

# A Geometric View to Optimal Transportation and Generative Model

David Xianfeng Gu<sup>1</sup>

<sup>1</sup>Computer Science & Applied Mathematics  
SUNY at Stony Brook University

Center of Mathematical Sciences and Applications  
Harvard University

Geometric Computation and Applications  
Trinity College, Dublin, Ireland

# Thanks

Thanks for the invitation.

These projects are collaborated with Shing-Tung Yau, Feng Luo, Zhongxuan Luo, Na Lei, Dimitris Samaras and so on.

- 1 Why does DL work?
- 2 How to quantify the learning capability of a DNN?
- 3 How does DL manipulate the probability distributions?



# Why dose DL work?

Deep learning is the mainstream technique for many machine learning tasks, including image recognition, machine translation, speech recognition, and so on. Despite its success, the theoretical understanding on how it works remains primitive.

# Manifold Assumption

We believe the great success of deep learning can be partially explained by the well accepted manifold assumption and the clustering assumption:

## Manifold Assumption

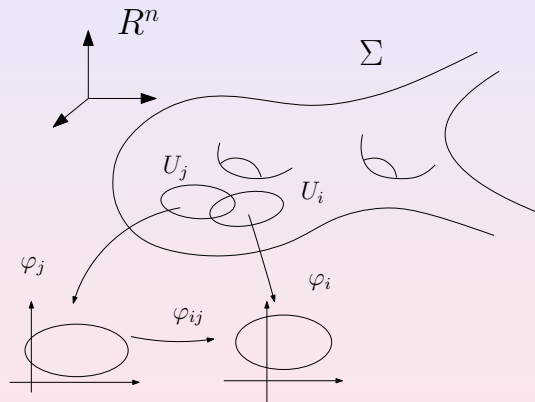
Natural high dimensional data concentrates close to a non-linear low-dimensional manifold.

## Clustering Assumption

The distances among the probability distributions of subclasses on the manifold are far enough to discriminate them.

Deep learning method can learn and represent the manifold structure, and transform the probability distributions.

# General Model



- Ambient Space-image space  $\mathbb{R}^n$
- manifold - Support of a distribution  $\mu$
- parameter domain - latent space  $\mathbb{R}^m$
- coordinates map  $\varphi_i$ -encoding/decoding maps
- $\varphi_{ij}$  controls the probability measure

## Definition (Manifold)

Suppose  $M$  is a topological space, covered by a set of open sets  $M \subset \bigcup_{\alpha} U_{\alpha}$ . For each open set  $U_{\alpha}$ , there is a homeomorphism  $\varphi_{\alpha} : U_{\alpha} \rightarrow \mathbb{R}^n$ , the pair  $(U_{\alpha}, \varphi_{\alpha})$  form a chart. The union of charts form an atlas  $\mathcal{A} = \{(U_{\alpha}, \varphi_{\alpha})\}$ . If  $U_{\alpha} \cap U_{\beta} \neq \emptyset$ , then the chart transition map is given by

$$\varphi_{\alpha\beta} : \varphi_{\alpha}(U_{\alpha} \cap U_{\beta}) \rightarrow \varphi_{\beta}(U_{\alpha} \cap U_{\beta}),$$

$$\varphi_{\alpha\beta} := \varphi_{\beta} \circ \varphi_{\alpha}^{-1}.$$

# Example



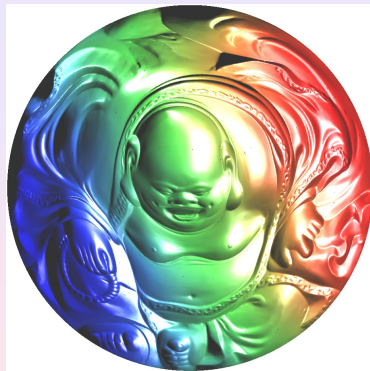
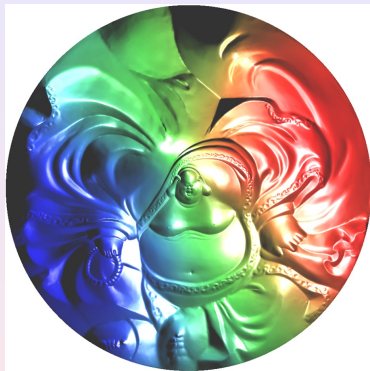
Image space  $\mathcal{X}$  is  $\mathbb{R}^3$ ; the data manifold  $\Sigma$  is the happy buddha.

# Example



The encoding map is  $\varphi_i : \Sigma \rightarrow \mathcal{Z}$ ; the decoding map is  $\varphi_i^{-1} : \mathcal{Z} \rightarrow \Sigma$ .

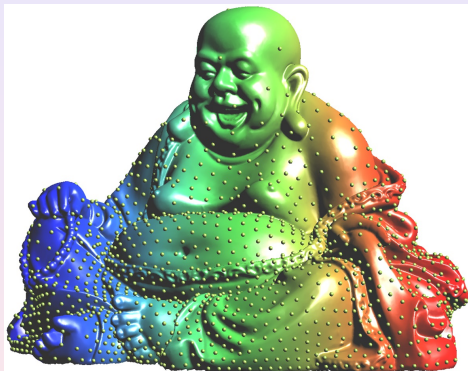
# Example



The automorphism of the latent space  $\varphi_{ij} : \mathcal{Z} \rightarrow \mathcal{Z}$  is the chart transition.

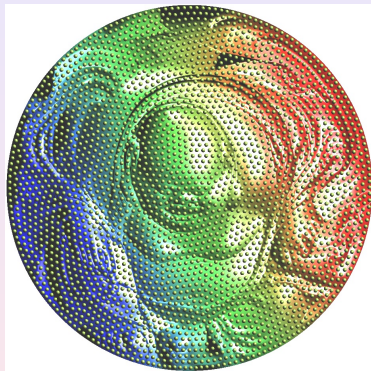
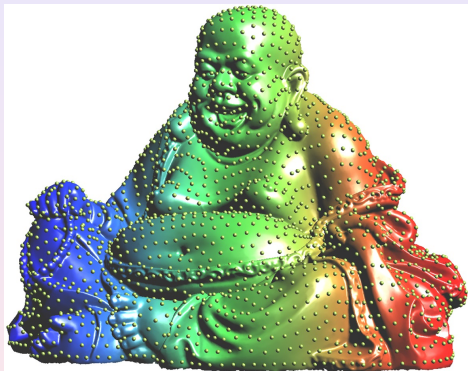


# Example



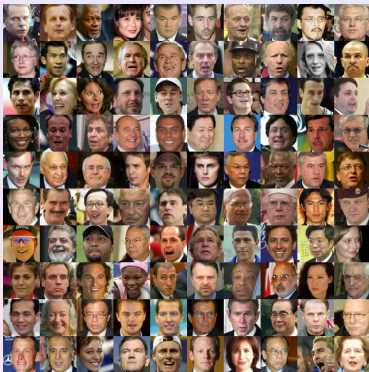
Uniform distribution  $\zeta$  on the latent space  $\mathcal{Z}$ , non-uniform distribution on  $\Sigma$  produced by a decoding map.

# Example



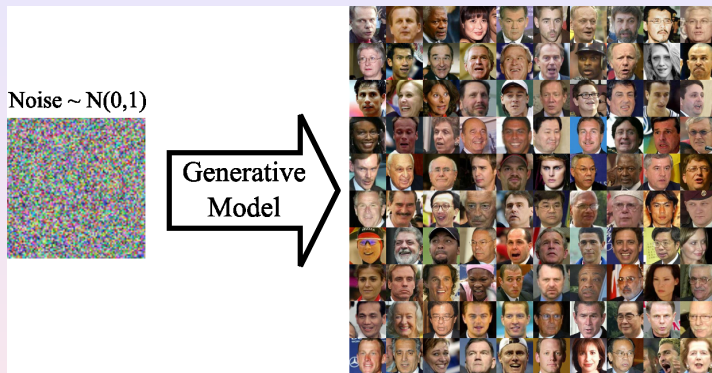
Uniform distribution  $\zeta$  on the latent space  $\mathcal{Z}$ , uniform distribution on  $\Sigma$  produced by another decoding map.

# Human Facial Image Manifold



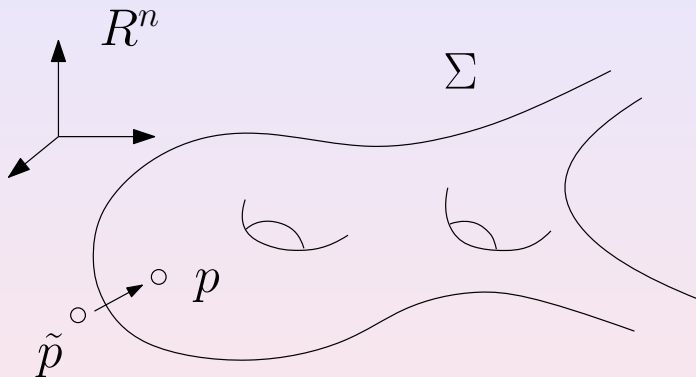
One facial image is determined by a finite number of genes, lighting conditions, camera parameters, therefore all facial images form a manifold.

# Manifold view of Generative Model



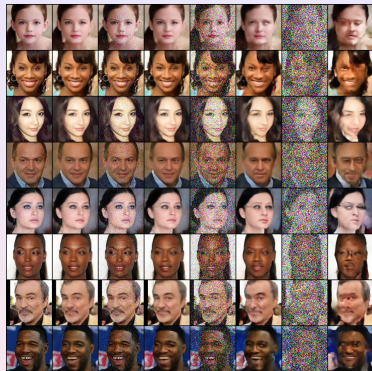
Given a parametric representation  $\phi : \mathcal{Z} \rightarrow \Sigma$ , randomly generate a parameter  $z \in \mathcal{Z}$  (white noise),  $\phi(z) \in \Sigma$  is a human facial image.

# Manifold view of Denoising



Suppose  $\tilde{p}$  is a point close to the manifold,  $p \in \Sigma$  is the closest point of  $\tilde{p}$ . The projection  $\tilde{p} \rightarrow p$  can be treated as denoising.

# Manifold view of Denoising



$\Sigma$  is the clean facial image manifold; noisy image  $\tilde{p}$  is a point close to  $\Sigma$ ; the closest point  $p \in \Sigma$  is the resulting denoised image.

# Manifold view of Denoising

## Traditional Method

Fourier transform the noisy image, filter out the high frequency component, inverse Fourier transform back to the denoised image.

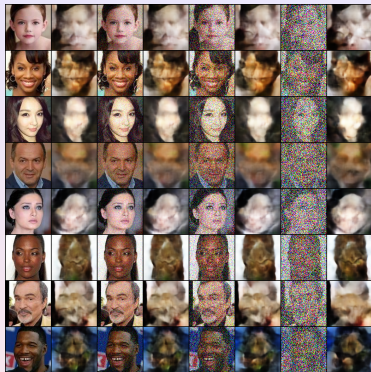
## ML Method

Use the clean facial images to train the neural network, obtain a representation of the manifold. Project the noisy image to the manifold, the projection point is the denoised image.

## Key Difference

Traditional method is independent of the content of the image; ML method heavily depends on the content of the image. The prior knowledge is encoded by the manifold.

# Manifold view of Denoising



If the wrong manifold is chosen, the denoising result is of non-sense. Here we use the cat face manifold to denoise a human face image, the result looks like a cat face.



# How dose DL learn a manifold?

The central tasks for Deep Learning are

- 1 Learn the manifold structure from the data;
- 2 Represent the manifold implicitly or explicitly.

# Autoencoder

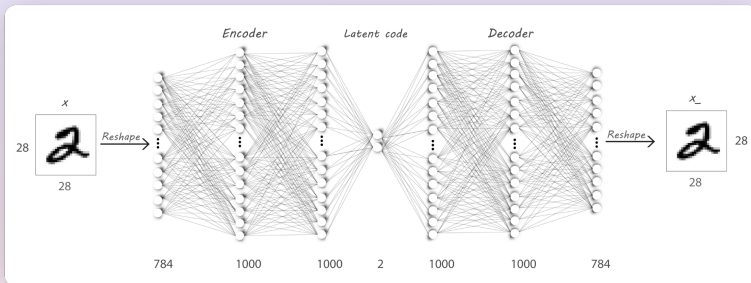


Figure: Auto-encoder architecture.

Ambient space  $\mathcal{X}$ , latent space  $\mathcal{Z}$ , encoding map  $\varphi_{\theta} : \mathcal{X} \rightarrow \mathcal{Z}$ , decoding map  $\psi_{\theta} : \mathcal{Z} \rightarrow \mathcal{X}$ .

# Autoencoder

The encoder takes a sample  $\mathbf{x} \in \mathcal{X}$  and maps it to  $\mathbf{z} \in \mathcal{F}$ ,  $\mathbf{z} = \varphi(\mathbf{x})$ . The decoder  $\psi : \mathcal{F} \rightarrow \mathcal{X}$  maps  $\mathbf{z}$  to the reconstruction  $\tilde{\mathbf{x}}$ .

$$\begin{array}{ccc} \{(\mathcal{X}, \mathbf{x}), \mu, M\} & \xrightarrow{\varphi} & \{(\mathcal{F}, \mathbf{z}), D\} \\ & \searrow \psi \circ \varphi & \downarrow \psi \\ & & \{(\mathcal{X}, \tilde{\mathbf{x}}), \tilde{M}\} \end{array}$$

An autoencoder is trained to minimise reconstruction errors:

$$\varphi, \psi = \operatorname{argmin}_{\varphi, \psi} \int_{\mathcal{X}} \mathcal{L}(\mathbf{x}, \psi \circ \varphi(\mathbf{x})) d\mu(\mathbf{x}),$$

where  $\mathcal{L}(\cdot, \cdot)$  is the loss function, such as squared errors. The reconstructed manifold  $\tilde{M} = \psi \circ \varphi(M)$  is used as an approximation of  $M$ .

## Definition (ReLU DNN)

For any number of hidden layers  $k \in \mathbb{N}$ , input and output dimensions  $w_0, w_{k+1} \in \mathbb{N}$ , a  $\mathbb{R}^{w_0} \rightarrow \mathbb{R}^{w_{k+1}}$  ReLU DNN is given by specifying a sequence of  $k$  natural numbers  $w_1, w_2, \dots, w_k$  representing widths of the hidden layers, a set of  $k$  affine transformations  $T_i : \mathbb{R}^{w_{i-1}} \rightarrow \mathbb{R}^{w_i}$  for  $i = 1, \dots, k$  and a linear transformation  $T_{k+1} : \mathbb{R}^{w_k} \rightarrow \mathbb{R}^{w_{k+1}}$  corresponding to weights of hidden layers.

The mapping  $\varphi_\theta : \mathbb{R}^{w_0} \rightarrow \mathbb{R}^{w_{k+1}}$  represented by this ReLU DNN is

$$\varphi = T_{k+1} \circ \sigma \circ T_k \circ \dots \circ T_2 \circ \sigma \circ T_1, \quad (1)$$

where  $\circ$  denotes mapping composition,  $\theta$  represent all the weight and bias parameters.

# Activated Path

Fix the encoding map  $\varphi_\theta$ , let the set of all neurons in the network is denoted as  $\mathcal{S}$ , all the subsets is denoted as  $2^{\mathcal{S}}$ .

## Definition (Activated Path)

Given a point  $\mathbf{x} \in \mathcal{X}$ , the *activated path* of  $\mathbf{x}$  consists all the activated neurons when  $\varphi_\theta(\mathbf{x})$  is evaluated, and denoted as  $\rho(\mathbf{x})$ . Then the activated path defines a set-valued function  $\rho : \mathcal{X} \rightarrow 2^{\mathcal{S}}$ .

# Cell Decomposition

## Definition (Cell Decomposition)

Fix an encoding map  $\varphi_\theta$  represented by a ReLU RNN, two data points  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$  are *equivalent*, denoted as  $\mathbf{x}_1 \sim \mathbf{x}_2$ , if they share the same activated path,  $\rho(\mathbf{x}_1) = \rho(\mathbf{x}_2)$ . Then each equivalence relation partitions the ambient space  $\mathcal{X}$  into cells,

$$\mathcal{D}(\varphi_\theta) : \mathcal{X} = \bigcup_{\alpha} U_{\alpha},$$

each equivalence class corresponds to a cell:  $\mathbf{x}_1, \mathbf{x}_2 \in U_{\alpha}$  if and only if  $\mathbf{x}_1 \sim \mathbf{x}_2$ .  $\mathcal{D}(\varphi_\theta)$  is called the cell decomposition induced by the encoding map  $\varphi_\theta$ .

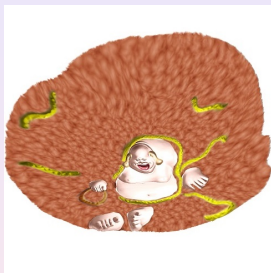
Furthermore,  $\varphi_\theta$  maps the cell decomposition in the ambient space  $\mathcal{D}(\varphi_\theta)$  to a cell decomposition in the latent space.

# Encoding/Decoding



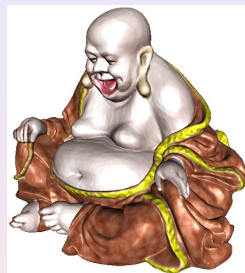
a. Input manifold

$$M \subset \mathcal{X}$$



b. latent representation

$$D = \varphi_{\theta}(M)$$



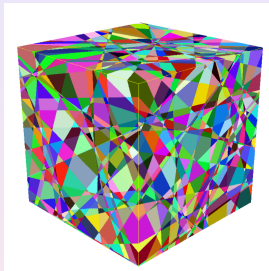
c. reconstructed manifold

$$\tilde{M} = \psi_{\theta}(D)$$

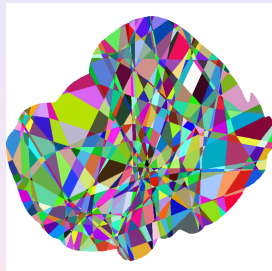
Figure: Auto-encoder pipeline.



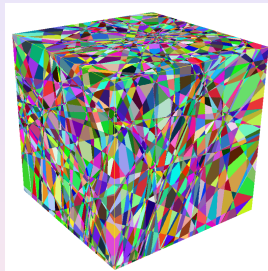
# Piecewise Linear Mapping



d. cell decomposition  
 $\mathcal{D}(\varphi_\theta)$



e. latent space  
cell decomposition



f. cell decomposition  
 $\mathcal{D}(\psi_\theta \circ \varphi_\theta)$

Piecewise linear encoding/decoding maps induce cell decompositions of the ambient space and the latent space.

## Definition (Rectified Linear Complexity of a ReLU DNN)

Given a ReLU DNN  $N(w_0, \dots, w_{k+1})$ , its rectified linear complexity is the upper bound of the number of pieces of all PL functions  $\varphi_\theta$  represented by  $N$ ,

$$\mathcal{N}(N) := \max_{\theta} \mathcal{N}(\varphi_\theta).$$

Rectified Linear complexity gives a measurement for the representation capability of a neural network.

# RL Complexity Estimate

## Lemma

*The maximum number of parts one can get when cutting  $d$ -dimensional space  $\mathbb{R}^d$  with  $n$  hyperplanes is denoted as  $\mathcal{C}(d, n)$ , then*

$$\mathcal{C}(d, n) = \binom{n}{0} + \binom{n}{1} + \binom{n}{2} + \cdots + \binom{n}{d}. \quad (2)$$

## Proof.

Suppose  $n$  hyperplanes cut  $\mathbb{R}^d$  into  $\mathcal{C}(d, n)$  cells, each cell is a convex polyhedron. The  $(n+1)$ -th hyperplane is  $\pi$ , then the first  $n$  hyperplanes intersection  $\pi$  and partition  $\pi$  into  $\mathcal{C}(d-1, n)$  cells, each cell on  $\pi$  partitions a polyhedron in  $\mathbb{R}^d$  into 2 cells, hence we get the formula

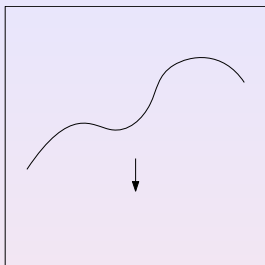
$$\mathcal{C}(d, n+1) = \mathcal{C}(d, n) + \mathcal{C}(d-1, n).$$

## Theorem (Rectified Linear Complexity of a ReLU DNN)

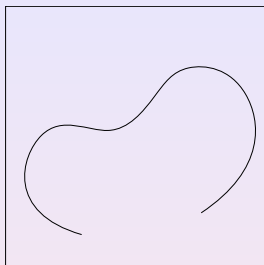
*Given a ReLU DNN  $N(w_0, \dots, w_{k+1})$ , representing PL mappings  $\varphi_\theta : \mathbb{R}^{w_0} \rightarrow \mathbb{R}^{w_{k+1}}$  with  $k$  hidden layers of widths  $\{w_i\}_{i=1}^k$ , then the linear rectified complexity of  $N$  has an upper bound,*

$$\mathcal{N}(N) \leq \prod_{i=1}^{k+1} \mathcal{C}(w_{i-1}, w_i). \quad (3)$$

# RL Complexity of Manifold



a. linear rectifiable



b. non-linear-rectifiable

## Definition (Linear Rectifiable Manifold)

Suppose  $M$  is a  $m$ -dimensional manifold, embedded in  $\mathbb{R}^n$ , we say  $M$  is linear rectifiable, if there exists an affine map  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , such that the restriction of  $\varphi$  on  $M$ ,  $\varphi|_M : M \rightarrow \varphi(M) \subset \mathbb{R}^m$ , is homeomorphic.  $\varphi$  is called the corresponding rectified linear map of  $M$ .

## Definition (Linear Rectifiable Atlas)

Suppose  $M$  is a  $m$ -dimensional manifold, embedded in  $\mathbb{R}^n$ ,  $\mathcal{A} = \{(U_\alpha, \varphi_\alpha)\}$  is an atlas of  $M$ . If each chart  $(U_\alpha, \varphi_\alpha)$  is linear rectifiable,  $\varphi_\alpha : U_\alpha \rightarrow \mathbb{R}^m$  is the rectified linear map of  $U_\alpha$ , then the atlas is called a linear rectifiable atlas of  $M$ .

## Definition (Rectified Linear Complexity of a Manifold)

Suppose  $M$  is a  $m$ -dimensional manifold embedded in  $\mathbb{R}^n$ , the rectified linear complexity of  $M$  is denoted as  $\mathcal{N}(\mathbb{R}^n, M)$  and defined as,

$$\mathcal{N}(\mathbb{R}^n, M) := \min \{ |\mathcal{A}| \mid \mathcal{A} \text{ is a linear rectifiable atlas of } M \}. \quad (4)$$

# Encodable Condition

## Definition (Encoding Map)

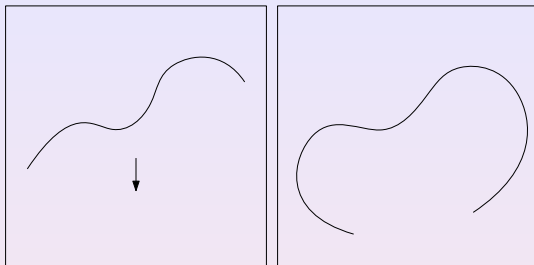
Suppose  $M$  is a  $m$ -dimensional manifold, embedded in  $\mathbb{R}^n$ , a continuous mapping  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is called an encoding map of  $(\mathbb{R}^n, M)$ , if restricted on  $M$ ,  $\varphi|_M : M \rightarrow \varphi(M) \subset \mathbb{R}^m$  is homeomorphic.

## Theorem (Encodable Condition)

*Suppose a ReLU DNN  $N(w_0, \dots, w_{k+1})$  represents a PL mapping  $\varphi_\theta : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $M$  is a  $m$ -dimensional manifold embedded in  $\mathbb{R}^n$ . If  $\varphi_\theta$  is an encoding mapping of  $(\mathbb{R}^n, M)$ , then the rectified linear complexity of  $N$  is no less than the rectified linear complexity of  $(\mathbb{R}^n, M)$ ,*

$$\mathcal{N}(\mathbb{R}^n, M) \leq \mathcal{N}(\varphi_\theta) \leq \mathcal{N}(N).$$

# Encodable Condition



## Lemma

Suppose a  $n$  dimensional manifold  $M$  is embedded in  $\mathbb{R}^{n+1}$ ,

$$M \xrightarrow{G} \mathbb{S}^n \xrightarrow{p} \mathbb{RP}^n$$

where  $G : M \rightarrow \mathbb{S}^n$  is the Gauss map,  $\mathbb{RP}^n$  is the real projective space, if  $p \circ G(M)$  covers the whole  $\mathbb{RP}^n$ , then  $M$  is not linear rectifiable.



# Representation Limitation Theorem

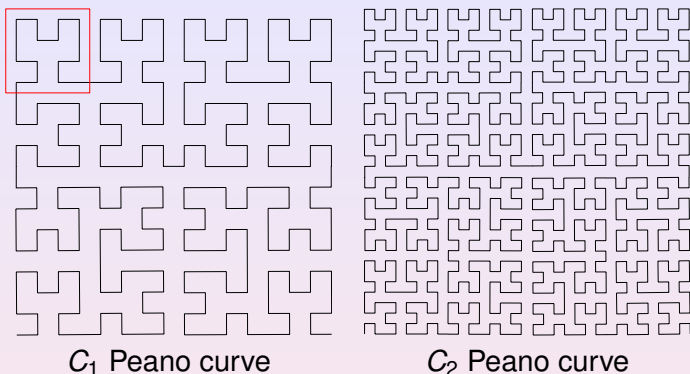


Figure:  $\mathcal{N}(\mathbb{R}^2, C_n) \geq 4^{n+1}$

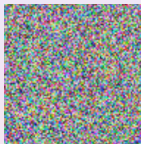
## Theorem

*Given any ReLU deep neural network  $N(w_0, w_1, \dots, w_k, w_{k+1})$ , there is a manifold  $M$  embedded in  $\mathbb{R}^{w_0}$ , such that  $M$  can not be encoded by  $N$ .*

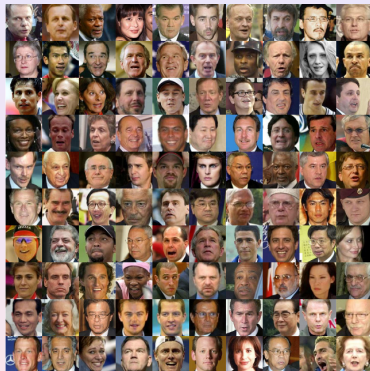
# How does DL control the probability distribution?

# Generative Model

Noise  $\sim N(0,1)$



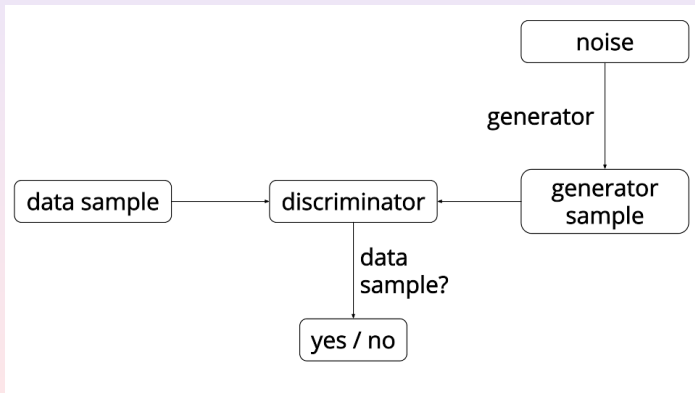
Generative  
Model



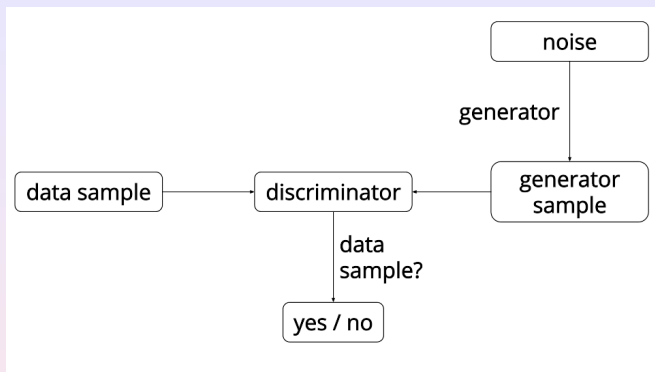
A generative model converts a white noise into a facial image.

# GAN Overview

The analogy that is often used here is that the generator is like a forger trying to produce some counterfeit material, and the discriminator is like the police trying to detect the forged items.



# GAN Overview



## Merits

- 1 Automatic generate samples, the requirement for the data samples is reduced;
- 2 Data sample distribution can be arbitrary, without closed form expression.

# GAN Overview

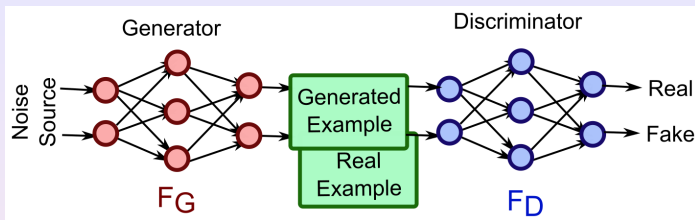


Figure: GAN DNN model.

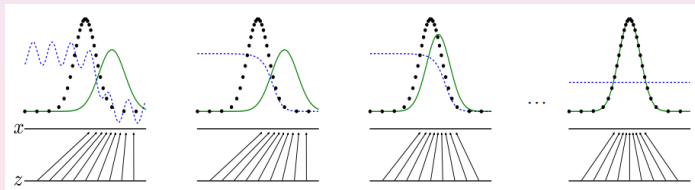
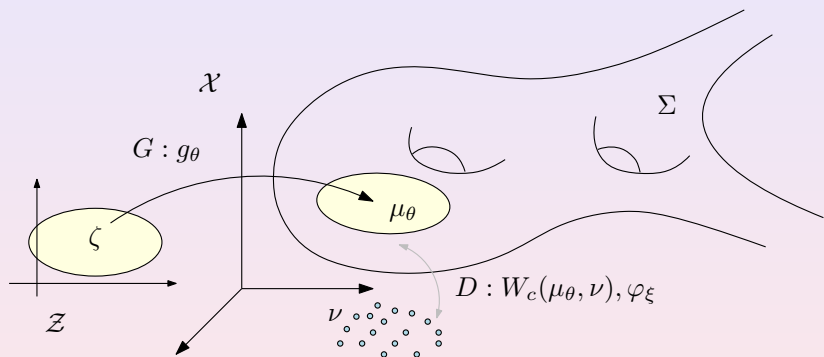


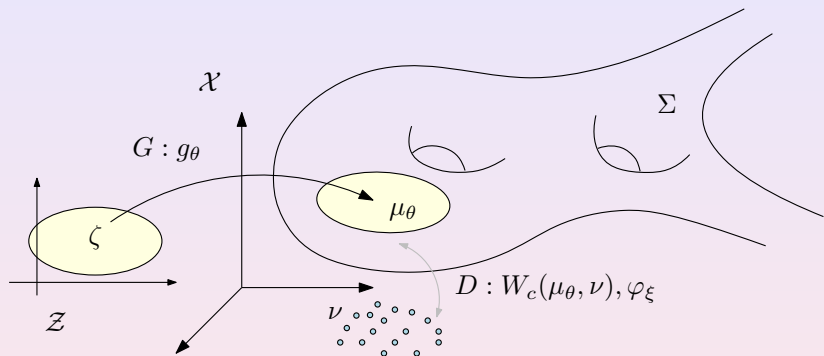
Figure: GAN learning process.

# Wasserstein GAN Model



$\mathcal{X}$ -image space;  $\Sigma$ -supporting manifold;  $\mathcal{Z}$ -latent space;

# Wasserstein GAN Model



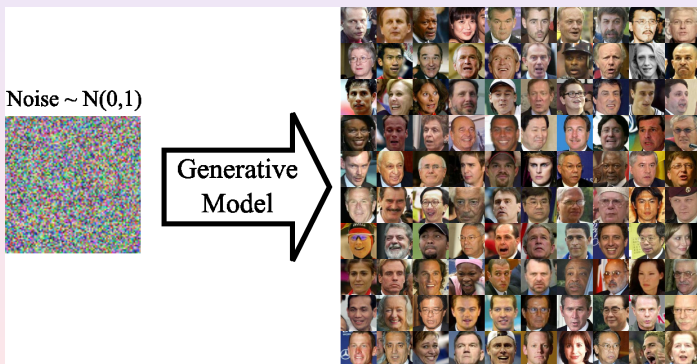
$\nu$ -training data distribution;  $\zeta$ -uniform distribution;  
 $\mu_\theta = g_{\theta\#} \zeta$ -generated distribution;  $G$  - generator computes  $g_\theta$ ;  
 $D$  - discriminator, measures the distance between  $\nu$  and  $\mu_\theta$ ,  
 $W_c(\mu_\theta, \nu)$ .



# Generative Model

## Generative Model

$G: \mathcal{Z} \rightarrow \mathcal{X}$  maps a fixed probability distribution  $\zeta$  to the training data probability distribution  $\nu$ .



## Wasserstein Space

Given a Riemannian manifold  $M$ , all the probability distributions on  $M$  form an infinite dimensional manifold Wasserstein space  $\mathcal{W}(M)$ , the distance between two probability distributions is given by the so-called Wasserstein distance.

## Optimal Mass Transportation

Given two probability measures  $\mu, \nu \in \mathcal{W}(M)$ , there is a unique optimal mass transportation map  $T : M \rightarrow M$ ,  $\phi$  maps  $\mu$  to  $\nu$  with the minimal transportation cost. The transportation cost of the optimal transportation map is the Wasserstein distance between  $\mu$  and  $\nu$ .

## Definition (Measure-Preserving Mapping)

Given two bounded domains in  $\mathbb{R}^n$  with probability measures  $(X, \mu)$  and  $(Y, \nu)$ , with equal total measure  $\mu(X) = \nu(Y)$ , a transportation mapping  $T : X \rightarrow Y$  is measure-preserving, if for any measurable set  $B \subset Y$ ,

$$\int_{T^{-1}(B)} d\mu(x) = \int_B d\nu(y),$$

and denoted as  $T_{\#}\mu = \nu$ .

Suppose  $T$  is a smooth map, then measure-preserving condition is equivalent to the Jacobian equation

$$\mu(x)dx = \nu(y)dy$$

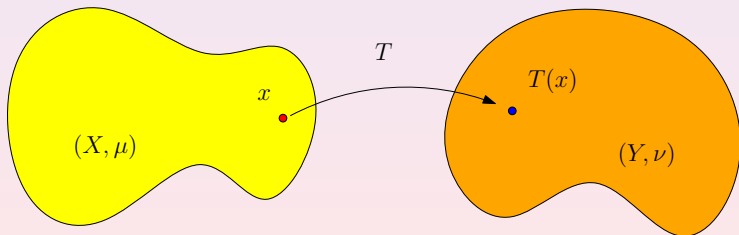
$$\det(DT) = \frac{\mu(x)}{\nu \circ T(x)}.$$

# Optimal Mass Transportation

## Definition (Transportation Cost)

Suppose the cost of moving a unit mass from point  $x$  to point  $y$  is  $c(x, y)$ , for a transportation map  $T : (X, \mu) \rightarrow (Y, \nu)$ , the total transportation cost is

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



# Cost Function $c(x, y)$

The cost of moving a unit mass from point  $x$  to point  $y$ .

$$\text{Monge}(1781) : c(x, y) = |x - y|.$$

This is the natural cost function. Other cost functions include

$$c(x, y) = |x - y|^p, p \neq 0$$

$$c(x, y) = -\log |x - y|$$

$$c(x, y) = \sqrt{\varepsilon + |x - y|^2}, \varepsilon > 0$$

Any function can be cost function. It can be negative.

# Monge Problem

## Problem (Monge)

*Find a measure-preserving transportation map  $T : (X, \mu) \rightarrow (Y, \nu)$  that minimizes the transportation cost,*

$$(MP) \quad \min_{T_{\#}\mu=\nu} \mathcal{C}(T) = \min_{T_{\#}\mu=\nu} \int_X c(x, T(x)) d\mu(x).$$

*such kind of map is called the optimal mass transportation map.*

## Definition (Wasserstein distance)

The transportation cost of the optimal transportation map  $T : (X, \mu) \rightarrow (Y, \nu)$  is called the Wasserstein distance between  $\mu$  and  $\nu$ , denoted as

$$W_c(\mu, \nu) := \min_{T_{\#}\mu=\nu} \mathcal{C}(T).$$

# Kantorovich Problem

Kantorovich relaxed transportation maps to transportation schemes.

## Problem (Kantorovich)

*Find an optimal transportation scheme, namely a joint probability measure  $\rho \in \mathcal{P}(X \times Y)$ , with marginal measures  $\rho_{x\#} = \mu$ ,  $\rho_{y\#} = \nu$ , that minimizes the transportation cost,*

$$(KP) \quad \min_{\rho} \left\{ \int_{X \times Y} c(x, y) d\rho(x, y) \mid \rho_{x\#} = \mu, \rho_{y\#} = \nu \right\}.$$

Kantorovich solved this problem by inventing linear programming, and won Nobel's prize in economics in 1975.

# Kantorovich Dual Problem

By the duality of linear programming, Kantorovich problem has the dual form:

## Problem (Kantorovich Dual)

*Find an functions  $\varphi : X \rightarrow \mathbb{R}$  and  $\psi : Y \rightarrow \mathbb{R}$ , such that*

$$(DP) \max_{\varphi, \psi} \left\{ \int_X \varphi(x) du(x) + \int_Y \psi(y) dv(y), \varphi(x) + \psi(y) \leq c(x, y) \right\}.$$



# Kantorovich Dual Problem

## Definition (c-transformation)

Given a function  $\varphi : X \rightarrow \mathbb{R}$ , and  $c(x, y) : X \times Y \rightarrow \mathbb{R}$ , its c-transform  $\varphi^c : Y \rightarrow \mathbb{R}$  is given by

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

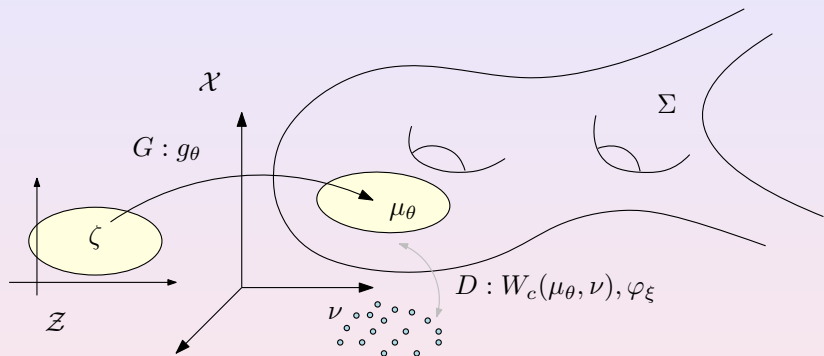
## Problem (Kantorovich Dual)

*The Kantorovich Dual problem can be reformulated as*

$$(DP) \quad \max_{\varphi} \left\{ \int_X \varphi(x) du(x) + \int_Y \varphi^c(y) dv(y) \right\}.$$

$\varphi$  is called Kantorovich potential.

# Wasserstein GAN Model



$\nu$ -training data distribution;  $\zeta$ -uniform distribution;  
 $\mu_\theta = g_{\theta\#} \zeta$ -generated distribution;  $G$  - generator computes  $g_\theta$ ;  
 $D$  -discriminator, measures the distance between  $\nu$  and  $\mu_\theta$ ,  
 $W_c(\mu_\theta, \nu)$ .

From the optimal transportation point of view, Wasserstein GAN performs the following tasks:

- The discriminator: computes the Wasserstein distance using Kantorovich Dual formula:

$$W_c(\mu_\theta, \nu) = \max_{\varphi_\xi} \int_X \varphi_\xi(x) d\mu_\theta(x) + \int_Y \varphi_\xi^c(y) d\nu(y),$$

namely computes the Kantorovich potential  $\varphi$ ;

- The generator: computes a measure-preserving transportation map  $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$ , s.t.  $g_{\theta\#}\zeta = \mu_\theta = \nu$ .
- The WGAN model: min-max optimization

$$\min_{\theta} \max_{\xi} \int_X \varphi_\xi \circ g_\theta(z) d\zeta(z) + \int_Y \varphi_\xi^c(y) d\nu(y)$$

## $L^1$ case

When  $c(x, y) = |x - y|$ ,  $\varphi_c = -\varphi$ , given  $\varphi$  is 1-Lipsitz, the WGAN model: min-max optimization

$$\min_{\theta} \max_{\xi} \int_X \varphi_{\xi} \circ g_{\theta}(z) d\zeta(z) - \int_Y \varphi_{\xi}(y) dv(y).$$

namely

$$\min_{\theta} \max_{\xi} \mathbb{E}_{z \sim \zeta}(\varphi_{\xi} \circ g_{\theta}(z)) - \mathbb{E}_{y \sim v}(\varphi_{\xi}(y)).$$

with the constraint that  $\varphi_{\xi}$  is 1-Lipsitz.

# Brenier's Approach

## Theorem (Brenier)

*If  $\mu, \nu > 0$  and  $X$  is convex, and the cost function is quadratic distance,*

$$c(\mathbf{x}, \mathbf{y}) = |\mathbf{x} - \mathbf{y}|^2$$

*then there exists a convex function  $u : X \rightarrow \mathbb{R}$  unique upto a constant, such that the unique optimal transportation map is given by the gradient map*

$$T : \mathbf{x} \rightarrow \nabla u(\mathbf{x}).$$

## Problem (Brenier)

*Find a convex function  $u : X \rightarrow \mathbb{R}$ , such that*

$$(BP) \quad (\nabla u)_{\#} \mu = \nu,$$

*$u$  is called the Brenier potential.*

# Brenier's Approach

From Jacobian equation, one can get the necessary condition for Brenier potential.

## Problem (Brenier)

*Find the Brenier potential  $u : X \rightarrow \mathbb{R}$  satisfies the Monge-Ampere equation*

$$(BP) \quad \det \left( \frac{\partial^2 u}{\partial x_i \partial x_j} \right) = \frac{\mu(\mathbf{x})}{v(\nabla f(\mathbf{x}))}.$$

# Kantorovich and Brenier potentials

## Theorem

*If the distance function  $c(x, y) = h(x - y)$ , where  $h : \mathbb{R} \rightarrow \mathbb{R}$  is a strictly convex function, the Kantorovich potential  $\phi : X \rightarrow \mathbb{R}$  gives the optimal mass transportation map directly:*

$$T(\mathbf{x}) = \mathbf{x} - (\nabla \phi)^{-1}(\nabla \phi(\mathbf{x}))$$

## Corollary

*Suppose  $c(x, y) = \frac{1}{2}|x - y|^2$ , then Kantorovich potential and Brenier potential satisfy the relation*

$$u(\mathbf{x}) = \frac{1}{2}|\mathbf{x}|^2 - \phi(\mathbf{x}).$$

## $L^2$ case

The discriminator computes the Kantorovich potential  $\varphi$ ; the generator  $G$  computes the optimal transportation map,  $T = \nabla u$ , where  $u$  is the Brenier potential; The Brenier potential equals to

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

Hence, in theory:

- $G$  can be obtained from the optimal  $D$  without training;
- $D$  can be obtained from the optimal  $G$  without training;
- The two deep neural networks are redundant;
- The competition between  $D$  and  $G$  is unnecessary.



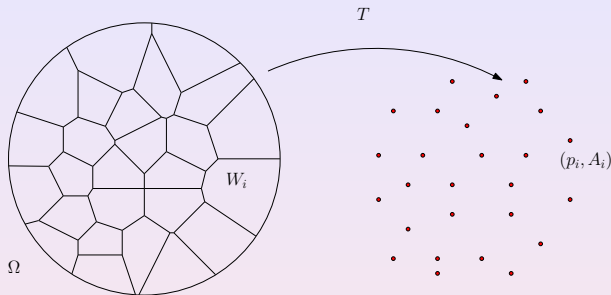
## Empirical Distribution

In practice, the target probability measure is approximated by empirical distribution:

$$\nu = \sum_{i=1}^n \delta(y - y_i) \nu_i,$$

in general  $\nu_i = 1/n$ .

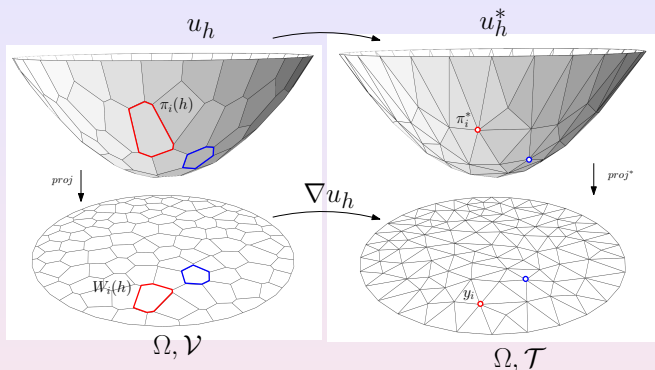
# Semi-discrete Optimal Transportation



Given a compact convex domain  $\Omega$  in  $\mathbb{R}^n$  and  $p_1, \dots, p_k$  in  $\mathbb{R}^n$  and  $A_1, \dots, A_k > 0$ , find a transport map  $T : U \rightarrow \{p_1, \dots, p_k\}$  with  $\text{vol}(T^{-1}(p_i)) = A_i$ , so that  $T$  minimizes the transport cost

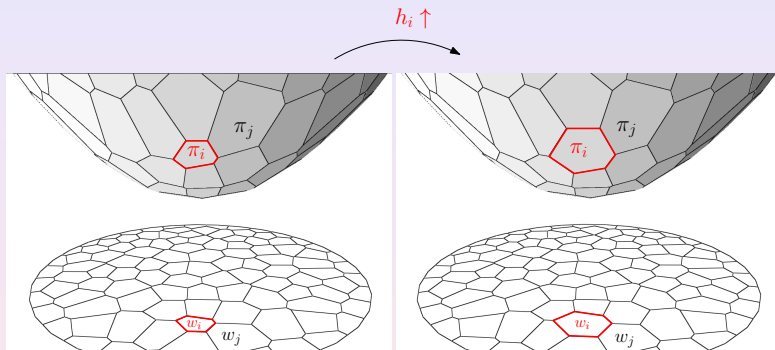
$$\frac{1}{2} \int_U |\mathbf{x} - T(\mathbf{x})|^2 d\mathbf{x}.$$

# Power Diagram vs Optimal Transport Map



- 1  $\forall y_i \in Y$ , construct a hyper-plane  $\pi_h^i(x) = \langle x, y_i \rangle - h_i$ ;
- 2 compute the upper envelope of the planes  
 $u_h(x) = \max_i \{ \pi_h^i(x) \}$
- 3 produce the power diagram of  $\Omega$ ,  $\mathcal{V}(h) = \cup_i W_i(h)$ ;
- 4 adjust the heights  $h$ , such that  $\mu(W_i(h)) = v_i$ .

# Power Diagram vs Optimal Transport Map



**Figure:** Variation of the  $\mu$ -volume of top-dimensional cells.

Adjust the height of each hyper-plane, such that  $\mu(W_i(h)) = v_i$ .

# Convex Geometry

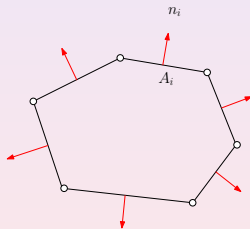
# Minkowski problem - 2D Case

## Example

A convex polygon  $P$  in  $\mathbb{R}^2$  is determined by its edge lengths  $A_i$  and the unit normal vectors  $\mathbf{n}_i$ .

Take any  $\mathbf{u} \in \mathbb{R}^2$  and project  $P$  to  $\mathbf{u}$ , then  $\langle \sum_i A_i \mathbf{n}_i, \mathbf{u} \rangle = 0$ , therefore

$$\sum_i A_i \mathbf{n}_i = \mathbf{0}.$$



# Minkowski problem - General Case

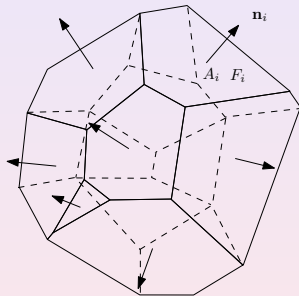
## Minkowski Problem

Given  $k$  unit vectors  $\mathbf{n}_1, \dots, \mathbf{n}_k$  not contained in a half-space in  $\mathbb{R}^n$  and  $A_1, \dots, A_k > 0$ , such that

$$\sum_i A_i \mathbf{n}_i = \mathbf{0},$$

find a compact convex polytope  $P$  with exactly  $k$  codimension-1 faces  $F_1, \dots, F_k$ , such that

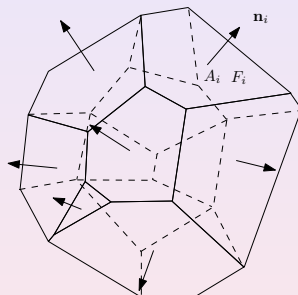
- 1  $\text{area}(F_i) = A_i$ ,
- 2  $\mathbf{n}_i \perp F_i$ .



# Minkowski problem - General Case

## Theorem (Minkowski)

*$P$  exists and is unique up to translations.*





# Brunn-Minkowski inequality

## Theorem (Brunn-Minkowski)

*For every pair of nonempty compact subsets  $A$  and  $B$  of  $\mathbb{R}^n$  and every  $0 \leq t \leq 1$ ,*

$$[\text{Vol}(tA \oplus (1-t)B)]^{\frac{1}{n}} \geq t[\text{vol}(A)]^{\frac{1}{n}} + (1-t)[\text{vol}(B)]^{\frac{1}{n}}.$$

*For convex sets  $A$  and  $B$ , the inequality is strict for  $0 < t < 1$  unless  $A$  and  $B$  are homothetic i.e. are equal up to translation and dilation.*

# Alexandrov Theorem

## Theorem (Alexandrov 1950)

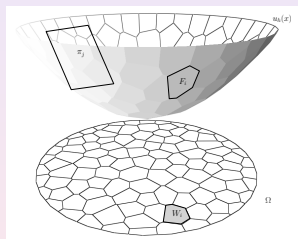
Given  $\Omega$  compact convex domain in  $\mathbb{R}^n$ ,  $p_1, \dots, p_k$  distinct in  $\mathbb{R}^n$ ,  $A_1, \dots, A_k > 0$ , such that  $\sum A_i = \text{Vol}(\Omega)$ , there exists PL convex function

$$f(\mathbf{x}) := \max\{\langle \mathbf{x}, \mathbf{p}_i \rangle + h_i \mid i = 1, \dots, k\}$$

unique up to translation such that

$$\text{Vol}(W_i) = \text{Vol}(\{\mathbf{x} \mid \nabla f(\mathbf{x}) = \mathbf{p}_i\}) = A_i.$$

Alexandrov's proof is topological, not variational. It has been open for years to find a constructive proof.



## Theorem (Gu-Luo-Sun-Yau 2013)

$\Omega$  is a compact convex domain in  $\mathbb{R}^n$ ,  $y_1, \dots, y_k$  distinct in  $\mathbb{R}^n$ ,  $\mu$  a positive continuous measure on  $\Omega$ . For any  $v_1, \dots, v_k > 0$  with  $\sum v_i = \mu(\Omega)$ , there exists a vector  $(h_1, \dots, h_k)$  so that

$$u(\mathbf{x}) = \max\{\langle \mathbf{x}, \mathbf{p}_i \rangle + h_i\}$$

satisfies  $\mu(W_i \cap \Omega) = v_i$ , where  $W_i = \{\mathbf{x} | \nabla f(\mathbf{x}) = \mathbf{p}_i\}$ .

Furthermore,  $\mathbf{h}$  is the maximum point of the convex function

$$E(\mathbf{h}) = \sum_{i=1}^k v_i h_i - \int_0^{\mathbf{h}} \sum_{i=1}^k w_i(\eta) d\eta_i,$$

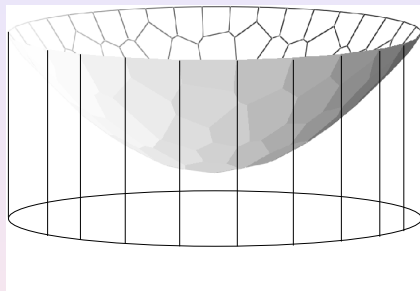
where  $w_i(\eta) = \mu(W_i(\eta) \cap \Omega)$  is the  $\mu$ -volume of the cell.

X. Gu, F. Luo, J. Sun and S.-T. Yau, “Variational Principles for Minkowski Type Problems, Discrete Optimal Transport, and Discrete Monge-Ampere Equations”, arXiv:1302.5472

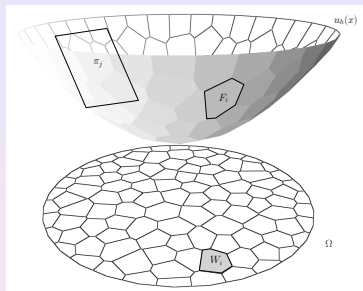


Accepted by Asian Journal of Mathematics (AJM)

# Geometric Interpretation



One can define a cylinder through  $\partial\Omega$ , the cylinder is truncated by the xy-plane and the convex polyhedron. The energy term  $\int^h \sum w_i(\eta) d\eta_i$  equals to the volume of the truncated cylinder.



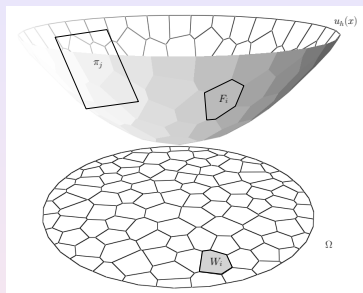
## Definition (Alexandrov Potential)

The concave energy is

$$E(h_1, h_2, \dots, h_k) = \sum_{i=1}^k v_i h_i - \int_0^h \sum_{j=1}^k w_j(\eta) d\eta_j,$$

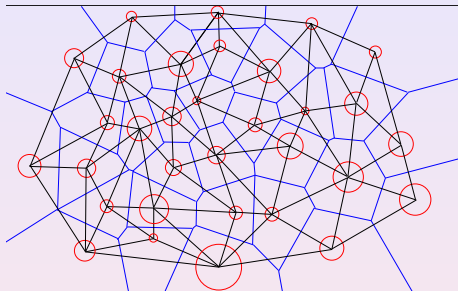
Geometrically, the energy is the volume beneath the parabola.

# Computational Algorithm



The gradient of the Alexanrov potential is the differences between the target measure and the current measure of each cell

$$\nabla E(h_1, h_2, \dots, h_k) = (v_1 - w_1, v_2 - w_2, \dots, v_k - w_k)$$



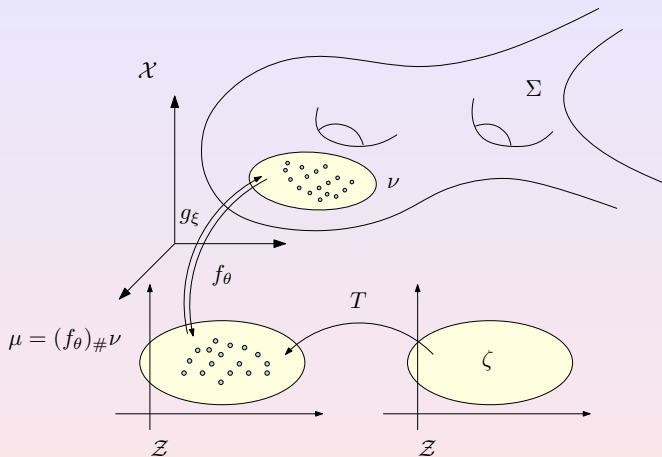
The Hessian of the energy is the length ratios of edge and dual edges,

$$\frac{\partial w_i}{\partial h_j} = \frac{|e_{ij}|}{|\bar{e}_{ij}|}$$



# Generative Model

# Autoencoder-OMT



Use autoencoder to realize encoder and decoder, use OMT in the latent space to realize probability transformation.

# Experiments



(a) real digits



(b) VAE



(c) WGAN

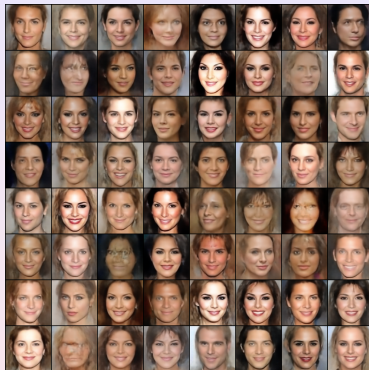


(d) AE-GMT

# Experiments



(a) VAE



(d) AE-OMT

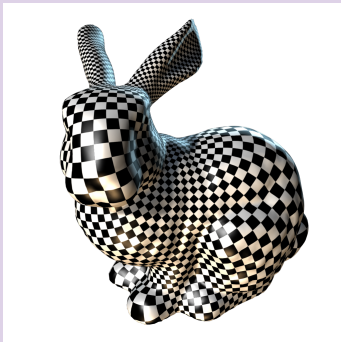
This work introduces a geometric understanding of deep learning:

- The manifold distribution assumption and the clustering assumption.
- Network complexity; manifold complexity.
- Geometric theory of optimal mass transportation.

- Na Lei, Zhongxuan Luo, Shing-Tung Yau and Xianfeng Gu, “Geometric Understanding of Deep Learning”, arXiv:1805.10451
- Na Lei, Kehua Su, Li Cui, Shing-Tung Yau and Xianfeng Gu, “A Geometric View of Optimal Transportation and Generative Model”, arXiv:1710.54888
- Xianfeng Gu, Feng Luo, Jian Sun and Shing-Tung Yau, Variational Principles for Minkowski Type Problems, Discrete Optimal Transport, and Discrete Monge-Ampere Equations, Vol. 20, No. 2, pp. 383-398, Asian Journal of Mathematics (AJM), April 2016.

# Thanks

For more information, please email to [gu@cs.stonybrook.edu](mailto:gu@cs.stonybrook.edu).



# Thank you!