates that once the optimal discriminator is known, the optimal generator is automatically obtained. The game between the discriminator and the generator is unnecessary.

The problem of finding a map that minimizes the inter-domain transportation cost while preserves measure quantities was first studied by Monge [4] in the 18th century. Let X and Y be two metric spaces with probability measures μ and ν respectively. Assume X and Y have equal total measure

$$\int_X d\mu = \int_Y d\nu.$$

Definition 3.1 (Measure-Preserving Map) A map $T : X \rightarrow Y$ is measure preserving if for any measurable set $B \subset Y$,

$$\mu(T^{-1}(B)) = \nu(B). \quad (5)$$

If this condition is satisfied, ν is said to be the push-forward of μ by T , and we write $\nu = T_{\#}\mu$.

If the mapping $T : X \rightarrow Y$ is differentiable, then measure-preserving condition can be formulated as the following Jacobian equation, $\mu(x)dx = \nu(T(x))dT(x)$,

$$\det(DT(x)) = \frac{\mu(x)}{\nu \circ T(x)}. \quad (6)$$

Let us denote the transportation cost for sending $x \in X$ to $y \in Y$ by $c(x, y)$, then the total transportation cost is given by

$$\mathcal{C}(T) := \int_X c(x, T(x))d\mu(x). \quad (7)$$

Problem 3.2 (Monge's Optimal Mass Transport[4]) Given a transportation cost function $c : X \times Y \rightarrow \mathbb{R}$, find the measure preserving map $T : X \rightarrow Y$ that minimizes the total transportation cost

$$(MP) \quad W_c(\mu, \nu) = \min_{T: X \rightarrow Y} \left\{ \int_X c(x, T(x))d\mu(x) : T_{\#}\mu = \nu \right\}. \quad (8)$$

The total transportation cost $W_c(\mu, \nu)$ is called the *Wasserstein distance* between the two measures μ and ν .

3.1 Kantorovich's Approach

In the 1940s, Kantorovich introduced the relaxation of Monge's problem [21]. Any strategy for sending μ onto ν can be represented by a joint measure ρ on $X \times Y$, such that

$$\rho(A \times Y) = \mu(A), \rho(X \times B) = \nu(B), \quad (9)$$

$\rho(A \times B)$ is called a *transportation plan*, which represents the share to be moved from A to B . We denote the projection to X and Y as π_x and π_y respectively, then $\pi_{x\#}\rho = \mu$ and $\pi_{y\#}\rho = \nu$. The total cost of the transportation plan ρ is

$$\mathcal{C}(\rho) := \int_{X \times Y} c(x, y)d\rho(x, y). \quad (10)$$

The Monge-Kantorovich problem consists in finding the ρ , among all the suitable transportation plans, minimizing $\mathcal{C}(\rho)$ in Eqn. 10

$$(KP) \quad W_c(\mu, \nu) := \min_{\rho} \left\{ \int_{X \times Y} c(x, y)d\rho(x, y) : \pi_{x\#}\rho = \mu, \pi_{y\#}\rho = \nu \right\} \quad (11)$$

3.2 Kantorovich Dual Formulation

Because Eqn. 11 is a linear program, it has a dual formulation, known as the Kantorovich problem [45]:

$$(DP) \quad W_c(\mu, \nu) := \max_{\varphi, \psi} \left\{ \int_X \varphi(x) d\mu(x) + \int_Y \psi(y) d\nu(y) : \varphi(x) + \psi(y) \leq c(x, y) \right\} \quad (12)$$

where $\varphi : X \rightarrow \mathbb{R}$ and $\psi : Y \rightarrow \mathbb{R}$ are real functions defined on X and Y . Equivalently, we can replace ψ by the c -transform of φ .

Definition 3.3 (c-transform) Given a real function $\varphi : X \rightarrow \mathbb{R}$, the c -transform of φ is defined by

$$\varphi^c(y) = \inf_{x \in X} (c(x, y) - \varphi(x)).$$

Then the Kantorovich problem can be reformulated as the following dual problem:

$$(DP) \quad W_c(\mu, \nu) := \max_{\varphi} \left\{ \int_X \varphi(x) d\mu(x) + \int_Y \varphi^c(y) d\nu(y) \right\}, \quad (13)$$

where $\varphi : X \rightarrow \mathbb{R}$ is called the *Kantorovich potential*.

For L^1 transportation cost $c(x, y) = |x - y|$ in \mathbb{R}^n , if the Kantorovich potential φ is 1-Lipsitz, then its c -transform has a special relation $\varphi^c = -\varphi$. The Wasserstein distance is given by

$$W_c(\mu, \nu) := \max_{\varphi} \left\{ \int_X \varphi(x) d\mu(x) - \int_Y \varphi(y) d\nu(y) \right\}, \quad (14)$$

For L^2 transportation cost $c(x, y) = 1/2|x - y|^2$ in \mathbb{R}^n , the c -transform and the classical Legendre transform has special relations.

Definition 3.4 Given a function $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$, its Legendre transform is defined as

$$\varphi^*(y) := \sup_x (\langle x, y \rangle - \varphi(x)). \quad (15)$$

Intuitively, Legendre transform has the following form:

$$\left(\int x dy \right)^* = \int y dx.$$

We can show the following relation holds when $c = 1/2|x - y|^2$,

$$\frac{1}{2}|y|^2 - \varphi^c = \left(\frac{1}{2}|x|^2 - \varphi \right)^*. \quad (16)$$

3.3 Brenier's Approach

At the end of 1980's, Brenier [5] discovered the intrinsic connection between optimal mass transport map and convex geometry. (see also for instance [44], Theorem 2.12(ii), and Theorem 2.32)

Suppose $u : X \rightarrow \mathbb{R}$ is a C^2 continuous convex function, namely its Hessian matrix is semi-positive definite. $(\partial^2 f / \partial x_i \partial x_j) \geq 0$. Its gradient map $\nabla u : X \rightarrow Y$ is defined as $x \mapsto \nabla u(x)$.

Theorem 3.5 (Brenier[5]) Suppose X and Y are the Euclidean space \mathbb{R}^n , and the transportation cost is the quadratic Euclidean distance $c(x, y) = |x - y|^2$. If μ is absolutely continuous and μ and ν have finite second order moments, then there exists a convex function $u : X \rightarrow \mathbb{R}$, its gradient map ∇u gives the solution to the Monge's problem, where u is called Brenier's potential. Furthermore, the optimal mass transportation map is unique.

This theorem converts the Monge's problem to solving the following Monge-Amperé partial differential equation:

$$\det \left(\frac{\partial^2 u}{\partial x_i \partial x_j} \right) (x) = \frac{\mu(x)}{\nu \circ \nabla u(x)}. \quad (17)$$

The function $u : X \rightarrow \mathbb{R}$ is called the *Brenier potential*. Brenier proved the polar factorization theorem.

Theorem 3.6 (Brenier Factorization[5]) *Suppose X and Y are the Euclidean space \mathbb{R}^n , $\varphi : X \rightarrow Y$ is measure preserving, $\varphi_{\#}\mu = \nu$. Then there exists a convex function $u : X \rightarrow \mathbb{R}$, such that*

$$\varphi = \nabla u \circ s,$$

where $s : X \rightarrow X$ preserves the measure μ , $s_{\#}\mu = \mu$. Furthermore, this factorization is unique.

Based on the generalized Brenier theorem we can obtain the following theorem.

Theorem 3.7 (Generator-Discriminator Equivalence) *Given μ and ν on a compact domain $\Omega \subset \mathbb{R}^n$ there exists an optimal transport plan γ for the cost $c(x, y) = h(x - y)$ with h strictly convex. It is unique and of the form $(\text{id}, T_{\#})\mu$, provided μ is absolutely continuous and $\partial\Omega$ is negligible. More over, there exists a Kantorovich potential φ , and T can be represented as*

$$T(x) = x - (\nabla h)^{-1}(\nabla \varphi(x)).$$

Proof: Assume ρ is the joint probability, satisfying the conditions $\pi_{x\#}\rho = \mu$, $\pi_{y\#}\rho = \nu$, (x_0, y_0) is a point in the support of ρ , by definition $\varphi^c(y_0) = \inf_x c(x, y_0) - \varphi(x)$, hence

$$\nabla \varphi(x_0) = \nabla_x c(x_0, y_0) = \nabla h(x_0 - y_0),$$

Because h is strictly convex, therefore ∇h is invertible,

$$x_0 - y_0 = (\nabla h)^{-1}(\nabla \varphi(x_0)),$$

hence $y_0 = x_0 - (\nabla h)^{-1}(\nabla \varphi(x_0))$. \square

When $c(x, y) = \frac{1}{2}|x - y|^2$, we have

$$T(x) = x - \nabla \varphi(x) = \nabla \left(\frac{x^2}{2} - \varphi(x) \right) = \nabla u(x).$$

In this case, the Brenier's potential u and the Kantorovich's potential φ is related by

$$u(x) = \frac{x^2}{2} - \varphi(x). \quad (18)$$

4 Convex Geometry

This section introduces Minkowski and Alexandrov problems in convex geometry, which can be described by Monge-Ampere equation as well. This intrinsic connection gives a geometric interpretation to optimal mass transportation map with L^2 transportation cost.

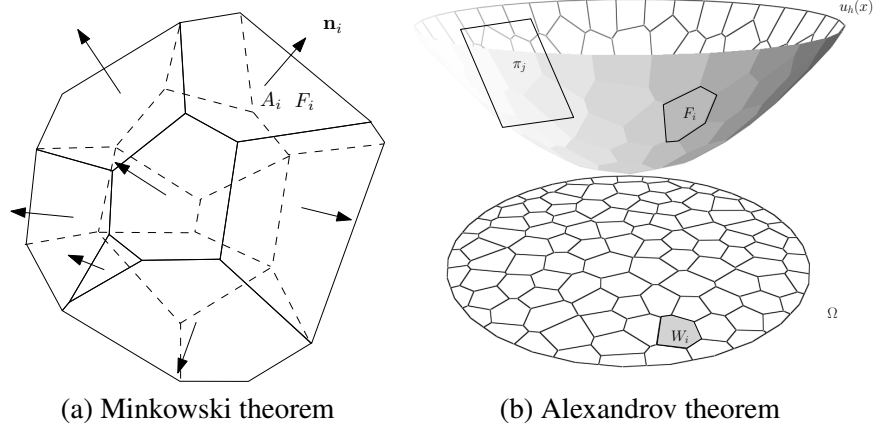


Figure 3: Minkowski and Alexandrov theorems for convex polytopes with prescribed normals and areas.

4.1 Alexandrov' Theorem

Minkowski proved the existence and the uniqueness of convex polytope with user prescribed face normals and the areas.

Theorem 4.1 (Minkowski) Suppose n_1, \dots, n_k are unit vectors which span \mathbb{R}^n and $\nu_1, \dots, \nu_k > 0$ so that $\sum_{i=1}^k \nu_i n_i = 0$. There exists a compact convex polytope $P \subset \mathbb{R}^n$ with exactly k codimension-1 faces F_1, \dots, F_k so that n_i is the outward normal vector to F_i and the volume of F_i is ν_i . Furthermore, such P is unique up to parallel translation.

Minkowski's proof is variational and suggests an algorithm to find the polytope. Minkowski theorem for unbounded convex polytopes was considered and solved by A.D. Alexandrov and his student A. Pogorelov. In his book on convex polyhedra [2], Alexandrov proved the following fundamental theorem (Theorem 7.3.2 and theorem 6.4.2)

Theorem 4.2 (Alexandrov[2]) Suppose Ω is a compact convex polytope with non-empty interior in \mathbb{R}^n , $n_1, \dots, n_k \subset \mathbb{R}^{n+1}$ are distinct k unit vectors, the $(n+1)$ -th coordinates are negative, and $\nu_1, \dots, \nu_k > 0$ so that $\sum_{i=1}^k \nu_i = \text{vol}(\Omega)$. Then there exists convex polytope $P \subset \mathbb{R}^{n+1}$ with exact k codimension-1 faces F_1, \dots, F_k so that n_i is the normal vector to F_i and the intersection between Ω and the projection of F_i is with volume ν_i . Furthermore, such P is unique up to vertical translation.

Alexandrov's proof is based on algebraic topology and non-constructive. Gu et al. [12] gave a variational proof for the generalized Alexandrov theorem stated in terms of convex functions.

Given $y_1, \dots, y_k \in \mathbb{R}^n$ and $h = (h_1, \dots, h_k) \in \mathbb{R}^k$, the piecewise linear convex function is defined as

$$u_h(x) = \max_i \{ \langle x, y_i \rangle + h_i \}.$$

The graph of u_h is a convex polytope in \mathbb{R}^{n+1} , the projection induces a cell decomposition of \mathbb{R}^n . Each cell is a closed convex polytope,

$$W_i(h) = \{x \in \mathbb{R}^n \mid \nabla u_h(x) = y_i\}.$$

Some cells may be empty or unbounded. Given a probability measure μ defined on Ω , the volume of $W_i(h)$ is defined as

$$w_i(h) := \mu(W_i(h) \cap \Omega) = \int_{W_i(h) \cap \Omega} d\mu.$$

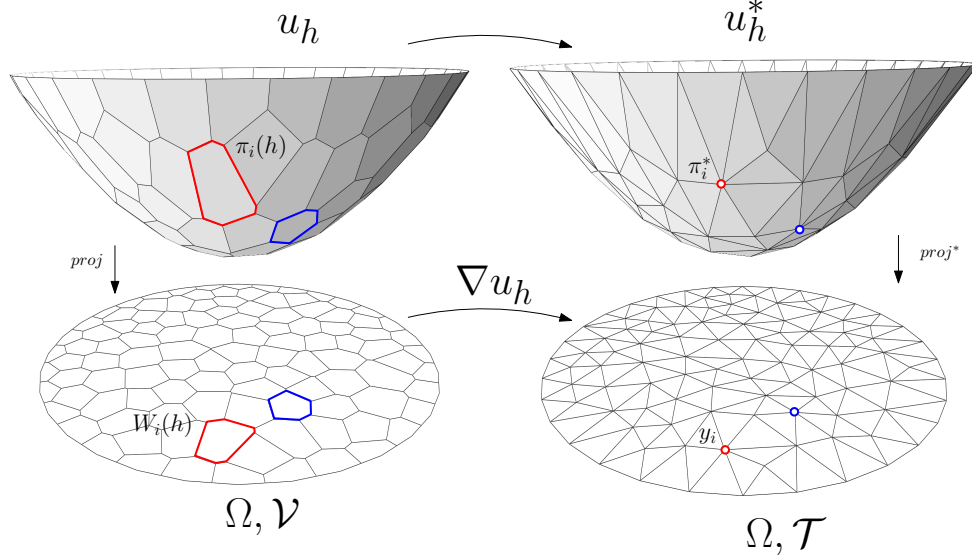


Figure 4: Geometric Interpretation to Optimal Transport Map: Brenier potential $u_h : \Omega \rightarrow \mathbb{R}$, Legendre dual u_h^* , optimal transportation map $\nabla u_h : W_i(h) \rightarrow y_i$, power diagram \mathcal{V} , weighted Delaunay triangulation \mathcal{T} .

Theorem 4.3 (Gu-Luo-Sun-Yau[12]) *Let Ω be a compact convex domain in \mathbb{R}^n , $\{y_1, \dots, y_k\}$ be a set of distinct points in \mathbb{R}^n and μ a probability measure on Ω . Then for any $\nu_1, \dots, \nu_k > 0$ with $\sum_{i=1}^k \nu_i = \mu(\Omega)$, there exists $h = (h_1, \dots, h_k) \in \mathbb{R}^k$, unique up to adding a constant (c, \dots, c) , so that $w_i(h) = \nu_i$, for all i . The vectors h are exactly maximum points of the concave function*

$$E(h) = \sum_{i=1}^k h_i \nu_i - \int_0^h \sum_{i=1}^k w_i(\eta) d\eta_i \quad (19)$$

on the open convex set

$$H = \{h \in \mathbb{R}^k | w_i(h) > 0, \forall i\}.$$

Furthermore, ∇u_h minimizes the quadratic cost

$$\int_{\Omega} |x - T(x)|^2 d\mu(x)$$

among all transport maps $T_{\#}\mu = \nu$, where the Dirac measure $\nu = \sum_{i=1}^k \nu_i \delta_{y_i}$.

For the convenience of discussion, we define the Alexandrov's potential as follows:

Definition 4.4 (Alexandrov Potential) *Under the above condition, the convex function*

$$\mathcal{A}(h) = \int_0^h \sum_{i=1}^k w_i(\eta) d\eta_i \quad (20)$$

is called the Alexandrov potential.

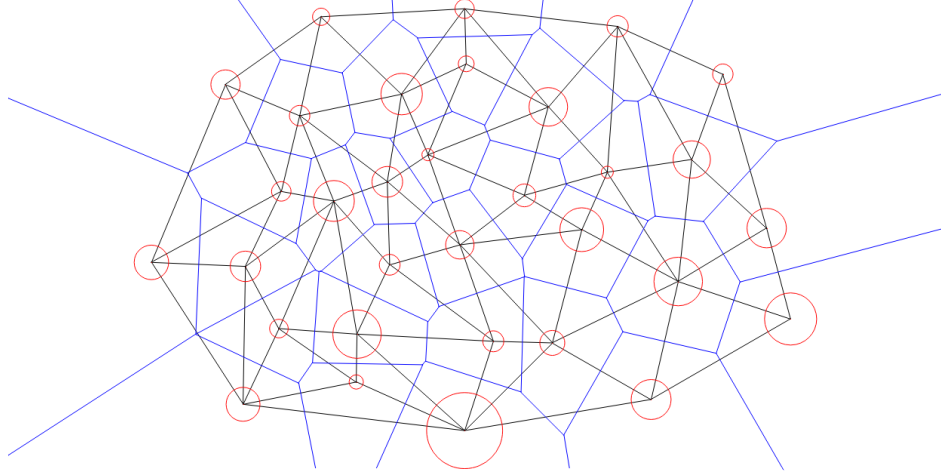


Figure 5: Power diagram (blue) and its dual weighted Delaunay triangulation (black), the power weight ψ_i equal to the square of radius r_i (red circle).

4.2 Power Diagram

Alexandrov's theorem has close relation with the conventional power diagram. We can use power diagram algorithm to solve the Alexandrov's problem.

Definition 4.5 (power distance) Given a point $y_i \in \mathbb{R}^n$ with a power weight ψ_i , the power distance is given by

$$\text{pow}(x, y_i) = |x - y_i|^2 - \psi_i.$$

Definition 4.6 (power diagram) Given weighted points $\{(y_1, \psi_1), (y_2, \psi_2), \dots, (y_k, \psi_k)\}$, the power diagram is the cell decomposition of \mathbb{R}^n , denoted as $\mathcal{V}(\psi)$,

$$\mathbb{R}^n = \bigcup_{i=1}^k W_i(\psi),$$

where each cell is a convex polytope

$$W_i(\psi) = \{x \in \mathbb{R}^n | \text{pow}(x, y_i) \leq \text{pow}(x, y_j), \forall j\}.$$

The weighted Delaunay triangulation, denoted as $\mathcal{T}(\psi)$, is the Poincaré dual to the power diagram, if $W_i(\psi) \cap W_j(\psi) \neq \emptyset$ then there is an edge connecting y_i and y_j in the weighted Delaunay triangulation.

Note that $\text{pow}(x, y_i) \leq \text{pow}(x, y_j)$ is equivalent to

$$\langle x, y_i \rangle + \frac{1}{2}(\psi_i - |y_i|^2) \geq \langle x, y_j \rangle + \frac{1}{2}(\psi_j - |y_j|^2).$$

let

$$h_i = 1/2(\psi_i - |y_i|^2), \quad (21)$$

we construct the convex function

$$u_h(x) = \max_i \{\langle x, y_i \rangle + h_i\}. \quad (22)$$

4.3 Convex Optimization

Now, we can use the power diagram to explain the gradient and the Hessian of the energy Eqn.19, by definition

$$\nabla E(h) = (\nu_1 - w_1(h), \nu_2 - w_2(h), \dots, \nu_k - w_k(h))^T. \quad (23)$$

The Hessian matrix is given by power diagram - weighted Delaunay triangulation, for adjacent cells in the power diagram,

$$\frac{\partial^2 E(h)}{\partial h_i \partial h_j} = \frac{\partial w_i(h)}{\partial h_j} = -\frac{\mu(W_i(h) \cap W_j(h) \cap \Omega)}{|y_j - y_i|} \quad (24)$$

Suppose edge e_{ij} is in the weighted Delaunay triangulation, connecting y_i and y_j . It has a unique dual cell D_{ij} in the power diagram, then

$$\frac{\partial w_i(h)}{\partial h_j} = -\frac{\mu(D_{ij})}{|e_{ij}|},$$

the volume ratio between the dual cells. The diagonal element in the Hessian is

$$\frac{\partial^2 E(h)}{\partial h_i^2} = \frac{\partial w_i(h)}{\partial h_i} = \sum_{j \neq i} \frac{\partial w_i(h)}{\partial h_j}. \quad (25)$$

Therefore, in order to solve Alexandrov's problem to construct the convex polytope with user prescribed normal and face volume, we can optimize the energy in Eqn. 19 using classical Newton's method directly.

Let's observe the convex function u_h^* , its graph is the convex hull $\mathcal{C}(h)$. Then the discrete Hessian determinant of u_h^* assigns each vertex v of $\mathcal{C}(h)$ the volume of the convex hull of the gradients of u_h^* at top-dimensional cells adjacent to v . Therefore, solving Alexandrov's problem is equivalent to solve a discrete Monge-Ampere equation.

5 Semi-discrete Optimal Mass Transport

In this section, we solve the semi-discrete optimal transportation problem from geometric point of view. This special case is useful in practice.

Suppose μ has compact support Ω on X , assume Ω is a convex domain in X ,

$$\Omega = \text{supp } \mu = \{x \in X | \mu(x) > 0\}.$$

The space Y is discretized to $Y = \{y_1, y_2, \dots, y_k\}$ with Dirac measure $\nu = \sum_{j=1}^k \nu_j \delta(y - y_j)$. The total mass are equal

$$\int_{\Omega} d\mu(x) = \sum_{i=1}^k \nu_i.$$

5.1 Kantorovich Dual Approach

We define the discrete Kantorovich potential $\psi : Y \rightarrow \mathbb{R}$, $\psi(y_j) = \psi_j$, then

$$\int_Y \psi d\nu = \sum_{j=1}^k \psi_j \nu_j. \quad (26)$$

The c-transformation of ψ is given by

$$\psi^c(x) = \min_{1 \leq j \leq k} \{c(x, y_j) - \psi_j\}. \quad (27)$$

This induces a cell decomposition of X ,

$$X = \bigcup_{i=1}^k W_i(\psi),$$

where each cell is given by

$$W_i(\psi) = \{x \in X | c(x, y_i) - \psi_i \leq c(x, y_j) - \psi_j, \forall 1 \leq j \leq k\}.$$

According to the dual formulation of the Wasserstein distance Eqn.13 and integration Eqn.26, we define the energy

$$E(\psi) = \int_X \psi^c d\mu + \int_Y \psi d\nu$$

then obtain the formula

$$E_D(\psi) = \sum_{i=1}^k \psi_i (\nu_i - w_i(\psi)) + \sum_{j=1}^k \int_{W_j(\psi)} c(x, y_j) d\mu. \quad (28)$$

where $w_i(\psi)$ is the measure of the cell $W_i(\psi)$,

$$w_i(\psi) = \mu(W_i(\psi)) = \int_{W_i(\psi)} d\mu(x). \quad (29)$$

Then the Wasserstein distance between μ and ν equals to

$$W_c(\mu, \nu) = \max_{\psi} E(\psi).$$

5.2 Brenier's Approach

Kantorovich's dual approach is for general cost functions. When the cost function is the L^2 distance $c(x, y) = |x - y|^2$, we can apply Brenier's approach directly.

We define a *height vector* $h = (h_1, h_2, \dots, h_k) \in \mathbb{R}^n$, consisting of k real numbers. For each $y_i \in Y$, we construct a hyperplane defined on X , $\pi_i(h) : \langle x, y_i \rangle + h_i = 0$. We define the Brenier potential function as

$$u_h(x) = \max_{i=1}^k \{\langle x, y_i \rangle + h_i\}, \quad (30)$$

then $u_h(x)$ is a convex function. The graph of $u_h(x)$ is an infinite convex polyhedron with supporting planes $\pi_i(h)$. The projection of the graph induces a polygonal partition of Ω ,

$$\Omega = \bigcup_{i=1}^k W_i(h), \quad (31)$$

where each cell $W_i(h)$ is the projection of a facet of the graph of u_h onto Ω ,

$$W_i(h) = \{x \in X | \nabla u_h(x) = y_i\} \cap \Omega. \quad (32)$$

The measure of $W_i(h)$ is given by

$$w_i(h) = \int_{W_i(h)} d\mu. \quad (33)$$

The convex function u_h on each cell $W_i(h)$ is a linear function $\pi_i(h)$, therefore, the gradient map

$$\nabla u_h : W_i(h) \rightarrow y_i, i = 1, 2, \dots, k. \quad (34)$$

maps each $W_i(h)$ to a single point y_i . According to Alexandrov's theorem, and the Gu-Luo-Yau theorem, we obtain the following corollary:

Corollary 5.1 *Let Ω be a compact convex domain in \mathbb{R}^n , $\{y_1, \dots, y_k\}$ be a set of distinct points in \mathbb{R}^n and μ a probability measure on Ω . Then for any $\nu = \sum_{i=1}^k \nu_i \delta_{y_i}$, with $\sum_{i=1}^k \nu_i = \mu(\Omega)$, there exists $h = (h_1, \dots, h_k) \in \mathbb{R}^k$, unique up to adding a constant (c, \dots, c) , so that $w_i(h) = \nu_i$, for all i . The vectors h are exactly maximum points of the concave function*

$$E_B(h) = \sum_{i=1}^k h_i \nu_i - \int_0^h \sum_{i=1}^k w_i(\eta) d\eta_i \quad (35)$$

Furthermore, ∇u_h minimizes the quadratic cost

$$\int_{\Omega} |x - T(x)|^2 d\mu$$

among all transport maps $T_{\#}\mu = \nu$.

5.3 Equivalence

For $c(x, y) = 1/2|x - y|^2$ cost cases, we have introduced two approaches: Kantorovich's dual approach and Brenier's approach. In the following, we show these two approaches are equivalent.

In Kantorovich's dual approach, finding the optimal mass transportation is equivalent to maximize the following energy:

$$E_D(\psi) = \sum_{i=1}^k \psi_i (\nu_i - w_i(\psi)) + \sum_{j=1}^k \int_{W_j(\psi)} c(x, y_j) d\mu.$$

In Brenier's approach, finding the optimal transportation map boils down to maximize

$$E_B(h) = \sum_{i=1}^k h_i \nu_i - \int_0^h \sum_{i=1}^k w_i(\eta) d\eta.$$

Lemma 5.2 *Let Ω be a compact convex domain in \mathbb{R}^n , $\{y_1, \dots, y_k\}$ be a set of distinct points in \mathbb{R}^n . Given μ a probability measure on Ω , $\nu = \sum_{i=1}^k \nu_i \delta_{y_i}$, with $\sum_{i=1}^k \nu_i = \mu(\Omega)$. If $c(x, y) = 1/2|x - y|^2$, then*

$$h_i = \psi_i - \frac{1}{2}|y_i|^2, \quad \forall i$$

and

$$E_D(\psi) - E_B(h) = \text{Const}$$

proof: Consider the power cell

$$c(x, y_i) - \psi_i \leq c(x, y_j) - \psi_j$$

is equivalent to

$$\langle x, y_i \rangle + \left(\psi_i - \frac{1}{2}|y_i|^2 \right) \geq \langle x, y_j \rangle + \left(\psi_j - \frac{1}{2}|y_j|^2 \right)$$

therefore $h_i = \psi_i - 1/2|y_i|^2$.

Let the transportation cost to be defined as

$$\mathcal{C}(\psi) = \sum_{j=1}^k \int_{W_j(\psi)} c(x, y_j) d\mu.$$

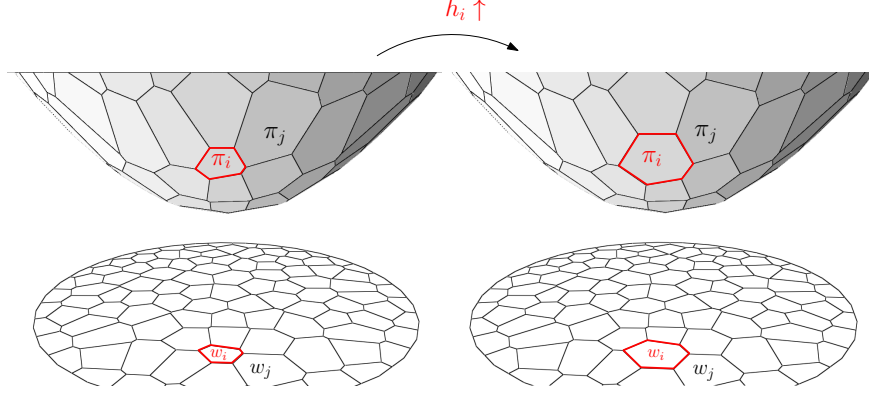


Figure 6: Variation of the volume of top-dimensional cells

Suppose we infinitesimally change h to $h + dh$, then we define

$$D_{ij} = W_j(h) \cap W_i(h + dh) \cap \Omega.$$

Then $\mu(D_{ij}) = dw_i$, also $\mu(D_{ij}) = -dw_j$. For each $x \in D_{ij}$, $c(x, y_i) - \psi_i = c(x, y_j) - \psi_j$, then $c(x, y_i) - c(x, y_j) = \psi_i - \psi_j$, hence

$$\int_{D_{ij}} (c(x, y_i) - c(x, y_j)) d\mu = \psi_i dw_i + \psi_j dw_j.$$

This shows $d\mathcal{C} = \sum_{i=1}^k \psi_i dw_i$, hence

$$\mathcal{C}(w) = \int^w \sum_{i=1}^k \psi_i dw_i.$$

The Legendre dual of \mathcal{C} is

$$\mathcal{C}^*(\psi) = \int^\psi \sum_{i=1}^k w_i d\psi_i.$$

Hence

$$\int^w \sum_{i=1}^k \psi_i dw_i + \int^\psi \sum_{i=1}^k w_i d\psi_i = \sum_{i=1}^k w_i \psi_i.$$

On the other hand, $\psi_i = h_i + 1/2|y_i|^2$, $d\psi_i = dh_i$,

$$\int^h \sum_{i=1}^k w_i dh_i = \int^\psi \sum_{i=1}^k w_i d\psi_i + \text{const.}$$

We put everything together

$$E_D(\psi) - E_B(h) = \sum_i (\psi_i - h_i) \nu_i - \left(\sum_i \psi_i \nu_i - \mathcal{C}(\psi) - \mathcal{C}^*(w) \right) - c_1 = c_2,$$

where C_1 and C_2 are two constants. \square

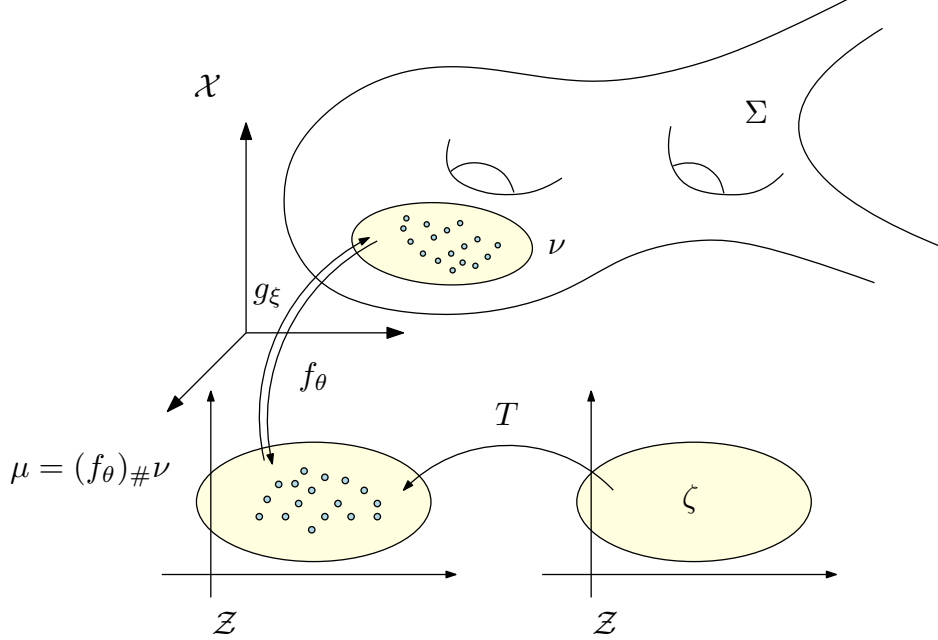


Figure 7: The framework for Geometric Generative Model.

This shows Kantorovich's dual approach and Brenier's approach are equivalent. At the optimal point, $\nu_i = w_i(\psi)$, therefore $E_D(\psi)$ equals to the transportation cost $\mathcal{C}(\psi)$. Furthermore, the Brenier's potential is

$$u_h(x) = \max_{i=1}^k \{ \langle x, p_i \rangle + h_i \},$$

where h_i is given by the power weight ψ_i . The Kantorovich's potential is the power distance

$$\varphi(x) = \psi^c(x) = \min_j \{ c(x, y_j) - \psi_j \} = \min_j \{ \text{pow}(x, y_j) \} = \frac{1}{2} |x|^2 - \max_j \{ \langle x, y_j \rangle + (\psi_j - \frac{1}{2} |y_j|^2) \}$$

hence at the optimum, the Brenier potential and the Kantorovich potential are related by

$$u_h(x) = \frac{1}{2} |x|^2 - \varphi(x). \quad (36)$$

6 Geometric Generative Model

In this section, we propose a novel generative framework, which combines the discriminator and the generator together. The model decouples the two processes

1. Encoding/decoding process: This step maps the samples between the image space \mathcal{X} and the latent (feature) space \mathcal{Z} by using deep neural networks, the encoding map is denoted as $f_\theta : \mathcal{X} \rightarrow \mathcal{Z}$, the decoding map is $g_\xi : \mathcal{Z} \rightarrow \mathcal{X}$. This step achieves the dimension deduction.
2. Probability measure transformation process: this step transform a fixed distribution $\zeta \in \mathcal{P}(\mathcal{Z})$ to any given distribution $\mu \in \mathcal{P}(\mathcal{Z})$. The mapping is denoted as $T : \mathcal{Z} \rightarrow \mathcal{Z}$, $T_\# \zeta = \mu$. This step can either use conventional deep neural network or use explicit geometric/numerical methods.

There are many existing methods to accomplish the encoding/decoding process, such as VAE model [23], therefore we focus on the second step.

As shown in Fig. 7, given an empirical distribution $\nu = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$ in the original ambient space \mathcal{X} , the support of ν is a sub-manifold $\Sigma \subset \mathcal{X}$. The encoding map $f_\theta : \Sigma \rightarrow \mathcal{Z}$ transform the support manifold Σ to the latent (or feature) space \mathcal{Z} , f_θ pushes forward the empirical distribution to μ defined on latent space

$$\mu = (f_\theta)_\# \nu = \frac{1}{n} \sum \delta_{z_i}. \quad (37)$$

where $z_i = f_\theta(y_i)$.

Let ζ be a fixed measure on the latent space, we would like to find an optimal transportation map $T : \mathcal{Z} \rightarrow \mathcal{Z}$, such that $T_\# \zeta = \mu$. This is equivalent to find the Brenier potential

$$u_h(z) = \max_{i=1}^n \{ \langle z, z_i \rangle + h_i \}.$$

Note that, u_h can be easily represented by linear combinations and ReLus. The height parameter can be obtained by optimizing the energy Eqn. 19

$$\frac{1}{n} \sum_{i=1}^k h_i - \int^h \sum_{i=1}^k w_i(\eta) d\eta_i.$$

The optimal transportation map $T = \nabla u_h$. This can be carried out as a power diagram with weighted points $\{(z_i, \psi_i)\}$, where

$$\psi_i = \frac{|z_i|^2}{2} - h_i.$$

The relation between the Kantorovich potential and the Brenier potential is

$$\varphi(x) = \frac{1}{2} |x|^2 - u_h(x).$$

The Wasserstein distance can be explicitly given by

$$W_c(\zeta, \mu) = \int_{\mathcal{Z}} \varphi(z) d\zeta(z) + \frac{1}{n} \sum_{j=1}^n \psi_j.$$

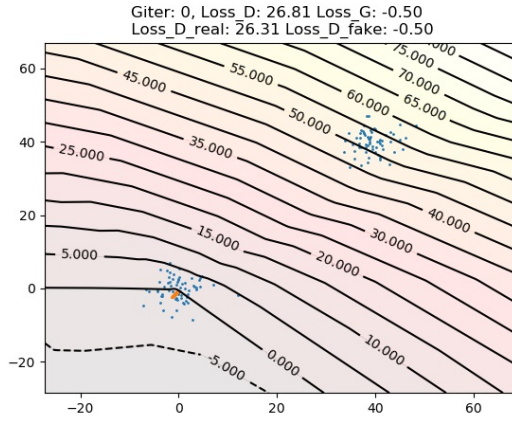
We use $g_\xi : \mathcal{Z} \rightarrow \mathcal{X}$ to denote the decoding map. Finally, the composition $g_\xi \circ T : \mathcal{Z} \rightarrow \mathcal{X}$ transforms ζ in the latent space to the original empirical distribution ν in the image space \mathcal{X} .

7 Experiments

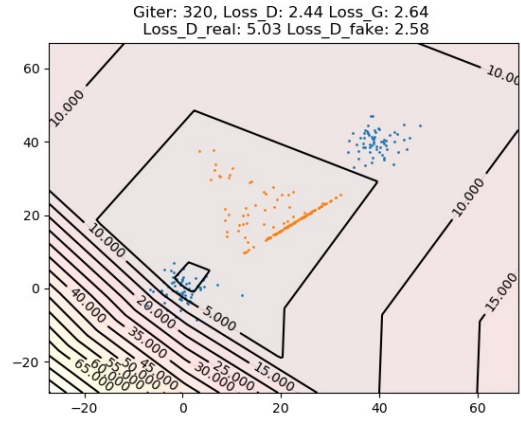
In order to demonstrate in principle the potential of our proposed method, we have designed and conducted the preliminary experiments.

7.1 Comparison with WGAN

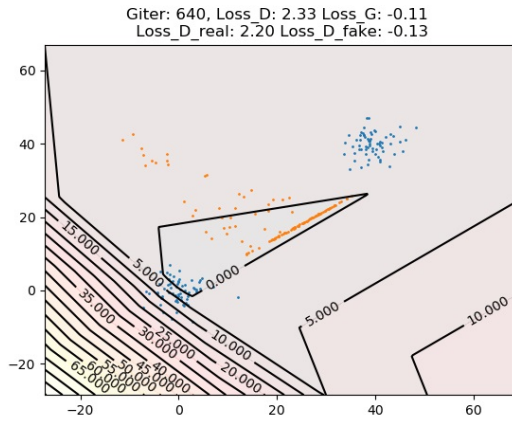
In the first experiment, we use Wasserstein Generative Adversarial Networks (WGANs) [3] to learn the mixed Gaussian distribution as shown in Fig. 8.



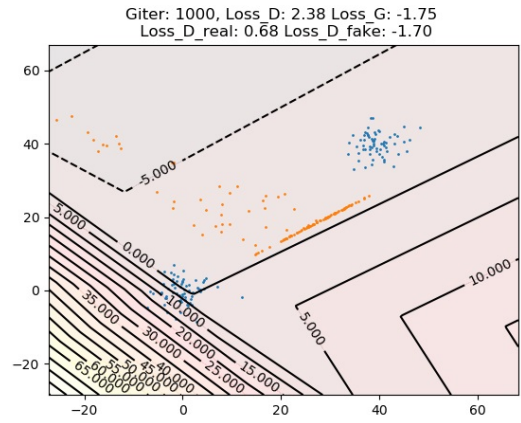
(a) initial stage



(b) after 320 iterations



(c) after 640 iterations



(d) final stage, after 1000 iterations

Figure 8: WGAN learns the Gaussian mixture distribution.

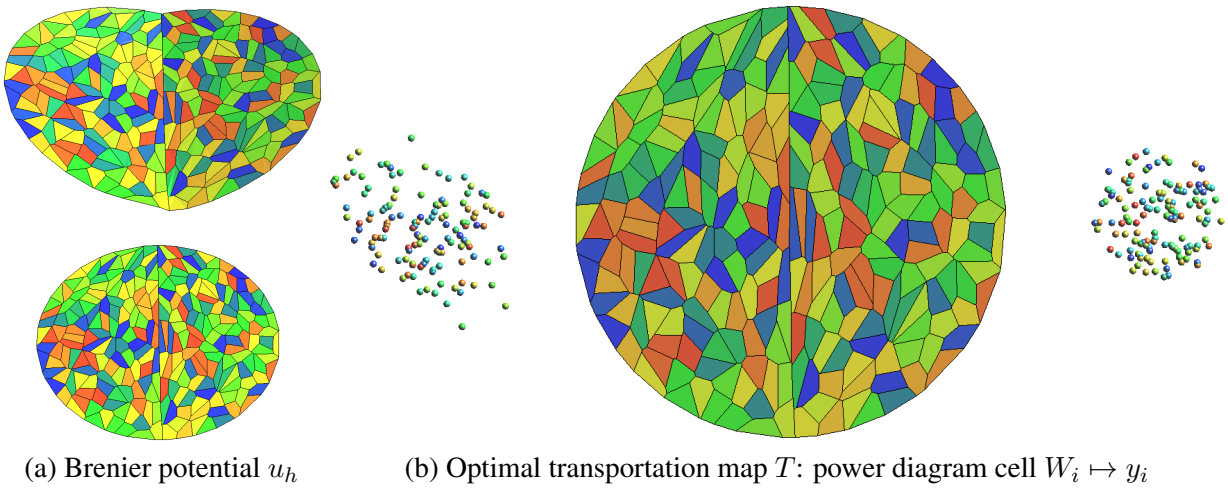


Figure 9: Geometric model learns the Gaussian mixture distribution .

Dataset The distribution of data ν is described by a point cloud on a $2d$ plane. We sample 128 data points as real data from two Gaussian distributions, $\mathcal{N}(p_k, \sigma_k^2)$, $k = 1, 2$, where $p_1 = (0, 0)$ and $\sigma_1 = 3$, $p_2 = (40, 40)$ and $\sigma_2 = 3$. The latent space \mathcal{Z} is a square on the $2d$ plane $[1k, 3k] \times [1k, 3k]$, the input distribution ζ is the uniform distribution on \mathcal{Z} . We use a generator to generate data from ζ to approximate the data distribution ν . We generate 128 samples in total.

Network Structure The structure of the discriminator is 2-layer (2×10 FC)-ReLU- $(10 \times 1$ FC) network, where FC denotes the fully connected layer. The number of inputs is 2 and the number of outputs is 1. The number of nodes of the hidden layer is 10.

The structure of the generator is a 6-layer (2×10 FC)-ReLU- $(10 \times 10$ FC)-ReLU- $(10 \times 10$ FC)-ReLU- $(10 \times 10$ FC)-ReLU- $(10 \times 2$ FC) network. The number of inputs is 2 and the number of outputs is 2. The number of nodes of all the hidden layer is 10.

Parameter Setting For WGAN, we clip all the weights to $[-0.5, 0.5]$. We use the RMSprop [16] as the optimizer for both discriminator and generator. The learning rate of both the discriminator and generator are set to $1e - 3$.

Deep learning framework and hardware We use the PyTorch [1] as our deep learning tool. Since the toy dataset is small, we do experiments on CPU. We perform experiments on a cluster with 48 cores and 193GB RAM. However, for this toy data, the running code only consumes 1 core with less than 500MB RAM, which means that it can run on a personal computer.

Results analysis In Fig. 8, the blue points represent the real data distribution and the orange points represent the generated distribution. The left frame shows the initial stage, the right frame illustrates the stage after 1000 iterations. It seems that WGAN cannot capture the Gaussian mixture distribution. Generated data tend to lie in the middle of the two Gaussians. One reason is the well known mode collapse problem in GAN, meaning that if the data distribution has multiple clusters or data is distributed in multiple isolated manifolds, then the generator is hard to learn multiple modes well. Although there are a couple of methods proposed to deal with this problem [15, 17], these methods require the number of clusters, which is still an open problem in the machine learning community.

Geometric OMT Figure 9 shows the geometric method to solve the same problem. The left frame shows the Brenier potential u_h , namely the upper envelope, which projects to the power diagram \mathcal{V} on a unit disk $\mathbb{D} \subset \mathcal{Z}$, $\mathcal{V} = \bigcup_k W_i(h)$. The right frame shows the discrete optimal transportation map $T : \mathbb{D} \rightarrow \{y_i\}$, which maps each cell $W_i(h) \cap \mathbb{D}$ to a sample y_i , the cell $W_i(h)$ and the sample y_i have the same color. All the cells have the same area, this demonstrates that T pushes the uniform distribution ζ to the exact empirical distribution $T_{\#}\zeta = 1/n \sum_i \delta_{y_i}$.

The samples $\{y_i\}$ are generated according to the same Gauss mixture distribution, therefore there are two clusters. This doesn't cause any difficulty for the geometric method. In the left frame, we can see the upper envelope has a sharp ridge, the gradients point to the two clusters. Hence, the geometric method outperforms the WGAN model in the current experiment.

7.2 Geometric Method

In this experiment, we use pure geometric method to generate uniform distribution on a surface Σ with complicated geometry. As shown in Fig. 10, the image space \mathcal{X} is the 3 dimensional Euclidean space \mathbb{R}^3 . The real data samples are distributed on a submanifold Σ , which is represented as a surface, as illustrated in

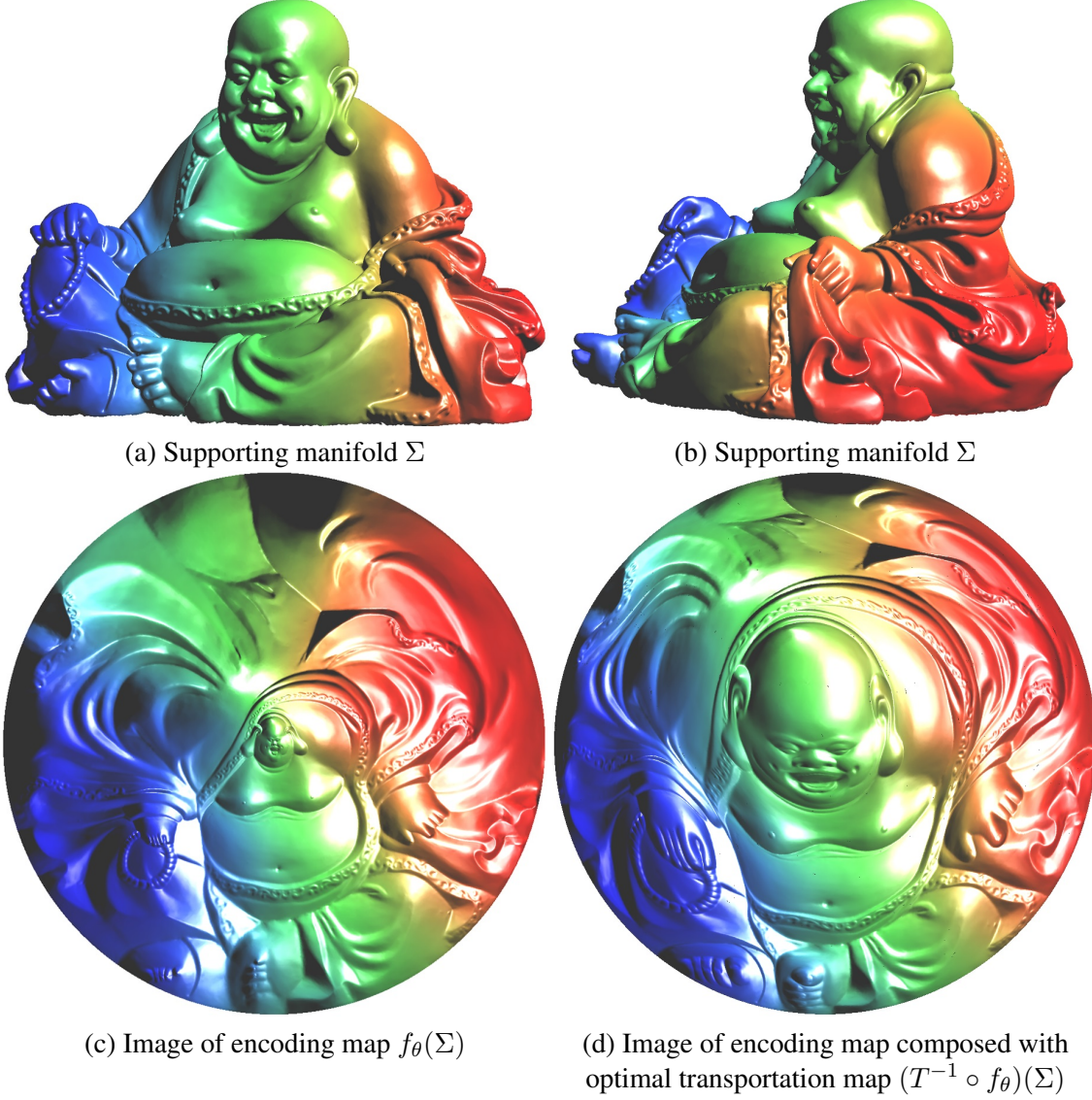


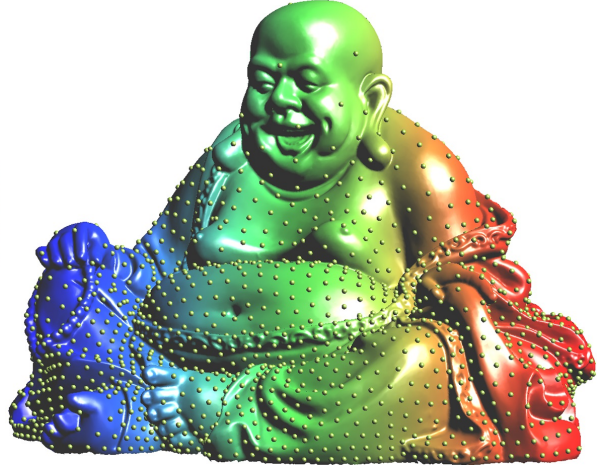
Figure 10: Illustration of geometric generative model.

(a) and (b). The encoding mapping $f_\theta : \Sigma \rightarrow \mathcal{Z}$ maps the supporting manifold to the latent space, which is a planar disk. The encoding map f_θ can be computed using discrete surface Ricci flow method [11]. We color-encode the normals to the surface, and push forward the color function from Σ to the latent space $f_\theta(\Sigma)$, therefore users can directly visualize the correspondence between Σ and its image in \mathcal{Z} as shown in (c). Then we construct the optimal mass transportation map $T : \mathcal{Z} \rightarrow \mathcal{Z}$, the image is shown in (d).

In Fig. 11, we demonstrate the generated distributions. In (a), we generate samples $\{z_1, \dots, z_k\}$ on the latent space $f_\theta(\Sigma)$ according to the uniform distribution ζ , the samples are pulled back to the surface Σ as $\{f_\theta^{-1}(z_1), \dots, f_\theta^{-1}(z_k)\}$ as shown in (b), which illustrate the distribution $(f_\theta^{-1})_\# \zeta$. It is obvious that the distribution generated this way is highly non-uniform on the surface. In frame (c), we uniformly generate samples on $(T^{-1} \circ f_\theta)(\Sigma)$, and map them back to the surface Σ as shown in (d). This demonstrates the generated distribution $(f_\theta^{-1} \circ T)_\# \zeta$ on Σ , which is uniform as desired.



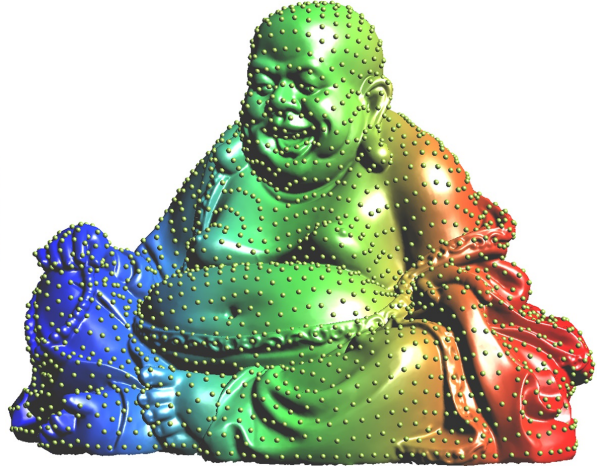
(a) sampling according to the uniform distribution ζ on \mathcal{Z}



(b) non-uniform sampling according to the distribution $(f_\theta^{-1})_\# \zeta$ on Σ



(c) sampling according to the uniform distribution ζ on \mathcal{Z}



(d) uniform sampling according to $(f_\theta^{-1} \circ T)_\# \zeta$ on Σ

Figure 11: Illustration of geometric generative model.

8 Discussion and Conclusion

In this work, we bridge convex geometry with optimal transportation, then use optimal transportation to analyze generative models. The basic view is that the discriminator computes the Wasserstein distance or equivalently the Kantorovich potential φ_ξ ; the generator calculates the transportation map g_θ . By selecting the transportation cost, such as L^p , $p > 1$ distance, φ_ξ and g_θ are related by a closed form, hence it is sufficient to train one of them.

For general transportation cost $c(x, y)$, the explicit relation between φ_ξ and g_θ may not exist, it seems that both training processes are necessary. In the following, we argue that it is still redundant. For a given cost function c , the optimal decoding map is g_1 ; by using L^2 cost function, the solution is g_0 . Both g_0 and g_1 induces the same measure,

$$(g_0)_\# \zeta = (g_1)_\# \zeta = \nu,$$

By Brenier polar factorization theorem 3.6, $g_1 = g_0 \circ s$, where $s : \mathcal{Z} \rightarrow \mathcal{Z}$, preserves the measure ζ , $s_{\#}\zeta = \zeta$. All such mappings form an infinite dimensional group, comparing to g_0 , the complexity of s increases the difficulty of finding g_1 . But both g_0 and g_1 generate the same distribution, there is no difference in terms of the performance of the whole system. It is much more efficient to use g_0 without the double training processes.

For high dimensional setting, rigorous computational geometric method to compute the optimal transportation map is intractable, due to the maintenance of the complex geometric data structures. Nevertheless, there exist different algorithms to handle high dimensional situation, such as socialistic method, sliced optimal transportation method, hierarchical optimal transportation method.

In the future, we will explore along this direction, and implement the proposed model in a large scale.

Acknowledgement

We thank the inspiring discussions with Dr. Dimitris Samaras from Stony Brook University, Dr. Shoucheng Zhang from Stanford University, Dr. Limin Chen from Ecole Centrale de Lyon and so on. The WGAN experiment is conducted by Mr. Huidong Liu from Stony Brook University. The project is partially supported by NSFC No. 61772105, 61772379, 61720106005, NSF DMS-1418255, AFOSR FA9550-14-1-0193 and the Fundamental Research Funds for the Central Universities No. 2015KJJC23.

Appendix

8.1 Commutative Diagram

The relations among geometric/functional objects are summarized in the following diagram:

$$\begin{array}{ccc}
 \mathcal{A} & \xrightarrow{\text{Legendre dual}} & \mathcal{C} \\
 \uparrow \text{integrate} & & \uparrow \\
 u_h & \xrightarrow{\text{Legendre dual}} & u_h^* \\
 \downarrow \text{graph} & & \downarrow \text{graph} \\
 \text{Env}(\{\pi_i\}) & \xrightarrow{\text{Poincare dual}} & \text{Conv}(\{\pi_i^*\}) \\
 \downarrow \text{proj} & & \downarrow \text{proj} \\
 \mathcal{V}(\psi) & \xrightarrow{\text{Poincare dual}} & \mathcal{T}(\psi)
 \end{array}$$

where each two adjacent layers are commutable. These relations can be observed from Fig. 4 as well.

8.2 Symbol list

The following is the symbol list

References

- [1] <http://pytorch.org/>.

Table 1: Symbol list.

\mathcal{X}	ambient space, image space	
Σ	support manifold for some distribution	
\mathcal{Z}	latent space, feature space	
ζ	a fixed probability measure on \mathcal{Z}	
g_θ	generating map	$g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$
φ_ξ	Kantorovich potential	
c	distance between two points	$c(x, y) = x - y ^p, p \geq 1$
W_c	Wasserstein distance	$W_c(\mu, \nu)$
X	source space	
Y	target space	
μ	source probability measure	
ν	target probability measure	
Ω	source domain	$\Omega \subset X$
y_i	the i -th sample in target	$\{y_1, \dots, y_k\} \in Y$
φ	Kantorovich potential	$\phi^c = \psi, \psi^c = \phi$
ψ	power weight	$\psi = (\psi_1, \dots, \psi_k)$
h	plane heights	$h = (h_1, \dots, h_k)$
π_i	hyper-plane	$\pi_i^h(x) = \langle y_i, x \rangle + h_i$
π_i^*	dual point of π_i	$\pi_i^* = (y_i, -h)$
pow	power distance	$\text{pow}(x, y_i) = c(x, y_i) - \psi_i$
W_i	power voronoi cell	$W_i(\psi) = \{x \in X \mid \text{pow}(x, y_i) \leq \text{pow}(x, y_j)\}$
w_i	the volume of W_i	$w_i(h) = \mu(W_i(h) \cap \Omega)$
u	Brenier potential	$u_h(x) = \max_i \{\langle x, y_i \rangle + h_i\}$
\mathcal{A}	Alexandrov potential	$\mathcal{A}(h) = \int^h \sum_i w_i dh_i$
T	transportation map	$T = \nabla u_h$
\mathcal{C}	transportation cost	$\mathcal{C}(T) = \int_X c(x, T(x)) d\mu$
Env	upper envelope of planes	$\text{Env}(\{\pi_i\})$ graph of u_h
Conv	convex hull of points	$\text{Conv}(\{\pi_i^*\})$ graph of u_h^*
\mathcal{V}	power diagram	$\mathcal{V}(\psi) : X = \bigcup_i W_i(\psi)$
\mathcal{T}	weighted Delaunay triangulation	

- [2] A. D. Alexandrov. *Convex polyhedra Translated from the 1950 Russian edition by N. S. Dairbekov, S. S. Kutateladze and A. B. Sossinsky*. Springer Monographs in Mathematics. Springer-Verlag, Berlin, 2005.
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pages 214–223, 2017.
- [4] Nicolas Bonnotte. From Knothe’s rearrangement to Brenier’s optimal transport map. *arXiv:1205.1099*, pages 1–29, 2012.
- [5] Yann Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Comm. Pure Appl. Math.*, 44(4):375–417, 1991.
- [6] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Neural photo editing with introspective adversarial networks. In *International Conference on Learning Representations*, 2017.

- [7] Hao Dong, Paarth Neekhara, Chao Wu, and Yike Guo. Unsupervised image-to-image translation with generative adversarial networks. *arXiv preprint arXiv:1701.02676*, 2017.
- [8] Herbert Edelsbrunner. *Voronoi Diagrams*, pages 293–333. Springer Berlin Heidelberg, Berlin, Heidelberg, 1987.
- [9] Aude Genevay, Marco Cuturi, Gabriel Peyr’e, and Fancis Bach. Stochastic optimization for large-scale optimal transport. In D.D. Lee, U.V. Luxburg, I. Guyon, and R.Garnett, editors, *Proceedings of NIPS’16*, pages 3432–3440, 2016.
- [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.
- [11] Xianfeng Gu, Feng Luo, Jian Sun, and Tianqi Wu. A discrete uniformization theorem for polyhedral surfaces. *Journal of Differential Geometry (JDG)*, 2017.
- [12] Xianfeng Gu, Feng Luo, Jian Sun, and Shing-Tung Yau. Variational principles for minkowski type problems, discrete optimal transport, and discrete monge-ampere equations. *Asian Journal of Mathematics (AJM)*, 20(2):383 C 398, 2016.
- [13] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of wasserstein gans. *arXiv preprint arXiv:1704.00028*, 2017.
- [14] Xin Guo, Johnny Hong, Tianyi Lin, and Nan Yang. Relaxed wasserstein with applications to gans. *arXiv preprint arXiv:1705.07164*, 2017.
- [15] Swaminathan Gurumurthy, Ravi Kiran Sarvadevabhatla, and Venkatesh Babu Radhakrishnan. Deligan: Generative adversarial networks for diverse and limited data. In *CVPR*, 2017.
- [16] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning-lecture 6a-overview of mini-batch gradient descent. lecture notes, 2012.
- [17] Quan Hoang, Tu Dinh Nguyen, Trung Le, and Dinh Phung. Multi-generator gernerative adversarial nets. *arXiv preprint arXiv:1708.02556*, 2017.
- [18] Rui Huang, Shu Zhang, Tianyu Li, and Ran He. Beyond face rotation: Global and local perception gan for photorealistic and identity preserving frontal view synthesis. In *International Conference on Computer Vision*, 2017.
- [19] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics*, 36(4):107, 2017.
- [20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *arXiv preprint arXiv:1611.07004*, 2016.
- [21] L. V. Kantorovich. On a problem of Monge. *Uspekhi Mat. Nauk.*, 3:225–226, 1948.
- [22] Taeksoo Kim, Moon-su Cha, Hyunsoo Kim, Jungkwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. *arXiv preprint arXiv:1703.05192*, 2017.
- [23] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *Preprint 1312.6114, ArXiv*, 2013.

- [24] Abhishek Kumar, Prasanna Sattigeri, and P Thomas Fletcher. Improved semi-supervised learning with gans using manifold invariances. *arXiv preprint arXiv:1705.08850*, 2017.
- [25] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. *arXiv preprint arXiv:1512.09300*, 2015.
- [26] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. *arXiv preprint arXiv:1609.04802*, 2016.
- [27] Jianan Li, Xiaodan Liang, Yunchao Wei, Tingfa Xu, Jiashi Feng, and Shuicheng Yan. Perceptual generative adversarial networks for small object detection. *arXiv preprint arXiv:1706.05274*, 2017.
- [28] Yijun Li, Sifei Liu, Jimei Yang, and Ming-Hsuan Yang. Generative face completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [29] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. *arXiv preprint arXiv:1703.00848*, 2017.
- [30] Ming-Yu Liu and Oncel Tuzel. Coupled generative adversarial networks. In *Advances in Neural Information Processing Systems*, pages 469–477, 2016.
- [31] Pauline Luc, Camille Couprie, Soumith Chintala, and Jakob Verbeek. Semantic segmentation using adversarial networks. *arXiv preprint arXiv:1611.08408*, 2016.
- [32] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015.
- [33] Augustus Odena. Semi-supervised learning with generative adversarial networks. *arXiv preprint arXiv:1606.01583*, 2016.
- [34] Eunbyung Park, Jimei Yang, Ersin Yumer, Duygu Ceylan, and Alexander C Berg. Transformation-grounded image generation network for novel 3d view synthesis. *arXiv preprint arXiv:1703.02921*, 2017.
- [35] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2536–2544, 2016.
- [36] Guim Perarnau, Joost van de Weijer, Bogdan Raducanu, and Jose M Álvarez. Invertible conditional gans for image editing. *arXiv preprint arXiv:1611.06355*, 2016.
- [37] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [38] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pages 2234–2242, 2016.
- [39] Wei Shen and Rujie Liu. Learning residual images for face attribute manipulation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

- [40] Zhixin Shu, Ersin Yumer, Sunil Hadap, Kalyan Sunkavalli, Eli Shechtman, and Dimitris Samaras. Neural face editing with intrinsic image disentangling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [41] Jost Tobias Springenberg. Unsupervised and semi-supervised learning with categorical generative adversarial networks. In *International Conference on Learning Representations*, 2016.
- [42] Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200*, 2016.
- [43] Ilya Tolstikhin, Sylvain Gelly, Olivier Bousquet, Carl-Johann Simon-Gabriel, and Bernhard Schölkopf. Adagan: Boosting generative models. *arXiv preprint arXiv:1701.02386*, 2017.
- [44] Cédric Villani. *Topics in optimal transportation*. Number 58 in Graduate Studies in Mathematics. American Mathematical Society, Providence, RI, 2003.
- [45] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- [46] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *Advances In Neural Information Processing Systems*, pages 613–621, 2016.
- [47] Xiaolong Wang, Abhinav Shrivastava, and Abhinav Gupta. A-fast-rcnn: Hard positive generation via adversary for object detection. *arXiv preprint arXiv:1704.03414*, 2017.
- [48] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Advances in Neural Information Processing Systems*, pages 82–90, 2016.
- [49] Raymond Yeh, Chen Chen, Teck Yian Lim, Mark Hasegawa-Johnson, and Minh N Do. Semantic image inpainting with perceptual and contextual losses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [50] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126*, 2016.
- [51] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint arXiv:1703.10593*, 2017.
- [52] Wentao Zhu and Xiaohui Xie. Adversarial deep structural networks for mammographic mass segmentation. *arXiv preprint arXiv:1612.05970*, 2016.