### A Geometric Understanding of Deep Learning

#### Na Lei

DUT-RU International School of Information Science and Engineering Dalian University of Technology

Nov. 6, 2022

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 1 / 76

< 6 b

- B- 6

### References

- Na Lei, Kehua Su, Li Cui, Shing-Tung Yau and Xianfeng Gu, "A Geometric View of Optimal Transportation and Generative Model", Computer Aided Geometric Design, 68(2019), 1-21.
- Na Lei, Dongsheng An, Yang Guo, Kehua Su, Shixia Liu, Zhongxuan Luo, Shing-Tung Yau, Xianfeng Gu, A Geometric Understanding of Deep Learning, Engineering, 6(2020), 361-374.
- Dongsheng An, Yang Guo, Na Lei, Zhongxuan Luo, Shing-Tung Yau, Xianfeng Gu, AE-OT: A New Generative Model Based on Extended Semi-Discrete Optimal Transport, ICLR 2020.
- Dongsheng An, Yang Guo, Min Zhang, Xin Qi, Na Lei, Xianfeng Gu, AE-OT-GAN: Training GANs from data specific latent distribution, ECCV 2020.
- Na Lei, Xianfeng Gu, FFT-OT: A Fast Algorithm for Optimal Transportation, ICCV 2021.

Nov. 6, 2022

2/76

- D. An, N. Lei and X. Gu, Approximate Discrete Optimal Transport Plan with Auxiliary Measure Method, ECCV 2022.
- D. An, N. Lei, X. Xu and X. Gu, Efficient Optimal Transport Algorithm by Accelerated Gradient Descent, AAAI, 2022.

# International Peer's Opinion



Fields Medalist, Cédric Villani, authoritative expert in optimal transpotation theory, introduced our work in the keynote speach at the New Frontiers in Mathematics International Conference.



• • • • • • • • • • • • •

### International Peer's Opinion

Alessio Figalli, who won the Fields Medal in 2018 for the theory of optimal transmission regularity, introduced our work in a public speech at the International Conference Dynamic, Educations and Applications.





DYNAMICS

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

### Outline

- Manifold Distribution Hypothesis
- Manifold Learning
- Oistribution Transformation
- Mode Collapse
- AE-OT framework
- AE-OT-GAN framework

A (10) A (10)

### **Manifold Distribution Hypothesis**

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

# Manifold Assumption

### Manifold Assumption[1,2]

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

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

[1] Roweis S T, Saul L K 2000. Nonlinear dimensionality reduction by locally linear embedding. Science 290(5500):2323-2326.

[2] Tenenbaum J B, De Silva V, Langford J C 2000. A global geometric framework for nonlinear dimensionality reduction. Science, 290(5500):2319-2323.

# MNIST tSNE Embedding





(a) LeCunn's MNIST handwritten (b) Hinton's t-SNE embemdding digits samples on manifold on latent space

- Each image 28 × 28 is treated as a point in the 784 dimensional image space;
- The hand-written digits image manifold *M* is of very low dimension;

### **General Model**



- Ambient Space- data space ℝ<sup>n</sup>
- manifold Support of a distribution μ
- parameter domain latent space ℝ<sup>m</sup>
- coordinates map φ<sub>i</sub>encoding/decoding maps
- φ<sub>ij</sub> controls the probability measure

# Encoding/Decoding



# a. Input manifold $M \subset \mathscr{X}$

b. latent representation c. reconstructed mfld  $D = \varphi_{\theta}(M) \subset \mathscr{Z}$   $\tilde{M} = \psi_{\theta}(D) \subset \mathscr{X}$ 

Figure: Auto-encoder pipeline.



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

Nov. 6, 2022 11 / 76

イロト イヨト イヨト イヨト



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

Nov. 6, 2022 12 / 76



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

Nov. 6, 2022 13 / 76

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))



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

Nov. 6, 2022 14 / 76

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

### **Generative Task**



The central tasks for a generative model are

- Learn the manifold structure from the data  $(f_{\theta}, g_{\xi})$ ;
- Oistribution transformation in the latent space (T).

### Learn the manifold structure

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 16 / 76

### Non-Adversarial Model: Autoencoder



#### Figure: Auto-encoder architecture.

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

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

# **ReLU DNN**

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

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

$$\varphi = T_{k+1} \circ \sigma \circ T_k \circ \cdots \circ T_2 \circ \sigma \circ T_1, \tag{1}$$

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

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

#### Definition (Activated Path)

Given a point  $\mathbf{x} \in \mathscr{X}$ , 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  $\rho : \mathscr{X} \to 2^{\mathscr{S}}$ .

# **Cell Decomposition**

### Definition (Cell Decomposition)

Fix a encoding map  $\varphi_{\theta}$  represented by a ReLU DNN, two data points  $\mathbf{x}_1, \mathbf{x}_2 \in \mathscr{X}$  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  $\mathscr{X}$  into cells,

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

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

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

# Encoding/Decoding



a. Input manifold  $M \subset \mathscr{X}$ 

b. latent representation  $D = \varphi_{\theta}(M)$ 

c. reconstructed mfld  $ilde{M} = \psi_{ heta}(D)$ 

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

Figure: Auto-encoder pipeline.

Nov. 6, 2022 21 / 76

# **Piecewise Linear Mapping**



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

e. latent space cell decomposition

f. cell decomposition  $\mathscr{D}(\psi_{ heta} \circ \varphi_{ heta})$ 

< D > < A.

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

Nov. 6, 2022 22 / 76

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

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

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

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  $\mathscr{C}(d, n)$ , then

$$\mathscr{C}(d,n) = \begin{pmatrix} n \\ 0 \end{pmatrix} + \begin{pmatrix} n \\ 1 \end{pmatrix} + \begin{pmatrix} n \\ 2 \end{pmatrix} + \dots + \begin{pmatrix} n \\ d \end{pmatrix}.$$
(2)

### Proof.

Suppose *n* hyperplanes cut  $\mathbb{R}^d$  into  $\mathscr{C}(d, n)$  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  $\mathscr{C}(d-1, n)$  cells, each cell on  $\pi$  partitions a polyhedron in  $\mathbb{R}^d$  into 2 cells, hence we get the formula

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

It is obvious that  $\mathscr{C}(2,1) = 2$ , the formula (2) can be easily obtained by induction. Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning Nov. 6, 2022 24/76

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

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

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

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

# **RL** Complexity of Manifold



#### 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  $\varphi : \mathbb{R}^n \to \mathbb{R}^m$ , such that the restriction of  $\varphi$  on *M*,  $\varphi|_M : M \to \varphi(M) \subset \mathbb{R}^m$ , is homeomorphic.  $\varphi$  is called the corresponding rectified linear map of *M*.

# Manifold RL Complexity

#### Definition (Linear Rectifiable Atlas)

Suppose *M* is a *m*-dimensional manifold, embedded in  $\mathbb{R}^n$ ,  $\mathscr{A} = \{(U_\alpha, \varphi_\alpha)\}$  is an atlas of *M*. If each chart  $(U_\alpha, \varphi_\alpha)$  is linear rectifiable,  $\varphi_\alpha : U_\alpha \to \mathbb{R}^m$  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  $\mathscr{N}(\mathbb{R}^n, M)$  and defined as,

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

### Definition (Encoding Map)

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

### Theorem (Encodable Condition)

Suppose a ReLU DNN  $N(w_0, ..., w_{k+1})$  represents a PL mapping  $\varphi_{\theta} : \mathbb{R}^n \to \mathbb{R}^m$ , 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 that the rectified linear complexity of  $(\mathbb{R}^n, M)$ ,

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

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

# **Representation Limitation Theorem**



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

#### Theorem

Given any ReLU deep neural network  $N(w_0, w_1, ..., w_k, w_{k+1})$ , 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?

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 30 / 76

# Concept

### Concept

Consider the total space  $\Omega$  consisting all possible *n* by *n* images, a concept is a subset *D* in the space,  $D \subset \Omega$ . A concept is represented as a probability measure

 $P(x) = Prob\{x \in D\}.$ 

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

# **Generative Model**

### **Generative Model**

 $G: \mathscr{Z} \to \mathscr{X}$  maps a fixed probability distribution  $\zeta$  to the training data probability distribution v.



### **Optimal Transport Problem**



Earth movement cost.

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 33 / 76

• • • • • • • • • • • •

# **Optimal Transportation**

### Definition (Transportation Cost)

Given two bounded domains in  $\mathbb{R}^n$  with probability measures  $(X, \mu)$ and (Y, v), with equal total measure  $\mu(X) = v(Y)$ . Suppose the cost of moving a unit mass from point *x* to point *y* is c(x, y), for a transportation map  $f: (X, \mu) \to (Y, v)$ , the total transportation cost is

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



Nov. 6, 2022 34 / 76

# **Optimal Mass Transportation**

### 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  $T : X \to Y$  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  $T_{\#}\mu = v$ .

Suppose *T* is a smooth map, then measure-preserving condition is equivalent to the Jacobian equation  $\mu(x)dx = v(y)dy$ 

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

# **Optimal Transport Problem**

### Problem (Monge)

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

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

such kind of map is called the optimal transportation map.

### Definition (Wasserstein distance)

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

$$W_c(\mu, v) := \min_{T_{\#}\mu=v} \mathscr{C}(T).$$
The cost of moving a unit mass from point *x* to point *y*.

*Monge*(1781) : 
$$c(x, y) = |x - y|$$
.

This is the natural cost function. Other cost functions include

$$\begin{array}{lll} c(x,y) &=& |x-y|^{\rho}, p \neq 0\\ c(x,y) &=& -\log |x-y|\\ c(x,y) &=& \sqrt{\varepsilon+|x-y|^2}, \varepsilon > 0 \end{array}$$

Nov. 6, 2022

37 / 76

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

Kantorovich relaxed transportation maps to transportation schemes.

#### Problem (Kantorovich)

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

$$(\mathcal{KP}) \qquad \min_{\rho} \left\{ \int_{X \times Y} c(x, y) d\rho(x, y) \left| \rho_{x\#} = \mu, \, \rho_{y\#} = v \right\} \right\}$$

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

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

Problem (Kantorovich Dual) Find an functions  $\varphi : X \to \mathbb{R}$  and  $\psi : Y \to \mathbb{R}$ , such that  $(DP) \max_{\varphi, \psi} \left\{ \int_X \varphi(x) du(x) + \int_Y \psi(y) dv(y), \varphi(x) + \psi(y) \le c(x, y) \right\}.$ 

## Kantorovich Dual Problem

#### Definition (c-transformation)

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

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

 $\varphi$  is called Kantorovich potential.

Nov. 6, 2022 40 / 76

3

# Brenier's Approach

### Theorem (Brenier)

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

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

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

$$T: \mathbf{X} \to \nabla u(\mathbf{X}).$$

Problem (Brenier)

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

$$(BP) \qquad (\nabla u)_{\#}\mu = v,$$

u is called the Brenier potential.

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

Problem (Brenier)

Find the C<sup>2</sup> Brenier potential  $u: X \to \mathbb{R}$  statisfies the Monge-Ampère equation

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

< ロ > < 同 > < 回 > < 回 >

### Wasserstein GAN Model



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

Nov. 6, 2022

43 / 76

# OMT view of WGAN

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

• The discriminator: computes the Wassersteind distance using Kantorovich Dual formula:

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

namely computes the Kantorovich potential  $\varphi$ ;

- The generator: computes a measure-preserving transportation map g<sub>θ</sub> : 𝔅 → 𝔅, s.t. g<sub>θ#</sub>ζ = μ<sub>θ</sub> = ν.
- The WGAN model: min-max optimization

$$\min_{\theta} \max_{\xi} \int_{X} \varphi_{\xi} \circ g_{\theta}(z) d\zeta(z) + \int_{Y} \varphi_{\xi}^{c}(y) dv(y)$$

- A TE N - A TE N

# OMT view of WGAN

### $L^1$ 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.

# OMT view of WGAN

### $L^2$ 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 competition between *D* and *G* should be replaced by cooperation.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

### **Reason of Mode Collapse**

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 47 / 76

- Generetive models are difficult to train and sensitive to hyper-parameters;
- Generetive models suffer from mode collapsing, the generated distributions miss some modes;
- Generetive models may generate unrealistic samples;

## Mode Collapse



Nov. 6, 2022 49 / 76

# Singularity Set of OT Maps



Figure: Singularity structure of an optimal transportation map.

We call  $\Sigma_{\Omega}$  as singular set of the optimal transportation map  $\nabla u : \Omega \to \Lambda$ .

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

```
Nov. 6, 2022 50 / 76
```

#### Intrinsic Conflict

Deep neural networks can only represent continuous mappings, but the transportation maps are discontinuous on singular sets. Namely, the target mappings are outside the functional space of DNNs. This conflict induces mode collapsing.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

# Discontinuity of Optimal Transportation Map



Nov. 6, 2022 52 / 76

# Discontinuity of Optimal Transportation Map



(c) concave support: multi modes

Nov. 6, 2022 53 / 76

# Singularity Set Detection



Nov. 6, 2022 54 / 76

# Mode Collapse



(a) a path inside the manifold



(b) a path through a singularity.

Figure: Facial images generated by an AE-OT model, the image in the center of (b) shows the transportation map is discontinuous.

# **AutoEncoder-Optimal Transportation Framework**

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 56 / 76

# **AE-OT Model**

Image Space  $\mathcal{X}$ 



Use autoencoder to realize encoder and decoder, use OT in the latent space to realize probability transformation. Image: A matrix

Nov. 6, 2022 57 / 76

### **Experiments - Mode Collapse**



Nov. 6, 2022 58/76

### **Experiments - Mode Collapse**



### **Experiments - Mnist**

З  $\boldsymbol{\varrho}$ t, З 9: a .5 (b) Lucic et al. 2018 (a) AE reconstruction n К -Ζ a D (c) Hoshen and Malik 2019 (d) AE-OT

Na Lei (Dalian University of Technology)

Geometric Understanding of Deep Learning

Nov. 6, 2022 60 / 76

# **Experiments - Mnist Fashion**



#### (a) AE reconstruction



#### (c) Hoshen and Malik 2019



### (b) Lucic et al. 2018



(d) AE-OT

Na Lei (Dalian University of Technology)

Geometric Understanding of Deep Learning

Nov. 6, 2022 61 / 76

### Experiments - Cifar 10



#### (a) AE reconstruction







#### (b) Lucic et al. 2018



(d) AE-OT

Na Lei (Dalian University of Technology)

Geometric Understanding of Deep Learning

Nov. 6, 2022 62 / 76

### **Experiments - CelebA**



#### (a) AE reconstruction







#### (b) Lucic et al. 2018



(d) AE-OT

Na Lei (Dalian University of Technology)

Geometric Understanding of Deep Learning

Nov. 6, 2022 63 / 76

### **Experiments - CelebA**



#### (a) WGAN-GP

(b) WGAN-div

Figure: Failure cases for WGAN-GP and WGAN-div.

Nov. 6, 2022 64 / 76

### **Experiments - CelebA**



#### (c) CRGAN

#### (d) AE-OT

#### Figure: Mode collapsing of CRGAN.

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 65 / 76

### **Experiments - AE-OT interpolation**



#### Figure: Curves on facial photo manifold.

Na Lei (Dalian University of Technology)

Geometric Understanding of Deep Learning

Nov. 6, 2022 66 / 76

# Quantitative Comparison with FID

	Adversarial						
Dataset	NS GAN	LSGAN	WGAN	BEGAN			
MNIST	6.8	7.8	6.7	13.1			
Fansion	26.5	30.7	21.5	22.9			
CIFAR-10	58.5	87.1	55.2	71.4			
CelebA	55.0	53.9	41.3	38.9			

	No	n-Adver	sarial	Reference		
Dataset	VAE	GLO	GLANN	AE Reconstruction	Ours	
MNIST	23.8	49.6	8.6	5.5	6.4	
Fansion	58.7	57.7	13.0	4.7	10.2	
CIFAR-10	155.7	65.4	46.5	28.2	38.1	
CelebA	85.7	52.4	46.3	67.5	68.4	

### **AE-OT-GAN Model**

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 68 / 76



For further improve the quality of the generated images, we proposed AE-OT-GAN model.

Nov. 6, 2022 69 / 76

< ロ > < 同 > < 回 > < 回 >



Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 70 / 76

ъ

• • • • • • • • • • • • •



CT-GAN [42] WGAN-GP [12] WGAN-div [39] WGAN-QC [24] Proposed method



Comparison between the proposed method and the SOTA methods.

Na Lei (Dalian University of Technology) Geometric L

Geometric Understanding of Deep Learning

Nov. 6, 2022 71 / 76

method

	CIFAR10			CelebA		
	Standard		Resnet		Standard	Resnet
	FID	IS	FID	IS	FID	FID
WGAN-GP $[12]$	40.2	6.68	19.6	7.86	21.2	18.4
PGGAN $[14]$	-	-	18.8	8.80	-	16.3
SNGAN [27]	25.5	7.58	21.7	8.22	-	-
WGAN-div [39]	-	-	18.1	-	17.5	15.2
WGAN-QC $[24]$	-	-	-	-	-	12.9
AE-OT $[2]$	34.2	6.62	28.5	7.67	24.3	28.6
AE-OT-GAN	25.2	7.62	17.1	8.24	11.2	7.6

イロト イヨト イヨト イヨト
## Table 1. Experiments on stacked MNIST.

	Stacked MNIST	
	Modes	KL
DCGAN	99.0	3.40
VEEGAN	150.0	2.95
PacDCGAN4	$1000.0\pm0.00$	$0.07\pm0.005$
WGAN	$314.3\pm38.54$	$2.44\pm0.170$
AE-OT	$1000.0\pm0.00$	$0.03 \pm 0.0008$
AE-OT-GAN	$1000.0\pm0.0$	$0.05\pm0.006$



## **AE-OT-GAN Interpolation**



The interpolation of the AE-OT-GAN model. The left column shows the generated images, and the right 5 images are the ones used to generate the left images in the latent space.

## Conclusion

This work introduces a geometric understanding of deep learning:

- The intrinsic pattern of natural data can be represented by manifold distribution hypothesis.
- The deep learning system has two major tasks: manifold learning and probability distribution transformation.
- Optimal transportation methods can be used to accomplish the second task.
- The regularity theory of Monge-Ampère equation explains mode collapse.
- The AE-OT framework can avoid mode collapse, and make half the blackbox transparent.



## Thank you!

For more information, please email to nalei@dlut.edu.cn.

Na Lei (Dalian University of Technology) Geometric Understanding of Deep Learning

Nov. 6, 2022 76 / 76

< ロ > < 同 > < 回 > < 回 >