

# Deep Learning Theory

**Taiji Suzuki**

The University of Tokyo

Deep Learning Theory Team/AIP-RIKEN

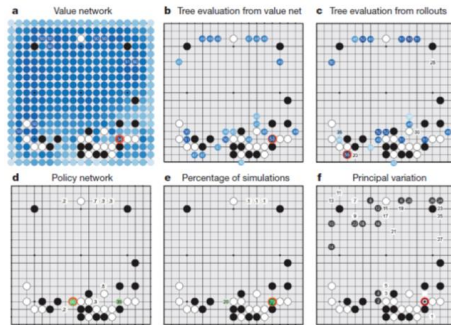
12<sup>nd</sup>/Mar/2024

MLSS2024@OIST

# Success of deep learning

Deep learning has shown great performances in the AI research field.  
→ Why?

## AlphaGo/Zero



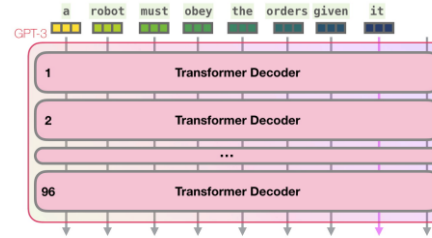
[Silver et al. (Google Deep Mind): Mastering the game of Go with deep neural networks and tree search, Nature, 529, 484—489, 2016]

## Image recognition



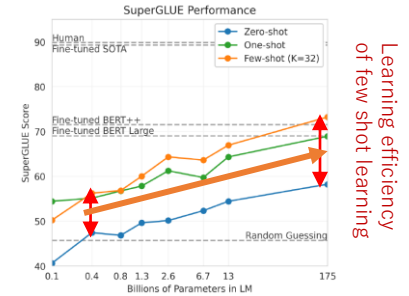
[He, Gkioxari, Dollár, Girshick: Mask R-CNN, ICCV2017]

## Large language model



[Alammar: How GPT3 Works - Visualizations and Animations, <https://jalammar.github.io/how-gpt3-works-visualizations-animations/>]

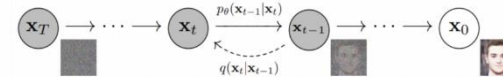
[Brown et al. “Language Models are Few-Shot Learners”, NeurIPS2020]



Performance of few-shot learning against model size

Learning efficiency of few shot learning

## Generative models (diffusion models)



[Ho, Jain, Abbeel: Denoising Diffusion Probabilistic Models. 2020]



Stable diffusion, 2022.



Jason Allen "Théâtre D'opéra Spatial" generated by **Midjourney**.  
Colorado State Fair generated by NovelAI  
1st prize in digital

# What we need to solve?

## Why does deep learning work well?

- Several theoretical work has been conducted.
- There are still many things that should be explored.

- Clarification of principle of deep learning
- What is essential to realize a “good” learning system?  
→ We may find a new method beyond DL.

### Issue in academic conference



Ali Rahimi's talk at NIPS(NIPS 2017 Test-of-time award presentation)

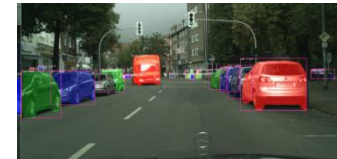


Ali Rahimi's talk at NIPS(NIPS 2017 Test-of-time award presentation)

Ali Rahimi's talk at NIPS2017 (test of time award).  
“Random features for large-scale kernel methods.”  
Criticism that DL is “alchemy.”

### Issue in industries

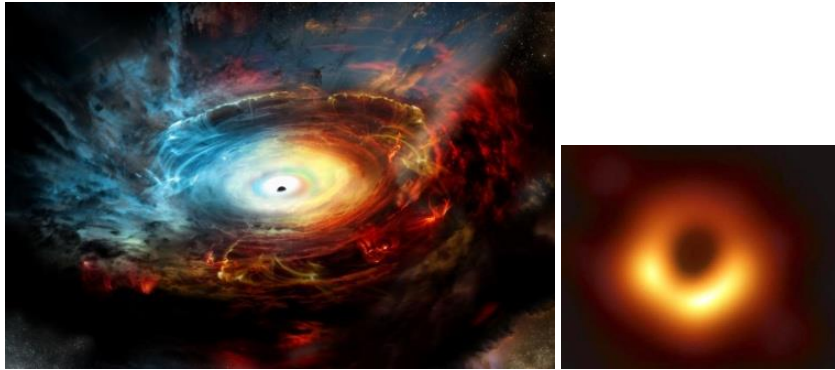
- We don't want to use black-box system.
- Accountabilities of companies.



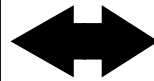
# Deep learning theory

# Role of mathematics

## Physics



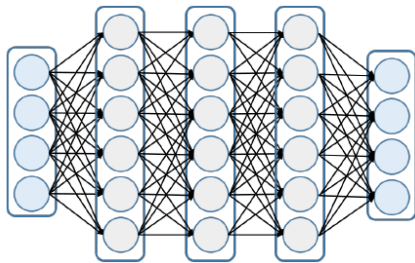
physical phenomenon


$$R_{\mu\nu} - \frac{1}{2}g_{\mu\nu}R = \frac{8\pi G}{c^4}T_{\mu\nu}$$

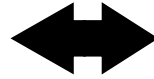
- Theory of Relative
  - Riemannian geometry
- Quantum mechanics
  - Functional analysis

## Mathematics

## Machine learning



Deep learning



Several mathematicians/physicists join the ML community.

- Prob. theory
- Functional anal.
- Wasserstein geom.
- Diffusion equation
- Statistics
- Optimization
- Numerical analysis

## Mathematics



# Layers of deep learning theory

Application

**Interpretability :**  
Accountability, visualization,  
easier maintenance

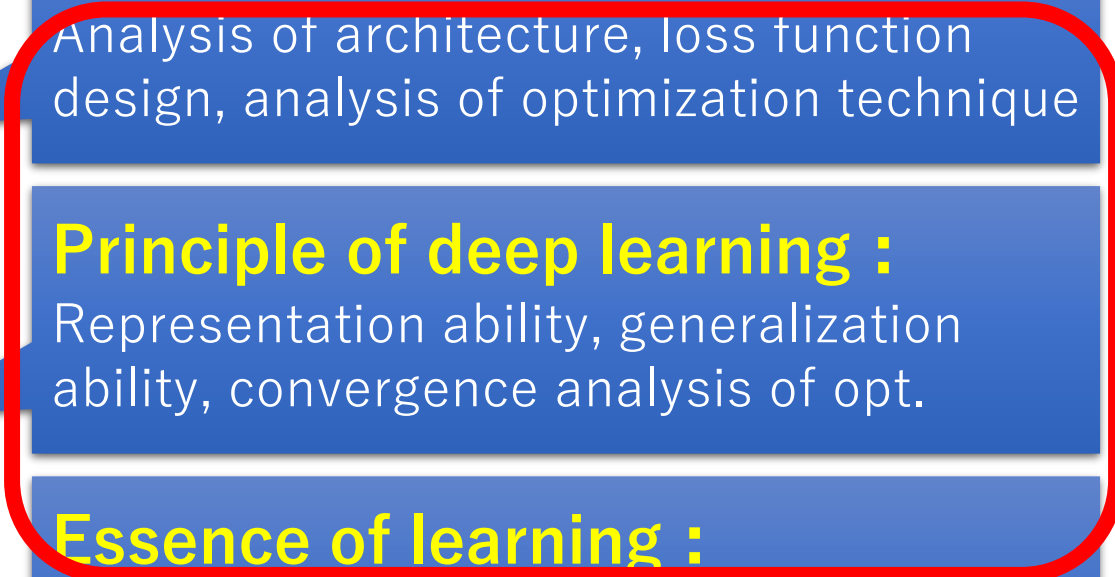
**Analysis of several techniques :**  
Analysis of architecture, loss function  
design, analysis of optimization technique

**Principle of deep learning :**  
Representation ability, generalization  
ability, convergence analysis of opt.

**Essence of learning :**  
Characterization of “good” learning  
methods, unified theory, beyond DL

Understanding behaviors of DL

- Accountability
- Clarifying possibility and limitations
- Guideline for designing a learning method



Today's topic

Foundation

# 3 issues of deep learning theory

## **Representation ability**

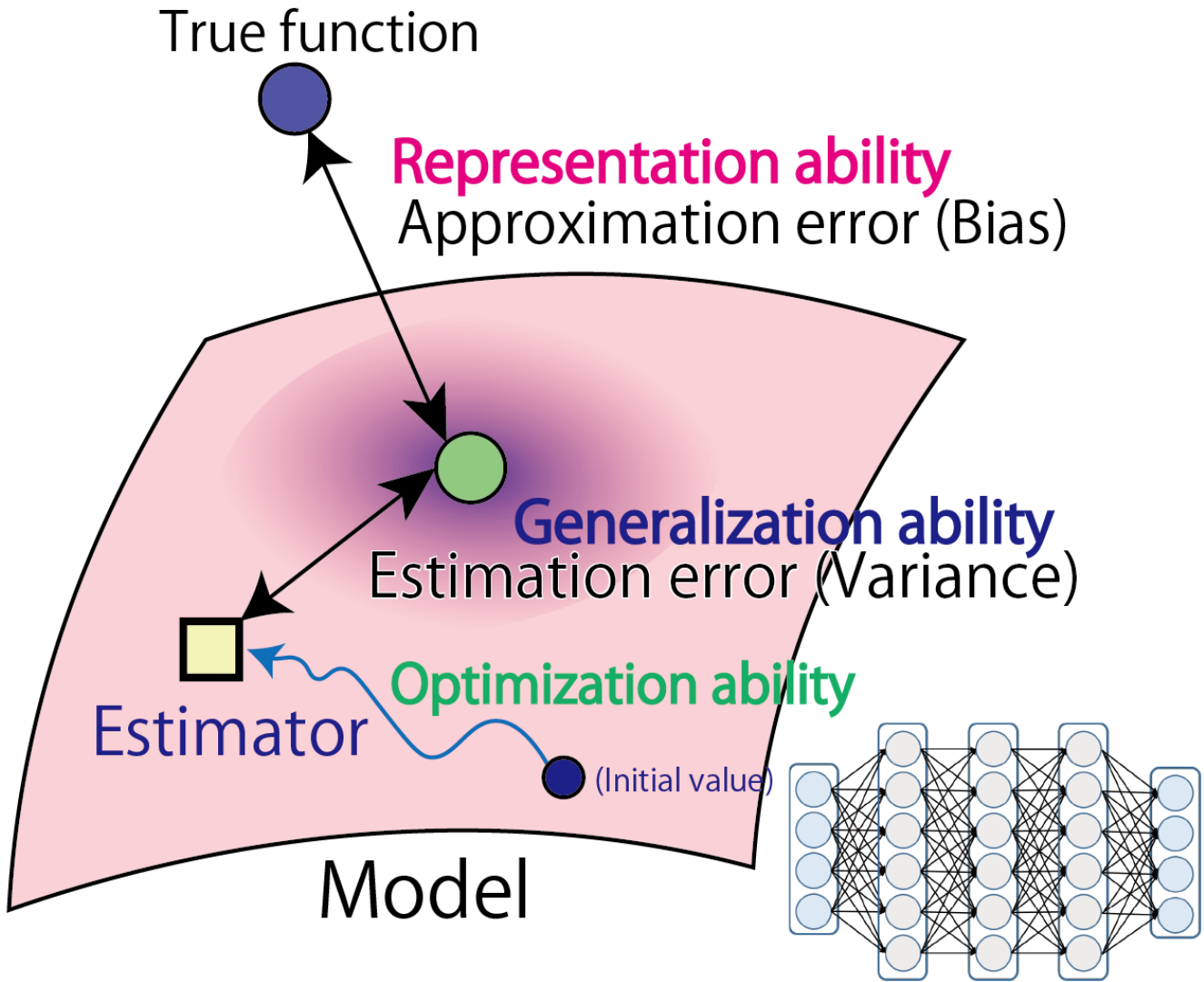
What kind of functions can DNN approximate?

## **Generalization ability**

How well can DL generalize from finite observations?

## **Optimization ability**

How fast can we find the optimal parameter?



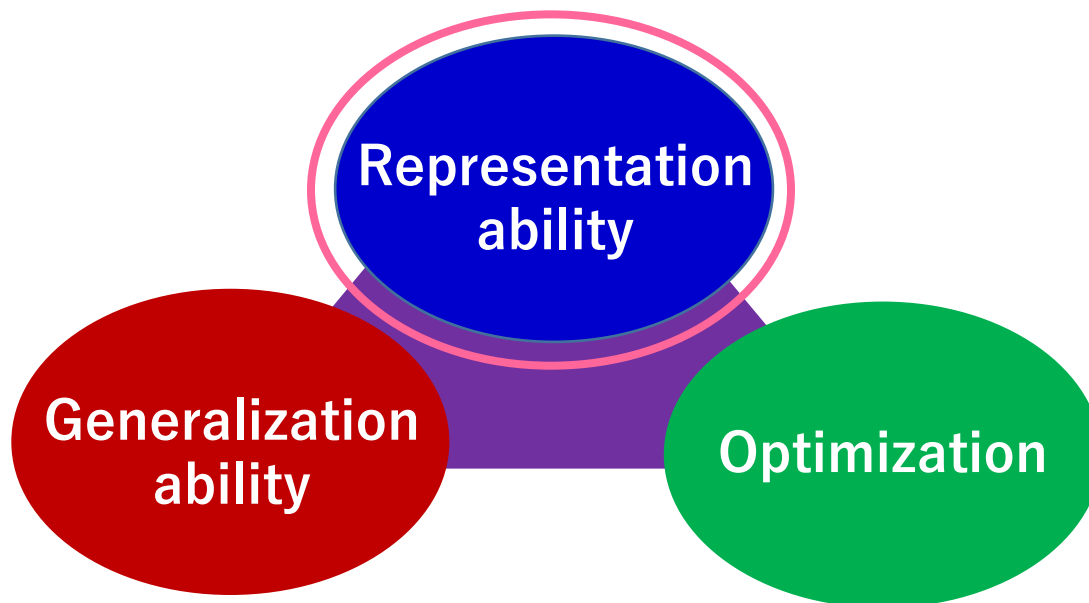
## 1. Representation ability + Generalization ability

- Universal approximator
- Depth separation
- Adaptivity of deep learning
  - Inhomogeneity of smoothness
  - Curse of dimensionality
- Foundation models
  - Diffusion model
  - Transformer

## 2. Optimization ability

- Noisy gradient descent
- Mean field Langevin
- CSQ lowerbound

# Representation ability of neural networks



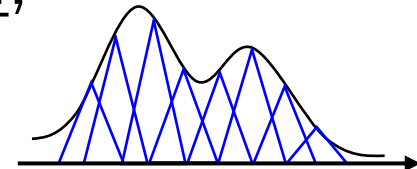
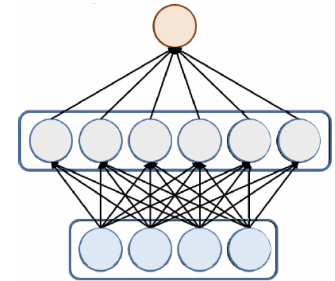


# Representation ability

**Universal approximator**  
Neural networks can approximate  
“any function” with “any precision”.

2-layer NNs can approximate any function,  
by increasing the number of neurons.

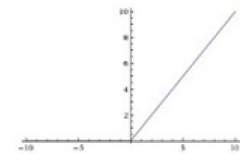
[Hecht-Nielsen,1987][Cybenko,1989]



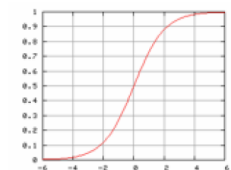
Year		Basis function	Space
1987	Hecht-Nielsen	Depending on the target	$C(R^d)$
1988	Gallant & White	Cos	$L_2(K)$
	Irie & Miyake	integrable	$L_2(R^d)$
1989	Carroll & Dickinson	Continuous sigmoidal	$L_2(K)$
	Cybenko	Continuous sigmoidal	$C(K)$
	Funahashi	Monotone & bounded	$C(K)$
1993	Mhaskar + Micchelli	Polynomial growth	$C(K)$
2015	Sonoda + Murata	<b>Unbounded</b> , admissible	$L_1(R^d), L_2(R^d)$

$K$  is any compact set.

**ReLU:**  $\eta(u) = \max\{u, 0\}$



**Sigmoid:**  $\eta(u) = \frac{1}{1+\exp(-u)}$

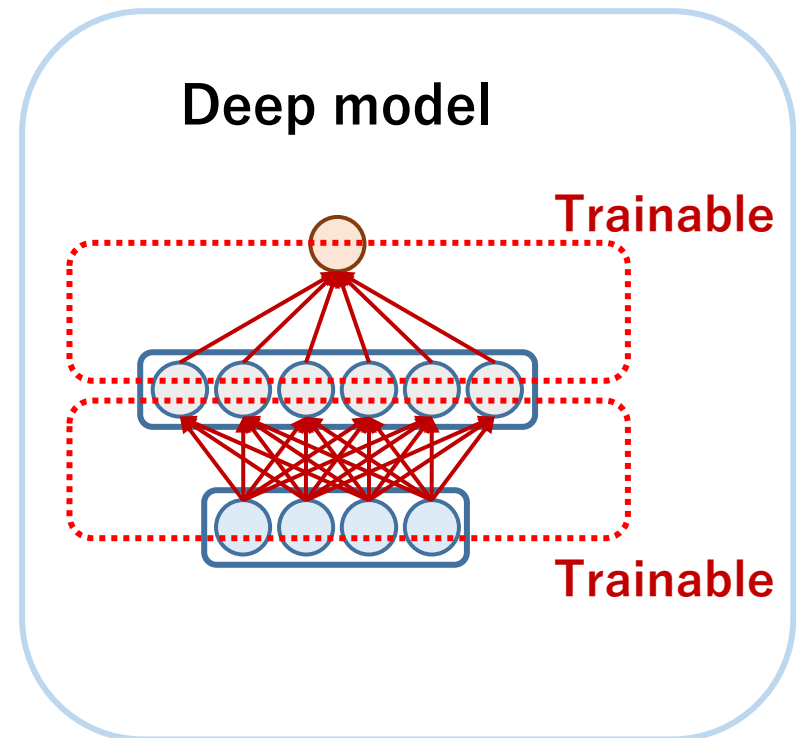
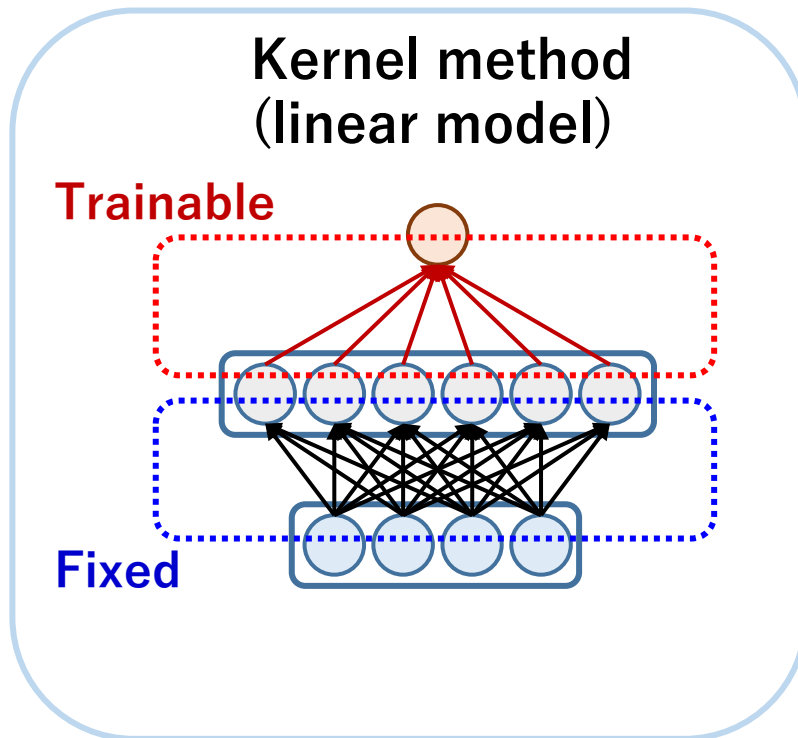
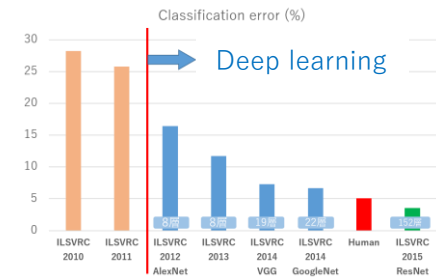


# What are we missing?

- **[Theory]** Kernel method is also a universal approximator.
- **[Practice]** DL performs better.

→ Why?

(in some case)

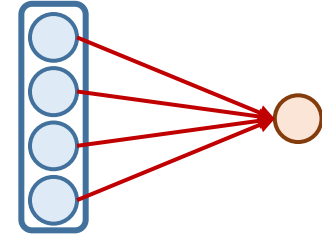


→ We compare “accuracy” of estimation/approximation.

# Feature learning

- Linear model

$$f(x) = \sum_{j=1}^d \alpha_j x_j$$

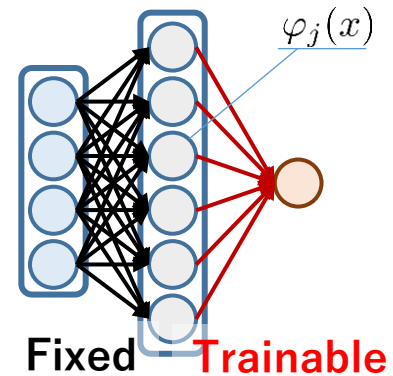


Nonlinear



- Kernel model

$$f(x) = \sum_{j=1}^M \alpha_j \underbrace{\varphi_j(x)}_{\text{Fixed}}$$

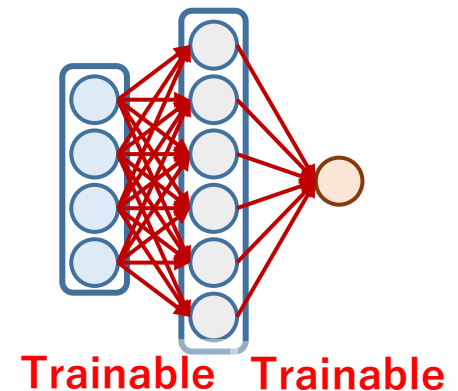


Adaptive  
basis learning

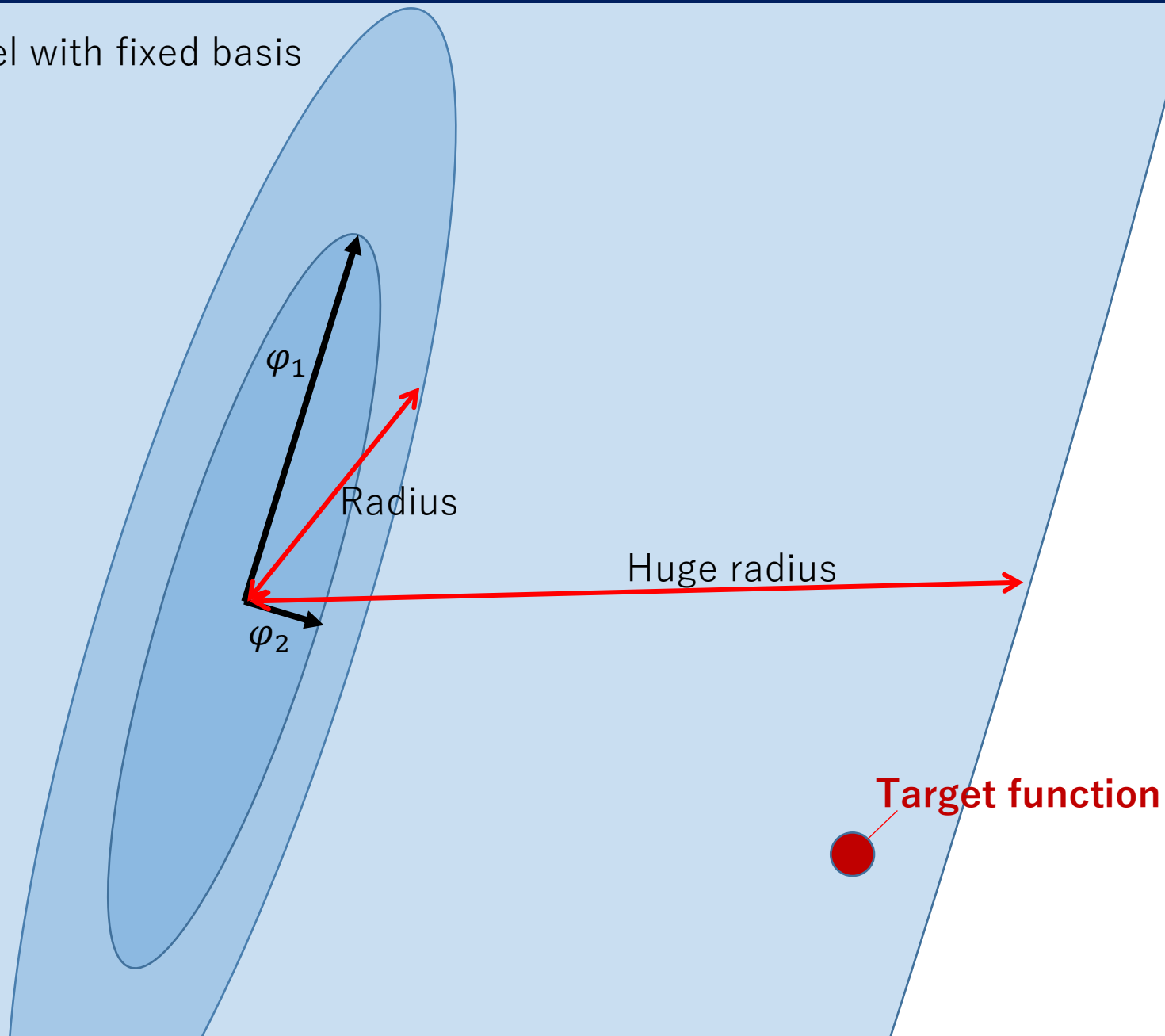


- Neural network

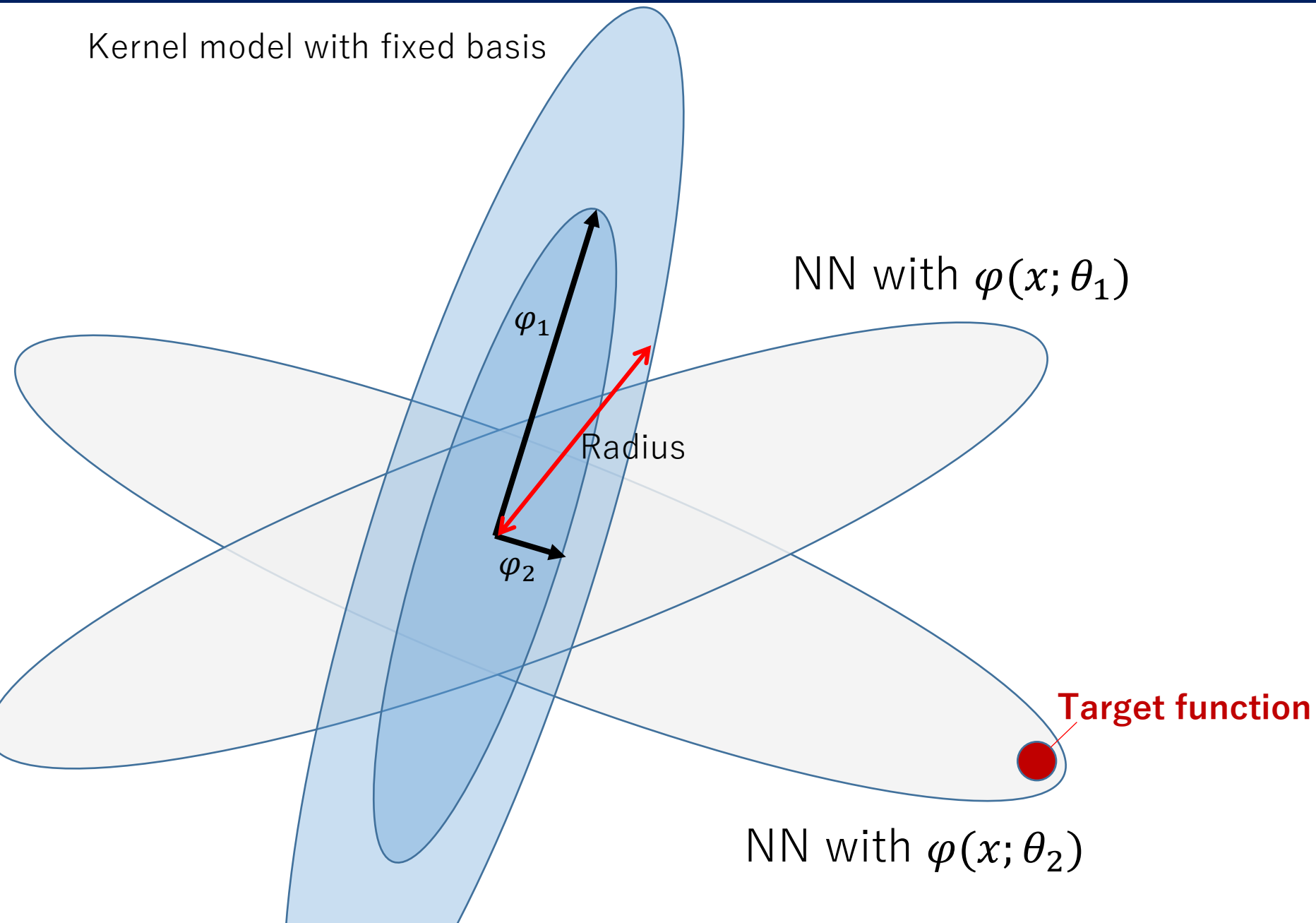
$$f(x) = \sum_{j=1}^M \alpha_j \underbrace{\varphi_j(x; \theta)}_{\text{Trainable}}$$



Kernel model with fixed basis



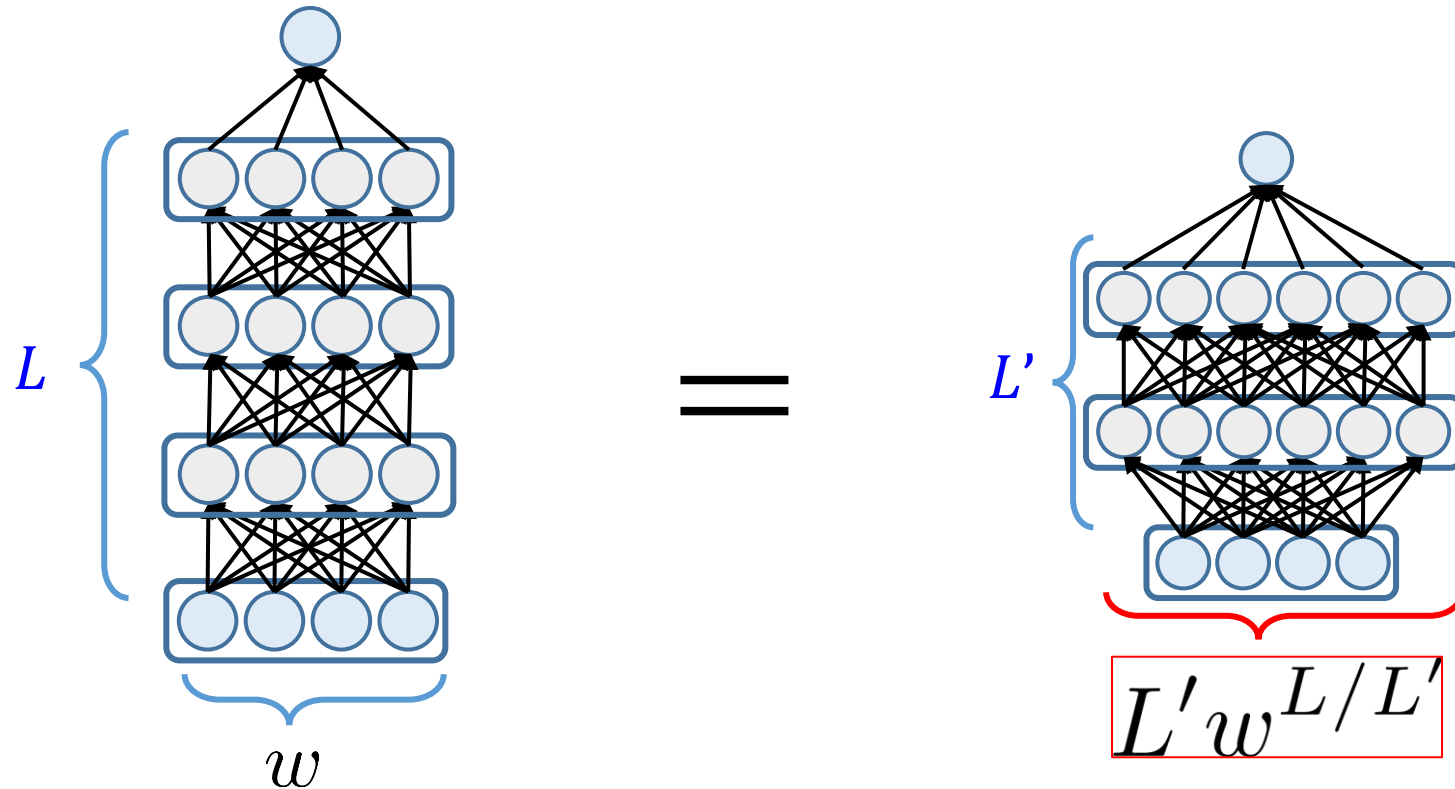
Kernel model with fixed basis





# Deep network is exponentially powerful<sup>15</sup>

Width vs Depth



Exponentially large width is required.

[Arora, Basu, Mianjy, Mukherjee: Understanding Deep Neural Networks with Rectified Linear Units. ICLR2018.]

# Curse of dimensionality/Barron class

$\pi$ : probability measure on  $\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d \mid \|x\| = 1\}$

- Mean field model defined by  $\pi$  ( $\sigma$ : ReLU)

$$\mathcal{H}_\pi = \left\{ \int_{\mathbb{S}^{d-1}} a(w) \sigma(w^\top x) \pi(dw) \mid \int_{\mathbb{S}^{d-1}} a(w)^2 \pi(dw) < \infty \right\}$$

$$\|f\|_{\mathcal{H}_\pi}^2 := \inf_a \mathbb{E}_{w \sim \pi} [|a(w)|^2] \quad \text{where } f = \int a(w) \sigma(w^\top x) \pi(dw)$$

Approx. error

- Neural network model with  $M$  neurons

$$\mathcal{H}_{\text{NN}}(M) = \left\{ \hat{f}(x) = \sum_{j=1}^M r_j \sigma(u_j^\top x) \mid r_j \in \mathbb{R}, u_j \in \mathbb{S}^{d-1} \right\} : \text{set of NNs}$$

The first layer can be tuned

$$\inf_{\hat{f} \in \mathcal{H}_{\text{NN}}} \|f - \hat{f}\|_{L_2(P_X)}^2 = O(1/M) \quad (\forall f \in \mathcal{H}_\pi)$$

- Random feature model with  $M$  neurons

$$\mathcal{H}_{\text{rand}}(M) = \left\{ \hat{f}(x) = \sum_{j=1}^M r_j \sigma(u_j^\top x) \mid r_j \in \mathbb{R} \right\}$$

**Generated randomly**  
(fixed independent of the target)

$$\inf_{\hat{f} \in \mathcal{H}_{\text{rand}}(M)} \|f - \hat{f}\|_{L_2(P_X)}^2 \gtrsim \frac{1}{d^2 M^{2/d}} \quad (\exists \pi, \exists f \in \mathcal{H}_\pi)$$



**Curse of dimensionality**

(To obtain  $\epsilon$  accuracy,  $M = \epsilon^{-\Omega(d)}$  is required)

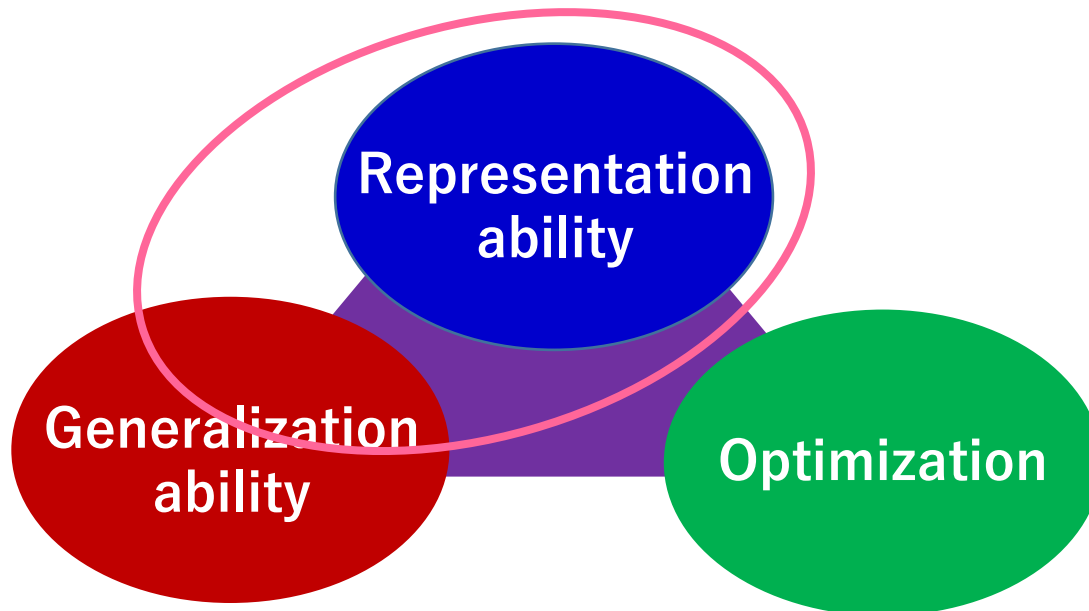
## 1. Representation ability + Generalization ability

- Universal approximator
- Depth separation
- Adaptivity of deep learning
  - Inhomogeneity of smoothness
  - Curse of dimensionality
- Foundation models
  - Diffusion model
  - Transformer

## 2. Optimization ability

- Noisy gradient descent
- Mean field Langevin
- CSQ lowerbound

# Analysis in nonparametric regression -Superiority of deep learning-

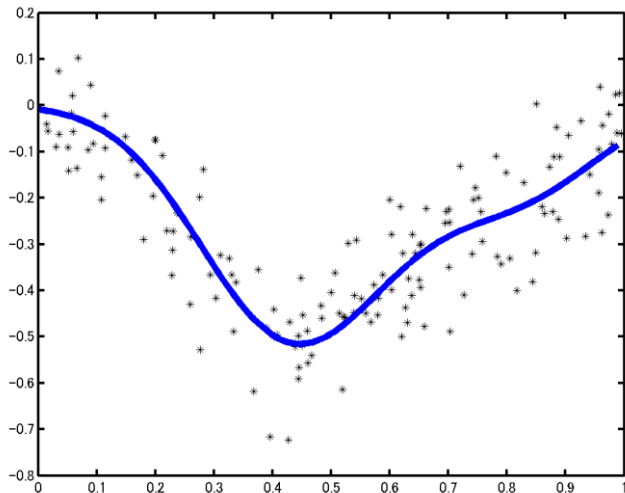


## Non-parametric regression

$$y_i = f^\circ(x_i) + \xi_i \quad (i = 1, \dots, n)$$

where  $\xi_i \sim N(0, \sigma^2)$  and  $x_i \in [0,1]^d \sim P(X)$  (i.i.d.).

We estimate  $f^\circ$  from  $(x_i, y_i)_{i=1}^n$ .



Estimation error:

$$\mathbb{E}[\|\hat{f} - f^\circ\|_{L_2(P)}^2] < ?$$

A similar argument can be applied to classification.



# Bias-Variance decomposition

Model  $\mathcal{F}$ :  $d$ -dimensional parameter

$n$ : data size

$$\hat{f} \leftarrow \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Model's degrees of freedom

$$\text{Predictive error} = \frac{\mathbf{M}}{n} + [\text{Approx. error}]$$

(mean squared error)

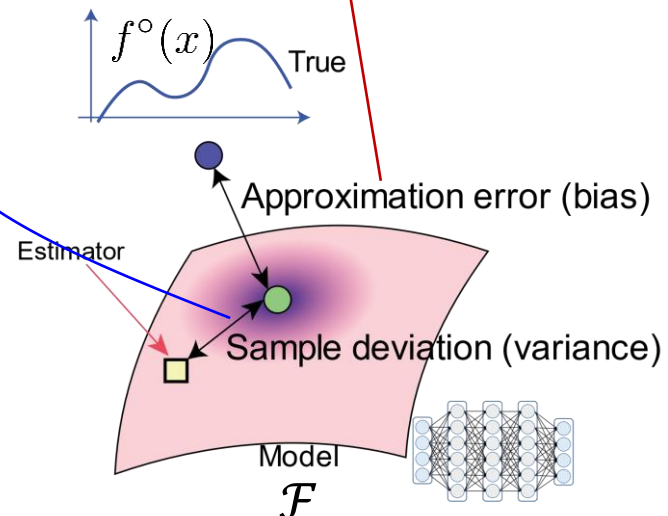
$$\mathbb{E}[(f^\circ(X) - \hat{f}(X))^2]$$

Variance

Model Bias

$$\inf_{f \in \mathcal{F}} \mathbb{E}[(f^\circ(X) - f(X))^2]$$

- **Bias-variance trade-off**
- **Large model is not necessarily good**
- $d$  can be replaced by an “intrinsic” dimensionality.
- Several learning theory is reduced to evaluate the bias and variance terms.



# Bias-Variance decomposition

Model  $\mathcal{F}$ :  $d$ -dimensional parameter

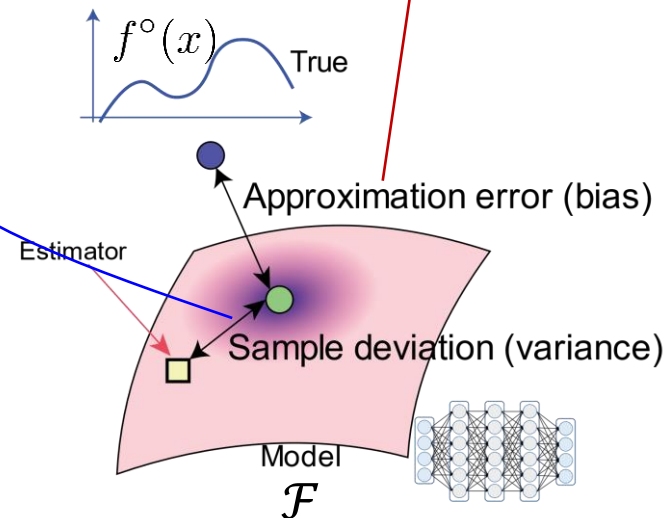
$n$ : data size

$$\hat{f} \leftarrow \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\mathbb{E}[(f^\circ(X) - \hat{f}(X))^2] = \frac{\log(\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, 1/n))}{n} + \inf_{f \in \mathcal{F}} \mathbb{E}[(f^\circ(X) - f(X))^2]$$

(mean squared error) Variance Model Bias

- **Bias-variance trade-off**
- **Large model is not necessarily good**
- $d$  can be replaced by an “intrinsic” dimensionality.
- Several learning theory is reduced to evaluate the bias and variance terms.



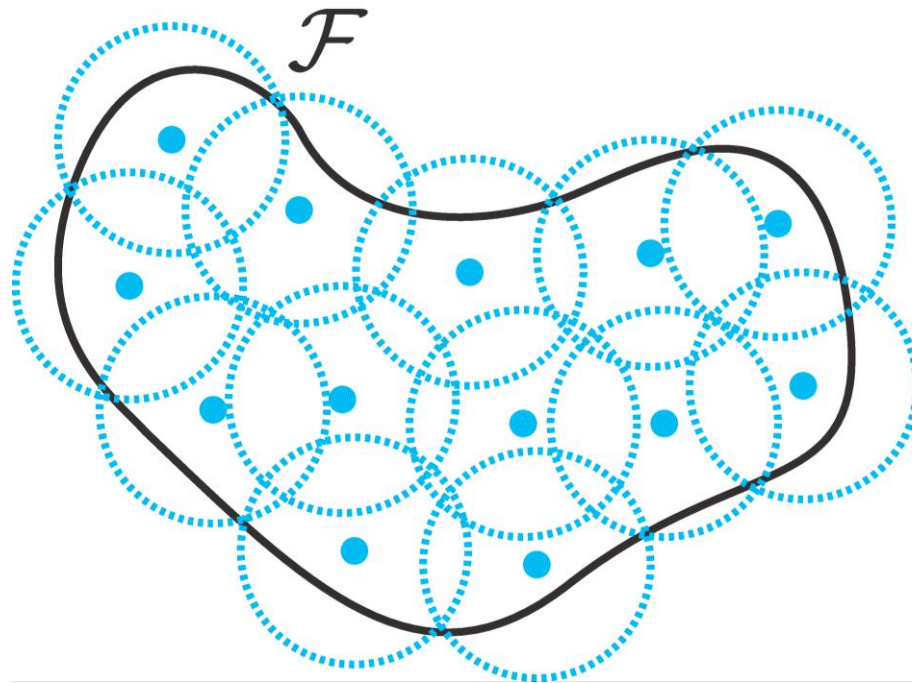
[Gine&Nickl: Mathematical Foundations of Infinite-Dimensional Statistical Models. 2015]

[Schmidt-Hieber, 2017; Hayakawa&Suzuki, 2020]

# Covering number

$$\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \epsilon)$$

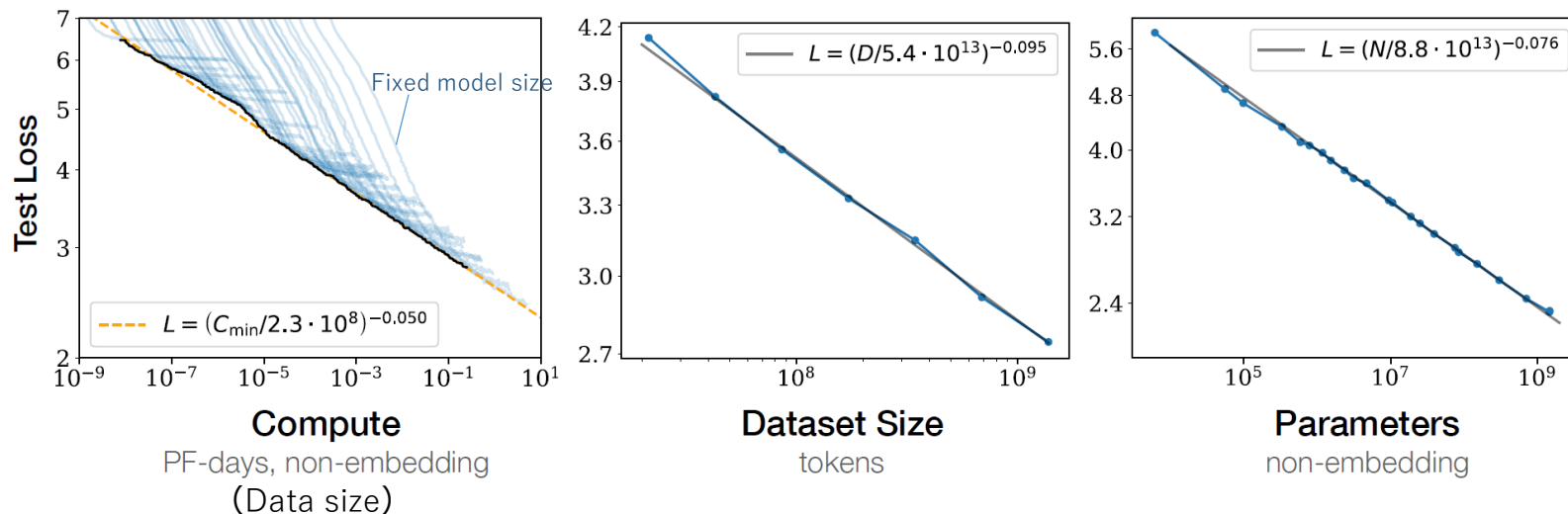
The smallest number of balls with radius  $\epsilon$  measured by the norm  $\|\cdot\|_\infty$  to cover the function class  $\mathcal{F}$ .



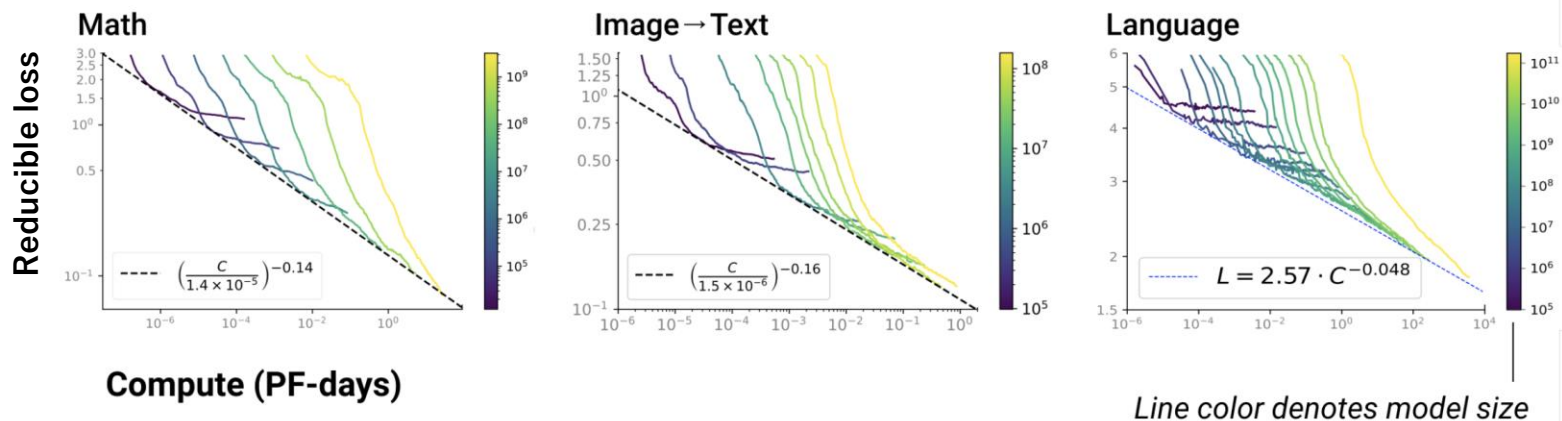
- The covering number measures how large the space  $\mathcal{F}$  is.
- In other words, it represents “complexity” of the model.

# Scaling law

[Kaplan et al.: Scaling Laws for Neural Language Models, 2020]



[Henighan et al.: Scaling Laws for Autoregressive Generative Modeling, 2020]



$$\log(\text{Predictive error}) = -\alpha \log(n) + \log(C)$$

[Brown et al.: Language Models are Few-Shot Learners, 2020]  $\leftarrow$  Analysis of GPT-3

# Analysis of kernel model

## Model

$$f^\circ(x) = \sum_{j=1}^{\infty} \alpha_j \varphi_j(x) \quad (\text{ONS in } L_2)$$

Observation:  $y_i = f^\circ(x_i) + \epsilon_i$

Assumption

$$\mu_j \sim j^{-a}$$

$$\sum_{j=1}^{\infty} \alpha_j^2 / \mu_j < \infty$$

## Training method

Least squares estimator

$M$  dimensional model

$$\hat{f}(x) = \sum_{j=1}^M \hat{\alpha}_j \varphi_j(x)$$



$$\min_{(\hat{\alpha}_j)_{j=1}^M} \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^M \hat{\alpha}_j \varphi_j(x_i) \right)^2$$

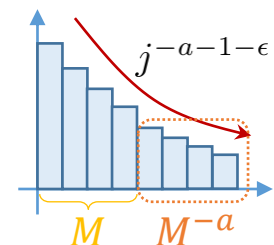
## Predictive error

$$\|f^\circ - \hat{f}\|_{L_2(P_X)}^2 \leq C \left( \overbrace{\frac{M}{n}}^{\text{Variance}} + \overbrace{M^{-a}}^{\text{Bias}} \right)$$

Optimal rate



$$n \sim \frac{a}{1+a} \quad (M^* = n^{\frac{1}{1+a}})$$



$$\log(\text{Pred. error}) = -\frac{a}{1+a} \log(n) + \log(C)$$

Variance=(dim of model)/n  
Bias=L2-norm of residual



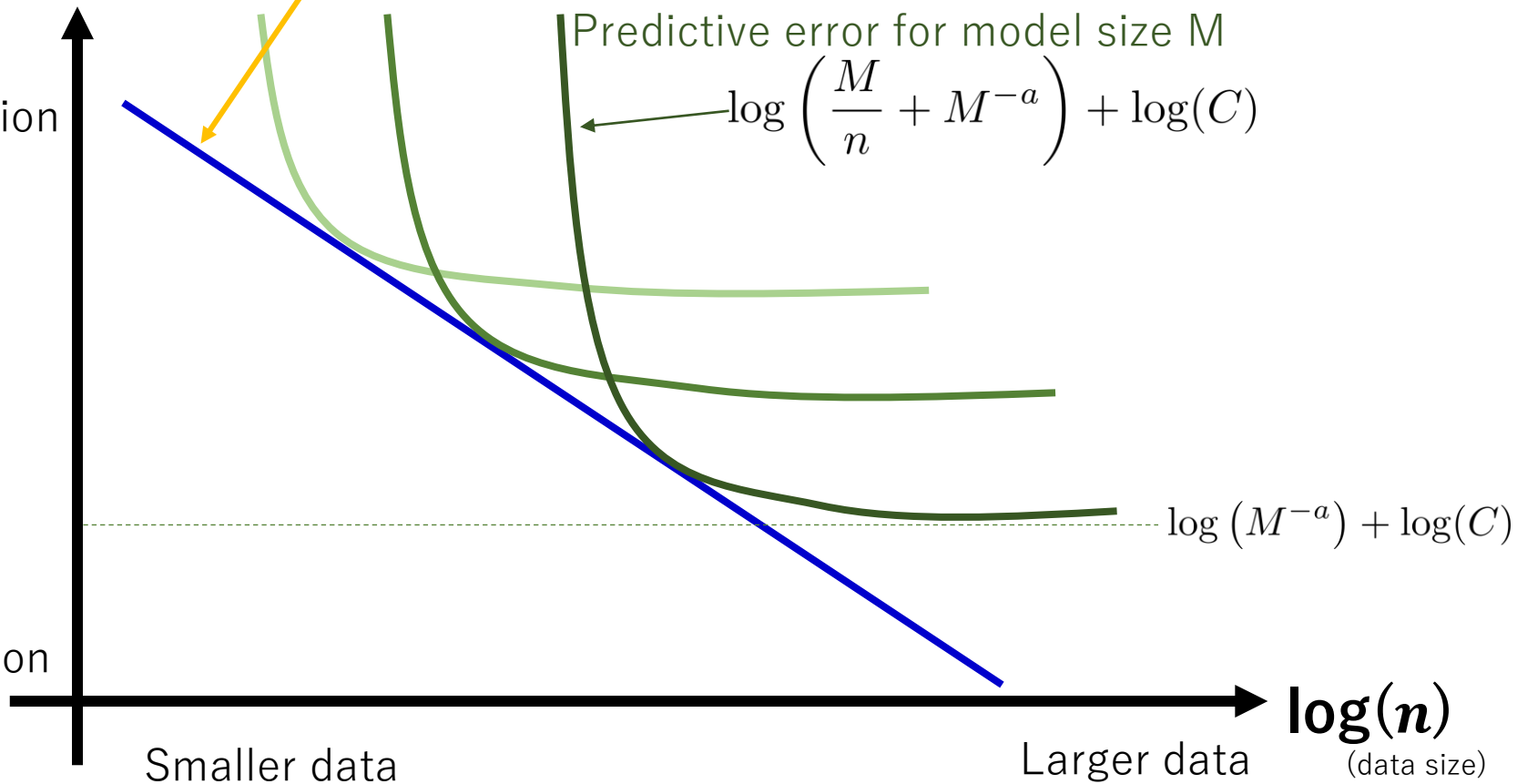
Predictive error with the optimal model size

$$\log(\text{Predictive error}) = -\frac{a}{1+a}\log(n) + \log(C)$$

**log(Pred. error)**

Worse prediction

Better prediction



## Statistical learning theory of kernel methods

- Caponnetto and De Vito. Optimal Rates for the Regularized Least-Squares Algorithm. *Foundations of Computational Mathematics*, volume 7, pp.331–368 (2007).
- Steinwart and Christmann. *Support Vector Machines*. 2008.

## Related recent papers

- Mei, Misiakiewicz, Montanari. Generalization error of random features and kernel methods: hypercontractivity and kernel matrix concentration. arXiv:2101.10588.
- Bordelon, Canatar, Pehlevan. Spectrum Dependent Learning Curves in Kernel Regression and Wide Neural Networks. ICML, 1024-1034, 2020.
- Canatar, Bordelon, Pehlevan. Spectral Bias and Task-Model Alignment Explain Generalization in Kernel Regression and Infinitely Wide Neural Networks. Nature Communications, volume 12, Article number: 2914 (2021).



# Why deep?

- There are many theories...

## Reduced rank regression

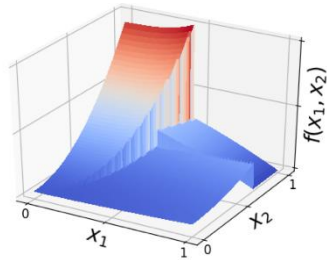
If there is low dimensional representation, deep is better.

$$Y_i = U V X_i$$

## Piece-wise smooth function

[Imaizumi&Fukumizu, 2019]

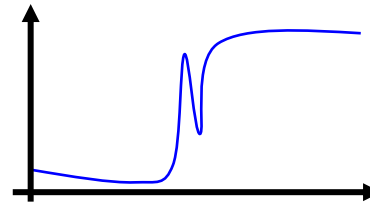
Deep is better to estimate a discontinuous function.



## Besov space

[Suzuki, 2019]

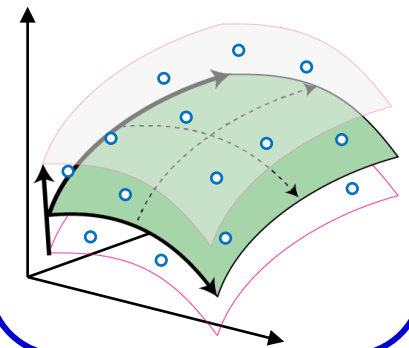
To estimate functions with non-uniform smoothness, DL is better.



## Low dimensional data structure

[Schmidt-Hieber, 2019] [Nakada&Imaizumi, 2019][Chen et al., 2019][Suzuki&Nitanda, 2019]

If data are distributed on low dim-manifold, DL is better.



Deep

$$\frac{r(M + N)}{n}$$

$$n^{-\frac{2s}{2s+d}} \vee n^{-\frac{\alpha}{\alpha+D-1}}$$

$$n^{-\frac{2s}{2s+d}}$$

$$n^{-\frac{2s}{2s+D}}$$

Kernel

$$\frac{MN}{n}$$

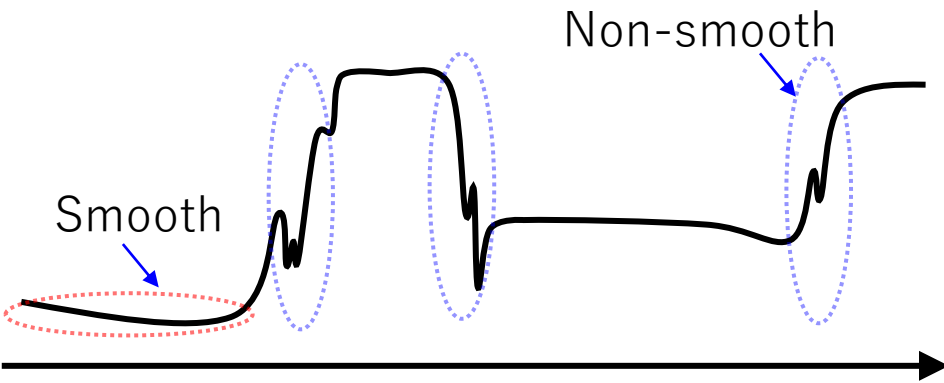
$$\frac{1}{\sqrt{n}}$$

$$n^{-\frac{2s-2d(1/p-1/2)_+}{2s+d-2d(1/p-1/2)_+}}$$

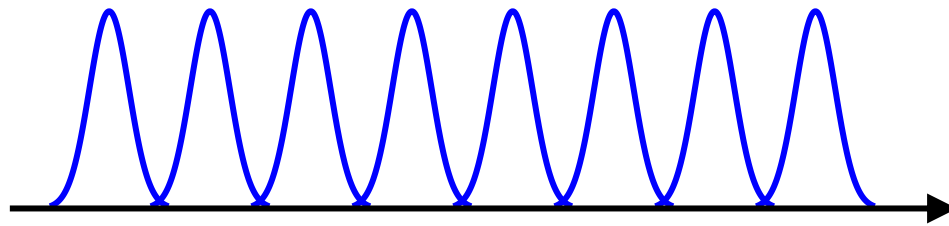
$$n^{-\frac{2(s-D/p+d/2)}{2(s-D/p+d/2)+d}} \vee n^{-\frac{2s}{2s+D}}$$

Est. error

# Typical situation

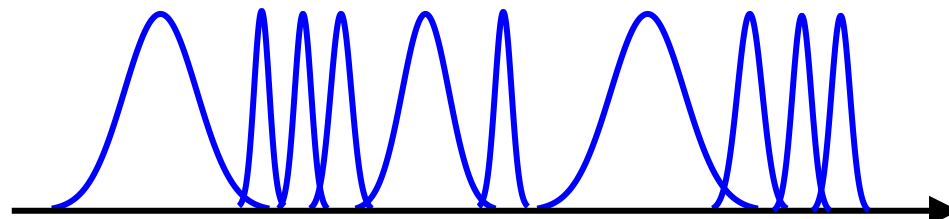


- Approximation by Gaussian kernel

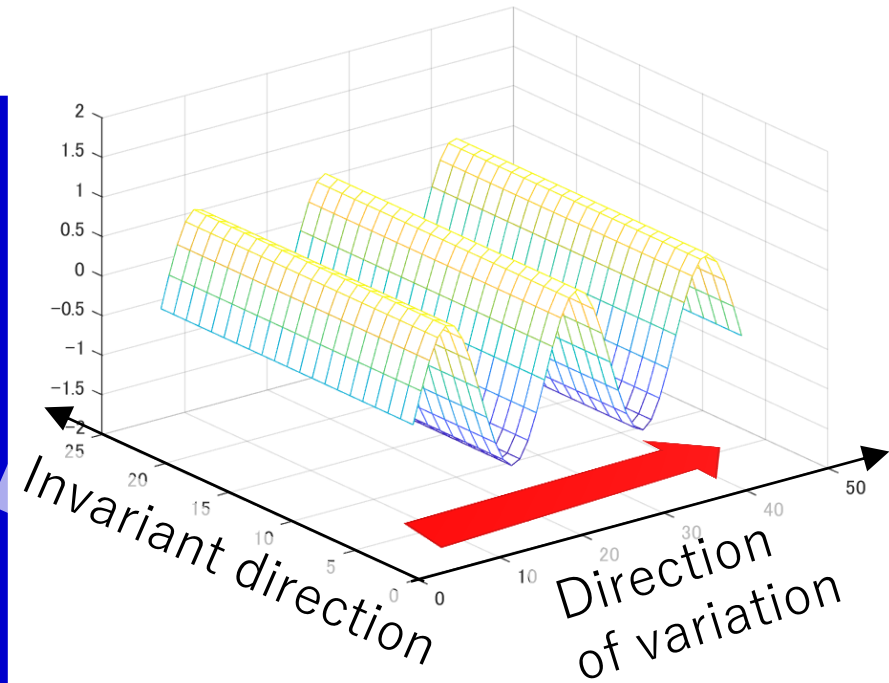


**Same Gaussian width → Less efficient**

- Approximation by NN



**Different resolution depending on location  
→ More efficient**



- Approximation by Gaussian kernel  
Kernel method cannot specify informative directions.  
→ Curse of dimensionality
- Approximation by NN  
Informative directions is extracted in internal layers.  
→ Avoids curse of dimensionality

# Hölder, Sobolev, Besov space

$$\Omega = [0, 1]^d \subset \mathbb{R}^d$$

- Hölder space ( $\mathcal{C}^\beta(\Omega)$ )

$$\|f\|_{\mathcal{C}^\beta} = \max_{|\alpha| \leq m} \|\partial^\alpha f\|_\infty + \max_{|\alpha|=m} \sup_{x \in \Omega} \frac{|\partial^\alpha f(x) - \partial^\alpha f(y)|}{|x - y|^{\beta-m}}$$

- Sobolev space ( $W_p^k(\Omega)$ )

$$\|f\|_{W_p^k} = \left( \sum_{|\alpha| \leq k} \|D^\alpha f\|_{L^p(\Omega)}^p \right)^{\frac{1}{p}}$$

- Besov space ( $B_{p,q}^s(\Omega)$ ) ( $0 < p, q \leq \infty, 0 < s \leq m$ )

Spatial inhomogeneity

$$\omega_m(f, t)_p := \sup_{\|h\| \leq t} \left\| \sum_{j=0}^m (-1)^{m-j} \binom{m}{j} f(\cdot + jh) \right\|_{L^p(\Omega)},$$

$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \left( \int_0^\infty [t^{-s} \omega_m(f, t)_p]^q \frac{dt}{t} \right)^{1/q}.$$

Smoothness

# Hölder, Sobolev, Besov space

$$\Omega = [0, 1]^d \subset \mathbb{R}^d$$

- Hölder space ( $C^\beta(\Omega)$ )

**Intuition**

$$\|f\|_{C^\beta} = \max_{|\alpha| \leq m} \|\partial^\alpha f\|_\infty + \max_{|\alpha|=m} \sup_{x, y \in \Omega} \frac{|\partial^\alpha f(x) - \partial^\alpha f(y)|}{|x - y|^\beta}$$

**Smoothness**

$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \|D^s f\|_{L^p(\Omega)}$$

**Uniformity of smoothness**

- Besov space ( $B_{p,q}^s(\Omega)$ ) ( $0 < p, q \leq \infty, 0 < s \leq m$ )

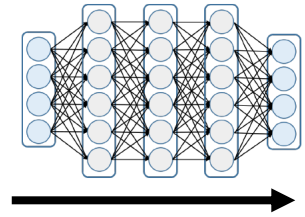
Spatial inhomogeneity

$$\omega_m(f, t)_p := \sup_{\|h\| \leq t} \left\| \sum_{j=0}^m (-1)^{m-j} \binom{m}{j} f(\cdot + jh) \right\|_{L^p(\Omega)}$$

$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \left( \int_0^\infty [t^{-s} \omega_m(f, t)_p]^q \frac{dt}{t} \right)^{1/q}$$

**Smoothness**

# Deep learning has adaptivity



- DL constructs basis function “adaptively”.  
→ Efficient learning
- Shallow learning should use “redundant” model.  
→ Inefficient learning (affected by redundant noise)

$f^\circ \in B_{p,q}^s([0,1]^d)$ : “Besov space”

[Suzuki: Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality, ICLR2019]

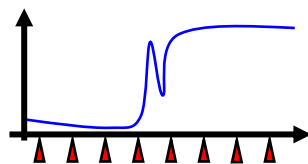
## Linear estimator (Shallow method)

e.g., kernel ridge regression:  $\hat{f}(x) = K_{x,X}(K_{X,X} + \lambda I)^{-1}Y$

$$n^{-\frac{2s - 2d(1/p - 1/2)_+}{2s + d - 2d(1/p - 1/2)_+}}$$

**Suboptimal**

( $n$ : sample size,  $p$ : uniformity of smoothness,  $s$ : smoothness)

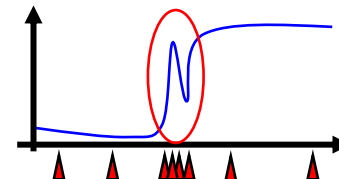


**Uniform resolution**

## Deep learning

$$n^{-\frac{2s}{2s + d}}$$

**Optimal**



**Adaptive resolution**

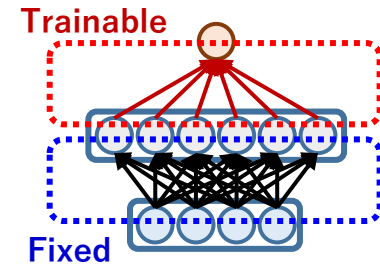
Convergence rate of estimation error (mean squared error  $\mathbb{E}[\|\hat{f} - f^*\|^2]$ )

- The rate of error decrease as the sample size  $n \rightarrow \infty$ .



## “Shallow” learning methods Kernel ridge regression:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^\infty} \sum_{i=1}^n (y_i - \beta^\top \psi(x_i))^2 + \lambda \beta^\top \beta$$



$\psi : \mathcal{X} \rightarrow \mathbb{R}^\infty$  (feature map)  
**Fixed**

$K_{X,X} = (\psi(x_i)^\top \psi(x_j))_{i,j=1}^{n,n}$   
Gram matrix (kernel function)

$$\hat{f}(x) = K_{x,X} (K_{X,X} + \lambda I)^{-1} \underline{Y}$$

(see also [Imaizumi&Fukumizu, 2019])

Nadaraya-Watson estimator:

$$\hat{f}(x) = \frac{\sum_{i=1}^n k(x_i, x) \underline{y}_i}{\sum_{i=1}^n k(x_i, x)}$$

Linear estimator: linear to the observation  $Y = (y_i)_{i=1}^n$ .

$$X_n = (x_1, \dots, x_n)$$

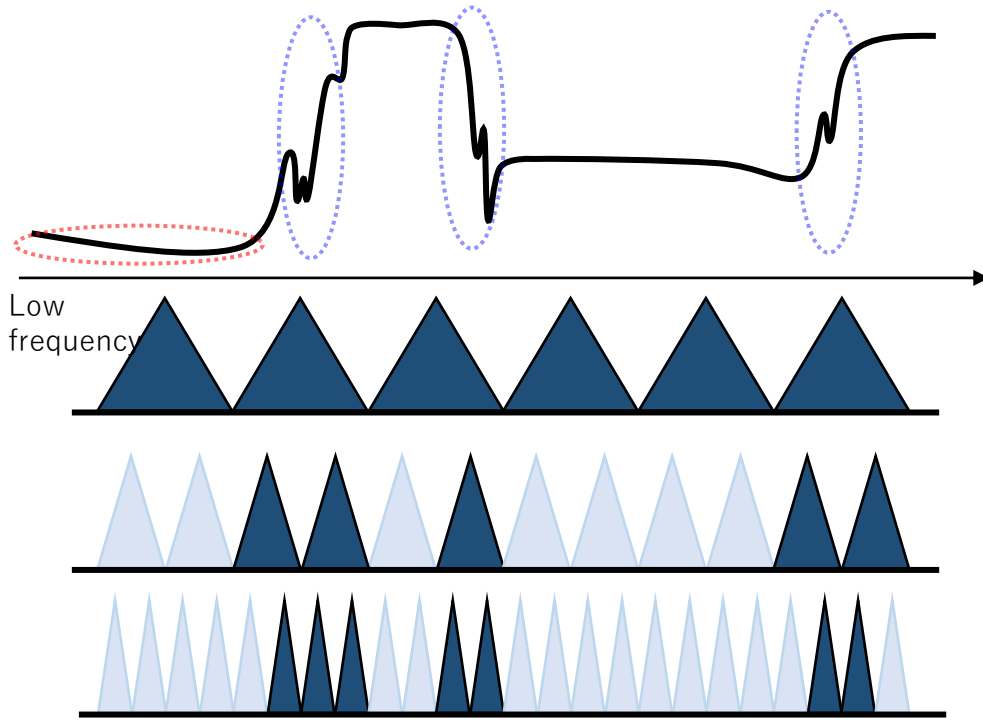
$$\hat{f}(x) = \sum_{i=1}^n \varphi_i(x; X_n) \underline{y}_i$$

linear

## Example

- Kernel ridge estimator
- Sieve estimator
- Nadaraya-Watson estimator
- k-NN estimator

# Relation to sparse estimation



**Sparse linear combination of wavelet basis functions is effective to capture input-dependent smoothness**

$$f = \sum_{k \in \mathbb{N}_+} \alpha_k \phi_k$$

Wavelet basis decomposition

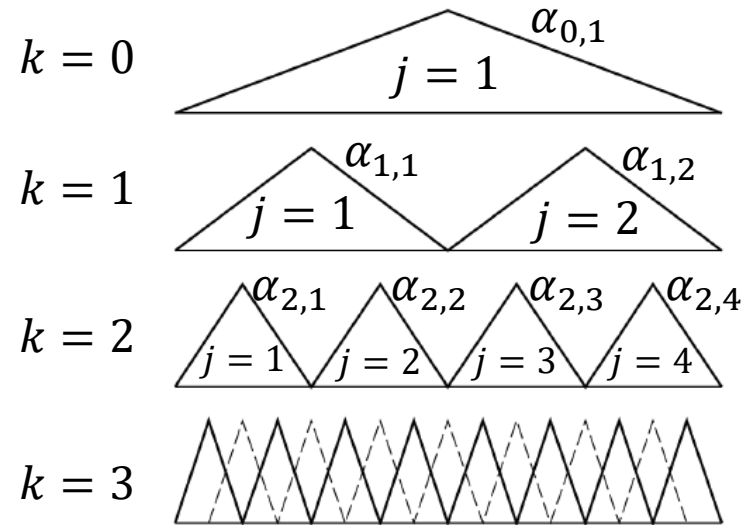
Small  $p$  = Sparse coefficient



Non-uniform smoothness over the space

## Wavelet basis

Resolution



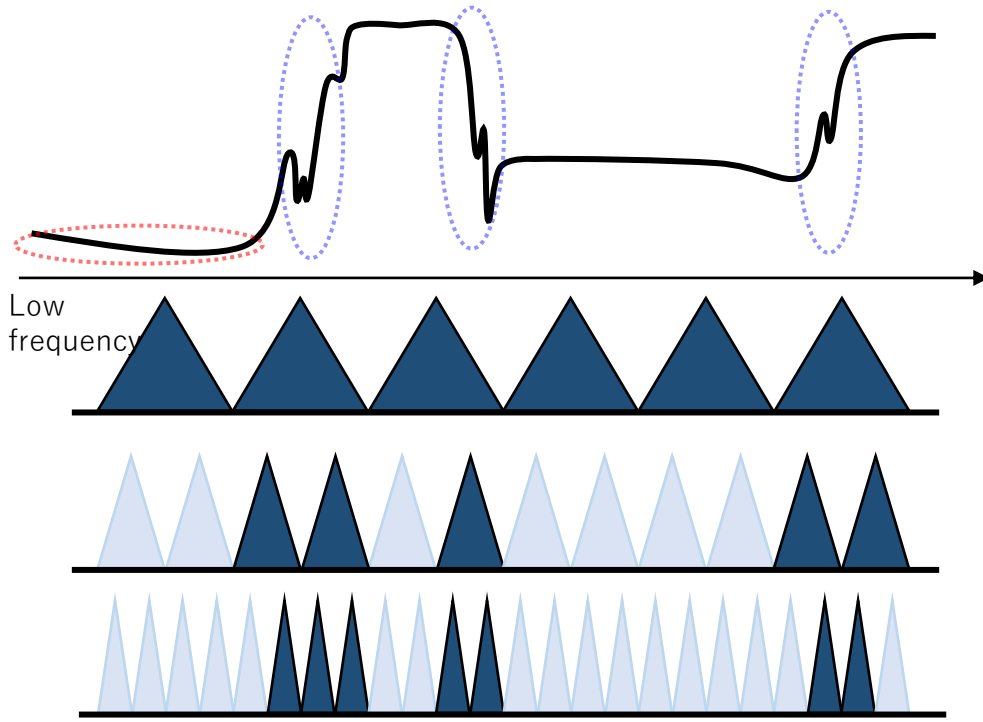
Multiresolution expansion

$$\|f\|_{B_{p,q}^s} = \left( \sum_{k \in \mathbb{N}_+} |\alpha_k|^p \right)^{1/p} \quad (0 < p)$$

(informal)

$$\|f\|_{B_{p,q}^s} \simeq \left[ \sum_{k=0}^{\infty} \{2^{sk} (2^{-kd} \sum_{j \in J(k)} |\alpha_{k,j}|^p)^{1/p}\}^q \right]^{1/q}$$

# Relation to sparse estimation



High frequency

**Sparse linear combination of wavelet basis functions is effective to capture input-dependent smoothness**

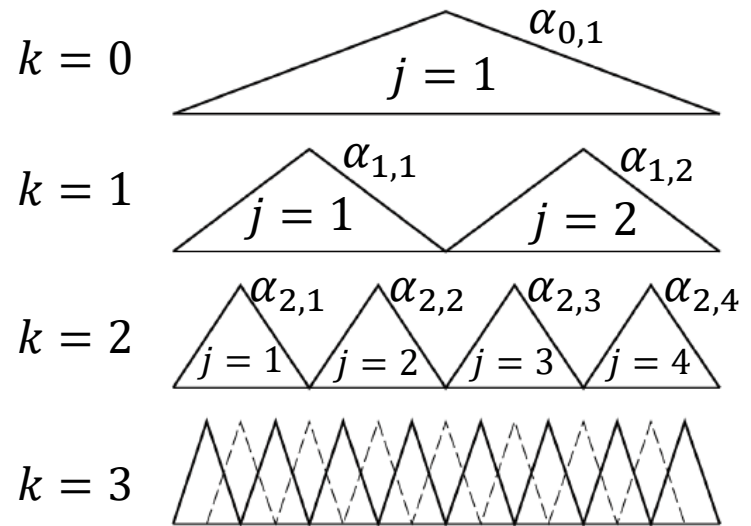
$$f = \sum_{k \in \mathbb{N}_+} \alpha_k \phi_k$$

Wavelet basis decomposition

Small  $p$  = Sparse coefficient

## Wavelet basis

Resolution



resolution expansion

$$\|f\|_{B_{p,q}^s} = \left( \sum_{k \in \mathbb{N}_+} |\alpha_k|^p \right)^{1/p} \quad (0 < p)$$

$$\|f\|_{B_{p,q}^s} \simeq \left[ \sum_{k=0}^{\infty} \{2^{sk} (2^{-kd} \sum_{j \in J(k)} |\alpha_{k,j}|^p)^{1/p}\}^q \right]^{1/q}$$

Wavelet-shrinkage



Non-uniform smoothness over the space

# Proof strategy

- **Step 1** : Find basis function expansion.

$$f^\circ \in \mathcal{F} \quad \Rightarrow \quad f^\circ(x) = \sum_{i=1}^{\infty} \alpha_i \psi_i(x)$$

$$f^\circ = \sum_{i=1}^N \alpha_i \psi_i + \underbrace{\sum_{i=N+1}^{\infty} \alpha_i \psi_i}_{\text{Adaptive approximation by B-Spline}}$$

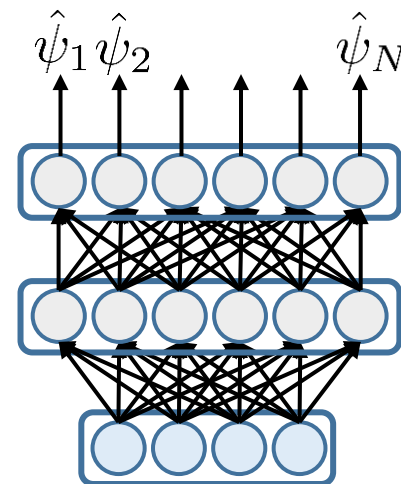
∴ Adaptive approximation by B-Spline [DeVore & Popov, 1988; Dung, 2011]

$$\|\cdot\|_{L^2} \leq N^{-s/d}$$

- **Step 2** : Approximate basis functions.

$\psi_i \simeq \hat{\psi}_i$  : Approximation by DNN.

$$\Rightarrow \quad \check{f} = \sum_{i=1}^N \alpha_i \hat{\psi}_i \quad : \text{linear combination.}$$



- **Step 3**: Combine the bounds of Step 1 and 2.

$$\begin{aligned} \|f^\circ - \check{f}\|_{L^2} &\leq \sum_{i=1}^N |\alpha_i| \underbrace{\|\psi_i - \hat{\psi}_i\|_{L^2}}_{\leq O(e^{-L})} + \underbrace{\left\| \sum_{i=N+1}^{\infty} \alpha_i \psi_i \right\|_{L^2}}_{\leq N^{-s/d}} \\ &\lesssim N^{-2s/d} \end{aligned}$$

# Bias variance decomposition

[Local Rademacher complexity bound]

$$\mathbb{E}[\|f^\circ - \hat{f}\|_{L^2(P_X)}^2] \lesssim \underbrace{\frac{S[L \log(BW) + \log(Ln)]}{n}}_{\text{Variance}} + \underbrace{\inf_{f \in \mathcal{F}(L, W, S, B)} \|f - f^\circ\|_{L^2(P_X)}^2}_{\text{Bias}}$$

Classic analysis on nonparametric regression.

See [Schmidt-Hieber, 2019; Hayakawa&Suzuki,2020] for the analysis of DNN.

Depth

Width

Sparseness

(# of non-zero parameters)

Upper bound of absolute value  
of each parameter

$$L = O(\log(N)), W = O(N), S = O(N \log(N)), B = O(N^{(d/p-s)_+})$$

$$\text{Variance} = \frac{N \log(N)^3}{n}$$

$$\text{Bias} = N^{-2s/d}$$

Choose the network size  $N$  to balance bias and variance:

$$N = \tilde{O}(n^{2d/(2s+d)}).$$

$$\text{Predictive error} = n^{-\frac{2s}{2s+d}} \log(n)^3$$

# Hardness: Convex hull argument

A function with a property that is destroyed by convex combination is hard to estimate by linear estimators.

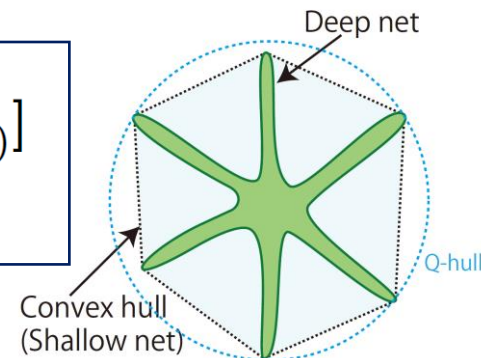
e.g., “Spatial inhomogeneity of smoothness”

## Theorem

$$\inf_{\hat{f}: \text{Linear}} \sup_{f^\circ \in \mathcal{F}} \mathbb{E}[\|\hat{f} - f^\circ\|_{L_2(P)}^2] = \inf_{\hat{f}: \text{Linear}} \sup_{f^\circ \in \text{conv}(\mathcal{F})} \mathbb{E}[\|\hat{f} - f^\circ\|_{L_2(P)}^2]$$

(This can be extended to Q-hull for the fixed design setting)

[Donoho & Johnstone, 1994] [Satoshi Hayakawa and Taiji Suzuki: 2020]



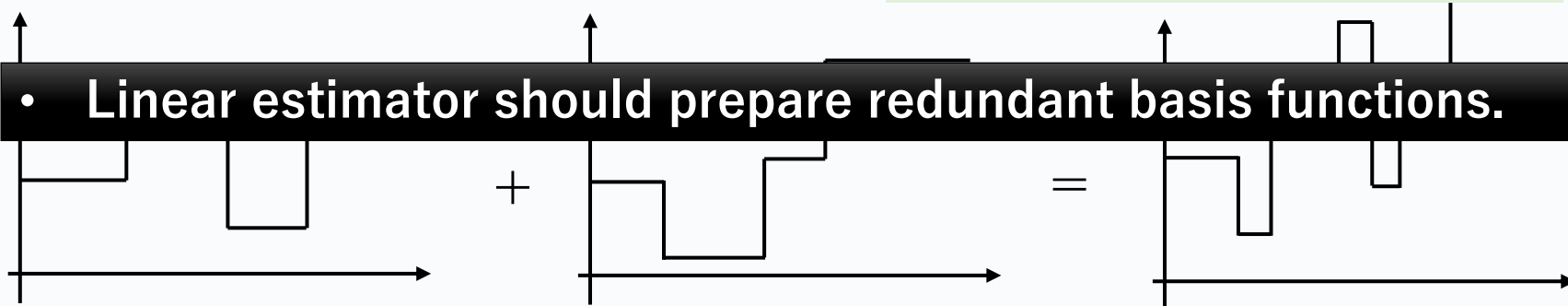
### Ex. Piecewise constant function with 3 jumps.

0.5x

0.5x

Deep:  $1/n$ , Kernel:  $1/\sqrt{n}$

• Linear estimator should prepare redundant basis functions.



3 jumps

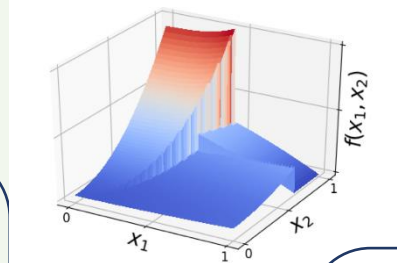
3 jumps

6 jumps

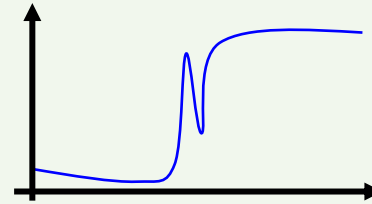
Reduced rank regression

$$Y_i = U V X_i$$

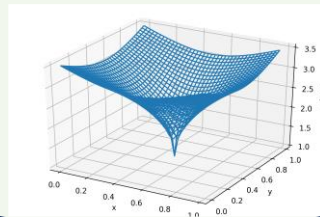
Piece-wise smooth



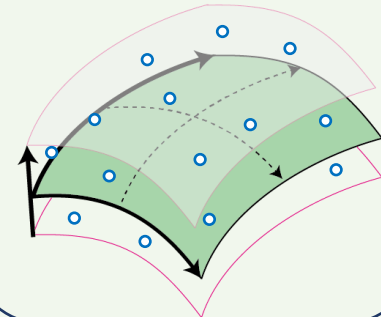
Besov space



Variable smooth Besov space



Low dim. data

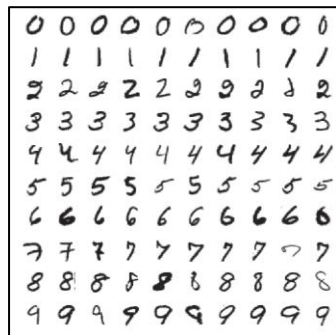


**Non-convexity  
sparsity**

Estimation error bound :

$$n^{-\frac{2s}{2s+d}}$$

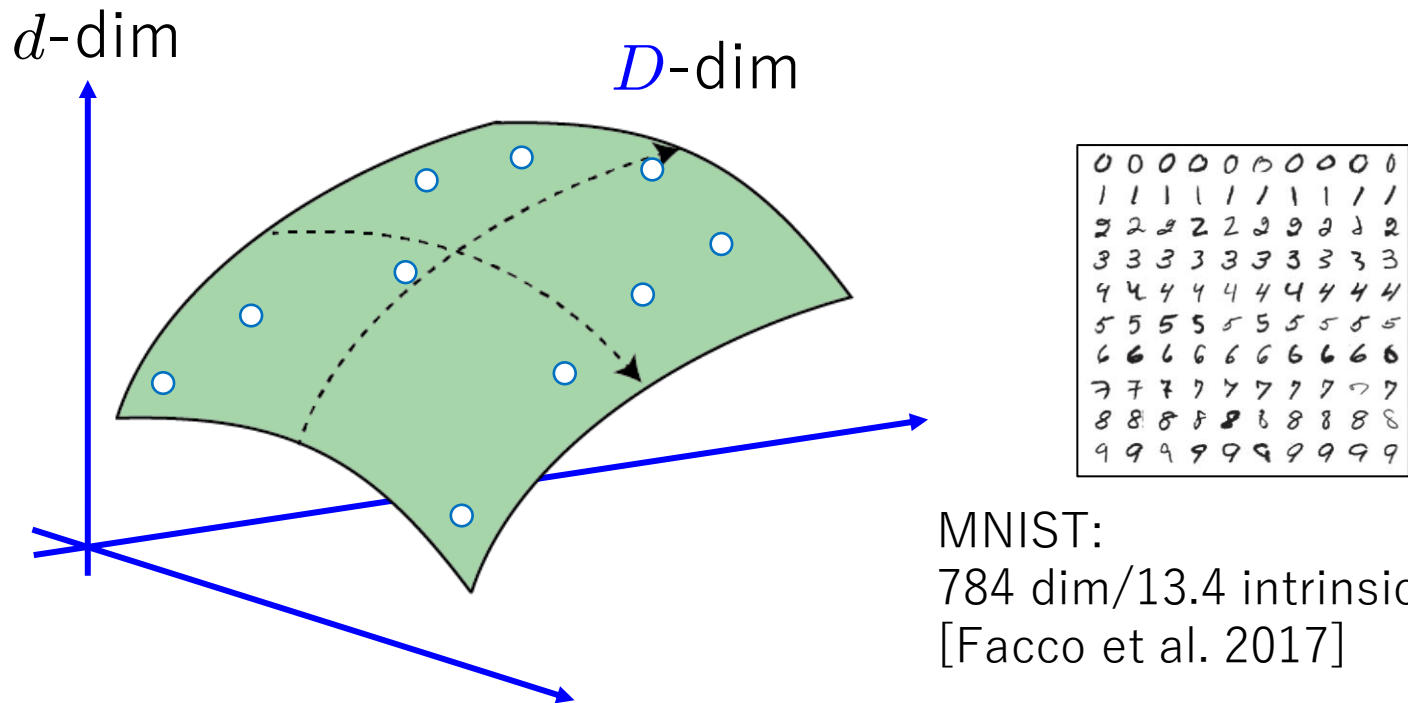
→ Curse of dimensionality



MNIST: 784 dim/13.4 intrinsic-dim  
[Facco et al. 2017]



# Dimensionality: Manifold regression <sup>41</sup>

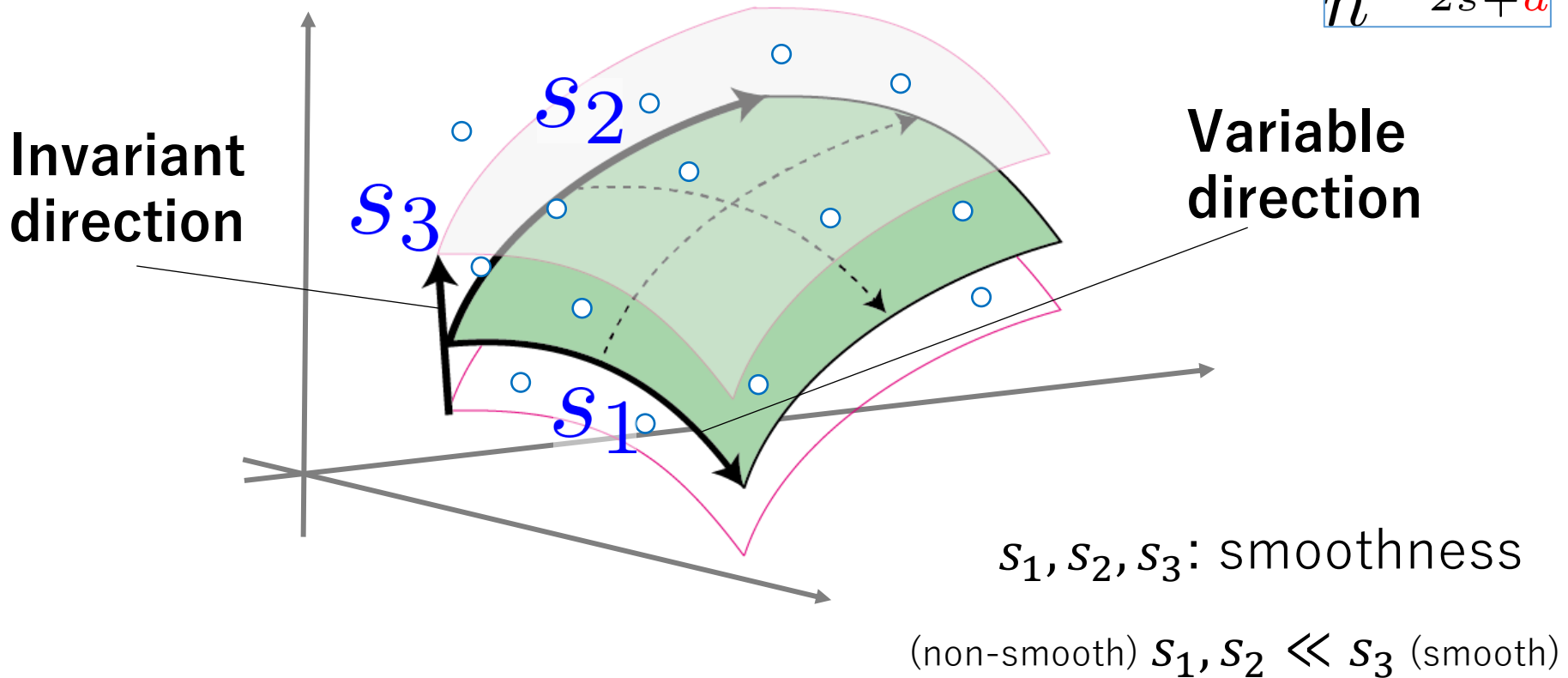


- **Classic nonparametric method:** Bickel & Li (2007); Yang & Tokdar (2015); Yang & Dunson (2016).
- **Deep learning:** Nakada & Imaizumi (2019); Schmidt-Hieber (2019); Chen et al. (2019).

$$n^{-\frac{2s}{2s+D}}$$

# More realistic setting

$$n^{-\frac{2s}{2s+d}}$$



Data are hardly distributed **exactly** on low-dim manifold.

- Smoothness could depend on directions.
- Local coordinate.

$$n^{-\frac{2\tilde{s}}{2\tilde{s}+1}}$$

$$\tilde{s} = (s_1^{-1} + \dots + s_d^{-1})^{-1}$$

# Anisotropic Besov space

Def. (Anisotropic Besov space)

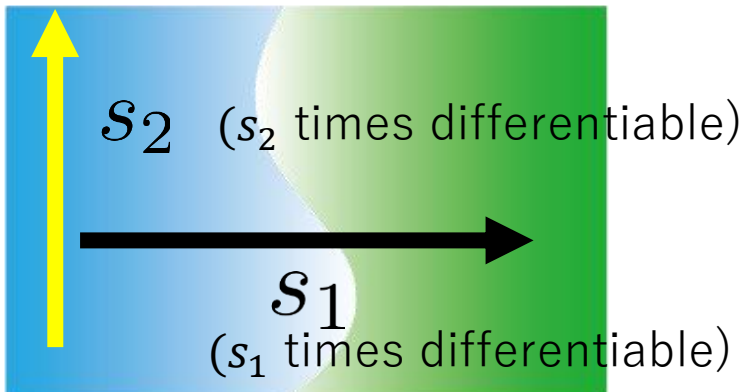
$$\Delta_h^r(f)(x) := \Delta_h^{r-1}(f)(x+h) - \Delta_h^{r-1}(f)(x), \quad (\text{finite difference})$$
$$\Delta_h^0(f)(x) := f(x) \quad (h \in \mathbb{R}^d)$$

$$w_{r,p}(f, t) = \sup_{h \in \mathbb{R}^d: |h_i| \leq t_i} \|\Delta_h^r(f)\|_p \quad (\text{modulus of smoothness})$$

$$s = (s_1, \dots, s_d) \in \mathbb{R}_{++}^d$$

$$|f|_{B_{p,q}^s} := \begin{cases} \left( \sum_{k=0}^{\infty} [2^k w_{r,p}(f, (2^{-k/s_1}, \dots, 2^{-k/s_d}))]^q \right)^{1/q} & (q < \infty), \\ \sup_{k \geq 0} 2^k w_{r,p}(f, (2^{-k/s_1}, \dots, 2^{-k/s_d})) & (q = \infty). \end{cases}$$

$$\|f\|_{B_{p,q}^s} := \|f\|_p + |f|_{B_{p,q}^s}$$



$L_p$ -norm of  $s_i$ -times derivative.

- $s_i$ : smoothness to the  $i$ -th coordinate
- $p$ : Uniformity of smoothness over the input space.

# Anisotropic Besov space

Def. (Anisotropic Besov space)

$$\Delta_h^r(f)(x) := \Delta_h^{r-1}(f)(x+h) - \Delta_h^{r-1}(f)(x), \quad (\text{finite difference})$$

$$\Delta_h^0(f)(x) := f(x) \quad (h \in \mathbb{R}^d)$$

$$w_{r,p}(f, t) = \sup_{h \in \mathbb{R}^d: |h_i| \leq t_i} \|\Delta_h^r(f)\|_p \quad (\text{modulus of smoothness})$$

$$s = (s_1, \dots, s_d) \in \mathbb{R}_{++}^d$$

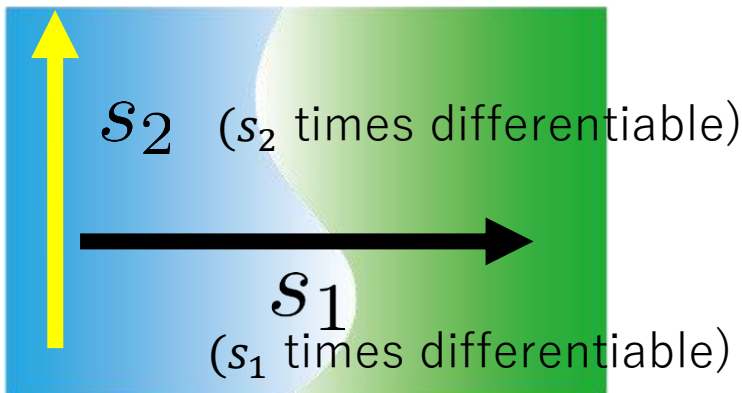
$$|f|_{B_{p,q}^s} := \begin{cases} \left( \sum_{k=0}^{\infty} [2^k w_{r,p}(f, (2^{-k/s_1}, \dots, 2^{-k/s_d}))]^q \right)^{1/q} & (q < \infty), \\ \sup_{k \geq 0} 2^k w_{r,p}(f, (2^{-k/s_1}, \dots, 2^{-k/s_d})) & (q = \infty). \end{cases}$$

$$\|f\|_{B_{p,q}^s} := \|f\|_p + |f|_{B_{p,q}^s}$$

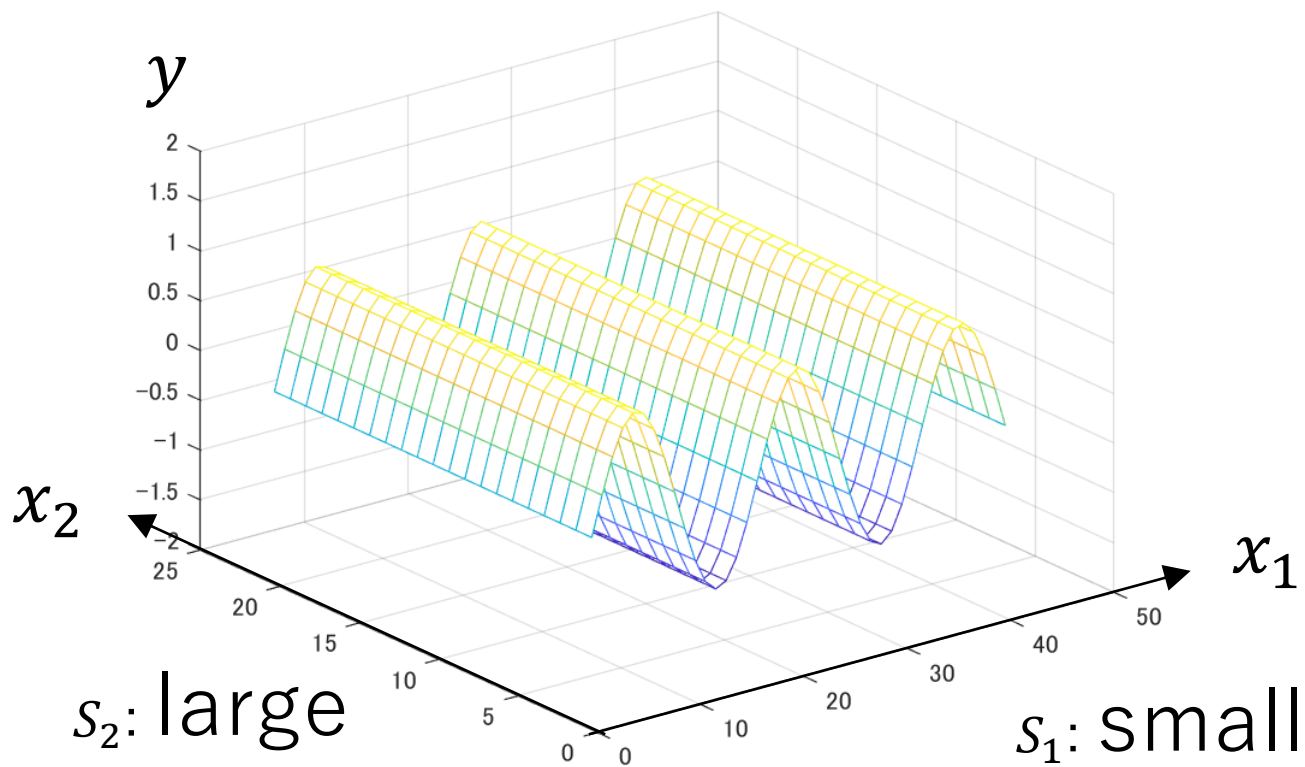
$$\|f\|_{B_{p,q}^s} = \|f\|_{L^p} + \sum_{i=1}^d \left\| \frac{\partial^{s_i} f}{\partial x_i^{s_i}} \right\|_{L^p}$$

$L_p$ -norm of  $s_i$ -times derivative.

- $s_i$ : smoothness to the  $i$ -th coordinate
- $p$ : Uniformity of smoothness over the input space.



$$\|f\|_{B_{p,q}^s} = \|f\|_{L^p} + \sum_{i=1}^d \left\| \frac{\partial^{s_i} f}{\partial x_i^{s_i}} \right\|_{L^p}$$



# Estimation error bound

$$f^\circ(x) = h_H \circ \cdots \circ h_1(x) \quad h_\ell \in B_{p,q}^{(s_1^{(\ell)}, \dots, s_{m_\ell}^{(\ell)})}([0, 1]^{m_\ell}) \quad h_\ell : [0, 1]^{m_\ell} \rightarrow [0, 1]^{m_{\ell+1}}$$

$$\hat{f} = \arg \min_{f: \text{deep neural-net}} \sum_{i=1}^n (y_i - f(x_i))^2$$

(least squares estimator)

✧ Here, we do not discuss optimization ability.

## Theorem

$$\text{Let } \tilde{s}^{(\ell)} := \left( \frac{1}{s_1^{(\ell)}} + \cdots + \frac{1}{s_{m_\ell}^{(\ell)}} \right)^{-1}, \quad \tilde{s}^{*(\ell)} := \tilde{s}^{(\ell)} \prod_{k=\ell+1}^H [(\min_j s_j^{(k)} - 1/p) \wedge 1]$$

$$\mathbb{E}[\|\hat{f} - f^\circ\|_{L^2(P_X)}^2] \lesssim \max_{\ell \in [H]} n^{-\frac{2\tilde{s}^{*(\ell)}}{2\tilde{s}^{*(\ell)}+1}} \log(n)^3$$

The rate of convergence is determined by smoothness parameters.

When  $H = 1$ ,


$$n^{-\frac{2\tilde{s}}{2\tilde{s}+1}} \\ \tilde{s} = \left( s_1^{-1} + \cdots + s_d^{-1} \right)^{-1}$$

If  $s_i$ s are small (non-smooth) toward small numbers of directions and large toward other directions, DNN can avoid the curse of dimensionality.

If  $s_2 = \cdots = s_d = \infty$ , then  $n^{-\frac{2s_1}{2s_1+1}}$ .

Comparison to linear estimator :


**DNN**

$$n^{-\frac{2\tilde{s}}{2\tilde{s}+1}} \\ \tilde{s} = (s_1^{-1} + s_2^{-1} + s_3^{-1})^{-1}$$


**Linear**

$$n^{-\frac{2s_1}{2s_1+d}}$$

(Affected by dimension)



- **Linear estimator cannot find important directions. Then, the rate of convergence is strongly affected by the most non-smooth ( $s_1$ ) parameter.**
- **We used a “convex hull argument” to show the rate of convergence.**

# Comparison to linear estimator

$$f^\circ(x) = g(Wx) \quad (W \in \mathbb{R}^{D \times d}, g \in B_{p,q}^s([0,1]^D))$$

$f^\circ$  depends only  $D$ -dimensional subspace.

**Deep**

$$n^{-\frac{2s}{2s+D}}$$

$$(n^{-\frac{2s}{2s+1}} \text{ when } D = 1)$$

**Optimal**

**Linear estimator**

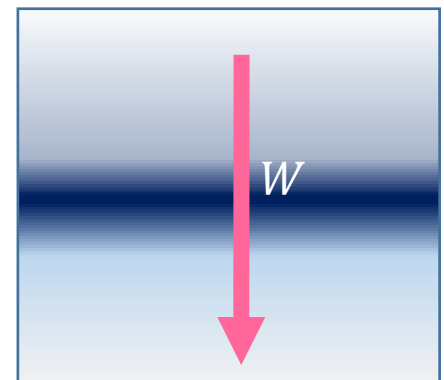
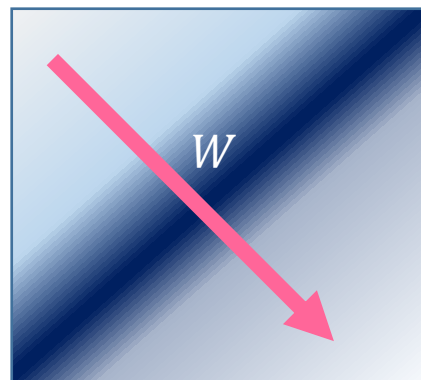
$$n^{-\frac{2(s-D/p+d/2+c)}{2(s-D/p+d/2+c)+d}}$$

$$c = 1 \text{ if } D < d/2, c = 0 \text{ if } D \geq d/2.$$

$$(n^{-\frac{2s+d}{2s+2d}} \text{ when } D = 1 \text{ and } p = 1)$$

**Suboptimal**

$\ll$



Deep can ease curse of dim.,  
but linear estimators directly  
suffers from curse of dim.

## 1. Representation ability + Generalization ability

- Universal approximator
- Depth separation
- Adaptivity of deep learning
  - Inhomogeneity of smoothness
  - Curse of dimensionality
- Foundation models
  - Diffusion model
  - Transformer

## 2. Optimization ability

- Noisy gradient descent
- Mean field Langevin
- CSQ lowerbound

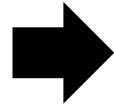


# Analysis of diffusion model

- Kazusato Oko, Shunta Akiyama, Taiji Suzuki: Diffusion Models are Minimax Optimal Distribution Estimators. ICML2023.

# Diffusion model

「An astronaut riding a horse in a photorealistic style」



DALL·E: [Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, Ilya Sutskever: Zero-Shot Text-to-Image Generation. ICML2021.]  
DALL·E2: [Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, Mark Chen: Hierarchical Text-Conditional Image Generation with CLIP Latents. arXiv:2204.06125]



Stable diffusion, 2022.



Jason Allen "Théâtre D'opéra Spatial" generated by **Midjourney**. Colorado State Fair's fine art competition, 1<sup>st</sup> prize in digital art category



Generated by NovelAI

# Movie generation

- SORA (OpenAI, 2024)
  - Diffusion Transformer



A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp and reflective, creating a mirror effect of the colorful lights. Many pedestrians walk about.

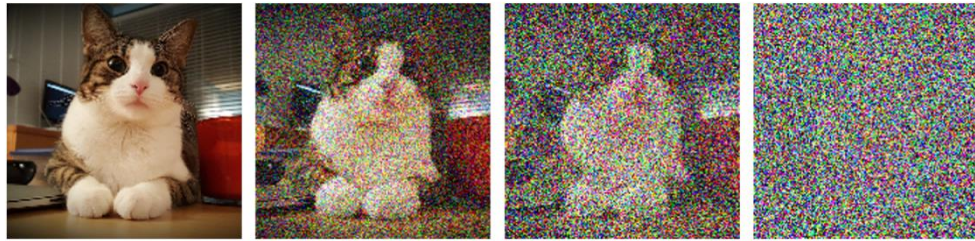


# Diffusion model

[Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Song et al., 2020; Ho et al., 2020; Vahdat et al., 2021]

**Forward process** : Convert the target distribution to a noise distribution (e.g., Gaussian)

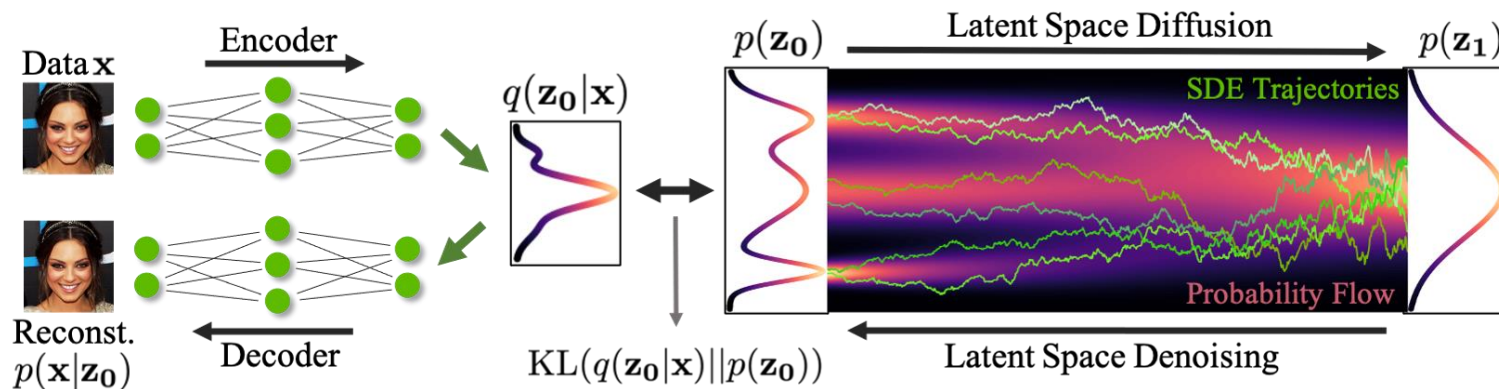
$$dX_t = -X_t dt + \sqrt{2}dB_t$$



$$dY_t = (Y_t + 2\nabla \log(p_{\bar{T}-t}(Y_t)))dt + \sqrt{2}dB_t$$

$(Y_t \sim X_{\bar{T}-t})$

**Reverse process** : Convert the noise distribution to the target distribution



[Vahdat, Kreis, Kautz: Score-based Generative Modeling in Latent Space. arXiv:2106.05931]

# Forward process

Forward process:

$$(X_{t+\eta} = X_t - \eta X_t + \sqrt{2\eta}\xi_t)$$

$$dX_t = -X_t dt + \sqrt{2}dB_t$$

OU process

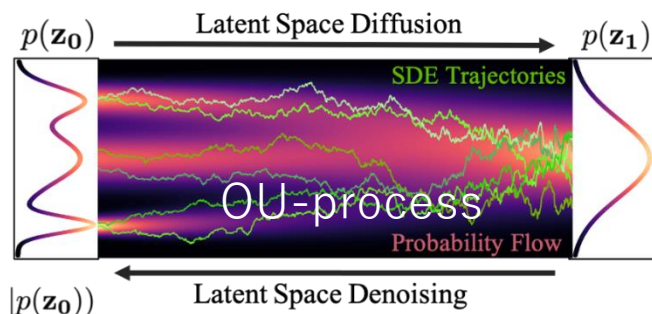
$$p_t = \text{Law}(X_t) \longrightarrow p_t = \int N(\mu_t X_0, \sigma_t^2 I) p_0(dX_0)$$

where  $\mu_t = \exp(-t)$ ,  $\sigma_t^2 = 1 - \exp(-2t)$ .

$$p_t(x) = \int p_0(y) \frac{1}{\sigma_t^d (2\pi)^{\frac{d}{2}}} \exp\left(-\frac{\|x - \mu_t y\|^2}{2\sigma_t^2}\right) dy$$

The forward process converges to the noise distribution (standard normal) exponentially:

$$\text{KL}(p_t || N(0, I)) \leq \exp(-2t) \text{KL}(p_0 || N(0, I))$$



Standard normal

[Vahdat, Kreis, Kautz: Score-based Generative Modeling in Latent Space. arXiv:2106.05931]

## Reverse process:

$$Y_0 \sim p_{\bar{T}}$$

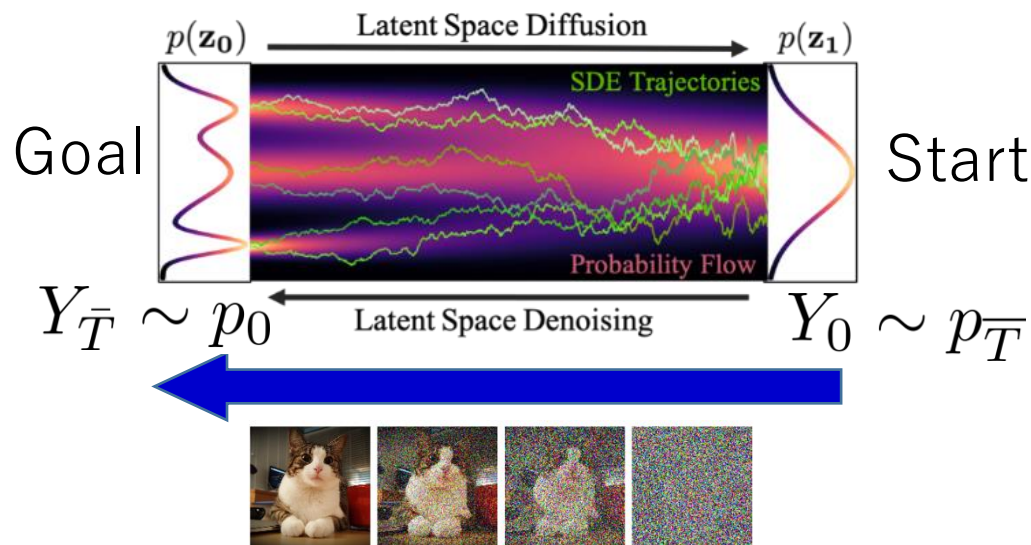
[Haussmann & Pardoux, 1986]

$$dY_t = (Y_t + 2\nabla \log(p_{\bar{T}-t}(Y_t)))dt + \sqrt{2}dB_t \quad (t \in [0, \bar{T}])$$

Fact :  $Y_t$ 's distribution =  $X_{\bar{T}-t}$ 's distribution

That is,  $Y_t \sim p_{\bar{T}-t}$

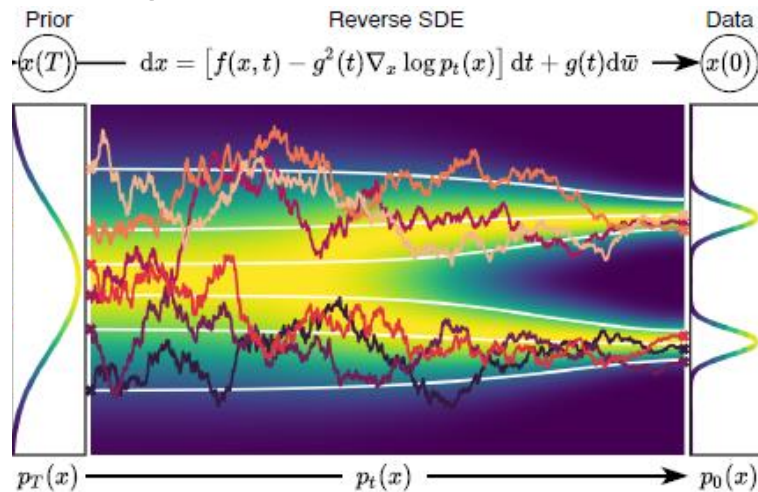
By following the forward process in reverse, noise that follows a (nearly) normal distribution can be gradually modified to reproduce the original distribution of the images.



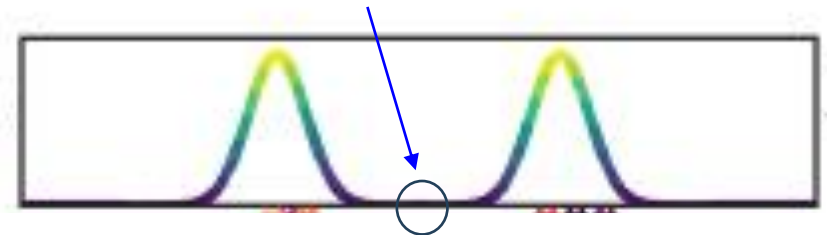
# Benefit of diffusion model

We can sample from multimodal distribution efficiently.

➤ “Easy distribution” → “Difficult distribution”



If we try to sample directly from the original distribution, then it could not get over the “gap”.



- Even though the score of the original distribution is complex, the distribution of the diffused  $X_t$  is smooth → easy to estimate → easy to generalize.
- Learning is more stable because it uses information of the intermediate distribution  $p_t$  instead of directly learning end-to-end mapping from the noise to the source distribution.



[<https://github.com/Kei18/tiny-tiny-diffusion>]

This is not generating an image of a dinosaur but the shape of the density function looks like a dinosaur. Each point corresponds to each image ( $X_t$ ).



## Reverse process:

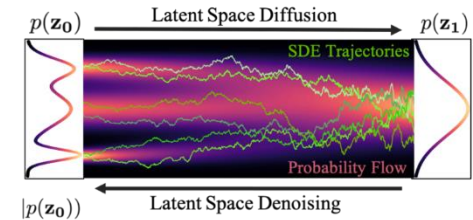
$$Y_0 \sim p_{\bar{T}} \quad (\text{unknown})$$

$$dY_t = (Y_t + 2 \nabla \log(p_{\bar{T}-t}(Y_t)))dt + \sqrt{2}dB_t \quad (\text{unknown})$$

$$\Rightarrow Y_t \sim p_{\bar{T}-t}$$

$$Y_{\bar{T}} \sim p_0$$

[Haußmann & Pardoux, 1986]



Approximated process (generative model):

$$\hat{Y}_0 \sim N(0, I)$$

$(N(0, I)$  is close to  $p_{\bar{T}}$ )

$$d\hat{Y}_t = (\hat{Y}_t + 2\hat{s}(\hat{Y}_t, \bar{T} - t))dt + \sqrt{2}dB_t$$

### Theorem (Girsanov's theorem)

If  $\hat{Y}_0 \sim p_{\bar{T}}$ , then

$$\text{KL}(p_0 || p_{\hat{Y}_{\bar{T}}}) \leq \frac{1}{4} \int_0^{\bar{T}} \mathbb{E}_{Y_t} [\|\nabla \log(p_{\bar{T}-t}(Y_t)) - \hat{s}(Y_t, \bar{T} - t)\|^2] dt$$

$\Rightarrow$  It suffices to estimate the score function  $\nabla \log(p_t)$  as accurate as possible.

$$\begin{aligned}
 & \int_0^{\bar{T}} \mathbb{E}_{Y_t} [\|\nabla \log(p_{\bar{T}-t}(Y_t)) - \hat{s}(Y_t, \bar{T} - t)\|^2] dt \\
 & \quad \text{Unknown. We cannot calculate.} \\
 = & \int_0^{\bar{T}} \mathbb{E}_{X_t} [\|\nabla \log(p_t(X_t)) - \hat{s}(X_t, t)\|^2] dt \quad (X_{\bar{T}-t} \text{ と } Y_t \text{ は同じ分布}) \\
 = & \int_0^{\bar{T}} \mathbb{E}_{X_t} [\|\nabla \log(p_t(X_t))\|^2 - 2\langle \nabla \log(p_t(X_t)), \hat{s}(X_t, t) \rangle + \|\hat{s}(X_t, t)\|^2] dt \\
 = & \int_0^{\bar{T}} \mathbb{E}_{X_t} \left[ \underbrace{-2 \left\langle \frac{\nabla p_t(X_t)}{p_t(X_t)}, \hat{s}(X_t, t) \right\rangle + \|\hat{s}(X_t, t)\|^2}_{\text{score matching}} \right] dt + (\text{const})
 \end{aligned}$$

$$\begin{aligned}
 & \int -2 \left\langle \frac{\nabla p_t(x_t)}{p_t(x_t)}, \hat{s}(x_t, t) \right\rangle p_t(x_t) dx_t + \mathbb{E} [\|\hat{s}(X_t, t)\|^2] \\
 = & \int -2 \langle \nabla p_t(x_t), \hat{s}(x_t, t) \rangle dx_t + \mathbb{E} [\|\hat{s}(X_t, t)\|^2] \\
 = & \int -2 \left\langle \nabla \int p_t(x_t|x_0) p_0(x_0) dx_0, \hat{s}(x_t, t) \right\rangle dx_t + \mathbb{E} [\|\hat{s}(X_t, t)\|^2] \\
 = & \int \int -2 \langle \nabla p_t(x_t|x_0), \hat{s}(x_t, t) \rangle p_0(x_0) dx_0 dx_t + \mathbb{E} [\|\hat{s}(X_t, t)\|^2] \\
 = & \int \int -2 \left\langle \frac{\nabla p_t(x_t|x_0)}{p_t(x_t|x_0)}, \hat{s}(x_t, t) \right\rangle p_t(x_t|x_0) p_0(x_0) dx_0 dx_t + \mathbb{E} [\|\hat{s}(X_t, t)\|^2] \\
 = & \mathbb{E}_{X_0, X_t} [-2 \langle \nabla \log(p_t(X_t|X_0)), \hat{s}(X_t, t) \rangle + \|\hat{s}(X_t, t)\|^2]
 \end{aligned}$$

$$\begin{aligned} & \int_0^{\bar{T}} \mathbb{E}_{Y_t} [\|\nabla \log(p_{\bar{T}-t}(Y_t)) - \hat{s}(Y_t, \bar{T} - t)\|^2] dt \\ &= \int_0^{\bar{T}} \mathbb{E}_{X_t} [\|\nabla \log(p_t(X_t)) - \hat{s}(X_t, t)\|^2] dt \\ &= \int_0^{\bar{T}} \mathbb{E}_{X_t} [\|\nabla \log(p_t(X_t))\|^2 - 2\langle \nabla \log(p_t(X_t)), \hat{s}(X_t, t) \rangle + \|\hat{s}(X_t, t)\|^2] dt \\ &= \int_0^{\bar{T}} \mathbb{E}_{X_t} \left[ \underbrace{-2 \left\langle \frac{\nabla p_t(X_t)}{p_t(X_t)}, \hat{s}(X_t, t) \right\rangle + \|\hat{s}(X_t, t)\|^2}_{\mathbb{E}_{X_0, X_t} [-2 \langle \nabla \log(p_t(X_t|X_0)), \hat{s}(X_t, t) \rangle + \|\hat{s}(X_t, t)\|^2]} \right] dt + (\text{const}) \\ &= \int_0^{\bar{T}} \mathbb{E}_{X_t, X_0} [\|\nabla \log(p_t(X_t|X_0)) - \hat{s}(Y_t, t)\|^2] dt + (\text{const}) \end{aligned}$$

# Score matching

$$\begin{aligned} & \int_0^{\bar{T}} \mathbb{E}_{Y_t} [\|\nabla \log(p_{\bar{T}-t}(Y_t)) - \hat{s}(Y_t, \bar{T} - t)\|^2] dt \\ &= \int_0^{\bar{T}} \mathbb{E}_{X_t} [\|\nabla \log(p_t(X_t)) - \hat{s}(X_t, t)\|^2] dt \\ &= \int_0^{\bar{T}} \mathbb{E}_{X_t, X_0} [\|\nabla \log(p_t(X_t|X_0)) - \hat{s}(Y_t, t)\|^2] dt + (\text{const}) \end{aligned}$$

Observation ( $n$  data points  $D_n = \{x_i\}_{i=1}^n$ ):

$$x_i \sim p_0 \quad (i = 1, \dots, n)$$

**Empirical score matching loss:**

$$\min_{s \in \text{DNN}} \frac{1}{n} \sum_{i=1}^n \int_{\underline{T}}^{\bar{T}} \mathbb{E}_{X_t|X_0=x_i} [\|s(X_t, t) - \nabla \log p_t(X_t|x_i)\|^2] dt$$

Can be sampled via OU process

$$N(x_i e^{-t}, 1 - e^{-2t})$$

Explicit form is available

$$-\frac{(X_t - e^{-t}x_i)}{1 - e^{-2t}}$$

# Error analysis of diffusion models 61

- Reverse SDE characterization: Song et al. (2021)

[Approximation error analysis]

- KL-divergence bound via Girsanov's theorem: Chen et al. (2022)
- Error bound with LSI: Lee et al. (2022a)
  - With smoothness: Chen et al. (2022) and Lee et al. (2022b)
- Error propagation with manifold assumption: Pidstrigach (2022)

[Generalization analysis]

- Wasserstein distance bound:  $O(n^{-\frac{1}{d}})$  with manifold assumption: De Bortoli (2022)

**Q:**

1. How accurately can we estimate the score functions?
2. How strongly does the estimation error of score functions affect the final result?

# Problem setting

## Assumption 1

The true distribution  $p_0$  is supported on  $[-1,1]^d$  and

$$p_0 \in B_{p,q}^s$$

with  $s > (1/p - 1/2)_+$  as a density function on  $[-1,1]^d$ .

## Assumption 2

$p_0$  is sufficiently smooth on the edge of the support  $[-1,1]^d \setminus [-1 + n^{-\frac{1-\delta}{d}}, 1 - n^{-\frac{1-\delta}{d}}]^d$ .

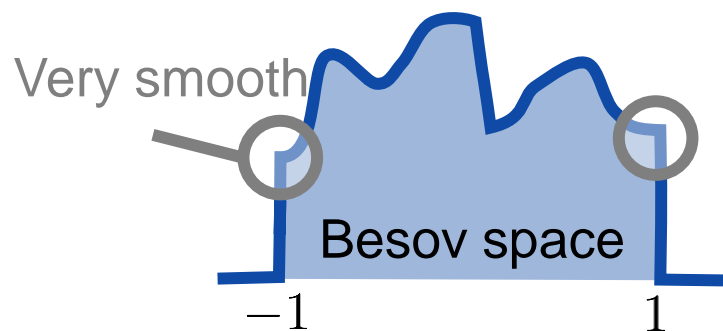
## Besov space ( $B_{p,q}^s(\Omega)$ )

$$\omega_m(f, t)_p := \sup_{\|h\| \leq t} \left\| \sum_{j=1}^m (-1)^{m-j} \binom{m}{j} f(\cdot + jh) \right\|_{L^p(\Omega)},$$

$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \left( \int_0^\infty [t^{-s} \omega_m(f, t)_p]^q \frac{dt}{t} \right)^{1/q}.$$

Smoothness

Spatial inhomogeneity



# Problem setting

## Assumption 1

The true distribution  $p_0$  is supported on  $[-1,1]^d$  and

$$p_0 \in B_{p,q}^s$$

with  $s > (1/p - 1/2)_+$  as a density function on  $[-1,1]^d$ .

## Assumption 2

$p_0$  is sufficiently smooth on the edge of the support  $[-1,1]^d \setminus [-1 + n^{-\frac{1-\delta}{d}}, 1 - n^{-\frac{1-\delta}{d}}]^d$ .

## Besov space ( $B_{p,q}^s(\Omega)$ )

Intuition

Smoothness

$\omega_m(f, t)_p$

$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \|D^s f\|_{L^p(\Omega)}$$

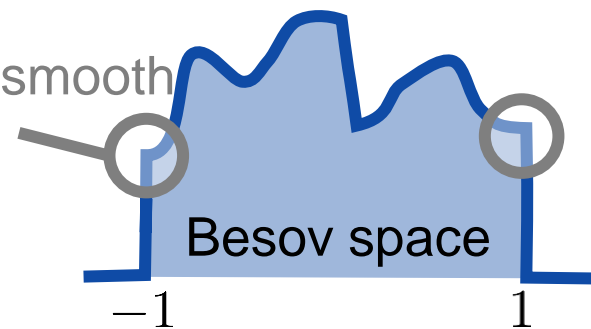
Very smooth

Uniformity of smoothness

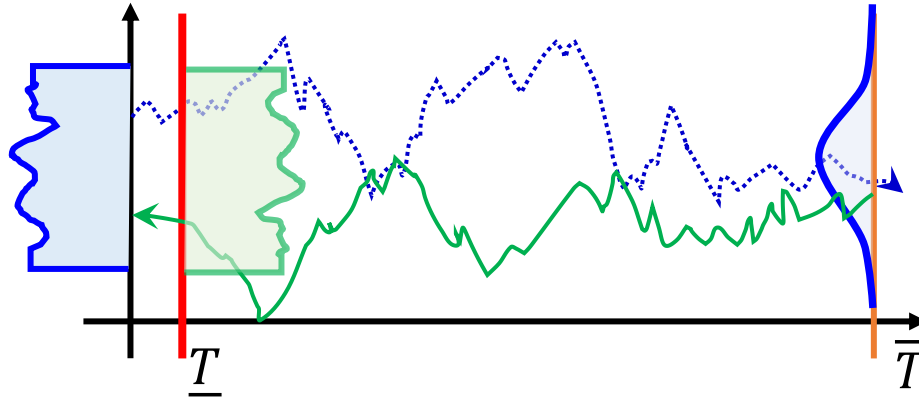
$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \left( \int_0^\infty [t^{-s} \omega_m(f, t)_p]^q \frac{dt}{t} \right)^{1/q}$$

Smoothness

Spatial inhomogeneity



# Convergence rate result



## Theorem (Estimation error in TV-distance)

Let  $\underline{T} = n^{-o(1)}$ ,  $\bar{T} = O(\log(n))$ . Then, the empirical risk minimizer  $\hat{s}$  in DNN satisfies

$$\mathbb{E}_{D_n} \left[ \text{TV}(\hat{Y}_{\bar{T}-\underline{T}}, X_0) \right] \lesssim n^{-\frac{s}{2s+d}} \log^9(n).$$

This is **minimax optimal**, that is, it holds

$$n^{-\frac{s}{2s+d}} \lesssim \inf_{\hat{\mu}:\text{estimator}} \sup_{p_0} \mathbb{E}_{D_n} [\text{TV}(\hat{\mu}, X_0)]$$

Although  $\hat{s}(x, t)$  is a function with  $d + 1$ -dimensional input, there appears “ $d$ ” in the bound instead of  $d + 1$ . This is because Gaussian convolution induces smoothness.



# B-spline basis decomposition

$$\nabla \log(p_t(x)) = \frac{\nabla p_t(x)}{p_t(x)}$$

Approximate each term by DNNs

- B-spline decomposition of a Besov function  $p_0$

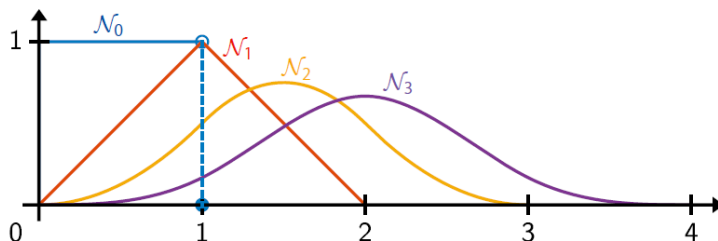
$$p_0(x) \approx \sum_{j=1}^N \alpha_j M_{a^j, b^j}^d(x)$$

$$\mathcal{N}(x) = \begin{cases} 1 & (x \in [0, 1]), \\ 0 & (\text{otherwise}) \end{cases}$$

Cardinal B-spline of order  $m$ :

$$\mathcal{N}_m(x) = \underbrace{(\mathcal{N} * \mathcal{N} * \dots * \mathcal{N})}_{m+1 \text{ times}}(x)$$

→ Piece-wise polynomial of order  $m$ .



Tensor product B-spline:

$$M_{a,b}^d(x) = \prod_{j=1}^d \mathcal{N}_m(2^{a_j} - b_j)$$

# Cardinal B-spline interpolation (DeVore & Popov, 1988) 66

Reference

- Atomic decomposition:

$$\mathcal{N}_{k,j}^{(d)}(x_1, \dots, x_d) = \prod_{i=1}^d \mathcal{N}_m(2^k x_i - j_i)$$

$f \in B_{p,q}^s$  can be decomposed into

$$f = \sum_{k \in \mathbb{N}} \sum_{j \in J(k)} \alpha_{k,j} \mathcal{N}_{k,j}^{(d)}$$

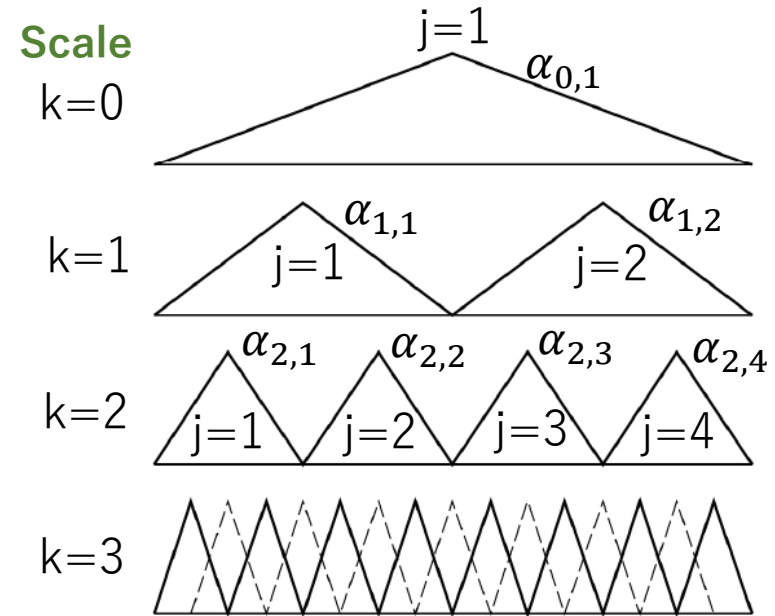
such that

(where  $J(k) = \{j \in \mathbb{Z}^d \mid -m < j_i < 2^{k_i+1} + m\}$ )

$$N(f) = \left[ \sum_{k=0}^{\infty} \left\{ 2^{sk} \left( 2^{-kd} \sum_{j \in J(k)} |\alpha_{k,j}|^p \right)^{1/p} \right\}^q \right]^{1/q} < \infty$$

$$\|f\|_{B_{p,q}^s} \simeq N(f) \quad (\text{Norm equivalence})$$

Wavelet/multi-resolution expansion



**DNN can approximate each B-spline basis efficiently.**

$$f = \sum_{k,j \in I_N} \alpha_{k,j} \mathcal{N}_{k,j}^{(d)} + O(N^{-s/d})$$

$N$  terms (should be appropriately chosen depending on  $f$ )

# Proof outline (1)

$$\nabla \log(p_t(x)) = \frac{\nabla p_t(x)}{p_t(x)}$$

Approximate each term by DNNs

- B-spline decomposition of a Besov function  $p_0$

$$p_0(x) \approx \sum_{j=1}^N \alpha_j M_{a^j, b^j}^d(x)$$

Approximation error  
 $O(N^{-s/d})$

- Diffused B-spline basis expansion of  $p_t$

$$p_t(x) = \int p_0(y) \underbrace{\frac{1}{\sigma_t^d (2\pi)^{\frac{d}{2}}} \exp\left(-\frac{\|x - \mu_t y\|^2}{2\sigma_t^2}\right)}_{=: K_t(x|y)} dy$$

Decompose

$$p_t(x) \approx \sum_{j=1}^N \alpha_j \underbrace{\int M_{a^j, b^j}^d(y) K_t(x|y) dy}_{=: E_{a^j, b^j}(x, t)}$$

**Diffused B-spline**

- We approximate Diffused B-splines by DNNs.

# Approximation error of Diffused B-spline 68

## Lemma (Approximation error of diffused B-spline)

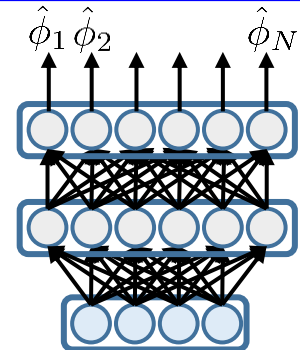
There exists a deep neural network  $\hat{\phi}: \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}^d$  such that

$$\left\| \hat{\phi}(x, t) - E_{a^j, b^j}(x, t) \right\|_{\infty} \leq \epsilon$$

with depth  $L = O(\log^4(\epsilon^{-1}))$ , width  $W_i = O(\log^6(\epsilon^{-1}))$ , sparsity (# of non-zero parameters)  $S = O(\log(\epsilon^{-1}))$ , and  $\ell^\infty$ -norm bound  $B = O(\exp(O(\log^2(\epsilon^{-1}))))$  on parameters.

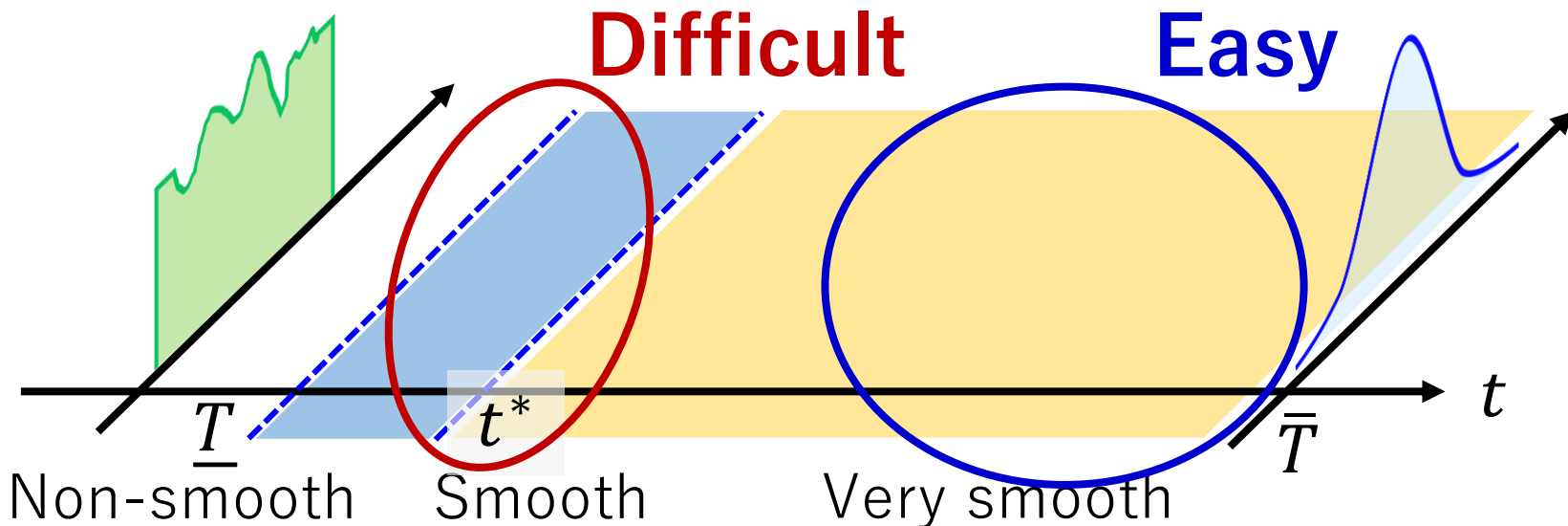
$$\check{f}_N(x, t) = \sum_{i=1}^N \alpha_i \hat{\phi}_i(x, t): \text{Deep neural network}$$

# of non-zero parameters:  $N \text{polylog}(N)$



$$\|p_t(\cdot) - \check{f}_N(\cdot, t)\|_{L^r} \leq \sum_{i=1}^N |\alpha_i| \underbrace{\|\phi_i(\cdot, t) - \hat{\phi}_i(\cdot, t)\|_{L^r}}_{\leq O(e^{-L})} + \underbrace{\left\| \sum_{i=N+1}^{\infty} \alpha_i \phi_i(\cdot, t) \right\|_{L^r}}_{\leq N^{-s/d}}$$

# Error bound of score



- Bound by diffused B-spline approximation

$$\|p_t - \check{f}_N(\cdot, t)\|_{L^r(X_t)} \lesssim N^{-s/d} \|p_0\|_{B_{p,q}^s}$$

$$p_t(x) \approx \sum_{j=1}^N \alpha_j E_{a^j, b^j}(x, t)$$

- Similar argument is applied to  $\nabla p_t$ :

$$\|\nabla \log p_t - \dot{f}_N(\cdot, t)\|_{L^2}^2 \lesssim \frac{N^{-2s/d} \log(N)}{\sigma_t^2}$$

$$(\sigma_t^2 = 1 - \exp(-2t))$$

- A tighter bound on the smooth part ( $t > t_*$ )

$$\|p_t\|_{W_p^k} = \sum_{|\alpha| \leq k} \left\| \frac{\partial^\alpha p_t}{\partial x^\alpha} \right\|_{L^p} \lesssim \sigma_t^{-k} \left( \leq t_*^{-\frac{k}{2}} \right)$$

- Useful for W1 bound.
- Smoothness around the edge (A2) is not required.

$$\Rightarrow \|p_t - \check{f}_{N'}\|_{L^2(X_t)}^2 \lesssim N'^{-2k/d} t_*^{-k}$$

(take  $k = s + 1$ )

# Error decomposition

$$\text{TV}(X_0, \hat{Y}_{\bar{T}-\underline{T}}) \leq \sqrt{\frac{1}{2} \text{KL}(X_0 \| \hat{Y}_{\bar{T}-\underline{T}})} \quad (\text{Pinsker's inequality})$$

Score matching loss

$$\text{TV}(X_0, \hat{Y}_{\bar{T}-\underline{T}}) \lesssim \left[ \int_{t=\underline{T}}^{\bar{T}} \mathbb{E}_{X_t \sim p_t} [\|\hat{s}(X_t, t) - \nabla \log p_t(X_t)\|^2] dt \right]^{\frac{1}{2}} + n^{O(1)} \sqrt{\underline{T}} + \exp(-O(\bar{T})) \lesssim n^{-\frac{s}{d+2s}} \log^9 n$$

Truncation loss at  $\underline{T}$ .
Truncation loss at  $\bar{T}$ .

$$t_* = N^{-(2-\delta)/d}$$

$$\int_{t=\underline{T}}^{\bar{T}} \mathbb{E}_{X_t} [\|\nabla \log p_t - \hat{s}(\cdot, t)\|^2] dt$$

Bias

$$\frac{\log(\text{covering num})}{n}$$

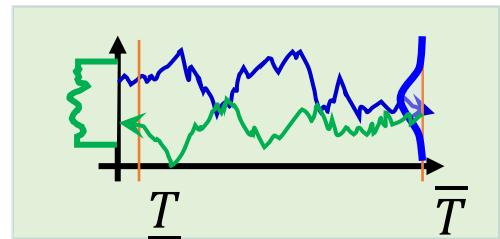
Variance

$$\lesssim \int_{\underline{T}}^{\bar{T}} \frac{N^{-2s/d}}{\sigma_t^2} \log(N) dt + \frac{N \text{polylog}(N)}{n}$$

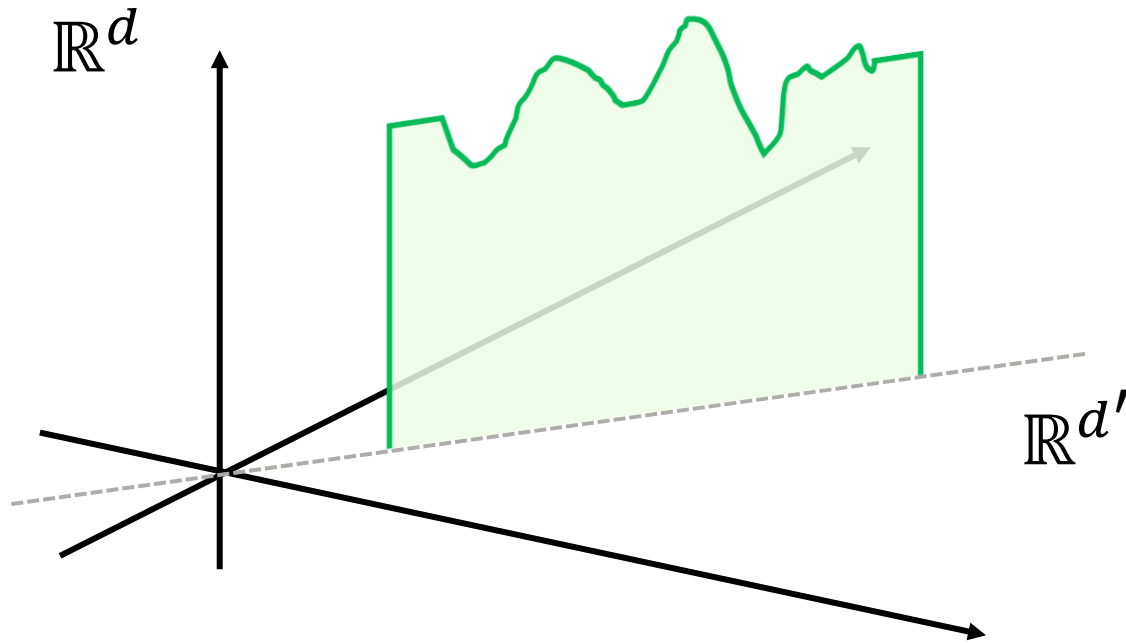
$$\lesssim \left( N^{-2s/d} + \frac{N}{n} \right) \text{polylog}(N)$$

$$N \simeq n^{d/(2s+d)}$$

$$\lesssim n^{-2s/(2s+d)} \text{polylog}(n)$$



# Low dimensional structure



The support of the target distribution is in a low dimensional subspace.

The estimated distribution is never absolutely continuous to the target distribution.

→ **Wasserstein distance**

# $W_1$ -distance convergence rate

## Theorem (Estimation error in $W_1$ -distance)

For any fixed  $\delta > 0$ , by slightly changing the estimator, the empirical risk minimizer  $\hat{s}$  in DNN satisfies

$$\mathbb{E}_{D_n} \left[ W_1(\hat{Y}_{\bar{T}-\underline{T}}, X_0) \right] \lesssim n^{-\frac{s+1-\delta}{2s+d'}}.$$

This is also known as **minimax optimal** (up to  $\delta$ ) [Niles-Weed & Berthet (2022)].

- $d'$  appears instead of  $d$ : **Diffusion model can avoid curse of dimensionality.**
- The minimax rate of Wasserstein distance is faster than that of TV distance, which makes it difficult to establish the bound.
  - We need more precise estimate of the score around  $t = 0$ .

$$(TV) \quad n^{-\frac{s}{2s+d}} \quad \longrightarrow \quad n^{-\frac{s+1}{2s+d}} \quad (W1)$$



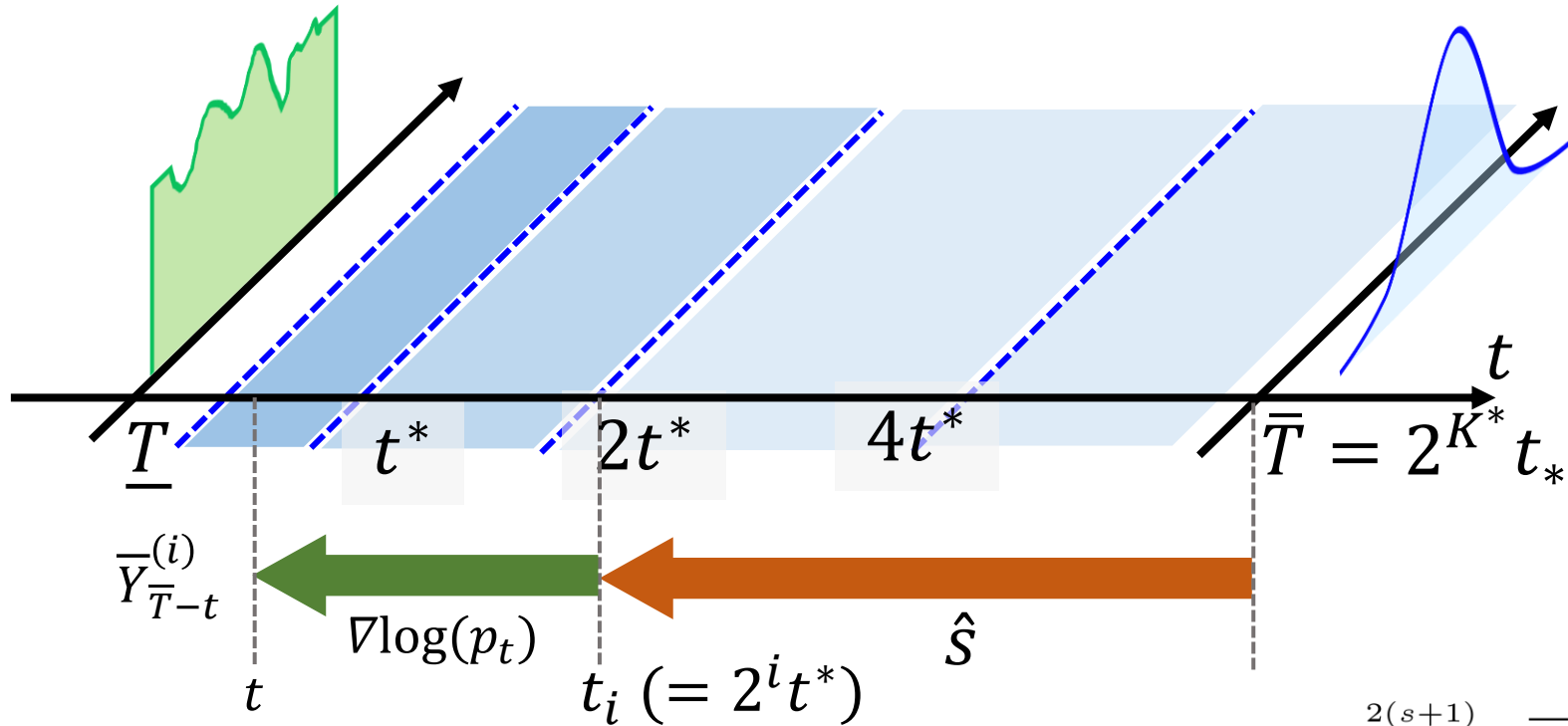
Lemma (tighter bound on W1 distance error)

$$W_1(X_0, \hat{Y}_{\bar{T}-\underline{T}}) \lesssim \sqrt{\int_{t=\underline{T}}^{\bar{T}} t \mathbb{E}_{X_t} [\|\hat{s}(X_t, t) - \nabla \log p_t(X_t)\|^2] dt} \\ + \sqrt{\underline{T}} + \exp(-O(\bar{T}))$$

- For large  $t$ , we can estimate the score more accurately.
- For small  $t$ , the error does not propagate so much due to the term  $t$ .

→ **Better rate.**

# Bound for W1 distance



$$t_* = n^{-\frac{2-\delta}{d+2s}}, \quad t_k = t_* 2^k$$

$$\underline{T} = n^{-\frac{2(s+1)}{2s+d}}, \quad \bar{T} \simeq \log(n)$$

$$W_1(X_0, \hat{Y}_{\bar{T}-\underline{T}}) \leq \underbrace{W_1(X_0, X_{\bar{T}})}_{\text{(negligible)}} + \underbrace{W_1(X_{\bar{T}}, \bar{Y}_{\bar{T}-\underline{T}}^{(K^*)})}_{\text{(exp(-\bar{T}))}} + \sum_{i=1}^{K^*} \underbrace{W_1(\bar{Y}_{\bar{T}-\underline{T}}^{(i-1)}, \bar{Y}_{\bar{T}-\underline{T}}^{(i)})}_{\text{(green bracket)}}$$

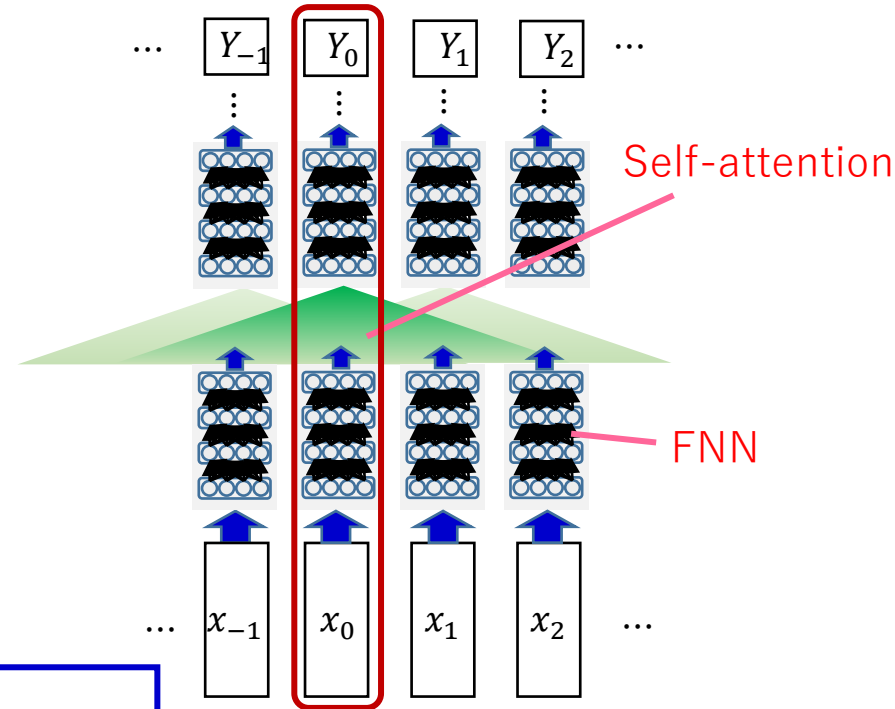
$$\sqrt{t_{i-1} \int_{t_{i-1}}^{t_i} \mathbb{E}_{X_t} [\|\hat{S}(X_t, t) - \nabla \log p_t(X_t)\|^2] dt} \lesssim n^{-\frac{s+1-\delta}{2s+d}}$$

[Shokichi Takakura, Taiji Suzuki: Approximation and Estimation Ability of Transformers for Sequence-to-Sequence Functions with Infinite Dimensional Input. ICML2023]

## Properties of Transformer

- It can output a value from wide range of tokens.
  - Curse of dimensionality?
- It can choose important tokens depending on input.
  - Can avoid curse of dim!

We showed minimax optimality to estimate a sequence-to-sequence function.

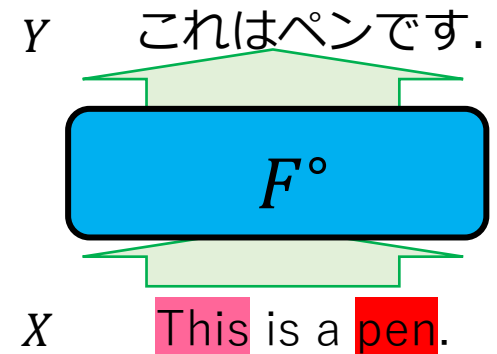


## Theorem (estimation error)

$$\frac{1}{r-l+1} \sum_{j=l}^r \mathbb{E}[\|\hat{F}_j - F_j^\circ\|_{L_2(P_X)}^2] \lesssim n^{-\frac{2a^\dagger}{2a^\dagger+1}} (\log n)^{2/\alpha+2+\max\{4/\alpha, 4\}}$$

(almost minimax optimal)

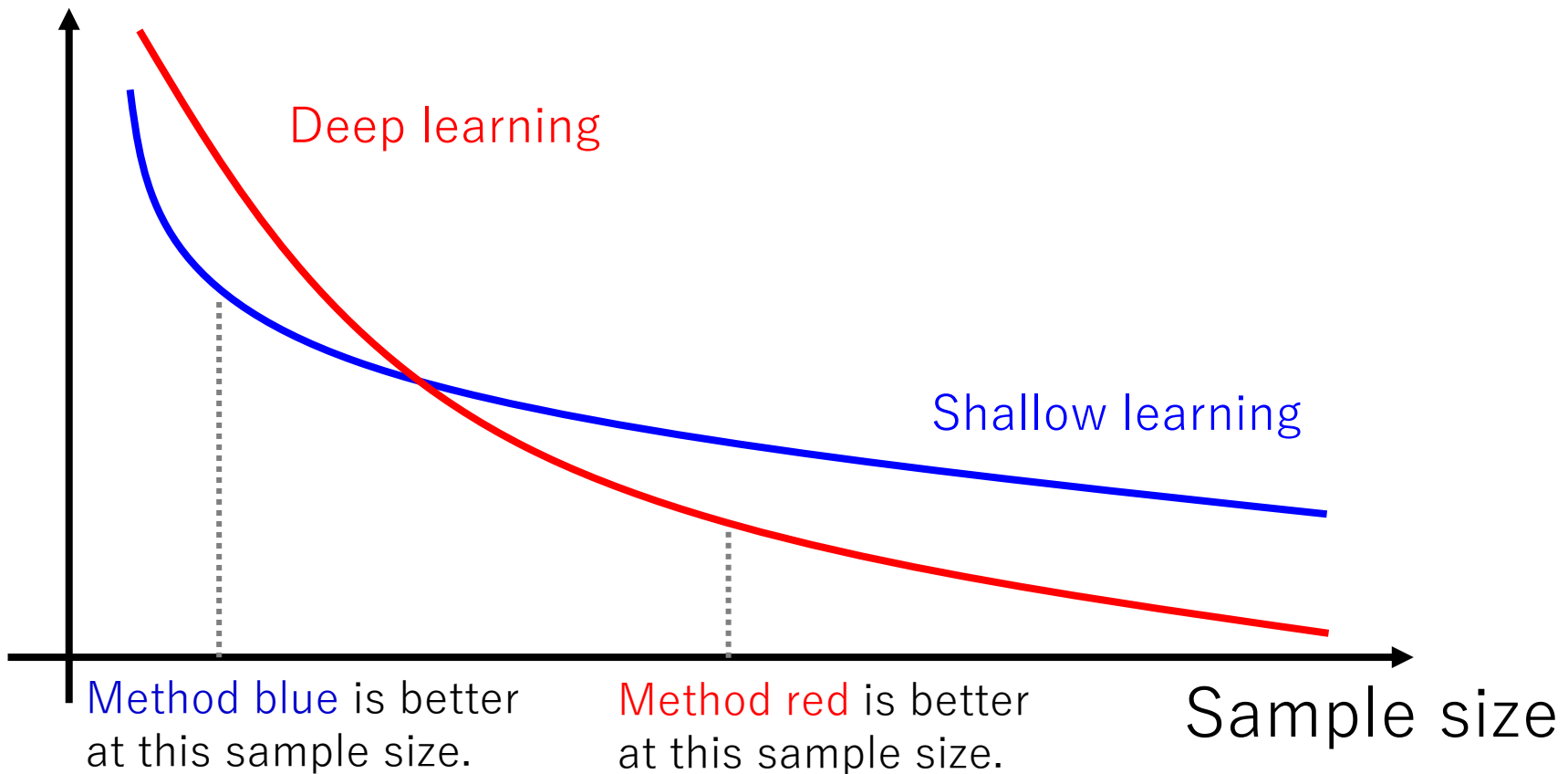
- It achieves polynomial order convergence even though input is infinite-dimensional.



# Remark on the rate of convergence<sup>76</sup>

Remark : Even if the rate is better, the method does not necessarily achieve better prediction.

Predictive error



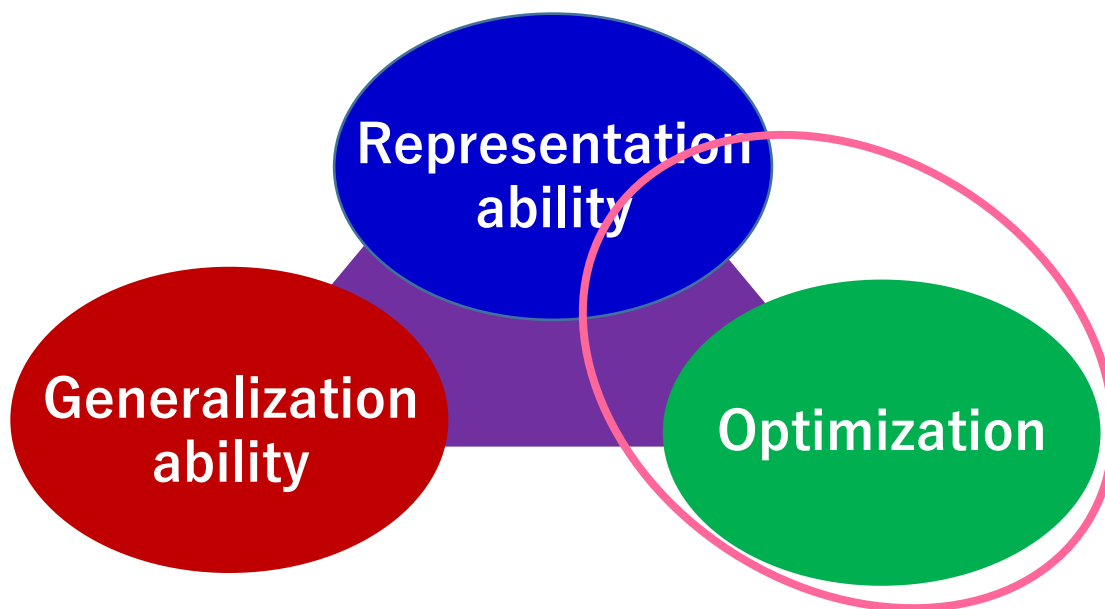
## 1. Representation ability + Generalization ability

- Universal approximator
- Depth separation
- Adaptivity of deep learning
  - Inhomogeneity of smoothness
  - Curse of dimensionality
- Foundation models
  - Diffusion model
  - Transformer

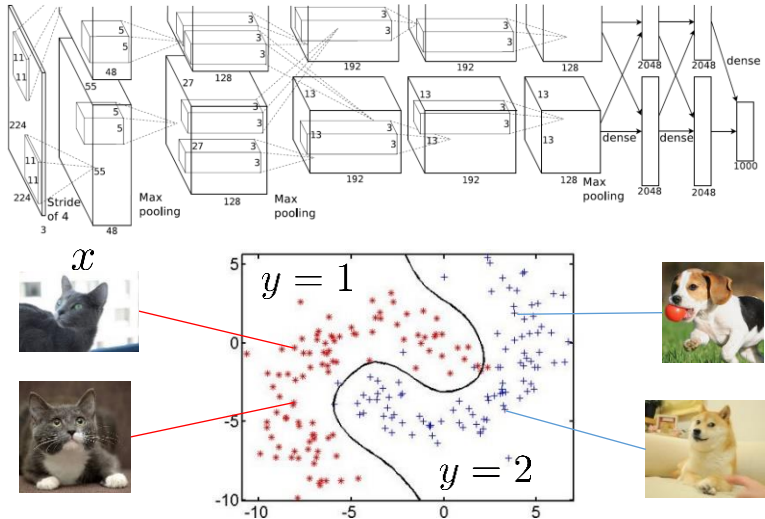
## 2. Optimization ability

- Noisy gradient descent
- Mean field Langevin
- CSQ lowerbound

# Optimization of NN



# Optimization of NN



We should “optimize” the parameters.

$$L(W) = \frac{1}{n} \sum_{i=1}^n \ell_i(W)$$

$W$ : parameter

Loss function for the  $i$ -th data point.

Loss function : degree of fit to the data

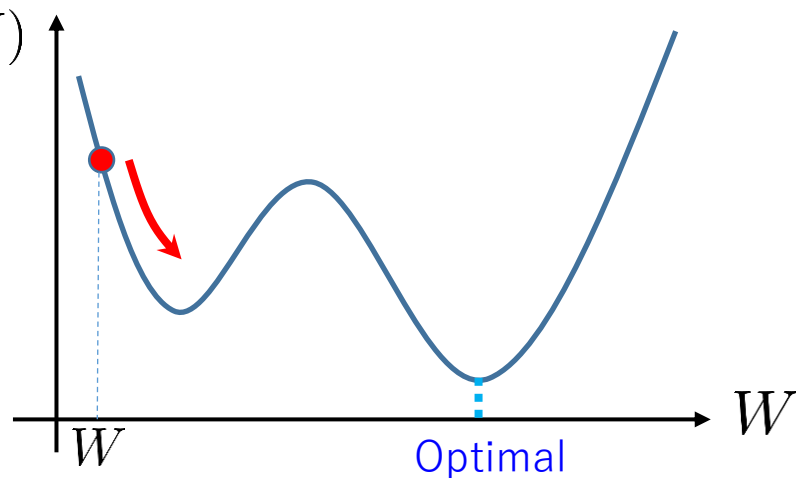
Loss function optimization

$$\min_W L(W)$$

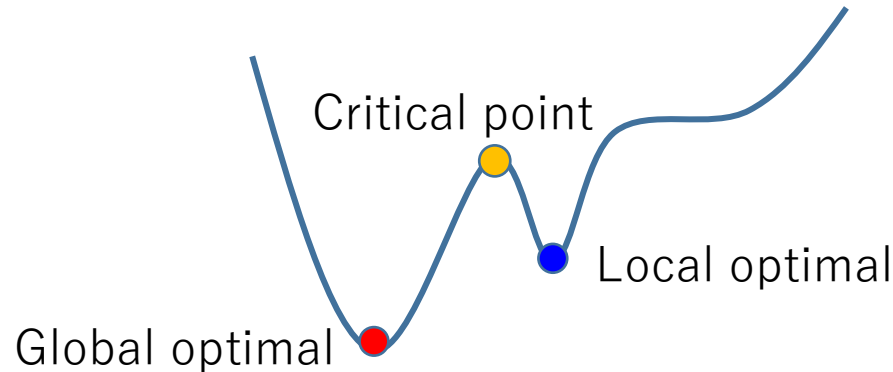
( $W$  could be billions dimensional)

Usually, **stochastic gradient descent** is used.

$L(W)$



Objective function of deep learning is non-convex.



- For linear deep NN, every local optimal is global optimal : Kawaguchi, 2016; Lu&Kawaguchi, 2017.

※**True only for linear NN.**

→ Sufficient conditions that a critical point is a global optimal was also derived by Yun, Sra&Jadbabaie (2018).

- Low rank matrix completion has no spurious local minimum : Ge, Lee&Ma, 2016; Bhojanapalli, Neyshabur&Srebro, 2016.

$$\min_{U \in \mathbb{R}^{M \times k}} \sum_{(i,j) \in E} (Y_{i,j} - (UU^T)_{i,j})^2$$



# Loss landscape

Wide neural network does not have any isolated local minimum.  
(a local optimal solution is connected to global optimal solutions)

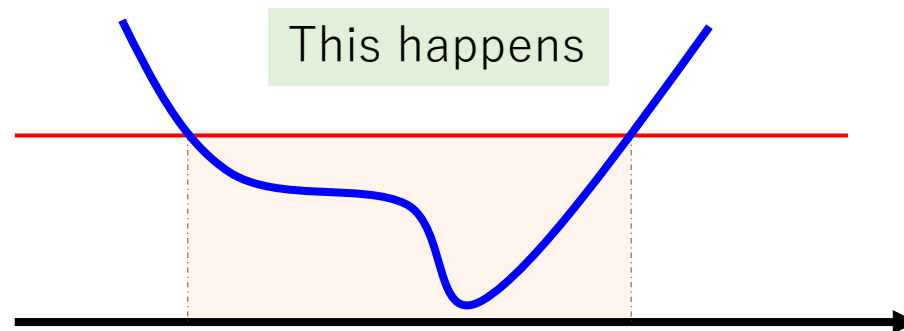
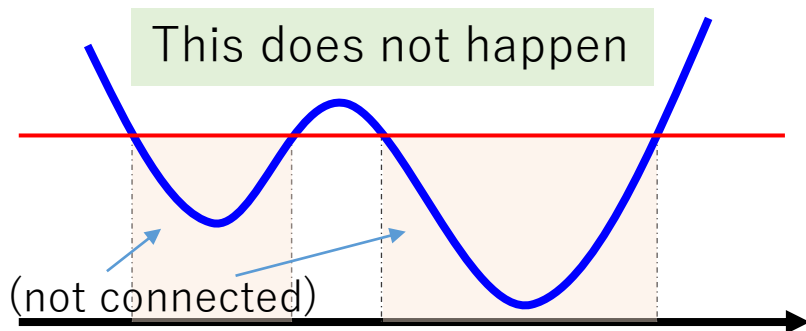
✖ This does not indicate GD can reach the global optimal.

## Theorem

Suppose that we are given  $n$ -training data  $(x_i, y_i)_{i=1}^n$ , and the loss function  $\ell$  is convex.

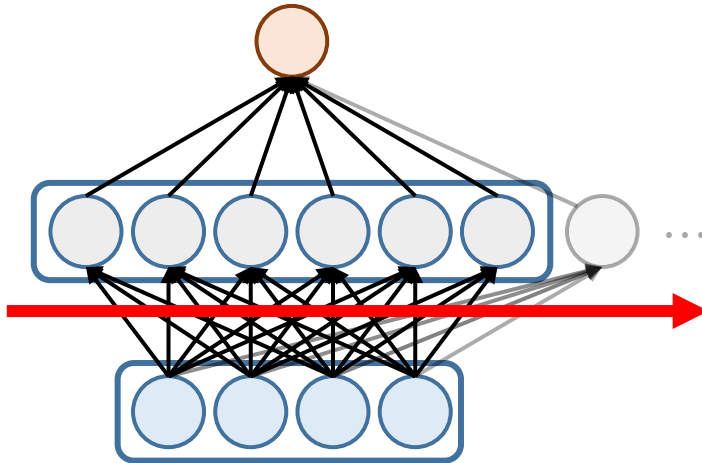
For two layer NN model  $f_{(a,W)}(x) = \sum_{m=1}^M a_m \eta(w_m^\top x)$  with continuous activation function, if the width is not smaller than the data size ( $M \geq n$ ), every arcwise connected component of a level set of the empirical loss  $\hat{L}(a, W) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{(a,W)}(x_i))$  contains the global optimal solution.

[Venturi, Bandeira, Bruna: Spurious Valleys in One-hidden-layer Neural Network Optimization Landscapes. JMLR, 20:1-34, 2019.]



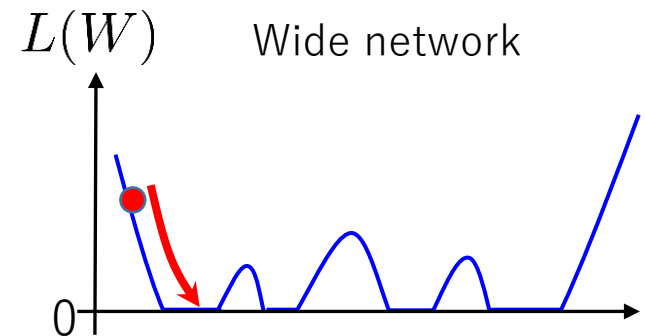
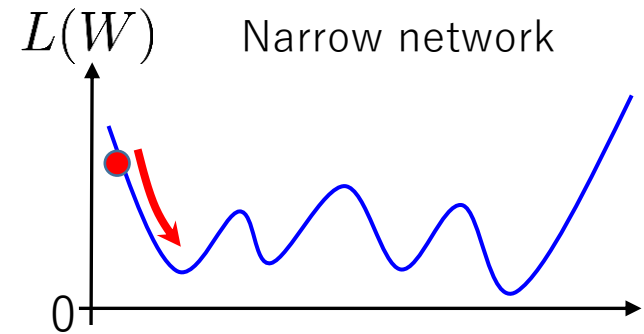
# Overparameterization

Wide neural network does not have spurious local minima.



Since the model complexity is increased, the initial solution is already close to the global optimal.

- Two types of analysis
  - Neural Tangent Kernel (NTK)
  - Mean-field analysis



$$f_W(x) = \sum_{j=1}^M a_j \eta(w_j^\top x)$$

- Neural Tangent Kernel regime (lazy learning )

- $a_j = \mathbf{O}(1/\sqrt{M})$

[Jacot+ 2018][Du+ 2019][Arora+ 2019]  
(Xavier initialization/He initialization)

- Mean field regime

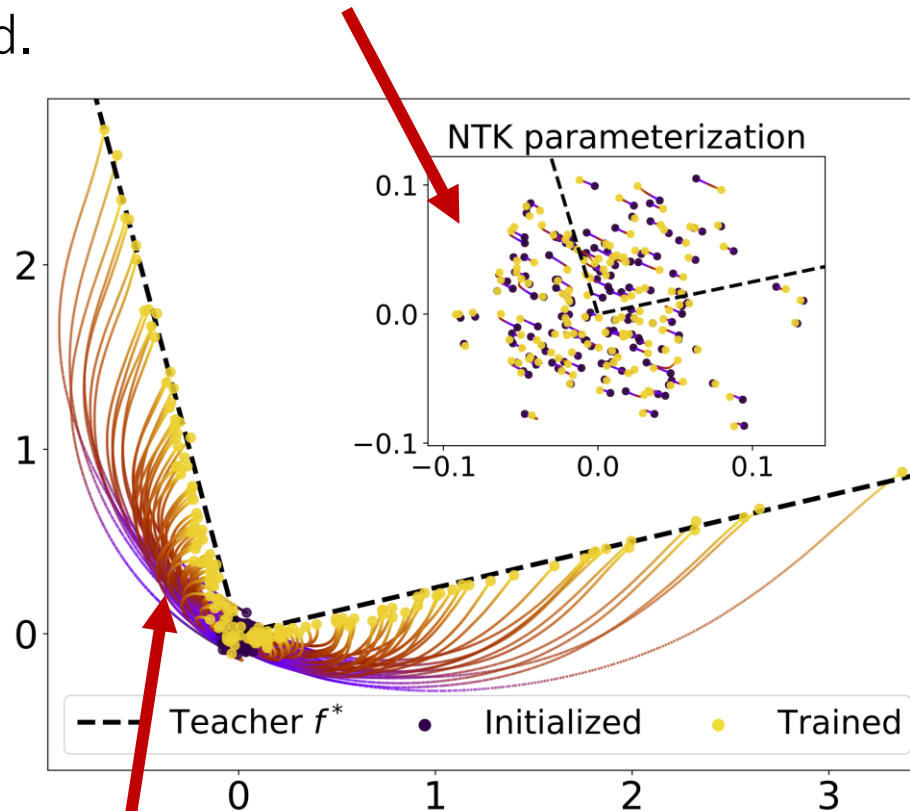
- $a_j = \mathbf{O}(1/M)$

[Nitanda & Suzuki (2017), Chizat & Bach (2018), Mei, Montanari, & Nguyen (2018)]

Different scaling of initial solution yields different behavior.

$$f(x) = \frac{1}{\sqrt{M}} \sum_{j=1}^M r_j \sigma(w_j^\top x)$$

**NTK:** Large scale initialization  $\rightarrow$  features are (almost) frozen.



Optimization trajectory of first layer parameters in a 2-layer NN:

$$f(x) = \sum_{j=1}^M a_j \sigma(w_j^\top x)$$

True function:

$$f^\circ(x) = \sum_{j=1}^2 \sigma(w_j^\top x)$$

[Ba et al., 2022]

**Mean field:** Small scale initialization  $\rightarrow$  features need to move significantly.

$$f(x) = \frac{1}{M} \sum_{j=1}^M r_j \sigma(w_j^\top x)$$

# ABC-parameterization

- ABC-parameterization [Yang&Hu, 2021]

$$x^l(\xi) = \phi(h^l(\xi)) \in \mathbb{R}^n, \quad h^{l+1}(\xi) = W^{l+1}x^l(\xi) \in \mathbb{R}^n, \quad \text{for } l = 1, \dots, L-1, \quad \underline{n: \text{width}}$$

(1) Parameterization

$$W^l = n^{-a_l} w^l$$

( $w^l$  is the actual trainable parameter)

(2) Initialization

$$w_{\alpha\beta}^l \sim \mathcal{N}(0, n^{-2b_l})$$

(3) Learning rate

$$\eta n^{-c}$$

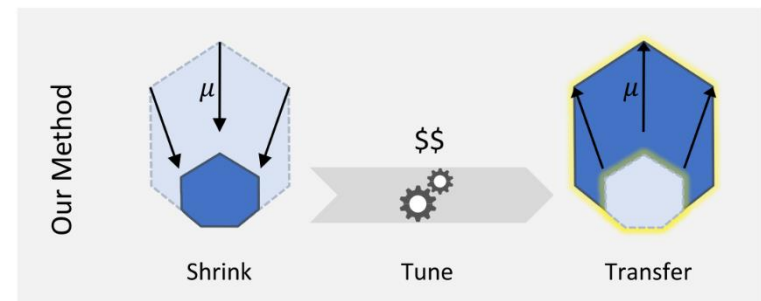
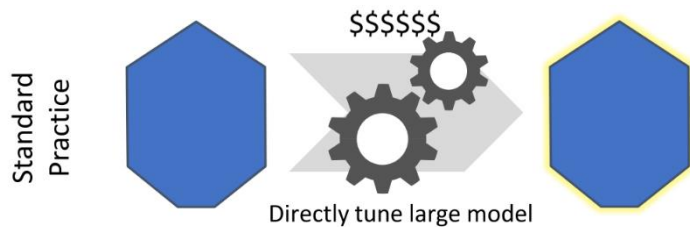
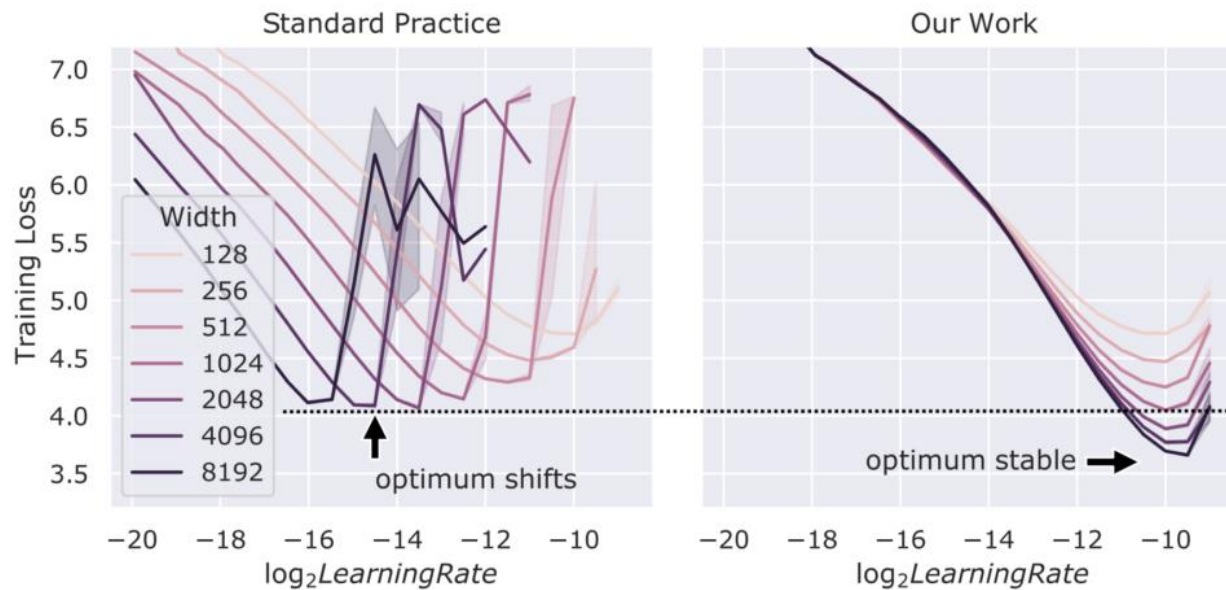
$$h^1 = W^1 \xi \in \mathbb{R}^n, x^l = \phi(h^l) \in \mathbb{R}^n, h^{l+1} = W^{l+1} x^l \in \mathbb{R}^n, f(\xi) = W^{L+1} x^L$$

Definition		SP (w/ LR $\frac{1}{n}$ )	NTP	MFP ( $L = 1$ )	$\mu$ P (ours)
$a_l$	$W^l = n^{-a_l} w^l$	0	$\begin{cases} 0 & l = 1 \\ 1/2 & l \geq 2 \end{cases}$	$\begin{cases} 0 & l = 1 \\ 1 & l = 2 \end{cases}$	$\begin{cases} -1/2 & l = 1 \\ 0 & 2 \leq l \leq L \\ 1/2 & l = L + 1 \end{cases}$
$b_l$	$w_{\alpha\beta}^l \sim \mathcal{N}(0, n^{-2b_l})$	$\begin{cases} 0 & l = 1 \\ 1/2 & l \geq 2 \end{cases}$	0	0	1/2
$c$	$LR = \eta n^{-c}$	1	0	-1	0
$r$	<b>Definition 3.2</b>	1/2	1/2	0	0
$2a_{L+1} + c$		1	1	1	1
$a_{L+1} + b_{L+1} + r$		1	1	1	1
Nontrivial?		✓	✓	✓	✓
Stable?		✓	✓	✓	✓
Feature Learning?				✓	✓
Kernel Regime?		✓	✓		

(Appropriate scaling)

The optimal hyper-parameter in a small size network can be transferred to huge model.

It is used to train GPT-3.5. Billions of dollars cost could be saved.



[Yang et al.:Tensor Programs V: Tuning Large Neural Networks via Zero-Shot Hyperparameter Transfer. arXiv:2203.03466]

<https://github.com/microsoft/mutransformers> ( $\mu$ P for Transformers)

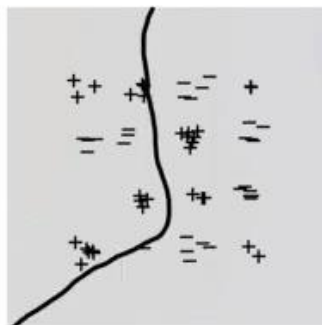
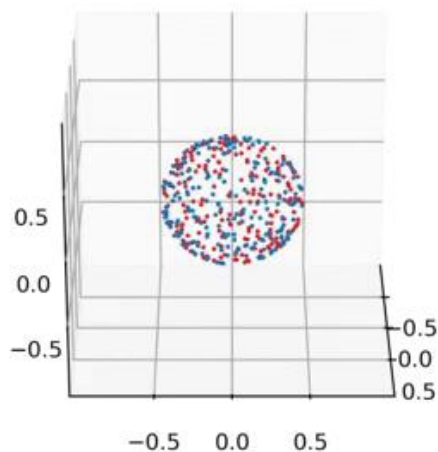
# Implicit regularization

Binary classification with exp-loss:

$$\min_{\rho} \sum_{i=1}^n \exp(-y_i f_{\rho}(x_i)) \quad \text{where} \quad f_{\rho}(x) = \int \eta(w^{\top} x) d\rho(w)$$

Optimization in the space of signed measures.

If we start from small initialization, only neurons that are necessary for classification “grow up.”



[Chizat&Bach: Implicit Bias of Gradient Descent for Wide Two-layer Neural Networks Trained with the Logistic Loss. COLT2020.]

Optimization dynamics implicitly regularize the solution.

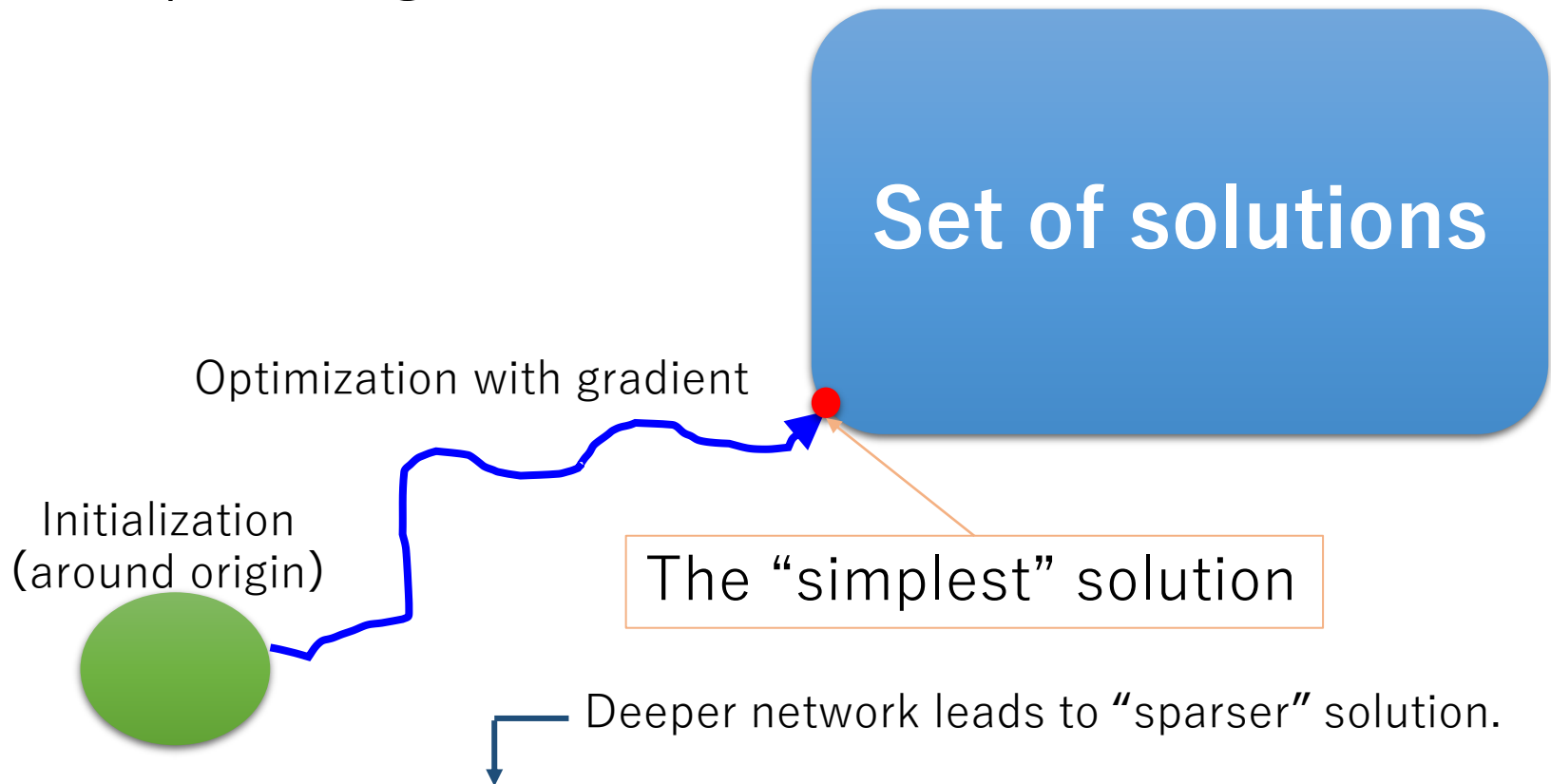
→ **Sparse solution : implicit regularization**

The solution converges to the max-margin solution under L1-norm constraint (if the sequence converges to a “global minimizer” direction):

$$\max_{\rho: \|\rho\|_{\mathcal{F}_1} \leq 1} \min_{i \in \{1, \dots, n\}} y_i f_{\rho}(x_i) \quad \|\rho\|_{\mathcal{F}_1} = |\rho|(\mathbb{R}^d)$$

# Gradient descent and implicit regularization<sup>88</sup>

- Dynamics starting from a small initialization converges to the minimum norm solution.  
→ Implicit regularization



[Gunasekar et al.: Implicit Regularization in Matrix Factorization, NIPS2017]

[Soudry et al.: The implicit bias of gradient descent on separable data. JMLR2018]

[Gunasekar et al.: Implicit Bias of Gradient Descent on Linear Convolutional Networks, NIPS2018]

[Moroshko et al.: Implicit Bias in Deep Linear Classification: Initialization Scale vs Training Accuracy, arXiv:2007.06738]



# Implicit regularization in each regime <sup>89</sup>

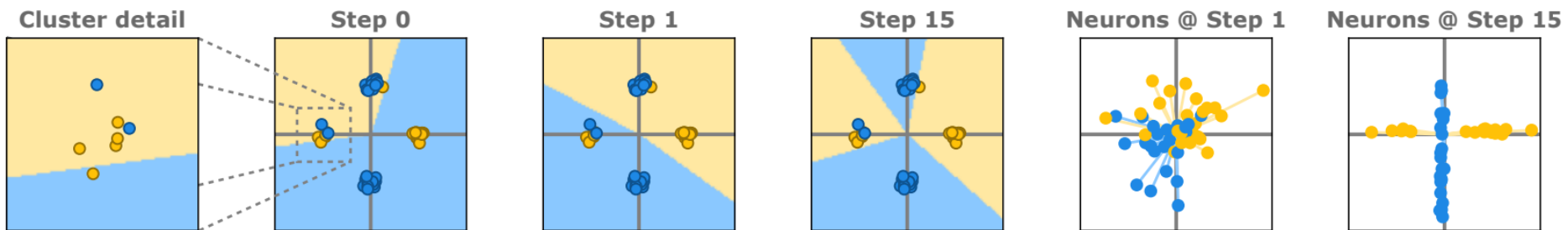
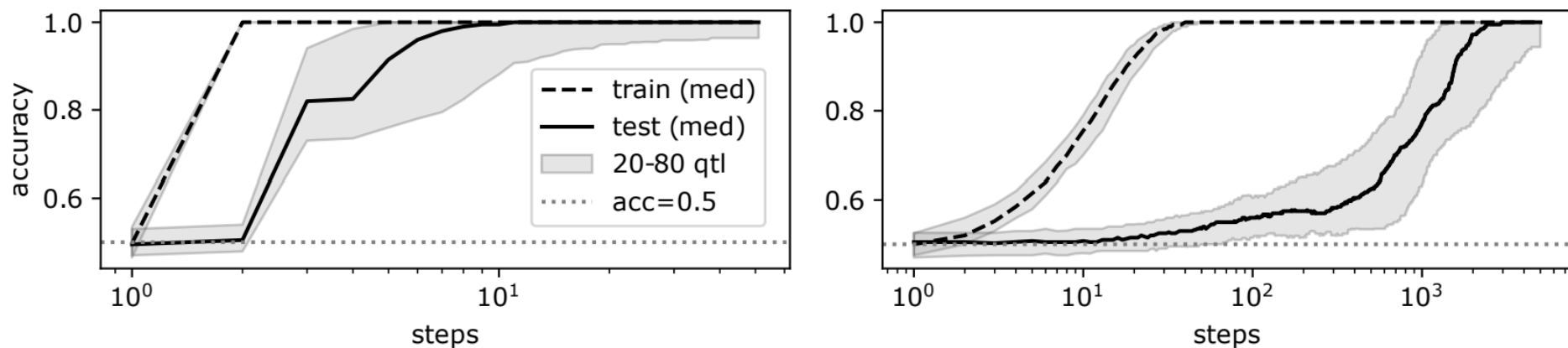
Regime	Implicit regularization
NTK, kernel method with early stopping	L2-regularization
Mean-field	L1-regularization

- Deep learning uses several **“explicit regularization”**.  
→ Batch normalization, Dropout, Weight decay, MixUp, ...
- On the other hand, the **“implicit regularization”** induced by the deep structure and optimization dynamics is also very important.  
→ Benign overfitting, Grokking, Flat-minimum, ...

# Grokking / Benign-overfitting

Grokking: [Power et al.: Grokking: Generalization beyond overfitting on small algorithmic datasets. arXiv:2201.02177]

[Xu et al.: Benign Overfitting and Grokking in ReLU Networks for XOR Cluster Data. arXiv2310.02541]



It is also called “hidden progress” [Barak et al. 2022].

See also [Meng et al.: Benign Overfitting in Two-Layer ReLU Convolutional Neural Networks for XOR Data. arXiv2310.01975]

## 1. Representation ability

- Universal approximator
- Adaptivity of deep learning
  - Inhomogeneity of smoothness
  - Curse of dimensionality

## 2. Generalization ability

- Double descent, Benign overfitting for overparameterized model
- Generalization gap analysis
  - Norm based bound
  - Compression based bound

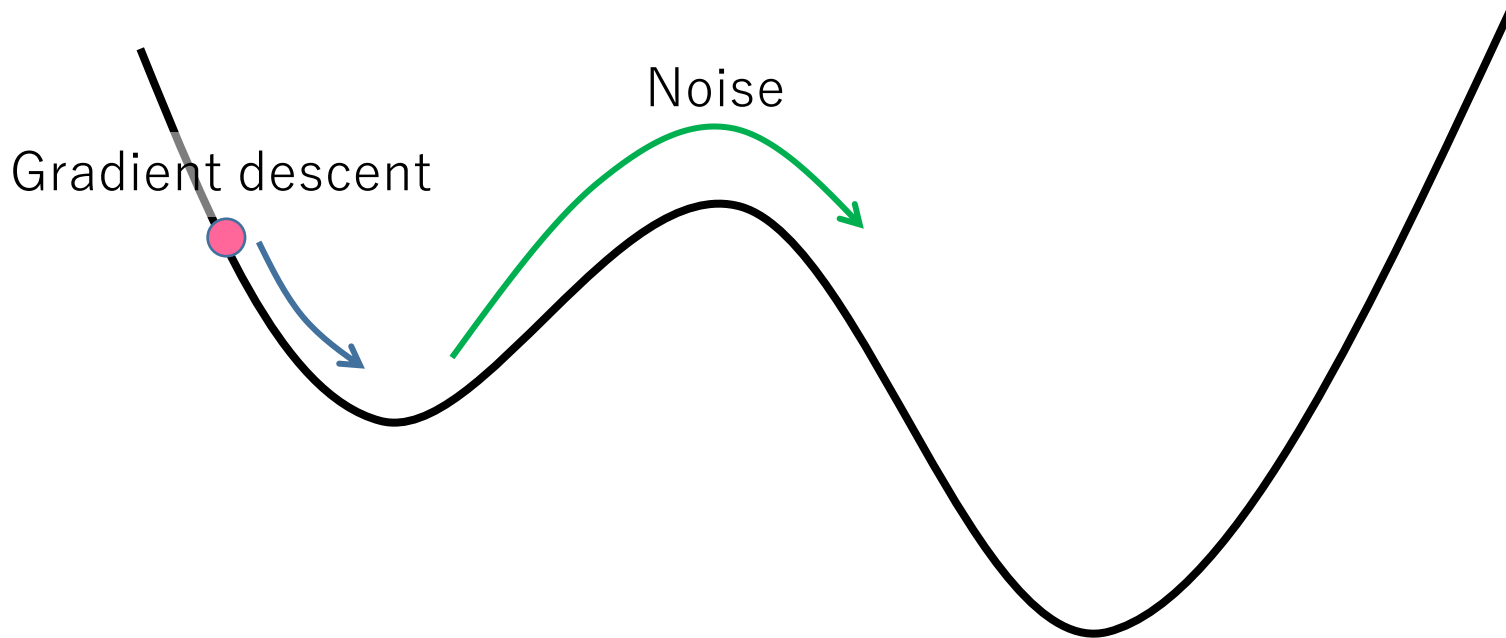
## 3. Optimization ability

- Neural Tangent Kernel
- Dynamics in a feature learning regime
- Mean field analysis

# Noisy gradient descent and its global optimality

# Noisy gradient descent

The model is not linearly approximated.  
We need to solve “non-convex” optimization.

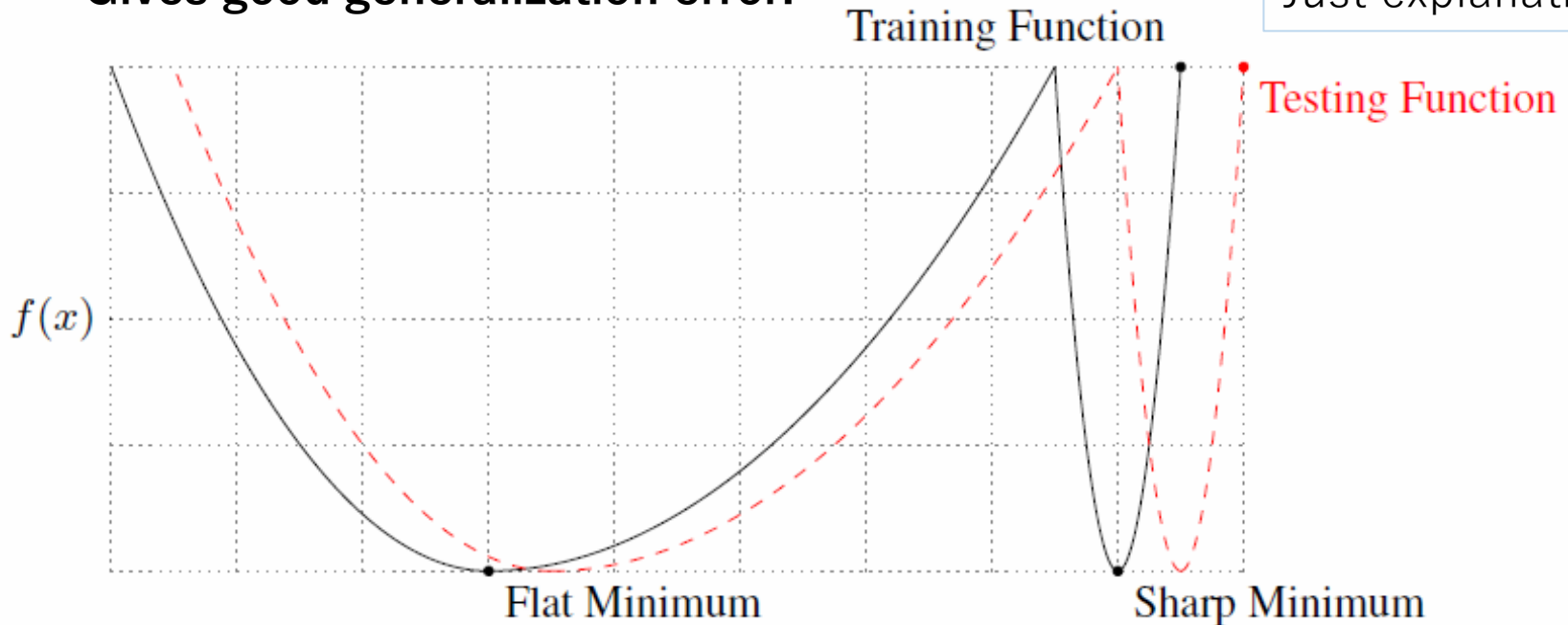


SGD is a noisy gradient descent.  
Noisy perturbation is helpful to escape local minimum.

# Sharp minima vs flat minima

It is said that SGD likely stay in “flat local minimum”  
→ Gives good generalization error.

This is not theory,  
Just explanation.



Keskar, Mudigere, Nocedal, Smelyanskiy, Tang (2017):

On large-batch training for deep learning: generalization gap and sharp minima.

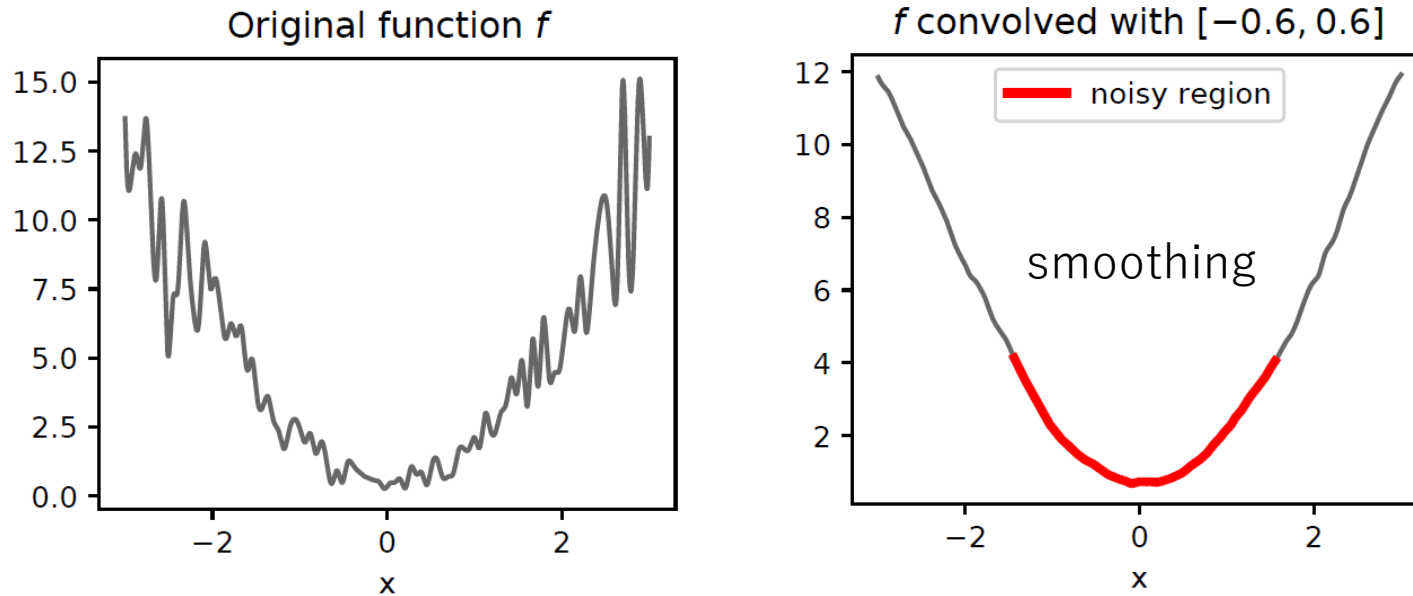
$$\theta_t = \theta_{t-1} - \alpha_b \underbrace{\left( \frac{1}{b} \sum_{j=1}^b \nabla_{\theta} \ell(z_{i_j}; \theta) \right)}_{\cong \text{Normal distribution}}$$

$\cong$  Normal distribution

→ Random walk is likely to be captured in a flat region.

- (criticism) The concept of “flat” depends on the choice of coordinate system. (Dinh et al., 2017)
- PAC-Bayesian analysis (Dziugaite, Roy, 2017)

# Smoothing by noisy gradient



[Kleinberg, Li, and Yuan, ICML2018]

**Stochastic gradient  $\Rightarrow$  Noise is added  $\Rightarrow$  Objective is smoothed**

$$x_t = x_{t-1} - \eta(\nabla L(x_{t-1}) + \xi_t) \quad (y_t = x_t + \eta\xi_t)$$

$$\Rightarrow y_t = y_{t-1} - \eta\xi_{t-1} - \eta\nabla L(y_{t-1} - \eta\xi_{t-1})$$

$$\Rightarrow \mathbb{E}_{\xi_{t-1}}[y_t] = y_{t-1} - \eta\nabla \mathbb{E}_{\xi_{t-1}}[L(y_{t-1} - \eta\xi_{t-1})]$$

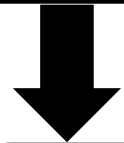
**SGD optimizes a “smoothed” objective:  $\bar{L}(y_t) = \mathbb{E}_{\xi_t}[L(y_t - \eta\xi_t)]$**

- Stochastic Gradient Langevin Dynamics (SGLD)

$$\min_{x \in \mathbb{R}^d} L(x) = \min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell_i(x) \quad (\text{Non-convex})$$

$$dX_t = -\nabla L(X_t)dt + \sqrt{2\lambda}dB_t \quad (\text{Gradient Langevin dynamics})$$

Stationary distribution :  $\pi \propto \exp(-L(X)/\lambda)$



**Discretization**

[Gelfand and Mitter (1991); Borkar and Mitter (1999); Welling and Teh (2011)]

**GLD:**  $X_{t+1} = X_t - \eta \nabla L(X_t) + \sqrt{2\eta\lambda}\xi_t$  (Euler-Maruyama scheme)  
 $\xi_t \sim N(0, I)$

**SGLD:**  $X_{t+1} = X_t - \eta \frac{1}{|I_B|} \sum_{i \in I_B} \nabla \ell_i(X_t) + \sqrt{2\eta\lambda}\xi_t$   
Stochastic gradient



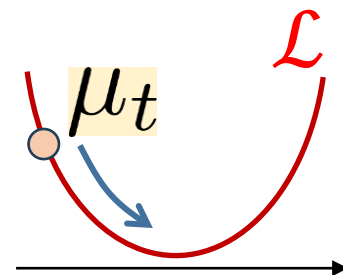
# GLD as a Wasserstein gradient flow<sup>97</sup>

$$dX_t = -\nabla L(X_t)dt + \sqrt{2\lambda}dB_t$$

$\mu_t$  : Distribution of  $X_t$  (we can assume it has a density)

PDE that describes  $\mu_t$ 's dynamics [**Fokker-Planck equation**]:

$$\begin{aligned}\partial_t \mu_t &= \nabla \cdot [\mu_t \nabla L] + \lambda \Delta_x \mu_t \\ &= \nabla \cdot [\mu_t (\nabla L + \lambda \nabla \log(\mu_t))]\end{aligned}$$



This is the **Wasserstein gradient flow** to minimize the following objective:

$$\mu^* = \arg \min_{\mu \in \mathcal{P}} \int L(x) d\mu(x) + \lambda \text{Ent}(\mu) =: \mathcal{L}(\mu)$$

[linear w.r.t.  $\mu$ ]

$$(\text{Ent}(\mu) = \int \log(\mu) d\mu)$$

➔  $\mu_t \rightsquigarrow \mu^*(x) \propto \exp(-L(x)/\lambda) =$  Stationary distribution  
c.f., Donsker-Varadhan duality formula

# Continuity equation

Continuity equation:  $\frac{\partial \mu_t}{\partial t} = -\nabla \cdot (v_t \mu_t)$   $v_t(x) = -(\nabla L(x) + \lambda \nabla \log(\mu_t)(x))$

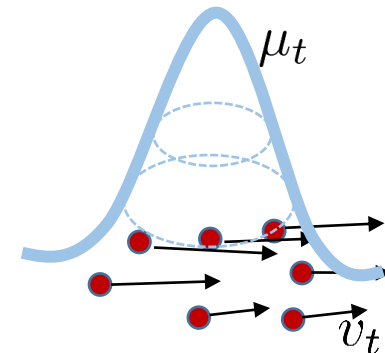
The meaning of this equation

$$\frac{d}{dt} \int f(x) \mu_t(x) dx = \int (\nabla f(x))^\top v_t(x) \mu_t(x) dx$$

$$\left( = - \int f(x) \nabla \cdot (v_t \mu_t) dx \right) \quad (\forall f: \text{compact support, } C^\infty\text{-class})$$

- Let  $T_t$  be a map generated by the vector field  $v_t$ :  $\frac{dT_t}{dt}(x) = v_t(T_t(x))$ .
- $\mu_t$  is the push-forward of  $\mu_0$  by a map  $T_t: \mathbb{R}^d \rightarrow \mathbb{R}^d$ :  $\mu_t = T_{t\#}\mu_0$ .  
That is,  $\mu_t$  is the distribution of  $T_t(x)$  where  $x \sim \mu_0$ .

$$\begin{aligned} \frac{d}{dt} \int f(x) \mu_t(x) dx &= \frac{d}{dt} \int f(T_t(x)) \mu_0(x) dx \\ &= \int \nabla f(T_t(x))^\top \frac{dT_t(x)}{dt} \mu_0(x) dx \\ &= \int \nabla f(T_t(x))^\top v_t(T_t(x)) \mu_0(x) dx \\ &= \int \nabla f(x)^\top v_t(x) \mu_t(x) dx. \end{aligned}$$



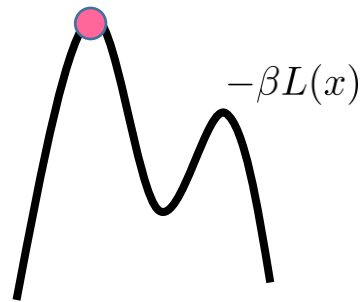
[continuity equation]

# Stationary distribution

$$dX_t = -\nabla L(X_t)dt + \sqrt{2\beta^{-1}}dB_t$$

Stationary distribution of the continuous time dynamics:

$$\mu^*(dx) \propto \exp(-\beta L(x))dx$$



The stationary distribution concentrates around the optimal solution.

# Wasserstein gradient flow

$$\begin{aligned}\lambda^{-1} \mathcal{L}(\mu) &= \int \lambda^{-1} L(x) d\mu(x) + \text{Ent}(\mu) && \boxed{\mu^*(x) \propto \exp(-\lambda^{-1} L(x))} \\ &= \int -\log(\mu^*) d\mu + \int \log(\mu) d\mu + (\text{const.}) \\ & && \text{We neglect this term below} \\ &= \int \log\left(\frac{\mu}{\mu^*}\right) d\mu = \text{KL}(\mu || \mu^*)\end{aligned}$$

By the continuity equation  $\mu_t = -\nabla \cdot [v_t \mu_t]$ , it holds that

$$\begin{aligned}\frac{d}{dt} \text{KL}(\mu_t || \mu^*) &= \frac{d}{dt} \int \log\left(\frac{\mu_t(x)}{\mu^*(x)}\right) \mu_t(x) dx \\ &= \int \log\left(\frac{\mu_t(x)}{\mu^*(x)}\right) \partial_t \mu_t(x) dx + \int \frac{\partial_t \mu_t(x)}{\mu_t(x)} \mu_t(x) dx \\ &= \int \log\left(\frac{\mu_t(x)}{\mu^*(x)}\right) \nabla \cdot (-v_t \mu_t(x)) dx && = 0 \\ &= \int \langle v_t, \nabla \log(\mu_t) - \nabla \log(\mu^*) \rangle \mu_t(x) dx\end{aligned}$$

$$\frac{d}{dt} \text{KL}(\mu_t || \mu^*) = \int \langle v_t, \nabla \log(\mu_t) - \nabla \log(\mu^*) \rangle \mu_t(x) dx$$

In particular, if

$$v_t = -(\lambda \nabla \log(\mu_t) + \nabla L) = -\lambda (\nabla \log(\mu_t) - \nabla \log(\mu^*)) \quad (\text{GLD})$$

then this is the steepest gradient descent direction such that

$$\partial_t \mu_t = \nabla \cdot \underbrace{[(\lambda \nabla \log(\mu_t) + \nabla L) \mu_t]}_{=:-v_t}$$

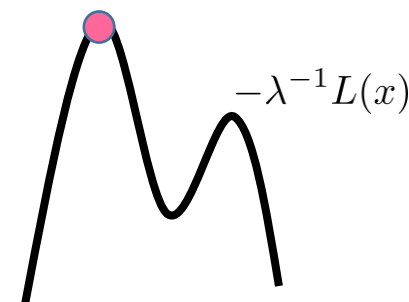
$$\begin{aligned} \frac{d}{dt} \text{KL}(\mu_t || \mu^*) &= -\lambda \int \|\nabla \log(\mu^*) - \nabla \log(\mu_t)\|^2 \mu_t dx \\ &= -\lambda I(\mu_t || \mu^*) \end{aligned}$$

**Fisher divergence:**

$$I(\mu || \nu) := \int \|\nabla \log(\nu) - \nabla \log(\mu)\|^2 \mu(x) dx$$

**GLD is the Wasserstein gradient flow to minimize the KL-div from  $\mu^*$ .**

Stationary distribution:  $\mu^*(x) \propto \exp(-\lambda^{-1}L(x))$



## Def (log-Sobolev inequality)

There exists a constant  $\alpha > 0$  such that for any probability measure  $\nu$  (absolutely-continuous w.r.t.  $\mu^*$ )

$$\text{KL}(\nu || \mu^*) \leq \frac{1}{2\alpha} I(\nu || \mu^*)$$

E.g. :

- Quadratic+Bounded
- Weak Morse function

KL-div

$$\text{KL}(\nu || \mu) = \int \log\left(\frac{\nu}{\mu}\right) \mu dx, \quad \text{Fisher-div} \quad I(\nu || \mu) = \int \left\| \nabla \log \frac{\nu}{\mu} \right\|^2 \nu dx$$

Fisher-div

➔ **Geometric ergodicity**  $\mu_t$ : distribution of  $X_t$

$$\frac{d}{dt} \text{KL}(\mu_t || \mu^*) = -\lambda I(\mu_t || \mu^*) \leq -2\alpha \text{KL}(\mu_t || \mu^*) \quad (\text{by log-Sobolev})$$

$$\text{KL}(\mu_t || \mu^*) \leq \exp(-2\alpha t) \text{KL}(\mu_0 || \mu^*)$$

**Linear convergence w.r.t. KL-div**

# Sufficient condition for log-Sobolev inequality

**Strongly convex (Bakry-Emery criterion):**

$$\mu^*(x) \propto \exp(-\lambda^{-1}L(x))$$

$$\nabla \nabla^\top L(x) \succeq \mu I \quad \Rightarrow \quad \alpha \geq \mu/\lambda$$

[Bakry and Émery, 1985]

Ex.: OU-process.  $L(x) = \frac{x^2}{2} \Rightarrow \mu = 1$

**Bounded perturbation lemma (Holley-Stroock):**

Suppose that  $\mu^*(x) = \mu(x) \exp(h(x))$  and  $\mu$  satisfies  $\alpha'$ -LSI, then

$$|h(x)| \leq B \quad (\forall x) \quad \longrightarrow \quad \mu^* \text{ satisfies } \alpha\text{-LSI with } \alpha \geq \alpha' \exp(-4B)$$

[R. Holley and D. Stroock. Logarithmic sobolev inequalities and stochastic Ising models. Journal of statistical physics, 46(5-6):1159-1194, 1987.]

Ex.:  $L(x) = \ell(x) + \lambda_1 x^2$  and  $|\ell(x)| \leq B$ , then  $\mu^*$  satisfies LSI with  $\alpha = \frac{2\lambda_1}{\lambda} \exp(-4B/\lambda)$ .

$$\mu^*(x) \propto \exp(-\lambda^{-1}L(x))$$

- **Finite dimensional Langevin dynamics:**

- [Convergence in low \(convex case\)](#): Dalalyan and Tsybakov, 2012; Dalalyan, 2016; Durmus and Moulines, 2015, ..
- [Non-convex Optimization](#): Raginsky et al., 2017; Xu et al., 2018; Erdogdu, Mackey and Shamir, 2018
- [Log-Sobolev inequality](#): Vempala and Wibisono, 2019.

- **Infinite dimensional Langevin dynamics:**

- Continuous time:
  - [Existence & Uniqueness of invariant measure](#): Da Prato and Zabczyk, 1992; Maslowski, 1989; Sowers, 1992.
  - [Geometric ergodicity](#): Jacquot and Royer, 1995; Shardlow, 1999; Hairer, 2002, Its explicit rate: Goldys and Maslowski, 2006.
- Discrete time:
  - [Weak approximation rate of discretized scheme](#): Hausenblas, 2003; Debussche, 2011; Bréhier, 2014; Bréhier and Kopec 2016.

Other topics (MCMC in Hilbert space):

- [preconditioned Crank–Nicolson \(pCN\)](#): Hairer et al., 2014; Eberle, 2014; Vollmer, 2015; Rudolf and Sprungk, 2018.
- [Metropolis-Adjusted Langevin Algorithm \(MALA\)](#): Durmus and Moulines, 2015; Beskos et al., 2017.



# Related work: Graduated optimization<sup>105</sup>

- Graduated non-convexity

Blake and Zisserman: *Visual reconstruction*, volume 2. MIT press Cambridge, 1987.

- Convolution with Gaussian kernel

Z. Wu. The effective energy transformation scheme as a special continuation approach to global optimization with application to molecular conformation. *SIAM Journal on Optimization*, 6(3):748-768, 1996.

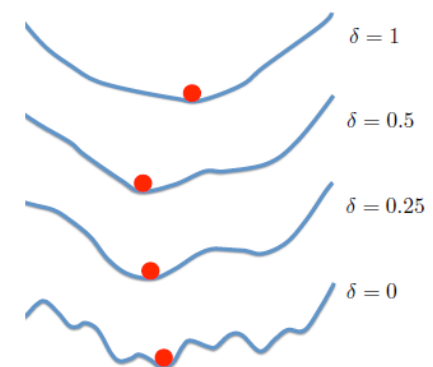
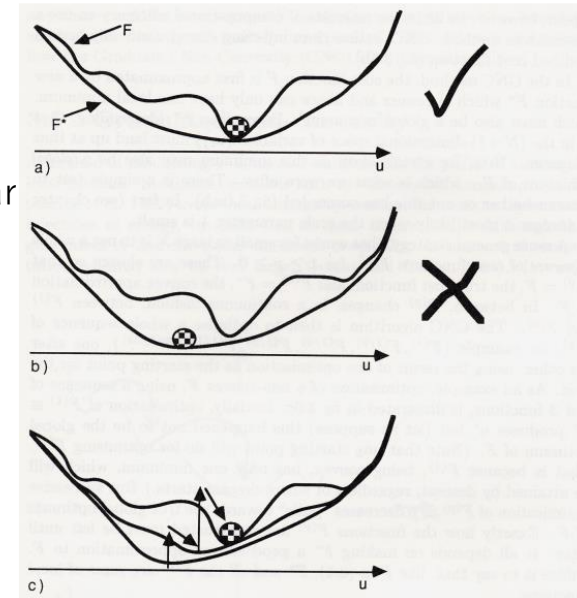
- Graduated optimization

Hazan, Levy, and Shalev-Shwartz: On graduated optimization for stochastic non-convex problems. *International conference on machine learning*, pp. 1833-1841, 2016.

- $\sigma$ -nice property:  $\hat{L}_\delta(x) = \mathbb{E}_{u \sim U(B(\mathbb{R}^d))} [L(x + \delta u)]$
- Polynomial time convergence.

Survey:

Mobahi and Fisher III. On the link between gaussian homotopy continuation and convex envelopes. *Energy Minimization Methods in Computer Vision and Pattern Recognition*, pp. 43-56, 2015.



# Optimization theory in mean field regime

# 2-layer NN in mean-field scaling 107

- 2-layer neural network:

$$f(z) = \frac{1}{M} \sum_{j=1}^M r_j \sigma(w_j^\top z)$$

Non-linear with respect to parameters  $(r_j, w_j)_{j=1}^M$ .

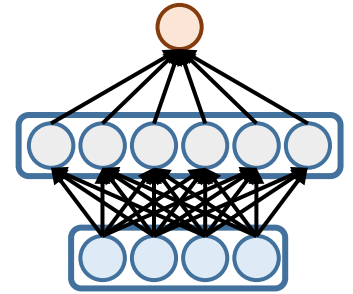
$$f_{\mathcal{X}}(z) = \frac{1}{M} \sum_{j=1}^M h_{X^{(j)}}(z)$$

where  $X^{(j)} = (r_j, w_j)$  and  $h_x(z) = r \sigma(w^\top z)$  for  $x = (r, w)$ .

**Loss function:**

$$F(\mathcal{X}) = \frac{1}{n} \sum_{i=1}^n \underset{\text{loss}}{\ell_i(f_{\mathcal{X}}(z_i))} + \lambda_1 \frac{1}{M} \sum_{j=1}^M \underset{\text{L2 regularization}}{\|X^{(j)}\|^2}$$

**Non-convex**



$$F(\mathcal{X}) = \frac{1}{n} \sum_{i=1}^n \ell_i(f_{\mathcal{X}}(z_i)) + \lambda_1 \frac{1}{M} \sum_{j=1}^M \|X^{(j)}\|^2$$

$$f_{\mathcal{X}}(z) = \frac{1}{M} \sum_{j=1}^M h_{X^{(j)}}(z)$$

Noisy gradient descent update:

$$X_{k+1}^{(j)} = X_k^{(j)} - \eta_k \nabla_{X_k^{(j)}} F(\mathcal{X}_k) + \sqrt{2\eta_k \lambda_2} \xi_k^{(j)} \quad \xi_k^{(j)} \sim \mathcal{N}(0, I)$$

$$\Leftrightarrow X_{k+1}^{(j)} = X_k^{(j)} - \frac{\eta_k}{M} \left( \frac{1}{n} \sum_{i=1}^n \ell'_i(f_{\mathcal{X}_k}(z_i)) \nabla_{X_k^{(j)}} h_{X_k^{(j)}}(z_i) + \lambda_1 X_k^{(j)} \right) + \sqrt{2\eta_k \lambda_2} \xi_k^{(j)}$$

## Does it converge?

Naïve application of existing theory in gradient Langevin dynamics yields

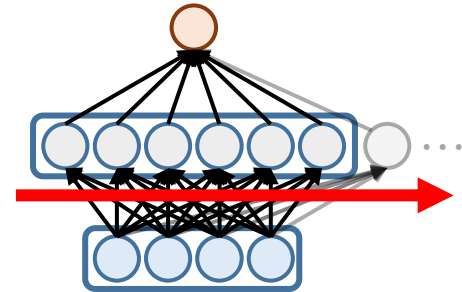
$$K = \exp(\mathcal{O}(Md)) \log(1/\epsilon)$$

iteration complexity to achieve  $\epsilon$  error.

→ Cannot be applied to wide neural network.

$$f(z) = \frac{1}{M} \sum_{j=1}^M r_j \sigma(w_j^\top z)$$

Non-linear with respect to the parameters  $(r_j, w_j)_{j=1}^M$ .



★ Mean field limit:

$$f(z) = \frac{1}{M} \sum_{j=1}^M r_j \sigma(w_j^\top z) \xrightarrow{M \rightarrow \infty} f_\mu(z) = \int r \sigma(w^\top z) d\mu(r, w)$$

Linear with respect to  $\mu$ .

[Nitanda&Suzuki, 2017][Chizat&Bach, 2018][Mei, Montanari&Nguyen, 2018][Rotskoff&Vanden-Eijnden, 2018]

**Loss function (empirical risk + regularization):**

$$F(\mathcal{X}) = \frac{1}{n} \sum_{i=1}^n \ell_i(f_{\mathcal{X}}(z_i)) + \lambda_1 \frac{1}{M} \sum_{j=1}^M \|X^{(j)}\|^2$$

$$\Rightarrow F(\mu) = \frac{1}{n} \sum_{i=1}^n \ell_i(f_\mu(z_i)) + \lambda_1 \mathbb{E}_\mu[\|x\|^2]$$

**Convex** w.r.t.  $\mu$  if the loss  $\ell_i$  is convex (e.g., squared / logistic loss).

$$\mathcal{L}(\mu) := \underbrace{F(\mu)}_{\text{convex}} + \lambda_2 \text{Ent}(\mu)$$

convex + strictly convex = strictly convex

$$F(\theta\mu + (1 - \theta)\nu) \leq \theta F(\mu) + (1 - \theta)F(\nu)$$

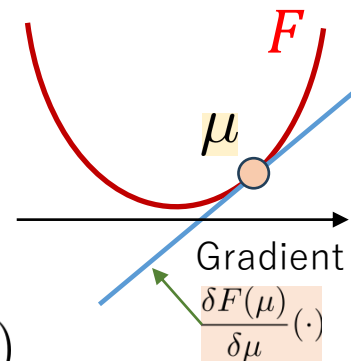
$$(\text{Ent}(\mu) = \int \log(\mu) d\mu)$$

## Mean field Langevin dynamics:

➤ SDE the Fokker-Planck equation of which corresponds to the Wasserstein GF:

$$dX_t = -\nabla \frac{\delta F(\mu_t)}{\delta \mu}(X_t) dt + \sqrt{2\lambda_2} dB_t$$

$$\mu_t = \text{Law}(X_t)$$



$$\text{GLD: } dX_t = -\nabla L(X_t) dt + \sqrt{2\lambda_2} dB_t, \quad \frac{\delta F(\mu)}{\delta \mu}(\cdot) = L(\cdot)$$

$$F(\mu) = \int L(x) d\mu$$

### Definition (first variation)

The first variation  $\frac{\delta F}{\delta \mu}: \mathcal{P} \times \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as a continuous functional such as

$$\lim_{\epsilon \rightarrow 0} \frac{F(\epsilon\nu + (1 - \epsilon)\mu) - F(\mu)}{\epsilon} = \int \frac{\delta F(\mu)}{\delta \mu}(x) d(\nu - \mu)(x)$$

# MF-LD to optimize mean field NN 111

Loss function:  $F(\mu) = \frac{1}{n} \sum_{i=1}^n \ell_i(f_\mu(z_i)) + \lambda_1 \mathbb{E}_\mu[\|X\|^2]$   $f_\mu(z) = \int h_x(z) d\mu(x)$

$$X_{k+1}^{(j)} = X_k^{(j)} - \frac{\eta_k}{M} \left( \frac{1}{n} \sum_{i=1}^n \ell'_i(\underline{f_{x_k}}(z_i)) \nabla_{X^{(j)}} h_{X_k^{(j)}}(z_i) + \lambda_1 X_k^{(j)} \right) + \sqrt{2\eta_k \lambda_2} \xi_k^{(j)}$$

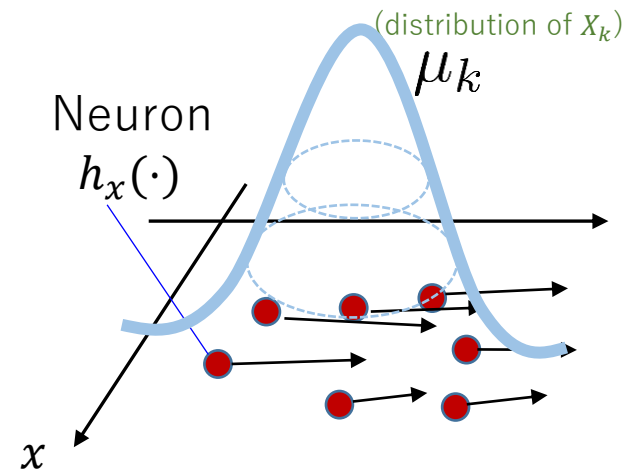
➔ 
$$X_{k+1} = X_k - \eta_k \left( \frac{1}{n} \sum_{i=1}^n \ell'_i(\underline{f_{\mu_k}}(z_i)) \nabla_X h_{X_k}(z_i) + \lambda_1 X_k \right) + \sqrt{2\eta_k \lambda_2} \xi_k$$

$$\mu_k = \text{Law}(X_k) \quad \nabla \frac{\delta F(\mu_k)}{\delta \mu}(X_k)$$

$$\frac{\delta F(\mu)}{\delta \mu}(X) = \frac{1}{n} \sum_{i=1}^n \ell'_i(f_\mu(z_i)) h_X(z_i) + \lambda_1 \|X\|^2$$

**Discrete time MFLD:**

$$X_{k+1} = X_k - \eta_k \nabla \frac{\delta F(\mu_k)}{\delta \mu}(X_k) + \sqrt{2\eta_k \lambda_2} \xi_k$$

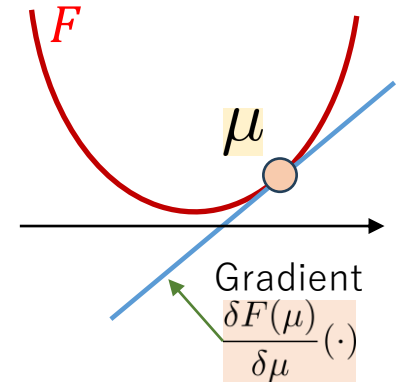


$$\mathcal{L}(\mu) = \underline{F(\mu)} + \lambda_2 \text{Ent}(\mu)$$

Linearized objective at  $\mu$ :

$$(\text{Ent}(\mu) = \int \log(\mu) d\mu)$$

$$\bar{\mathcal{L}}_{\mu}(\nu) = \int \frac{\delta F(\mu)}{\delta \mu}(x) d\nu(x) + \lambda_2 \text{Ent}(\nu)$$



Minimizer



$$p_{\mu}(x) \propto \exp\left(-\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x)\right)$$

$$F(\mu) = \int L(x) d\mu$$

$$\Rightarrow p_{\mu} \propto \exp(-\lambda_2^{-1} L(x))$$

**Proximal Gibbs measure**

- The proximal Gibbs measure is a kind of “tentative” target.
- It plays important role in the convergence analysis.



# Entropy sandwich

Proximal Gibbs measure:

$$p_\mu(x) \propto \exp\left(-\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x)\right) \quad p_\mu = \arg \min_{\nu \in \mathcal{P}} (\nu - \mu) \frac{\delta F(\mu)}{\delta \mu} + \lambda_2 \text{Ent}(\nu)$$

Theorem (Entropy sandwich) [Nitanda, Wu, Suzuki (AISTATS2022)][Chizat (2022)]

$$\mu^* = \arg \min_{\mu \in \mathcal{P}} \mathcal{L}(\mu)$$

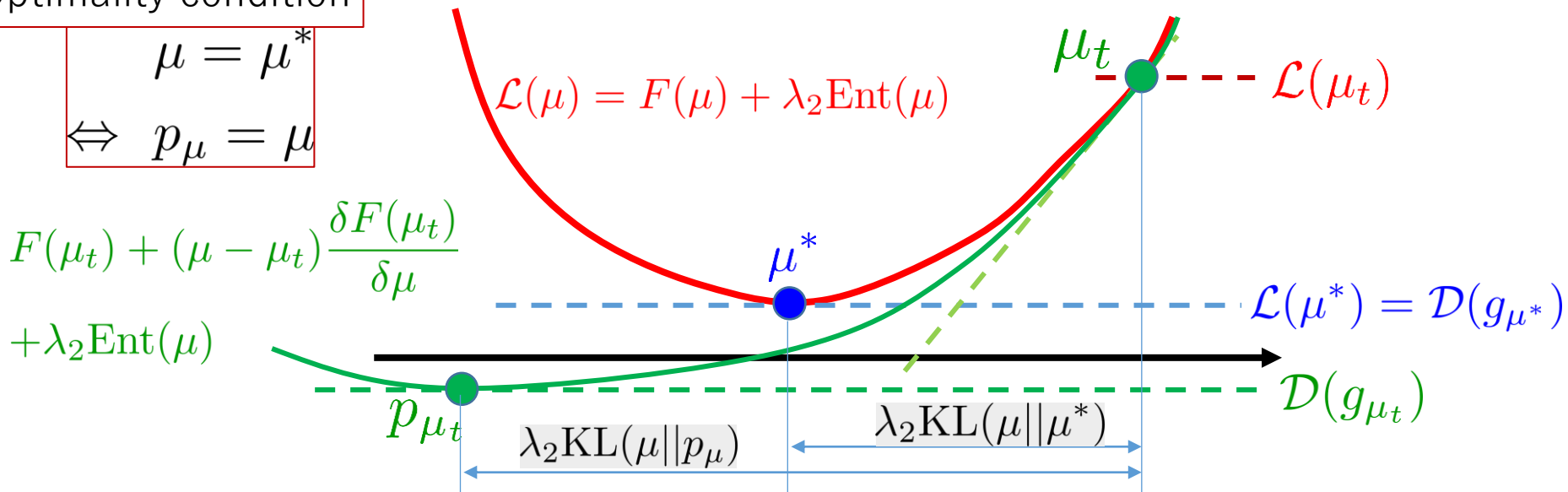
$$\lambda_2 \text{KL}(\mu || \mu^*) = \mathcal{L}(\mu) - \mathcal{L}(\mu^*) \leq \mathcal{L}(\mu) - \mathcal{D}(g_\mu) = \lambda_2 \text{KL}(\mu || p_\mu)$$

$$\mathcal{D}(g_{\mu^*})$$

Optimality condition

$$\mu = \mu^*$$

$$\Leftrightarrow p_\mu = \mu$$



# Duality (informal)

[Nitanda, Oko, Wu, Suzuki (ICML2023); Nitanda, Wu, Suzuki (AISTATS2022); Oko, Suzuki, Nitanda, Wu (ICLR2022)]

$$(\text{Ent}(\mu) = \int \log(\mu) d\mu)$$

Primal  $\min_{\mu \in \mathcal{P}} \mathcal{L}(\mu) = F(\mu) + \lambda_2 \text{Ent}(\mu)$

$$\parallel \min_{x \in \mathcal{X}} f(Ax) + g(x) = - \min_{g \in \mathcal{Y}^*} f^*(g) + g^*(-A^*g) \quad \text{(Fenchel's duality theorem)}$$

[Rockafellar (1967)]

$A : \mathcal{X} \rightarrow \mathcal{Y}$  (bounded linear)

Dual  $\max_{g: \mathbb{R}^d \rightarrow \mathbb{R}} \mathcal{D}(g) = -F^*(g) - \lambda_2 \log \left( \int \exp \left( -\frac{g(x)}{\lambda_2} \right) dx \right)$

$$F^*(g) := \sup_{\mu \in \mathcal{P}} \left\{ \int g(x) d\mu(x) - F(\mu) \right\}$$

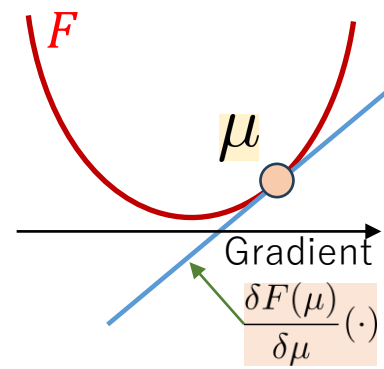
**Primal-Dual variable correspondence:**

$$\begin{array}{ccc} \text{(P)} & & \text{(D)} \\ \mu & \longrightarrow & g_\mu(x) := \frac{\delta F(\mu)}{\delta \mu}(x) \\ & & \longrightarrow \text{(P)} \end{array} \quad p_\mu(x) \propto \exp \left( -\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x) \right)$$

**Duality gap and divergence:**

$$\mu^* = \arg \min_{\mu \in \mathcal{P}} \mathcal{L}(\mu)$$

- $\mathcal{L}(\mu) - \mathcal{D}(g_\mu) = \lambda_2 \text{KL}(\mu || p_\mu) \geq 0$
- $\mathcal{L}(\mu^*) = \mathcal{D}(g_{\mu^*}) \Leftrightarrow \mu^* = p_{\mu^*}$   
(optimality condition)



Proximal Gibbs measure:

$$p_\mu(x) \propto \exp\left(-\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x)\right) \quad p_\mu = \arg \min_{\nu \in \mathcal{P}} (\nu - \mu) \frac{\delta F(\mu)}{\delta \mu} + \lambda_2 \text{Ent}(\nu)$$

## Assumption (Log-Sobolev inequality)

c.f., Polyak-Lojasiewicz condition  
 $f(x) - f(x^*) \leq C \|\nabla f(x)\|^2$

There exists  $\alpha > 0$  such that for any probability measure  $\nu$  (abs. cont. w.r.t.  $p_\mu$ )

$$\text{KL}(\nu || p_\mu) \leq \frac{1}{2\alpha} I(\nu || p_\mu)$$

KL-div

$$\text{KL}(\nu || \mu) = \int \log \left( \frac{d\nu}{d\mu} \right) d\nu$$

Fisher-div

$$I(\nu || \mu) = \int \left\| \nabla \log \frac{d\nu}{d\mu} \right\|^2 d\nu$$

Theorem (Linear convergence) [Nitanda, Wu, Suzuki (AISTATS2022)][Chizat (2022)]

If  $p_{\mu_t}$  satisfies the LSI condition for any  $t \geq 0$ , then

$$\mathcal{L}(\mu_t) - \mathcal{L}(\mu^*) \leq \exp(-2\alpha\lambda_2 t) (\mathcal{L}(\mu_0) - \mathcal{L}(\mu^*))$$

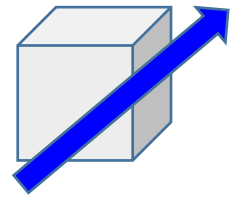
**The rate of convergence is characterized by LSI**

# Proof outline of convergence

- MF-LD obeys the following nonlinear Fokker-Planck equation:

$$\begin{aligned} \partial_t \mu_t &= \lambda_2 \Delta_x \mu_t + \nabla \cdot \left[ \mu_t \nabla \frac{\delta F(\mu_t)}{\delta \mu} \right] \\ &= \nabla \cdot \left[ \underbrace{\left( \lambda_2 \nabla \log(\mu_t) + \nabla \frac{\delta F(\mu_t)}{\delta \mu} \right)}_{=: -v_t} \mu_t \right] \\ &= -\nabla \cdot [v_t \mu_t] \quad \text{[Continuity equation]} \end{aligned}$$

Mass:  $\mu_t(x)$



Vector field:  $b(x, \mu_t)$

Then,

$$\mathcal{L}(\mu) := F(\mu) + \lambda_2 \text{Ent}(\mu)$$

$$\frac{d}{dt} \mathcal{L}(\mu_t) = \int \left\langle v_t, \nabla \frac{\delta \mathcal{L}(\mu_t)}{\delta \mu} \right\rangle d\mu_t \quad (\text{:continuity equation})$$

(Definition of  $p_{\mu_t}$ )

$$p_{\mu}(x) \propto \exp\left(-\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x)\right)$$

$$= \int \left\langle v_t, \nabla \frac{\delta F(\mu_t)}{\delta \mu} + \lambda_2 \nabla \log(\mu_t) \right\rangle d\mu_t$$

$$= - \int \|v_t\|^2 d\mu_t = -\lambda_2^2 I(\mu_t || p_{\mu_t})$$

**LSI & Entropy sandwich**

$$\leq -2\alpha \lambda_2^2 \text{KL}(\mu_t || p_{\mu_t}) \leq -2\alpha \lambda_2 (\mathcal{L}(\mu_t) - \mathcal{L}(\mu^*))$$

✘ Since  $\frac{\delta F(\mu_t)}{\delta \mu}$  nonlinearly depends on  $\mu_t$ , we say “nonlinear Fokker-Planck”.

$$\text{GLD: } F(\mu) = \int L(x) d\mu \Rightarrow \frac{\delta F(\mu)}{\delta \mu}(\cdot) = L(\cdot)$$

# Log-Sobolev inequality

L2-regularized loss function for mean field 2-layer NN:

$$f_\mu(z) = \int h_x(z) d\mu(x) \quad \text{where} \quad h_x(z) = r\sigma(w^\top z) \quad \text{for} \quad x = (r, w)$$

$$F(\mu) = \frac{1}{n} \sum_{i=1}^n \ell_i(f_\mu(z_i)) + \lambda_1 \mathbb{E}_\mu[\|X\|^2]$$

➔ Proximal Gibbs:

$$p_\mu(x) \propto \exp\left(-\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x)\right) = \exp\left[-\frac{1}{\lambda_2} \left(\underbrace{\frac{1}{n} \sum_{i=1}^n \ell'_i(f_\mu(z_i)) h_x(z_i)}_{\text{Bounded } (\leq B)} + \underbrace{\lambda_1 \|x\|^2}_{\text{Strongly convex}}\right)\right]$$

If  $\sup_z |\ell'_i(f_\mu(\cdot)) h_x(\cdot)| \leq B$ , the proximal Gibbs measure  $p_\mu$  satisfies the LSI with a constant  $\alpha$  with

$$\alpha \geq \frac{2\lambda_1}{\lambda_2} \exp(-4B/\lambda_2)$$

∴ **Bakry-Emery criterion** (1985) and **Holley-Strook bounded perturbation lemma** (1987)

Mean field Langevin dynamics can be applied to several problems where a distribution is optimized.

- **Nonparametric density estimation via MMD minimization**

$$F(\mu) = \text{MMD}^2(g * \mu, \hat{\mu}_n) + \lambda_1 \mathbb{E}_\mu[\|x\|^2]$$

$k$ : positive definite kernel

$$\text{MMD}^2(\nu_1, \nu_2) := \|k_{\nu_1} - k_{\nu_2}\|_{\mathcal{H}_k}^2$$

where  $k_\mu = \int k(x, \cdot) \mu(dx)$  (kernel embedding).

➤  $g(x) = \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \exp\left(-\frac{\|x\|^2}{2\sigma^2}\right)$

➤  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  : Empirical distribution (training data)

(see also Chizat (2022, TMLR))

- **Variational inference to approximate Bayesian posterior**

$$F(\mu) = \text{KSD}(\mu) + \lambda_1 \mathbb{E}_\mu[\|x\|^2]$$

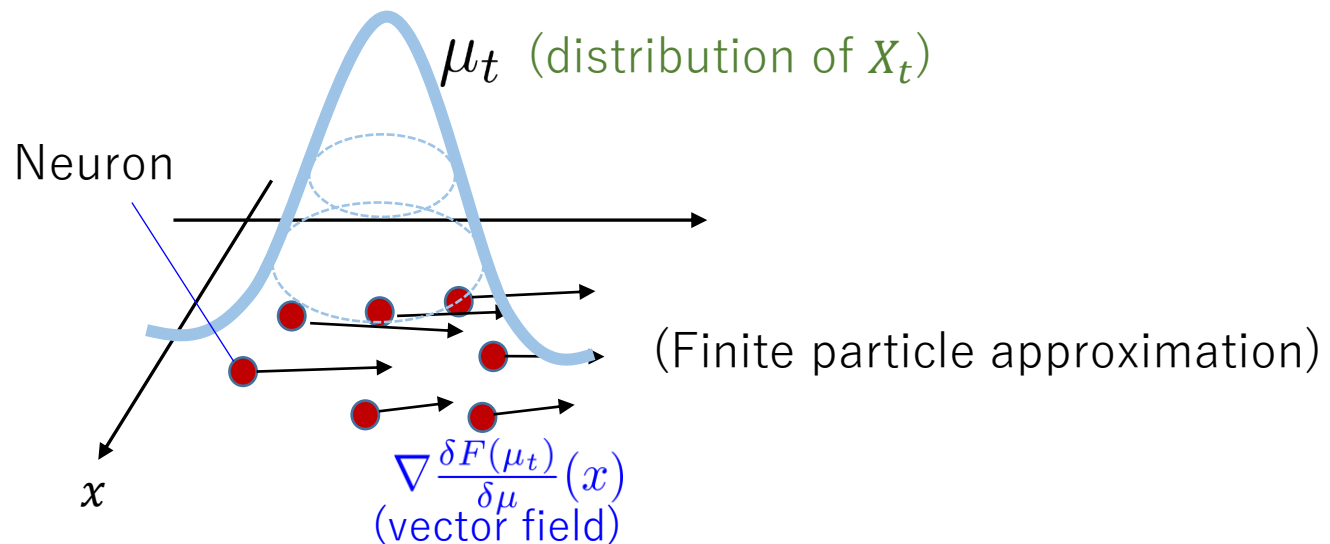
(KSD: Kernel Stein Discrepancy from a posterior distribution)

# Finite particles & discrete time algorithm

We have obtained a convergence of infinite width and continuous time dynamics.

## Question:

Can we evaluate a finite particles & discrete time approximation errors?



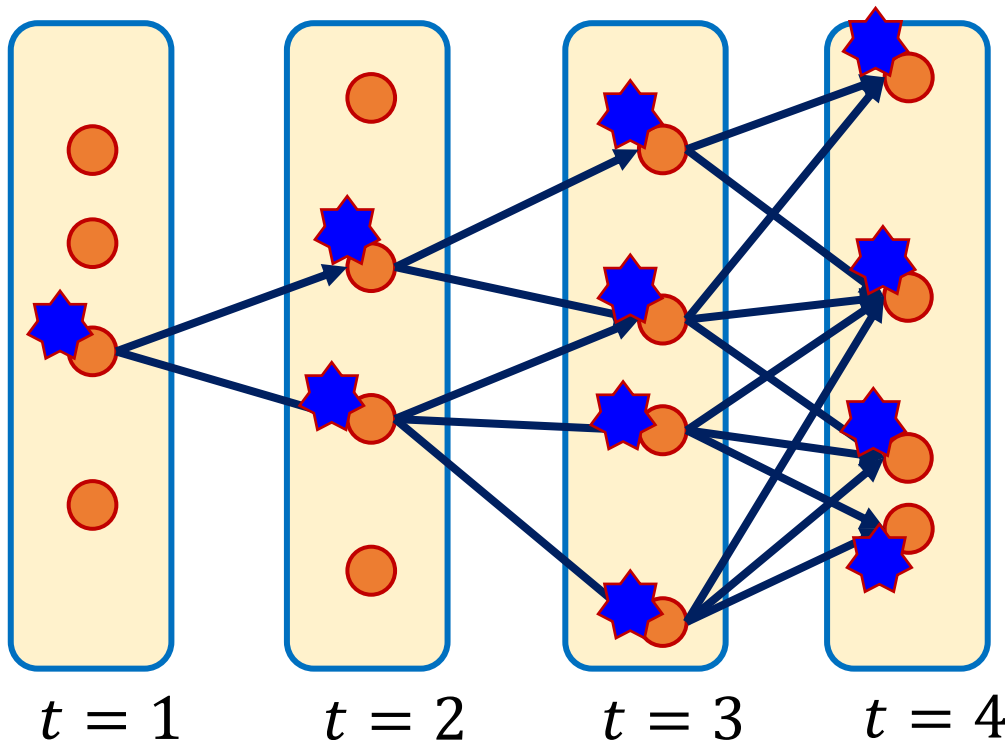


- SDE of interacting particles (McKean, Kac, ..., 60')

## Propagation of chaos [Sznitman, 1991; Lacker, 2021]:

The particles behave as if they are independent as the number of particles increases to infinity.

Finite particle approximation error can be amplified through time.  
→ It is difficult to bound the perturbation uniformly over time.



- A naïve evaluation gives exponential growth on time:

$$\exp(t) / M$$

[Mei et al. (2018, Theorem 3)]

- Weak interaction/Strong regularization in existing work

$$dX_t = -\nabla \frac{\delta F(\mu_t)}{\delta \mu}(X_t) dt + \sqrt{2\lambda_2} dB_t$$

$M$  particles  $(X_k^{(i)})_{i=1}^M$

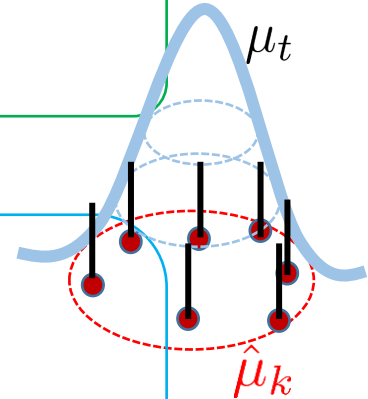
(time discretization)

$$X_{k+1}^{(i)} = X_k^{(i)} - \eta_k v_k^i + \sqrt{2\eta_k \lambda_2} \xi_k^{(i)}$$

where  $\mathbb{E}[v_k^i] = \nabla \frac{\delta F(\hat{\mu}_k)}{\delta \mu}(X_k^i)$  and  $\hat{\mu}_k = \frac{1}{M} \sum_{i=1}^M \delta_{X_k^{(i)}}$

(stochastic gradient)

(space discretization)



➤ Noisy gradient descent on 2-layer NN with finite width.

- **Time discretization:**  $t \rightarrow k\eta$  ( $\eta$ : step size,  $k$ : # of steps)
- **Space discretization:**  $\mu_t$  is approximated by  $M$  particles

$$\mu_t \rightarrow \hat{\mu}_k = \frac{1}{M} \sum \delta_{X_k^{(i)}}$$

- **Stochastic gradient:**  $\nabla \frac{\delta F(\mu)}{\delta \mu} \rightarrow v_k^i$

# Convergence analysis

$p_\mu(x) \propto \exp\left(-\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x)\right)$  : proximal Gibbs measure

Theorem (One-step update) [Suzuki, Wu, Nitanda (2023)]

Suppose that  $p_\mu$  satisfies log-Sobolev inequality with a constant  $\alpha$ . Under smoothness and boundedness of the loss function, it holds that

$$\mathcal{L}^{(M)}(\hat{\mu}_{k+1}) - \mathcal{L}(\mu^*)$$

$$\leq \exp(-\lambda_2 \eta_k \alpha) \left( \mathcal{L}^{(M)}(\hat{\mu}_k) - \mathcal{L}(\mu^*) \right)$$

$$+ C \left( \underbrace{\eta_k^3 + \lambda_2 \eta_k^2}_{\text{Time discr.}} + \underbrace{\frac{\eta_k}{M}}_{\text{Space discr.}} + \underbrace{\eta_k^{\frac{3}{2}} \lambda_2^{\frac{1}{2}} \sigma_k \tilde{\sigma}_k}_{\text{Stochastic approx.}} \right)$$

$\mathbf{O}(1/M)$

**Naïve bound:**

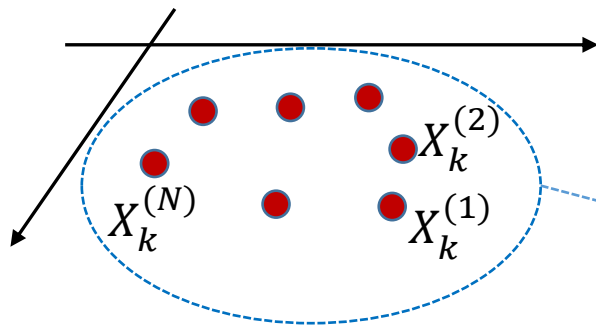
$$\eta_k \sigma_k^2$$

$$\sigma_k^2 = \max_i \mathbb{E} [\|v_k^i - \mathbb{E}[v_k^i]\|^2]$$

$$\tilde{\sigma}_k^2 = \max_i \mathbb{E} \left[ \left\| \nabla v_k^{i\top}(\mathcal{X}) - \nabla \nabla^\top \frac{\delta F(\mu, \mathcal{X})}{\delta \mu}(X^i) \right\|_{\text{op}}^2 \right]$$

## Assumption:

1.  $F: \mathcal{P} \rightarrow \mathbb{R}$  is convex and has a form of  $F(\mu) = L(\mu) + \lambda_1 \mathbb{E}_\mu [\|x\|^2]$ .
2. (smoothness)  $\left\| \nabla \frac{\delta L(\mu)}{\delta \mu}(x) - \nabla \frac{\delta L(\nu)}{\delta \mu}(y) \right\| \leq C(W_2(\mu, \nu) + \|x - y\|)$  and  
(boundedness)  $\left\| \nabla \frac{\delta L(\mu)}{\delta \mu}(x) \right\| \leq R$ . (+ second order differentiability)



$\mathcal{X}_k = \left( X_k^{(i)} \right)_{i=1}^M \sim \mu_k^{(M)}$ : Joint distribution of  $M$  particles.

Potential of the joint distribution  $\mu_k^{(M)}$  on  $\mathbb{R}^{d \times M}$ :

$$\mathcal{L}^M(\mu_k^{(M)}) = M \mathbb{E}_{\mathcal{X} \sim \mu_k^{(M)}} [F(\hat{\mu}_{\mathcal{X}})] + \lambda_2 \text{Ent}(\mu_k^{(M)}).$$

$$\text{where } \hat{\mu}_{\mathcal{X}} = \frac{1}{M} \sum_{i=1}^M \delta_{X^{(i)}} \quad (\mathcal{X} = (X^{(i)})_{i=1}^M)$$

➤ The finite particle dynamics is the Wasserstein gradient flow that minimizes  $\mathcal{L}^M$ .

**(Approximate) Uniform log-Sobolev inequality** [Chen et al. 2022]

**For any  $M$ ,**

$$\frac{1}{M} \mathcal{L}^M(\mu_k^{(M)}) - \mathcal{L}(\mu^*) \leq \frac{\lambda_2}{2\alpha} \left( \frac{1}{M} I(\mu_k^{(M)} \| p^{(M)}) \right) + \frac{C_{\alpha, \lambda_2}}{M}$$

(Fisher divergence)

$$\text{where } p^{(M)}(\mathcal{X}) \propto \exp\left(-\frac{M}{\lambda_2} F(\hat{\mu}_{\mathcal{X}})\right)$$

Recall  $\mathcal{L}(\mu) = F(\mu) + \lambda_2 \text{Ent}(\mu)$

[Chen, Ren, Wang. Uniform-in-time propagation of chaos for mean field Langevin dynamics. arXiv:2212.03050, 2022.]

## SG-MFLD

$$F(\mu) = \frac{1}{n} \sum_{j=1}^n \ell_j(\mu) + \lambda_1 \mathbb{E}[\|X\|^2] \quad (\text{finite sum}),$$

$$v_k^i = \frac{1}{B} \sum_{j \in I_k} \nabla \frac{\delta \ell_j(\hat{\mu}_k)}{\delta \mu} (X_k^{(i)}) + \lambda_1 X_k^{(i)} \quad (\text{stochastic gradient})$$

(Mini-batch size =  $B$ )

$$\mathcal{L}^{(N)}(\hat{\mu}_k) - \mathcal{L}(\mu^*) \lesssim \exp(-\lambda_2 \eta k \alpha) + \frac{1}{\alpha \lambda_2} \left( \underbrace{\eta^2 + \lambda_2 \eta}_{\text{Time discr.}} + \underbrace{\frac{1}{M}}_{\text{Space discr.}} + \underbrace{\frac{\eta + \sqrt{\eta \lambda_2}}{B}}_{\text{Stochastic approx.}} \right)$$

Iteration complexity:

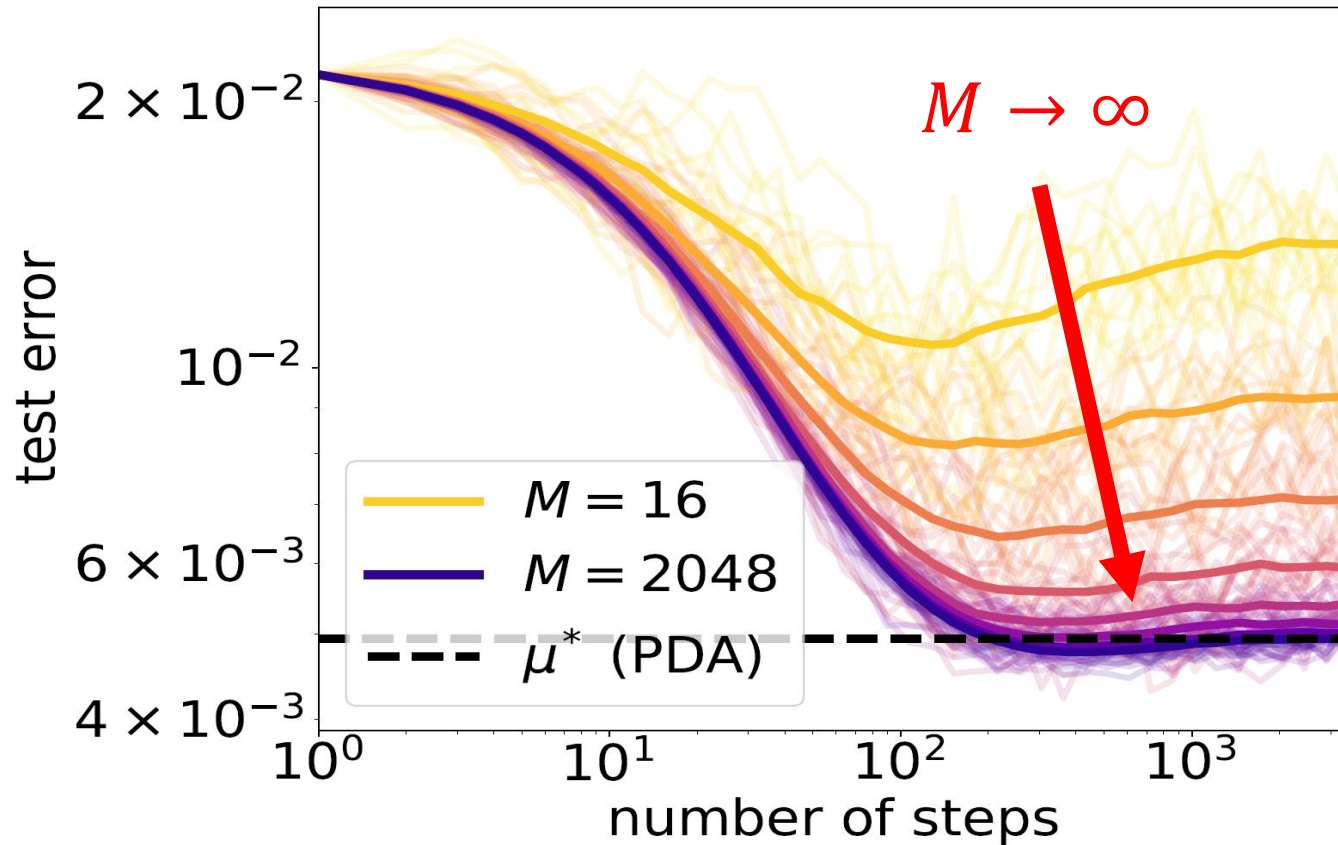
- Approximation errors are uniform in time.
- No exponential dependency on  $M$  (number of

$$k = O \left( \frac{1}{\epsilon \alpha} + \sqrt{\frac{1}{\lambda_2 \epsilon \alpha}} + \left( \frac{1}{\lambda_2 \epsilon \alpha} \right)^2 \frac{\lambda_2}{B^2} + \frac{1}{\lambda_2 \epsilon \alpha B} \right) \frac{1}{\lambda_2 \alpha} \log(\epsilon^{-1})$$

to achieve  $\epsilon + O(1/(\lambda_2 \alpha N))$  accuracy.

- $B = \sqrt{1/(\lambda_2 \alpha \epsilon)}$  is the optimal mini-batch size.  $\rightarrow k = O(\log(\epsilon^{-1})/\epsilon)$ .

# Numerical experiment



Test error v.s. Number of steps  
(regularization term:  $r(x) = \|x\|^2$ )

# Generalization error analysis

So far, we have obtained convergence of MFLD.

⇒ How effective is the feature learning of MFLD in terms of generalization error?

- Benefit of feature learning?

**Neural network** vs **Kernel method**

(NTK vs mean field)

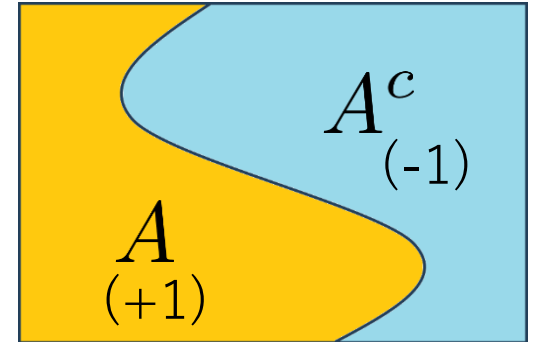


**Problem setting (classification):**

$$Y = \mathbf{1}_A(Z) - \mathbf{1}_{A^c}(Z) \in \{\pm 1\}$$

$$Z \in \mathbb{R}^d$$

Training data:  $\{(z_i, y_i)\}_{i=1}^n$



**Loss function and model:**

$$F(\mu) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f_\mu(z_i)) + \lambda_1 \mathbb{E}_\mu[\|x\|^2]$$

$$f_\mu(z) = \int h_x(z) d\mu(x)$$

➤ **Logistic loss:**  $\ell(yf) = \log(1 + \exp(-yf))$

➤ **Tanh activation:**  $h_x(z) = \bar{R} \cdot [\tanh(\langle x_1, z \rangle + x_2) + 2 \cdot \tanh(x_3)] / 3$

# Assumptions

Objective of MFLD:

$$\begin{aligned}\mathcal{L}(\mu) &= \frac{1}{n} \sum_{i=1}^n \ell(y_i f_\mu(z_i)) + \lambda_1 \mathbb{E}_\mu[\|x\|^2] + \lambda \text{Ent}(\mu) \\ &= \frac{1}{n} \sum_{i=1}^n \ell(y_i f_\mu(z_i)) + \lambda \text{KL}(\nu, \mu)\end{aligned}$$

where  $\nu = N(0, \lambda/(2\lambda_1))$ . KL-regularization

Assumption

There exists  $\mu^*$  such that

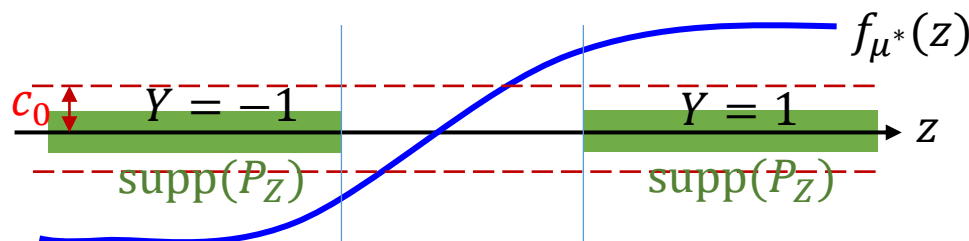
1.  $\text{KL}(\nu, \mu^*) \leq R$ ,

2.  $Y f_{\mu^*}(Z) \geq c_0$  (a.s.),

for some constants  $R, c_0 > 0$ .

(+ classification calibration condition)

The Bayes classifier is attained by  $\mu^*$  with a bounded KL-div from  $\nu$ .



# Main theorem

## Theorem 1

Suppose that  $\lambda = \Theta(1/R)$ , then it holds that

$$P(Y f_{\hat{\mu}}(Z) \leq 0) \lesssim \left[ \frac{\bar{R}^2 R}{n} (1 + t + \log(\log(n))) + \frac{\bar{R} + 1/(Rn)}{n} \right]$$

Class. error

with probability  $1 - \exp(-t)$ .

$$O\left(\frac{\bar{R}^2 R}{n}\right)$$

**Existing bound:** Chen et al. (2020); Nitanda, Wu, Suzuki (2021)

Class. Error  $\leq O\left(\frac{1}{\sqrt{n}}\right)$ . (Rademacher complexity bound)

- Our bound provides *fast learning rate* (faster than  $1/\sqrt{n}$ ).

$$O(R/n) \ll O(1/\sqrt{n})$$

$$\mathcal{L}(\mu) = \frac{1}{n} \sum_{j=1}^n \ell(y_j f_{\mu}(z_j)) + \lambda \text{KL}(\nu, \mu)$$

- $\mu^*: \text{KL}(\nu, \mu^*) \leq R, Y f_{\mu^*}(Z) \geq c_0$
- $h_x(z) = \bar{R} \cdot [\tanh(\langle x_1, z \rangle + x_2) + 2 \cdot \tanh(x_3)]/3$

## Theorem 2

Suppose that  $\lambda = \Theta(1/R)$  and

$$\frac{\bar{R}^3 R^2}{n} (\lambda + \bar{R}) \lesssim c_0$$

then it holds that

$$P(Y f_{\hat{\mu}}(Z) \leq 0) = 0 \quad \text{with probability } 1 - \exp\left(-\frac{c_0^3 n / R^2}{2\bar{R}^4}\right)$$

Theorem 1:  $\mathbb{E}[\text{Class. Error}] \leq O\left(\frac{R}{n}\right)$ .

Theorem 2:  $\mathbb{E}[\text{Class. Error}] \leq O(\exp(-O(n/R^2)))$  if  $n \geq R^2$ .

If we have sufficiently large training data, we have exponential convergence of test error.

We only need to evaluate  $R$  to obtain a test error bound.

$$\bar{F}(\mu) := \mathbb{E}_{X,Y}[\ell(Y f_{\mu}(X))]$$

**Expected loss**

$$\mu^{\circ} = \arg \min_{\mu \in \mathcal{P}} \{ \bar{F}(\mu) + \lambda \text{KL}(\nu, \mu) \}$$

Local Rademacher complexity yields the following bound:

$$\underbrace{\bar{F}(\hat{\mu}) - \bar{F}(\mu^{\circ}) - (\hat{\mu} - \mu^{\circ}) \frac{\delta \bar{F}(\mu^{\circ})}{\delta \mu}}_{\text{(I)}} + \underbrace{\lambda \text{KL}(\mu^{\circ}, \hat{\mu})}_{\text{(II)}} \lesssim \sqrt{\frac{\text{KL}(\mu^{\circ}, \hat{\mu})}{n}}$$

**Adaptive to the KL-divergence**  $\leq \frac{1}{2n\lambda} + \frac{\lambda}{2} \text{KL}(\mu^{\circ}, \hat{\mu})$

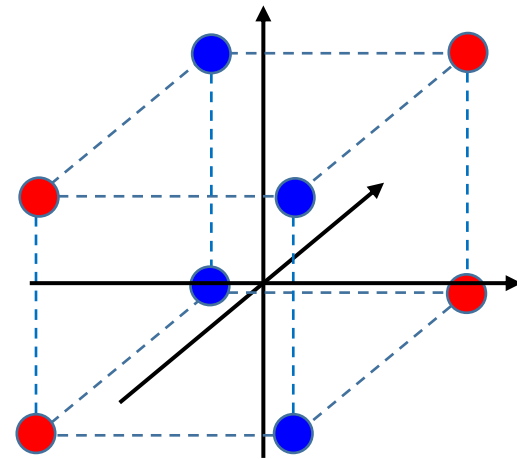
By setting  $\lambda = O(1/R)$ ,

- The term (I) gives the first bound:  $O(R/n)$ .
- The term (II) gives the second bound:  $O(\exp(-O(n/R^2)))$ .

# Example: $k$ -sparse parity problem<sup>134</sup>

- $k$ -sparse parity problem on high dimensional data
  - $Z \sim \text{Unif}(\{-1,1\}^d)$  (up to freedom of rotation)
  - $Y = \prod_{j=1}^k Z_j$
  - ※ Assume we don't know which coordinate corresponds to  $Z_j$ .

**Q: Can we learn sparse  $k$ -parity with GD?  
Is there any benefit of neural network?**



$k = 2$ : XOR problem  
 $d = 3, k = 2$

Complexity to learn XOR function ( $k = 2$ )

Reference	Algorithm	Technique	$m$	$n$	$t$
(Ji and Telgarsky, 2020b)	SGD	perceptron	$d^8$	$d^2/\epsilon$	$d^2/\epsilon$
Theorem 2.1	SGD	perceptron	$d^2$	$d^2/\epsilon$	$d^2/\epsilon$
(Barak et al., 2022)	2-phase SGD	correlation	$\mathcal{O}(1)$	$d^4/\epsilon^2$	$d^2/\epsilon^2$
(Wei et al., 2018)	WF+noise	margin	$\infty$	$d/\epsilon$	$\infty$
(Chizat and Bach, 2020)	WF	margin	$\infty$	$d/\epsilon$	$\infty$
Theorem 3.3	scalar GF	margin	$d^d$	$d/\epsilon$	$\infty$

Table 1 of [Telgarsky: Feature selection and low test error in shallow low-rotation ReLU networks, 2020]

## Reminder

Suppose that there exists  $\mu^*$  such that

$$\mu^*: \text{KL}(\nu, \mu^*) \leq R, Yf_{\mu^*}(Z) \geq c_0 \text{ (perfect classifier with margin } c_0)$$

Then,

$$\text{Theorem 1: } \mathbb{E}[\text{Class. Error}] \leq O\left(\frac{R}{n}\right).$$

$$\text{Theorem 2: } \mathbb{E}[\text{Class. Error}] \leq O(\exp(-O(n/R^2))) \text{ if } \underline{n \geq R^2}.$$

We can evaluate  $R$  required for the  $k$ -sparse parity problem:

## Lemma

For the  $k$ -parity problem, we may take

$$R = \mathcal{O}(k \log(k)d)$$

# Generalization

Corollary (Test accuracy of MFLD)

- Setting 1:  $n > d$

- Test error (classification error) =  $\mathbf{O}(d/n)$

- Setting 2:  $n > d^2$

- Test error (classification error) =  $\mathbf{O}(\exp(-n/d^2))$

Our analysis provides

- better sample complexity
- discrete-time/finite-width analysis
- $d$  and  $k$  are “decoupled.”

(Computational complexity is  $\exp(O(d))$ ) (But, can be relaxed to  $O(1)$  if  $X$  is anisotropic)

**These are better than NTK (kernel method);**

**Sample complexity of NTK  $n = \Omega(d^k)$  vs NN  $n = \mathbf{O}(d)$**

Trade-off between computational complexity and sample complexity.

Authors	regime/method	$k$ -parity	class error	width	# iterations
Ji and Telgarsky (2019)	NTK/SGD	No	$d^2/n$	$d^8$	$d^2/\epsilon$
Telgarsky (2023)	NTK/SGD	No	$d^2/n$	$d^2$	$d^2/\epsilon$
Barak et al. (2022)	Two phase SGD	Yes	$d^{(k+1)/2}/\sqrt{n}$	$O(1)$	$d/\epsilon^2$
Wei et al. (2019)	mean-field/GF	No	$d/n$	$\infty$	$\infty$
Telgarsky (2023)	mean-field/GF	No	$d/n$	$d^d$	$\infty$
Ours	mean-field/MFLD	Yes	$\exp(-O(n/d^2))$	$e^{O(d)}$	$e^{O(d)}$
Ours	mean-field/MFLD	Yes	$d/n$	$e^{O(d)}$	$e^{O(d)}$



# Discussion

- The CSQ lower bound states that  $\mathcal{O}(d^{k-1})$  sample complexity is optimal for methods with polynomial order computational complexity. [Abbe et al. (2023); Refinetti et al. (2021); Ben Arous et al. (2022); Damian et al. (2022)]
- On the other hand, our analysis is about full-batch GD.

	Minibatch size	# of iterations	Sample complexity
<b>Our analysis</b>	$n$	$e^d$	$d$
SGD (CSQ-lower bound)	1	$d^{k-1}$	$d^{k-1}$

We obtain a better sample complexity than  $\mathcal{O}(d^{k-1})$  with higher computational complexity.

→ We can obtain a polynomial order method with MFLD for anisotropic input.

## Def (Correlational Statistical Query (CSQ) algorithm)

[Ben-David, Itai, Kushilevitz, 1995; Kearns, 1998; Bshouty, Feldman, 2002]

A CSQ algorithm can access the data only via queries  $\phi: \mathbb{R}^d \rightarrow [-1,1]$  and returns  $g \in \mathbb{R}^d$  with tolerance  $\tau$  such that

$$g \in \mathbb{E}_{Z,Y}[\phi(X)Y] + [-\tau, \tau]$$

Ex. Online SGD for a squared loss:

$$\mathbb{E}_Z[(f^\circ(Z) - f_{\mathcal{X}}(Z))^2] \quad \rightarrow \quad f^\circ(z_i) \nabla f_{\mathcal{X}}(z_i) \quad (\text{CSQ})$$

• Boolean case:  $z \sim \text{Unif}(\{-1, +1\}^d)$

**$k$ -parity**:  $f^\circ(z) = \prod_{j \in S} z_j$  where  $|S| = k$ .

• Gaussian case:  $x \sim N(0, I)$

**Single index model**:  $f^\circ(x) = g(w^\top x)$

where  $g: \mathbb{R} \rightarrow \mathbb{R}$  and  $w \in \mathbb{R}^d$ .

For the Gaussian single index model, the information exponent plays an important role.

Hermite polynomial expansion of the link function:

$$g(z) = \sum_{k=1}^{\infty} \alpha_k h_k(z)$$

**Def (information exponent [Ben Arous, Gheissari, Jagannath, 2021])**

$$k^* := \arg \min_k \{k \mid \alpha_k \neq 0\}$$

The computational complexity of a CSQ algorithm is lower bounded as:

**Theorem (CSQ lower bound [Abbe, Boix-Adser`, Misiakiewicz, 2023])**

A CSQ algorithm with error tolerance  $\tau$  requires at least  $N$  queries to obtain an estimator  $\hat{f}$  s.t.  $\mathbb{E}[(\hat{f} - f^\circ)^2] \leq 0.1$  where

$$N/\tau^2 \geq \begin{cases} d^k & \text{(Boolean case),} \\ d^{k^*}/2 & \text{(Gaussian case).} \end{cases}$$

(we suppose  $k^* > 2$ )

Note that the gradient computation at each iteration consumes  $O(d)$  queries. Thus,  $d^{k-1}$  iterations are enough.

- SGD with smoothing operation achieves the Gaussian optimal rate:

Damien et al.: Smoothing the Landscape Boosts the Signal for SGD Optimal Sample Complexity for Learning Single Index Models. NeurIPS2023.

- Near optimal complexity of SGD to learn XOR problem:

Glasgow: SGD Finds then Tunes Features in Two-Layer Neural Networks with near-Optimal Sample Complexity: A Case Study in the XOR problem. ICML2024.

- Optimal SQ sample complexity to learn Gaussian single index model with the “generative” information exponent:

Damian, Pillaud-Vivien, Lee, Bruna: The Computational Complexity of Learning Gaussian Single-Index Models. arXiv:2403.05529.

The setting of  $k^*=1$ .

# Feature learning with one-step gradient descent

[Ba, Erdogdu, Suzuki, Wang, Wu, Yang: High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation. NeurIPS2022]



Jimmy Ba



Murat A. Erdogdu



Zhichao Wang



Denny Wu



Greg Yang

# Gradient descent and kernel alignment<sup>142</sup>

$$f_{\text{NN}}(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle x, w_i \rangle) = \frac{1}{\sqrt{N}} a^\top \sigma(W^\top x)$$

**Question** : Can we obtain “good” features from data by updating the first layer parameter  $W$  by gradient descent?

**Result** : GD with large step size can extract the leading term of the true function. Especially, for the single index model ( $f^*(x) = \sigma^*(\langle x, w^* \rangle)$ ), the **predictive risk provably outperforms random feature methods**.

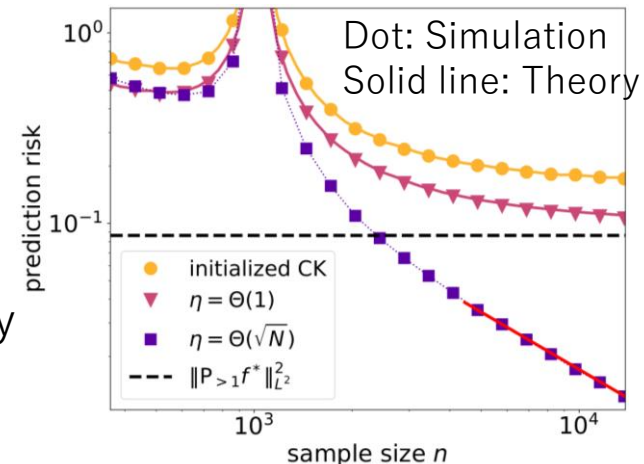
→ Kernel alignment, feature learning.

$$W_{k+1} = W_k - \eta \sqrt{N} \nabla_W L(f_{\text{NN}})$$

We consider the **proportional limit** ( $n, d, N \rightarrow \infty$ ), and evaluate predictive risk of **one-step GD**.

- $\eta = \Theta(\sqrt{N})$  can outperform random feature models.
- $\eta = \Theta(1)$  can outperform the initial setting of  $W$ .
- $\eta = o(1)$  does not improve the performance.

Gaussian equivalence property + Random matrix theory  
→ Exact risk evaluation.

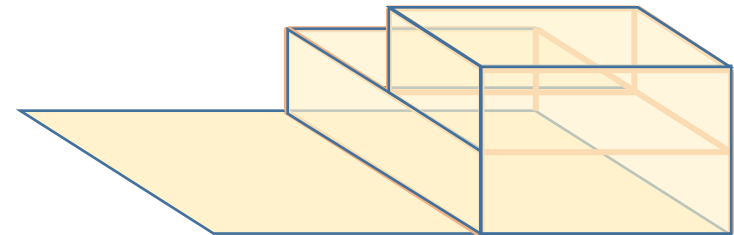


The first few step of GD with large learning rate can extract informative features.

- Staircase function

[Abbe et al., NeurIPS2021; Abbe et al., arXiv2202.08658]

Small number of gradient descent can extract nonlinear features to estimate “staircase” function. The trained features for GD can outperform random feature model.



- Benign overfitting with feature learning

[Cao et al., arXiv:2202.06526; Frei et al., arXiv:2202.05928]

Gradient descent in two-layer NN can yield benign overfitting and achieves almost the Bayes error in binary classification.

## Observation model:

$$y_i = f^*(x_i) + \epsilon_i \quad (i = 1, \dots, n)$$

where  $x_i \sim N(0, I)$ ,  $\epsilon_i \sim N(0, 1)$ , and  $x_i \in \mathbf{R}^d$ .

- We fit 2-layer NN of mean field scaling:  $(\because a_i = O_p(1/\sqrt{N}))$   
Mean field regime  $O(1/N)$

$$f_{\text{NN}}(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle x, w_i \rangle) = \frac{1}{\sqrt{N}} \overbrace{a^\top}^{\text{Mean field regime } O(1/N)} \sigma(W^\top x)$$

where  $a_i \sim N(0, 1/\underbrace{N}_{\text{var}})$  and  $W_{ij} \sim N(0, 1/\underbrace{d}_{\text{var}})$ .

## Empirical risk:

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

## Predictive risk:

$$\mathcal{R}(f) = \mathbb{E}[(f^*(X) - f(X))^2]$$

**Question: Can we provably improve the predictive risk by gradient descent?**

We analyze the risk especially for the single index model:

$$f^*(x) = \sigma^*(\langle x, \beta^* \rangle)$$



# Feature learning with optimization guarantee

$$f_{\text{NN}}(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle x, w_i \rangle) = \frac{1}{\sqrt{N}} a^\top \sigma(W^\top x)$$

$$W_{k+1} = W_k - \eta \sqrt{N} \nabla_W L(f_{\text{NN}})$$

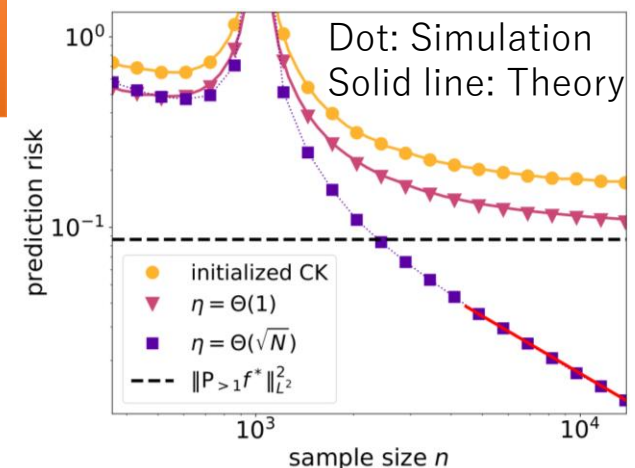
We consider the **proportional limit** ( $n, d, N \rightarrow \infty$  with  $n/d \rightarrow \psi_1, N/d \rightarrow \psi_2$ ).  
 It allows to derive precise risk.

We evaluate predictive risk of **one-step GD**.

Take home message:  
 GD with Large step-size can outperform **any** random feature model by only one-step update.

[Outline of our result]

- $\eta = \Theta(\sqrt{N})$  can get out of NTK regime and outperform random feature models.
- $\eta = \Theta(1)$  can outperform the initial setting of  $W$ .
- $\eta = o(1)$  does not improve the performance.



Feature learning vs Random feature

**Random features** (without feature learning):

- Conjugate kernel at initialization:

$$\phi_{\text{CK}}(x) = \frac{1}{\sqrt{N}} \sigma(W_0^\top x)$$

- NTK (Neural tangent kernel):

$$\phi_{\text{NTK}}(x) = \frac{1}{\sqrt{Nd}} \text{Vec}(\sigma'(W_0^\top x) x^\top)$$

Precise asymptotics has been extensively studied. (e.g., [Louart, Liao, and Couillet, 2018; Mei and Montanari, 2019])

$$\hat{a}_{\text{RF}} = \arg \min_{a \in \mathbb{R}^N} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \langle a, \phi_{\text{RF}}(x_i) \rangle)^2 + \frac{\lambda}{N} \|a\|^2 \right\} \quad \text{RF} \in \{\text{CK}, \text{NTK}\}$$

**Trained feature:**

$$\phi_{\text{CK}^{(t)}}(x) = \frac{1}{\sqrt{N}} \sigma(W_t^\top x)$$

- (1) Random feature models and
- (2) GD updates with small learning rate can learn only linear functions in the proportional

[El Karoui (2010); Ghorbani et al. (2019), Hu and Lu (2020),  $\mathcal{R}_{\mathcal{X}\mathcal{X}}(f) = \mathbb{E}[(f^*(X) - \hat{f}_{\mathcal{X}\mathcal{X}}(X))^2]$

## Theorem (Lower bound of predictive risk for RF)

If the step size is not large  $\eta = \Theta(1)$ , then for any finite number steps  $t$ , we have

$$\inf_{\lambda > 0} \min\{\mathcal{R}_{\text{CK}}(\lambda), \mathcal{R}_{\text{NTK}}(\lambda), \mathcal{R}_{\text{CK}^{(t)}}(\lambda)\} \geq \|P_{>1}f^*\|_{L^2(P_X)}^2 + o_{p,d}(1)$$

Nonlinear part cannot be trained!

$$P_{>1}f^* := (I - P_{\leq 1})f^*$$

where  $P_{\leq 1}$  is the projection operator in  $L^2(P_X)$  to the subspace consisting of linear functions and constants.

Remark: The same is true for “rotational invariant kernel” [El Karoui (2010)].

This is because in high dimensional setting, a central limit theorem yields

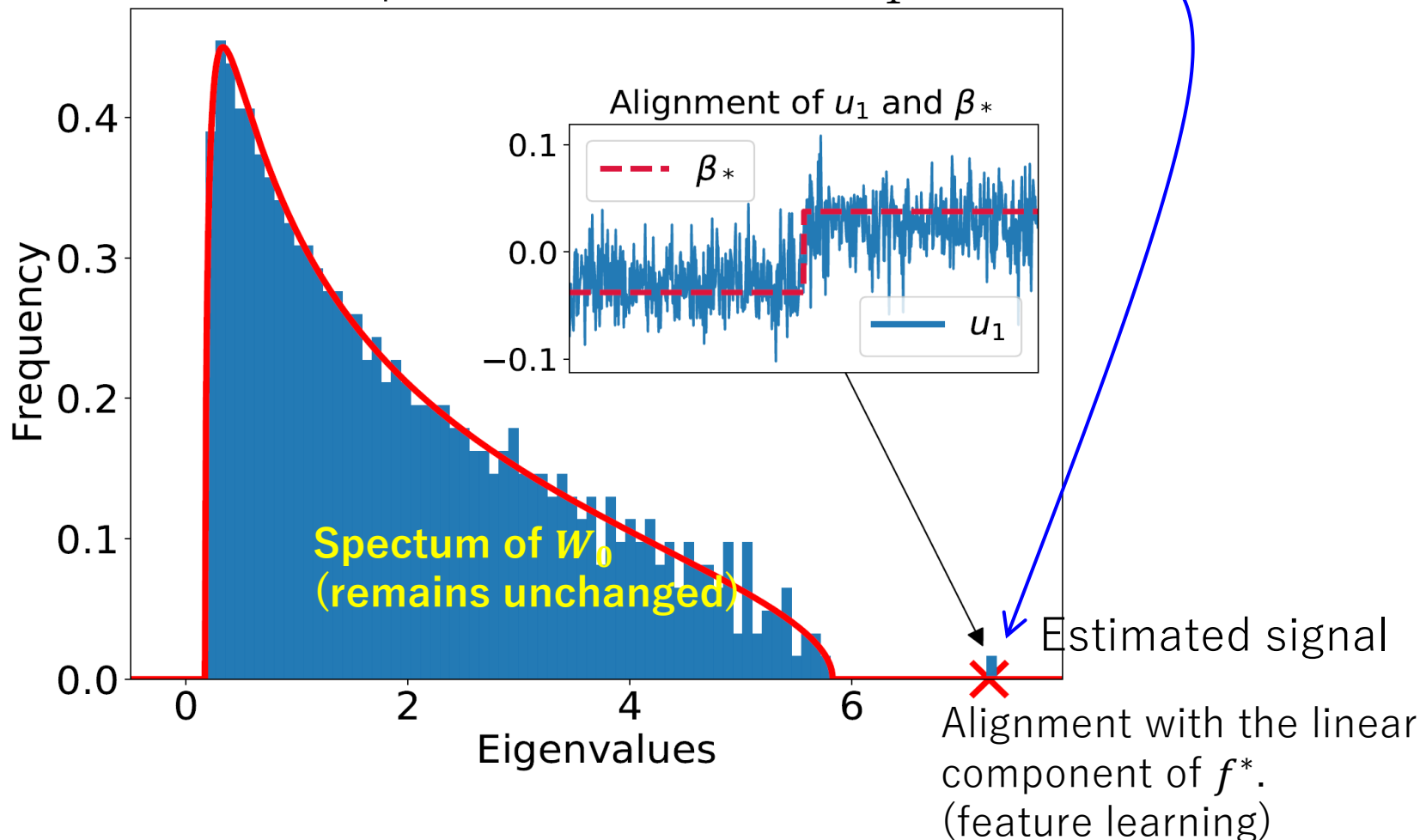
$$a^\top \phi_{\text{CK}}(x) = \frac{1}{\sqrt{N}} a^\top \sigma(W_0^\top x_i) \approx \frac{1}{\sqrt{N}} a^\top (\mu_1 W_0^\top x_i + \mu_2 z) \quad \begin{array}{l} \text{(linear function;} \\ \text{Gaussian equivalence)} \end{array}$$

# Effect of large step-size update

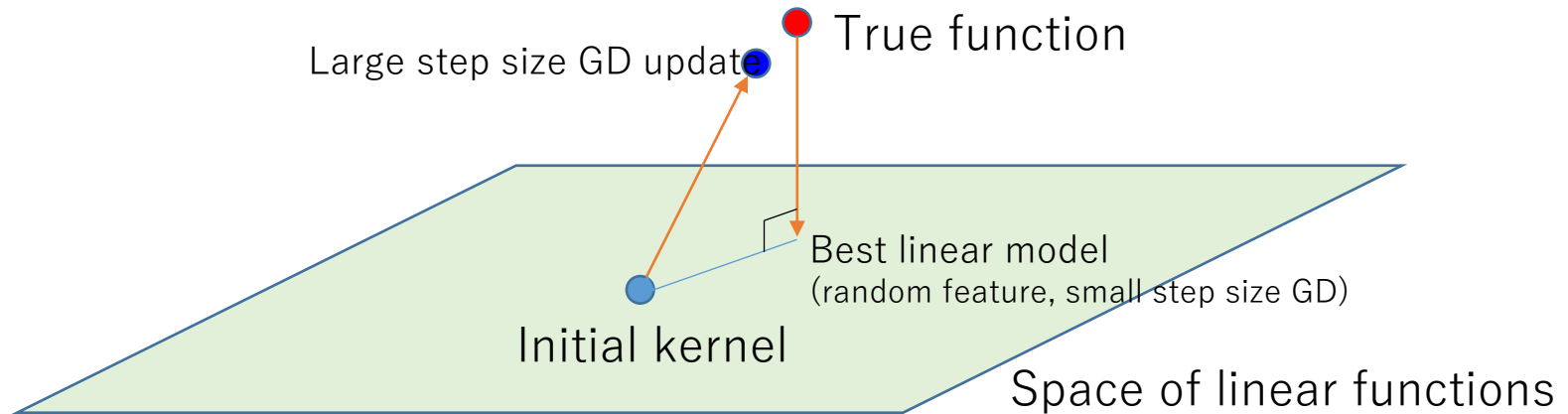
$$W_1 = W_0 + \eta\sqrt{N} \underbrace{\left(-\nabla_W \mathcal{L}(f_{\text{NN}}^{(0)})/2\right)}$$

**Theorem: almost rank 1**

Spectral distribution of  $W_1$



# Improvement over the Initial CK 149



- $\eta = \Theta(\sqrt{N})$  (large learning rate):

Known as **maximal update parameterization** ( $\mu P$ ) [Yang and Hu, 2020].

$$\tau^* = \inf_{\eta > 0} \mathbb{E}_{\xi_1 \sim N(0,1)} [\sigma^*(\xi_1) - \mathbb{E}_{\xi_2 \sim N(0,1)} [\sigma(\eta\xi_1 + \xi_2)]]$$

(measure for model misspecification)

$f^*(x) = \sigma^*(\langle x, \beta^* \rangle)$  is assumed.

- $\tau^* = 0$  if  $\sigma = \sigma^* = \text{erf}$ .
- $\tau^* \ll 1$  if  $\sigma = \sigma^* = \text{tanh}$ .

$$\mathcal{R}_{W_1}(\lambda) \leq 16\tau^* + C(\sqrt{\tau^*}\psi_1^{-1/2} + \psi_1^{-1}) + o_p(1)$$

$$n/d \rightarrow \psi_1, N/d \rightarrow \psi_2$$

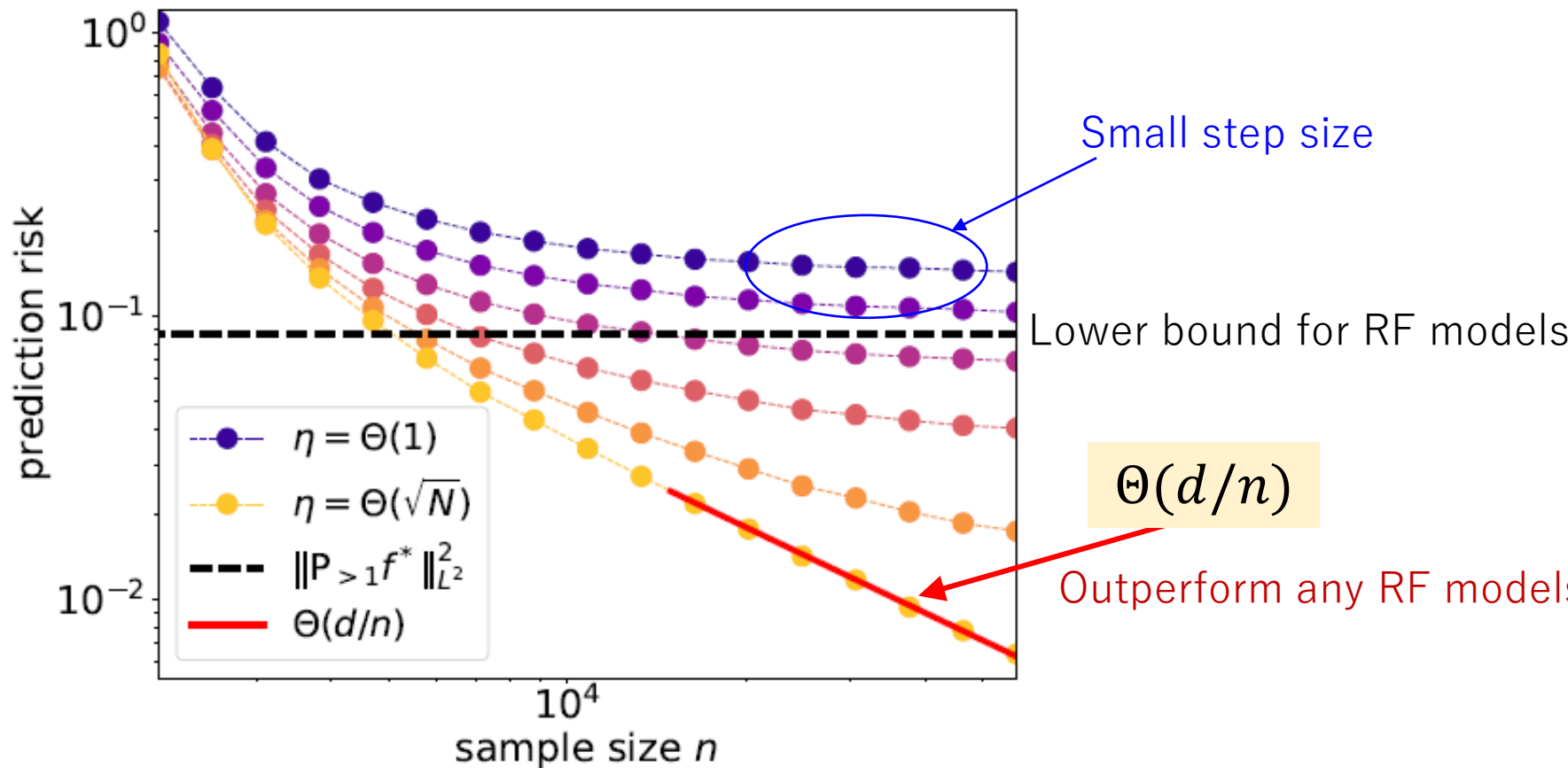
Large learning rate yields feature learning and can be better than the small step size regime if  $\tau^* \ll \|P_{>1} f^*\|^2$ .

# Implications

Corollary

If  $\sigma = \sigma^* = \text{erf}$ , then  $\tau^* = 0$ .

In particular, we have  $R_{W_1}(\lambda) = \Theta(\psi_1^{-1}) = \Theta(d/n)$ .



Predictive risk of ridge regression on CK obtained by one step GD (empirical simulation,  $d = 1024$ ): brighter color represents larger step size scaled as  $\eta = N^\alpha$  for  $\alpha \in [0, 1/2]$ . We chose  $\sigma = \sigma^* = \text{erf}$ ,  $\psi_2 = 2$ ,  $\lambda = 10^{-3}$ , and  $\sigma_\epsilon = 0.1$ .

- **Representation/Generalization ability**
  - Depth separation
  - Adaptivity of deep learning: separation between linear (shallow) and deep methods
- **Optimization ability**
  - Overparameterization
  - Noisy gradient descent: a near global optimum
    - ✓ Estimation error separation between kernel and deep learning
  - Mean field Langevin
  - CSQ lower bound

Deep learning theory that makes DL white box that can be *controllable*.

It would reveal the essence of “good learning system” which would be useful to develop methods *beyond* DL.