

Representation theory and optimization of neural networks

Taiji Suzuki

University of Tokyo / AIP-RIKEN
(Deep learning theory team)



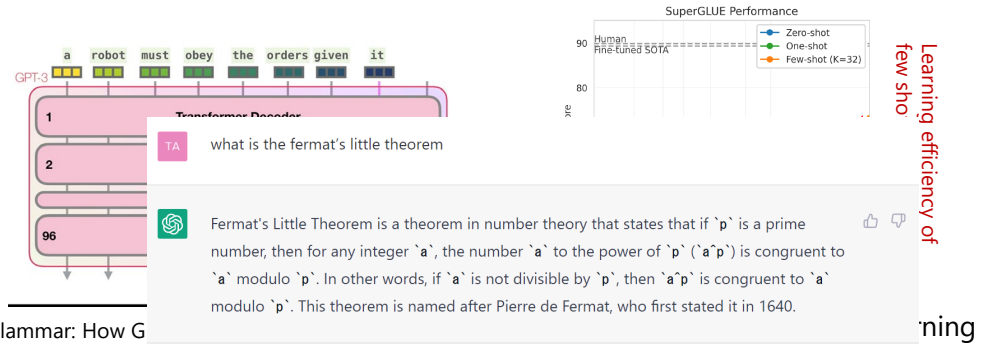
21th/Mar/2023

RIKEN-AIP & PRAIRIE Joint Workshop 2023

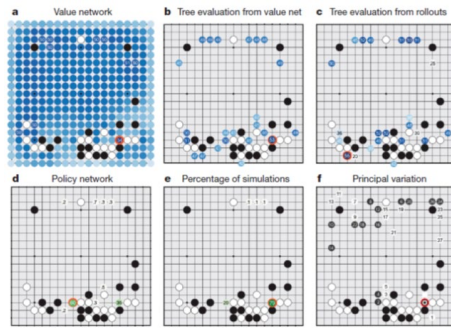
Success of deep learning

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

Large language model



AlphaGo/Zero

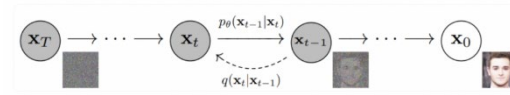


[Silver et al. (Google Deep Mind): Mastering the game of Go with deep neural networks and tree search, Nature, 529, 484—489, 2016]

[Alammar: How GPT works
<https://jalammar.com/gpt-3-animations/>]

[Brown et al. "Language Models are Few-Shot Learners" [ChatGPT. OpenAI2022]

Generative models (diffusion models)



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

Image recognition



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



Stable diffusion, 2022.



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

[Representation theory]

- Minimax optimality of diffusion model
 - Total variation distance and Wasserstein distance
 - Avoids curse of dimensionality

[Kazusato Oko, Shunta Akiyama, Taiji Suzuki: Diffusion Models are Minimax Optimal Distribution Estimators. arXiv:2303.01861, 2023]

[Optimization]

- Mean field Langevin dynamics
 - Unifying frame-work
 - (1) Time discretization, (2) Space discretization, (3) Stochastic gradient

[Taiji Suzuki, Atsushi Nitanda, Denny Wu: Convergence of mean-field Langevin dynamics: Time and space discretization, stochastic gradient, and variance reduction. 2023]

Minimax optimality of diffusion model

[Kazusato Oko, Shunta Akiyama, Taiji Suzuki: Diffusion Models are Minimax Optimal Distribution Estimators. arXiv:2303.01861, 2023]



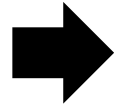
Kazusato Oko
(The University of Tokyo
/RIKEN-AIP)



Shunta Akiyama
(The University of Tokyo)

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, 1st prize in digital art category



Generated by NovelAI

Decoder : 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)

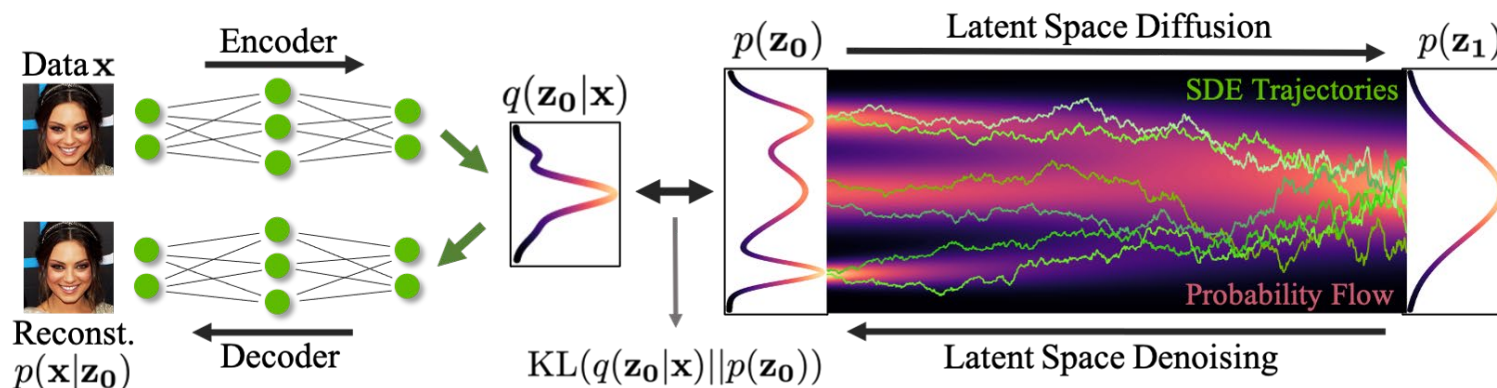
$$dX_t = -X_t dt + \sqrt{2}dB_t$$



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

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

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:

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

OU process

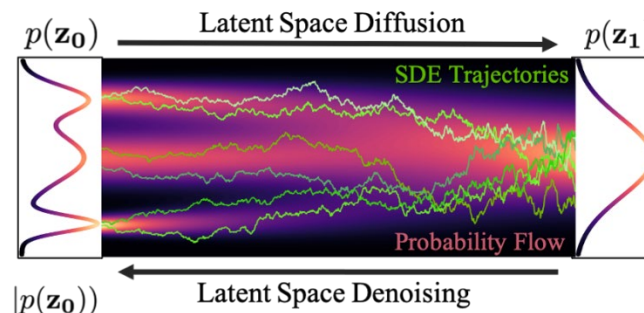
- Marginal distribution

$$p_t = \text{Law}(X_t) \longrightarrow p_t(x) = \int N_0(y) \mu_t \frac{1}{\sigma_t} \exp\left(-\frac{\|x - \mu_t y\|^2}{2\sigma_t^2}\right) dy$$

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

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

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



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

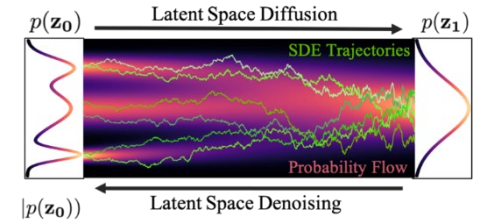
Reverse process

Reverse process:

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

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

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



[Hausmann & 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; Chen et al. (2023))

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.

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

Explicit form is available

- 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 dist bound ($n^{-1/d}$) with manifold assumption: De Bortoli (2022)

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 $[-1,1]^d \setminus [-1 + n^{-\frac{1-\delta}{d}}, 1 - n^{-\frac{1-\delta}{d}}]^d$.

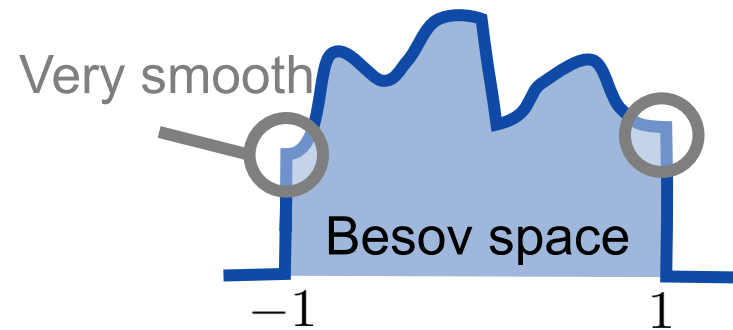
Besov space $(B_{p,q}^s(\Omega))$

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

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

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

$\omega_m(f, t)_p$

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

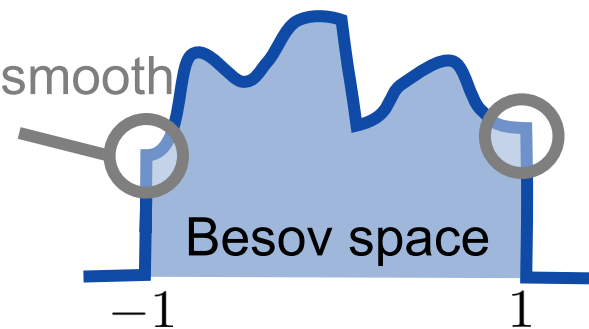
very smooth

Uniformity of smoothness

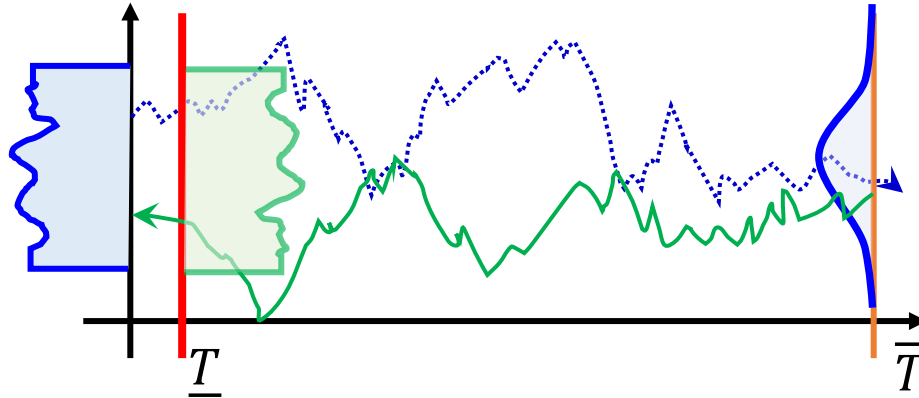
$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \left(\int_0^\infty [t^{-s} \omega_m(f, t)_p]^q \frac{dt}{t} \right)^{1/q}$$

Smoothness

Spatial inhomogeneity



Convergence rate result



Theorem (Estimation error in TV-distance)

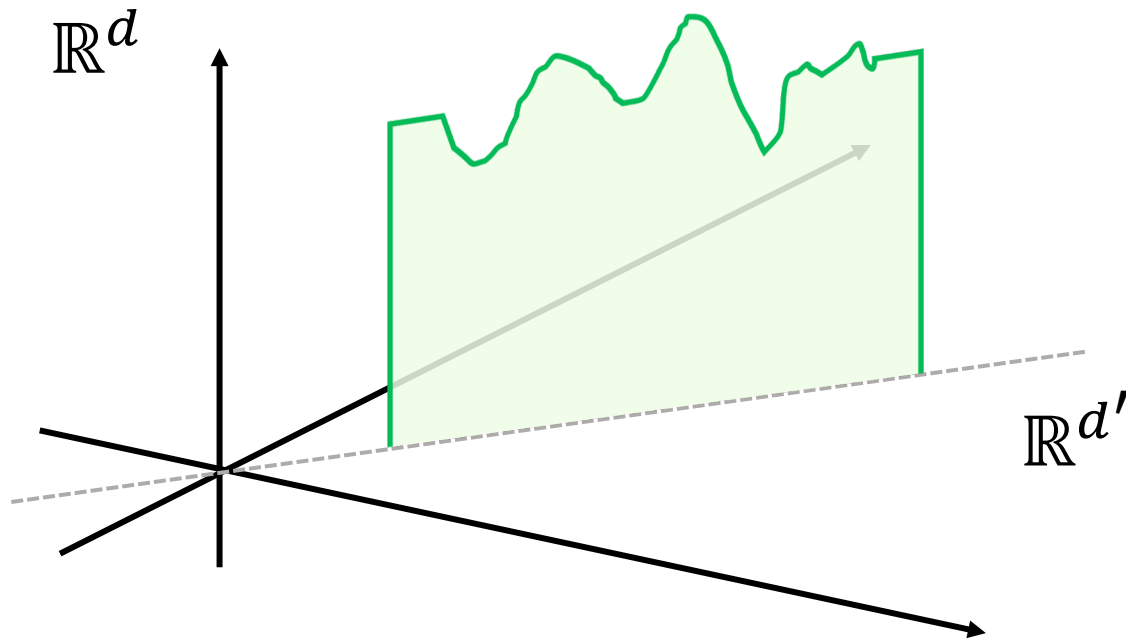
Let $\underline{T} = n^{-o(1)}$, $\bar{T} = O(\log(n))$. Then, the empirical risk minimizer \hat{s} in DNN satisfies

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

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.



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)

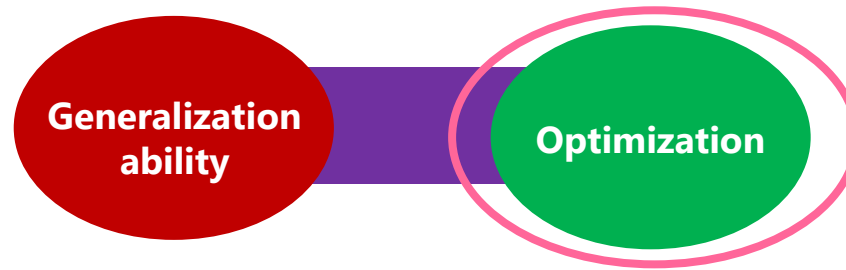
For any fixed $\delta > 0$, by slightly changing the estimator, the empirical risk minimizer \hat{s} in DNN satisfies

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

This is also known as **minimax optimal** (up to δ) [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 $t = 0$.

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



Mean field Langevin dynamics to optimize two-layer NN

[Suzuki, Nitanda, Wu: Convergence of mean-field Langevin dynamics: Time and space discretization, stochastic gradient, and variance reduction. 2023]



Atsushi Nitanda
(Kyusyu Institute of
Technology)



Denny Wu
(University of
Toronto)

Mean field Langevin dynamics:

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

convex

- Wasserstein gradient flow to minimize \mathcal{L} :

$$\partial_t \mu_t = \nabla \cdot \left[\left(\nabla \frac{\delta F(\mu_t)}{\delta \mu} + \lambda_2 \nabla \log(\mu_t) \right) \mu_t \right]$$

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

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

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

Vanilla GLD:

$$dX_t = -\nabla L(X_t) dt + \sqrt{2\lambda_2} dB_t$$

$$\mathcal{L}(\mu) = \int L(x) d\mu(x) + \lambda_2 \text{Ent}(\mu)$$

$$F(\mu) \Rightarrow \frac{\delta F}{\delta \mu} = L$$

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

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

MFLD for mean field NN

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

$$h_x(z) = r\sigma(w^\top z) \text{ for } x = (r, w)$$

Loss function:

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

$$dX_t = -\nabla_{X_t} \left(\frac{1}{n} \sum_{i=1}^n \ell'_i(f_{\mu_t}) h_{X_t}(z_i) + \lambda_1 \|X_t\|^2 \right) dt + \sqrt{2\lambda_2} dB_t$$

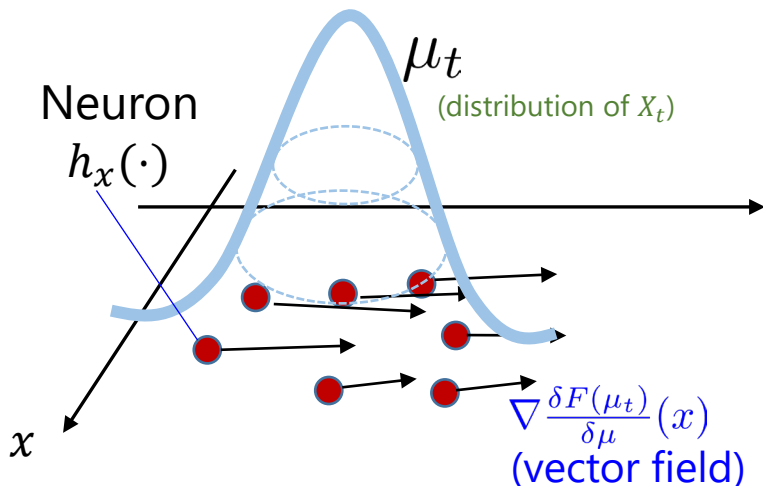
(escape from local min.)

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

$\nabla \frac{\delta F(\mu_t)}{\delta \mu}(X_t)$

[Noise perturbation]

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



Finite particle approximation:

$$d\hat{X}_t^i = -\nabla \frac{\delta F\left(\frac{1}{N} \sum_{j=1}^N \delta_{\hat{X}_t^j}\right)}{\delta \mu}(\hat{X}_t^i) dt + \sqrt{2\lambda_2} dB_t^i$$

(GLD to optimize the finite width neural network)

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**

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

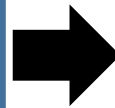
(KSD: Kernel Stein Discrepancy from a posterior distribution)

Summary of our research

Infinite particles / Continuous time

Linear convergence of mean field Langevin:

[Nitanda, Wu, Suzuki (AISTATS2022)]
[Chizat (TMLR2022)]



Finite particle / Discrete time

Double loop method:

- PDA [Nitanda, Wu, Suzuki: NeurIPS2021]
- P-SDCA [Oko, Suzuki, Wu, Nitanda: ICLR2022]
- Infinite-dim extension [Nishikawa, Suzuki, Nitanda: NeurIPS2022]

Difficult :

Propagation of chaos (McKean, Kac,..., 60's)

Finite particle / Continuous time

Uniform-in-time propagation of chaos:

- Super log-Sobolev ineq.
[Suzuki, Nitanda, Wu (ICLR2023)]
- Leave-one-out type evaluation/Uniform-log-Sobolev
[Chen, Ren, Wang (arXiv2022)]



Finite particle / Discrete time

Single loop method:

Time-space discretization,
stochastic gradient
[Suzuki, Nitanda, Wu (2023)]

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

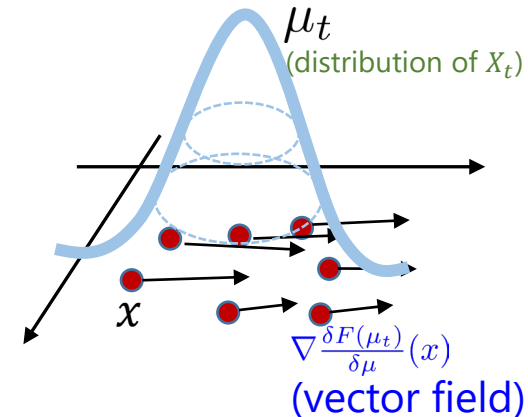
(time discretization)

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

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

(stochastic gradient)

(space discretization)

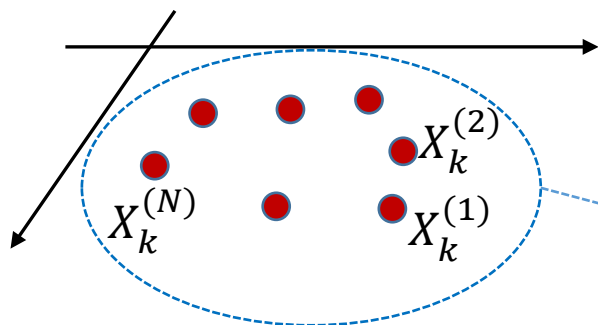


➤ Noisy gradient descent on 2-layer NN with finite width.

- **Time discretization:** $t \rightarrow k\eta$
- **Space discretization:** μ_t is approximated by N particles

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

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



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

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

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

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

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

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

For any N ,

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

(Fisher divergence)

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

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

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, Nitanda, Wu (2023)]

Suppose that p_μ satisfies log-Sobolev inequality with a constant α .

Under smoothness and boundedness of the loss function, it holds that

$$\begin{aligned} & \mathcal{L}^{(N)}(\hat{\mu}_{k+1}) - \mathcal{L}(\mu^*) \\ & \leq \exp(-\lambda_2 \eta_k / \alpha) \left(\mathcal{L}^{(N)}(\hat{\mu}_k) - \mathcal{L}(\mu^*) \right) \\ & \quad + C \left(\underbrace{\eta_k^3 + \lambda_2 \eta_k^2}_{\text{Time discr.}} + \underbrace{\frac{\eta_k}{N}}_{\text{Space discr.}} + \underbrace{\eta_k^{\frac{3}{2}} \lambda_2^{\frac{1}{2}} \max_i \text{Var}[v_k^i]}_{\text{Stochastic approx.}} \right) \end{aligned}$$

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) $\left\| \nabla \frac{\delta L(\mu)}{\delta \mu}(x) \right\| \leq R$.

Space discretization is shown through the uniform-log-Sobolev inequality shown by Chen et al. 2022.

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

- SG-MFLD

$$F(\mu) = \frac{1}{n} \sum_{j=1}^n f_j(\mu) \quad (\text{finite sum}),$$

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

(Mini-batch size = B)

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

Iteration complexity:

By setting $\eta = O\left(\frac{\lambda_2 \epsilon}{\alpha} \lambda_2^{-1} \wedge \left(\frac{\lambda_2 \epsilon}{\alpha}\right)^2 \frac{B^2(n-1)^2}{(n-B)^2 \lambda_2} \wedge \sqrt{\frac{\lambda_2 \epsilon}{\alpha}}\right)$,
the iteration complexity becomes

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

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

➤ $B = n \wedge \sqrt{\alpha/(\lambda_2 \epsilon)}$ is the optimal mini-batch size.

Deep learning theory Representation ability + Optimization

[Representation theory]

- Minimax optimality of diffusion model
 - Total variation distance and Wasserstein distance
 - Avoids curse of dimensionality

[Kazusato Oko, Shunta Akiyama, Taiji Suzuki: Diffusion Models are Minimax Optimal Distribution Estimators. arXiv:2303.01861, 2023]

[Optimization]

- Mean field Langevin dynamics
 - Unifying frame-work
 - (1) Time discretization, (2) Space discretization, (3) Stochastic gradient

[Taiji Suzuki, Atsushi Nitanda, Denny Wu: Convergence of mean-field Langevin dynamics: Time and space discretization, stochastic gradient, and variance reduction. 2023]

**We are still at a primitive stage.
Hope to have collaborations!**

Appendix

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

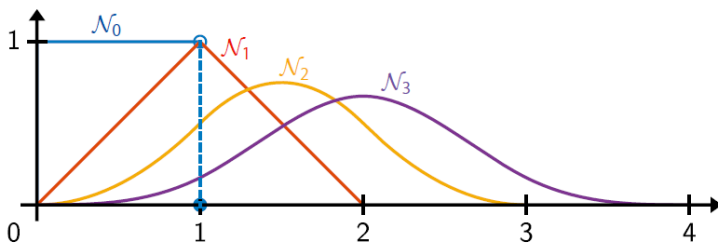
$$p_0(x) \approx \sum_{j=1}^N \alpha_j M_{a^j, b^j}^d(x)$$

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

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)²⁹

- 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

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

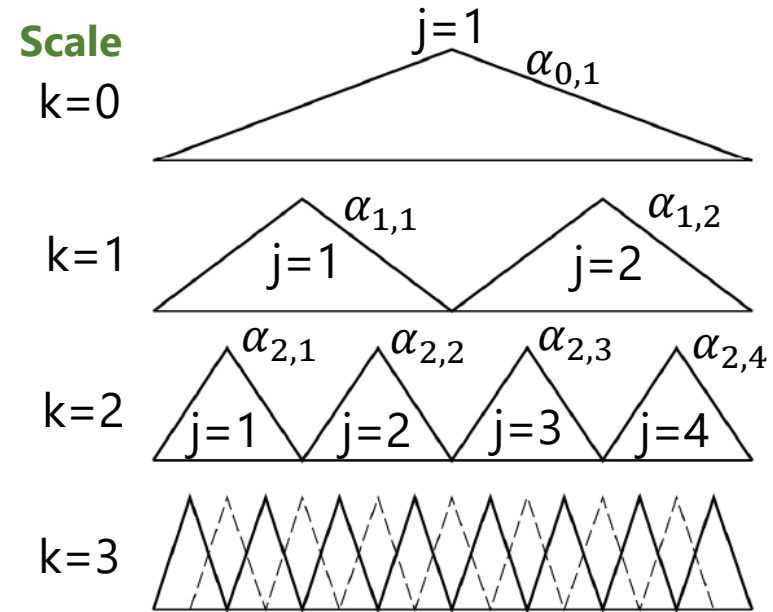
such that

(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.

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

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

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

- 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 ²¹

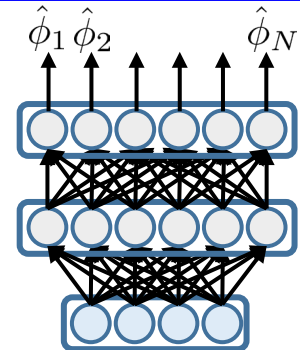
Lemma (Approximation error of diffused B-spline)

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

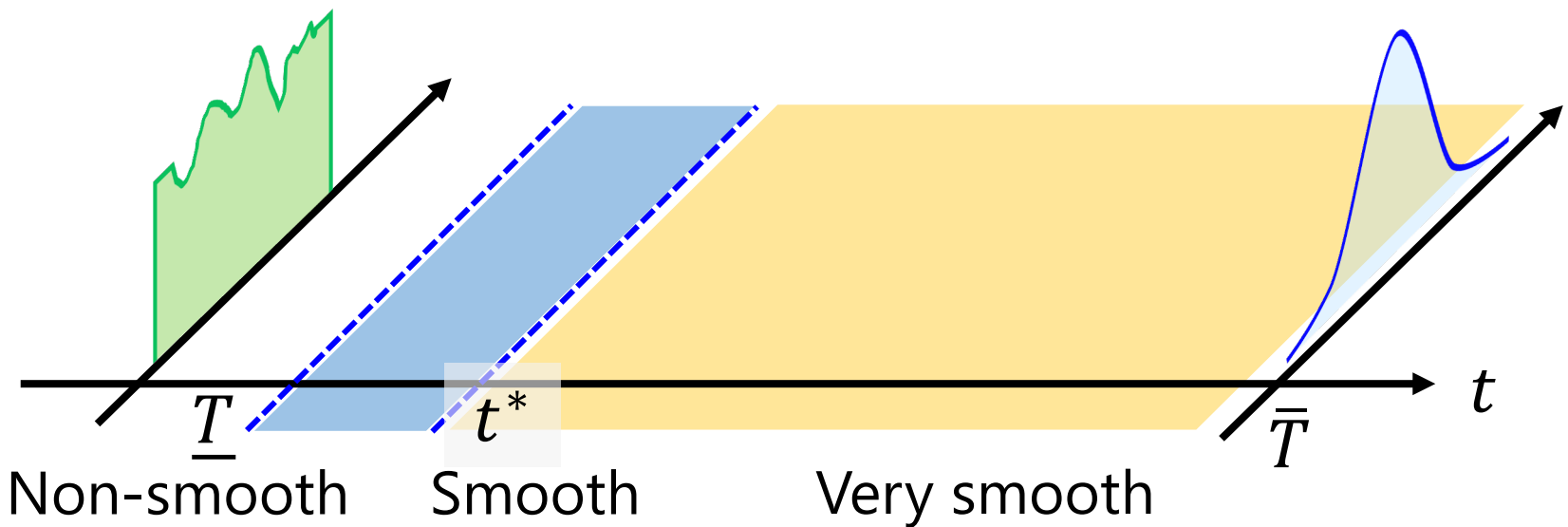
$$\left\| \hat{\phi}(x, t) - E_{a^j, b^j}(x, t) \right\|_{\infty} \leq \epsilon$$

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 ℓ^∞ -norm bound $B = O(\exp(O(\log^2(\epsilon^{-1}))))$ on parameters.

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



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



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

- 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.

➡ $\|p_t - \check{f}_{N'}\|_{L^2(X_t)}^2 \lesssim N'^{-2k/d} t_*^{-k}$ (take $k = s + 1$)

Error decomposition

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

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

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

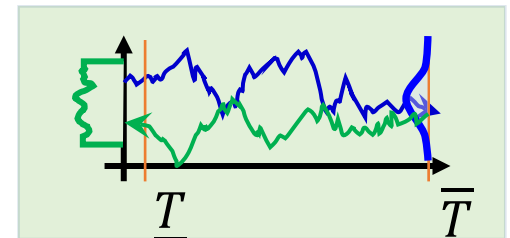
Variance

$$\frac{N \text{polylog}(N)}{n}$$

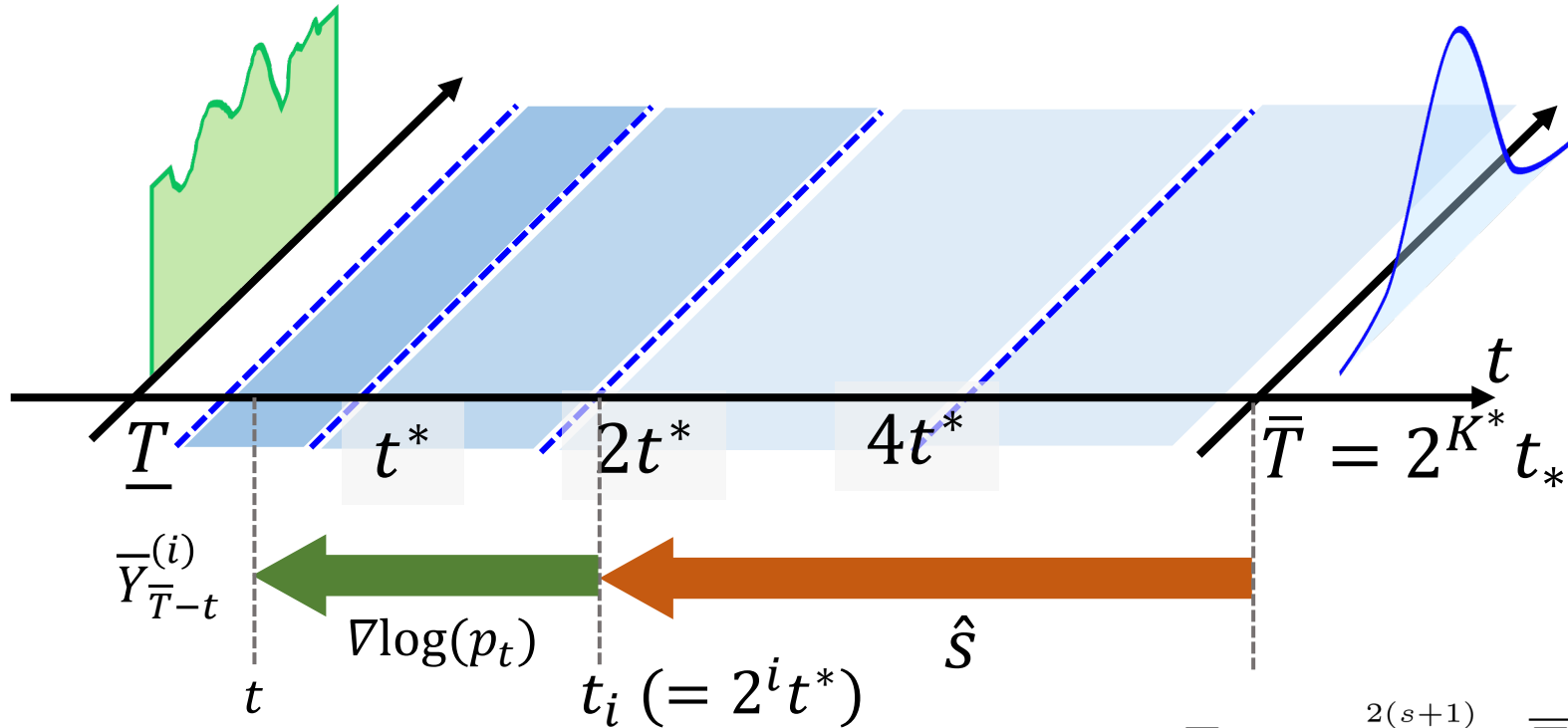
$$\lesssim \left(N^{-2s/d} + \frac{N}{n} \right) \text{polylog}(N)$$

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

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



Bound for W1 distance



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

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

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

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

Implementable discretization

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

Finite sample approximation

$$\min_{s \in \text{DNN}} \frac{1}{M} \sum_{j=1}^M \|s(x_{t_j, j}, t_j) - \nabla \log p_{t_j}(x_{t_j, j} | x_{i_j})\|^2$$

- $i_j \sim \text{Unif}(\{1, \dots, n\})$
- $t_j \sim \text{Unif}([\underline{T}, \bar{T}])$
- $x_{t_j, j} \sim p_{t_j}(\cdot | x_{i_j})$

Prop

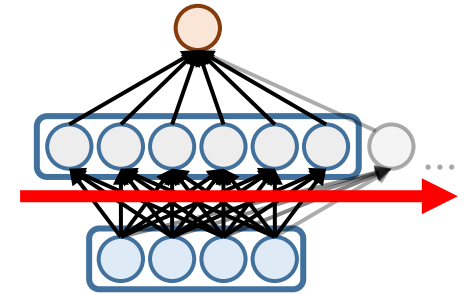
$$M \gtrsim n \cdot \underline{T}^{-1} = n^{1 + \frac{2(s+1)}{2s+d}}$$

is sufficient to attain the same convergence rate.

Mean field limit of 2-layer NN

- 2-layer neural network:

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

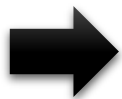


Non-linear with respect to parameters $(r_j, w_j)_{j=1}^M$.

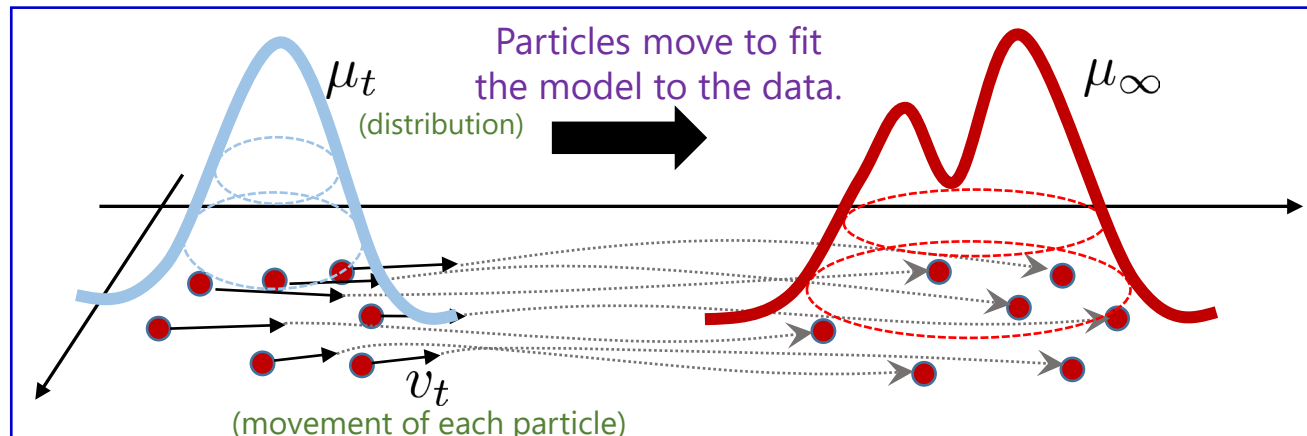
- Overparameterization (Mean field limit):

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

Linear with respect to the prob. measure μ .



Optimization of $f \Leftrightarrow$ Optimization of μ



GLD as a Wasserstein gradient flow²⁷

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

μ_t : Distribution of X_t (we can assume it has a density)

PDE that describes μ_t 's dynamics [Fokker-Planck equation]:

$$\partial_t \mu_t = \nabla \cdot \left[\mu_t \left(\nabla L + \frac{1}{\beta} \nabla \log(\mu_t) \right) \right]$$

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

$$\mu^* = \arg \min_{\mu \in \mathcal{P}} \int L(x) d\mu(x) + \frac{1}{\beta} \text{Ent}(\mu) =: \mathcal{L}(\mu)$$

[linear w.r.t. μ]

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

➔ $\mu^*(x) \propto \exp(-\beta L(x))$: Stationary distribution

Difficulty

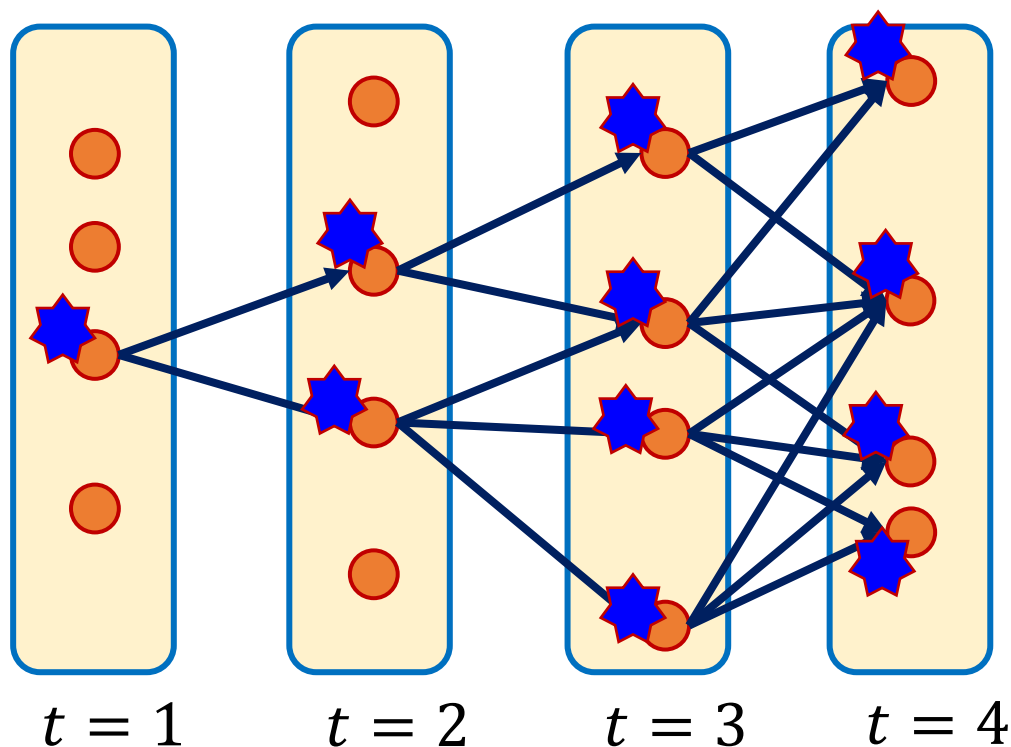
- SDE of interacting particles (McKean, Kac,..., 60')

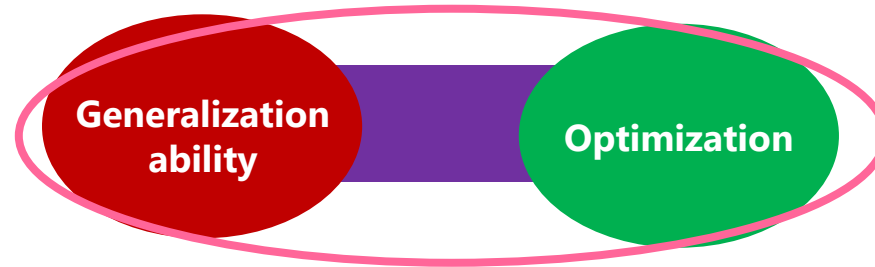
Propagation of chaos [Sznitman, 1991; Lacker, 2021]:

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

Finite particle approximation error can propagate through time.

→ It is difficult to bound the perturbation uniformly over time.





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

Problem setting

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 $O(1/N)$

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

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

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 guarantee⁴¹

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

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

We consider the **proportional limit** ($n, d, N \rightarrow \infty$ with $n/d \rightarrow \psi_1, N/d \rightarrow \psi_2$).

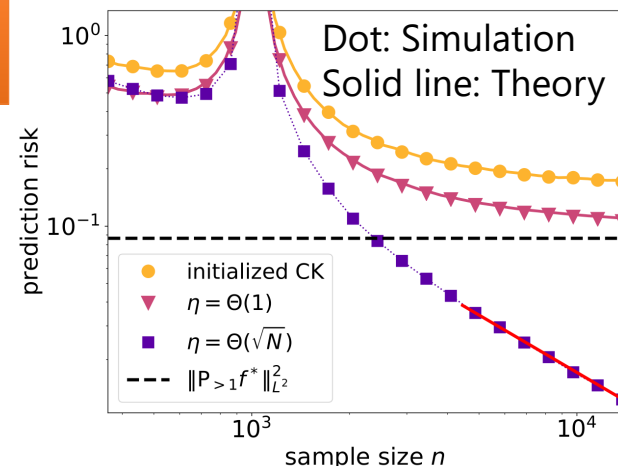
It allows to derive precise risk.

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:

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

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

- NTK (Neural tangent kernel):

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

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

Trained feature:

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

Rank 1 property of first gradient step 43

Reference

- The gradient G_t can be approximated by rank one matrix.
 \Rightarrow There appears "spike" in the spectral distribution of W_1 .

$$G_t = -\frac{1}{n} X^\top \left[\left(\frac{1}{\sqrt{N}} \left(\frac{1}{\sqrt{N}} \sigma(XW_t) a - y \right) a^\top \right) \odot \sigma'(XW_t) \right]$$

(generally, this is not low rank due to the nonlinearity of σ')

Theorem (Rank one approximation of gradient)

Remember that $G_0 = \frac{1}{\eta\sqrt{N}} (W_1 - W_0)$ ($\because W_1 = W_0 + \eta\sqrt{N}G_0$)

Let $\mu_1 = \mathbb{E}[z\sigma(z)]$, $\mu_2 = \sqrt{\mathbb{E}[\sigma(z)^2] - \mu_1^2}$, where $z \sim \mathcal{N}(0, 1)$.

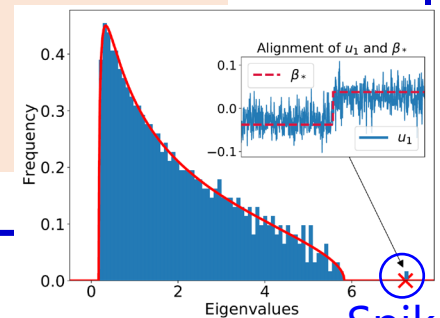
Define $A := \frac{\mu_1}{n\sqrt{N}} X^\top y a^\top$ (rank one matrix), then we have

$$\|G_0 - A\| \lesssim \frac{\log^2(n)}{\sqrt{n}} \cdot \|G_0\|$$

with high probability for sufficiently large n, d, N .

$W_1 = W_0 + \eta \times$ (rank one matrix).

\Rightarrow For large step size η , spike becomes more dominant.



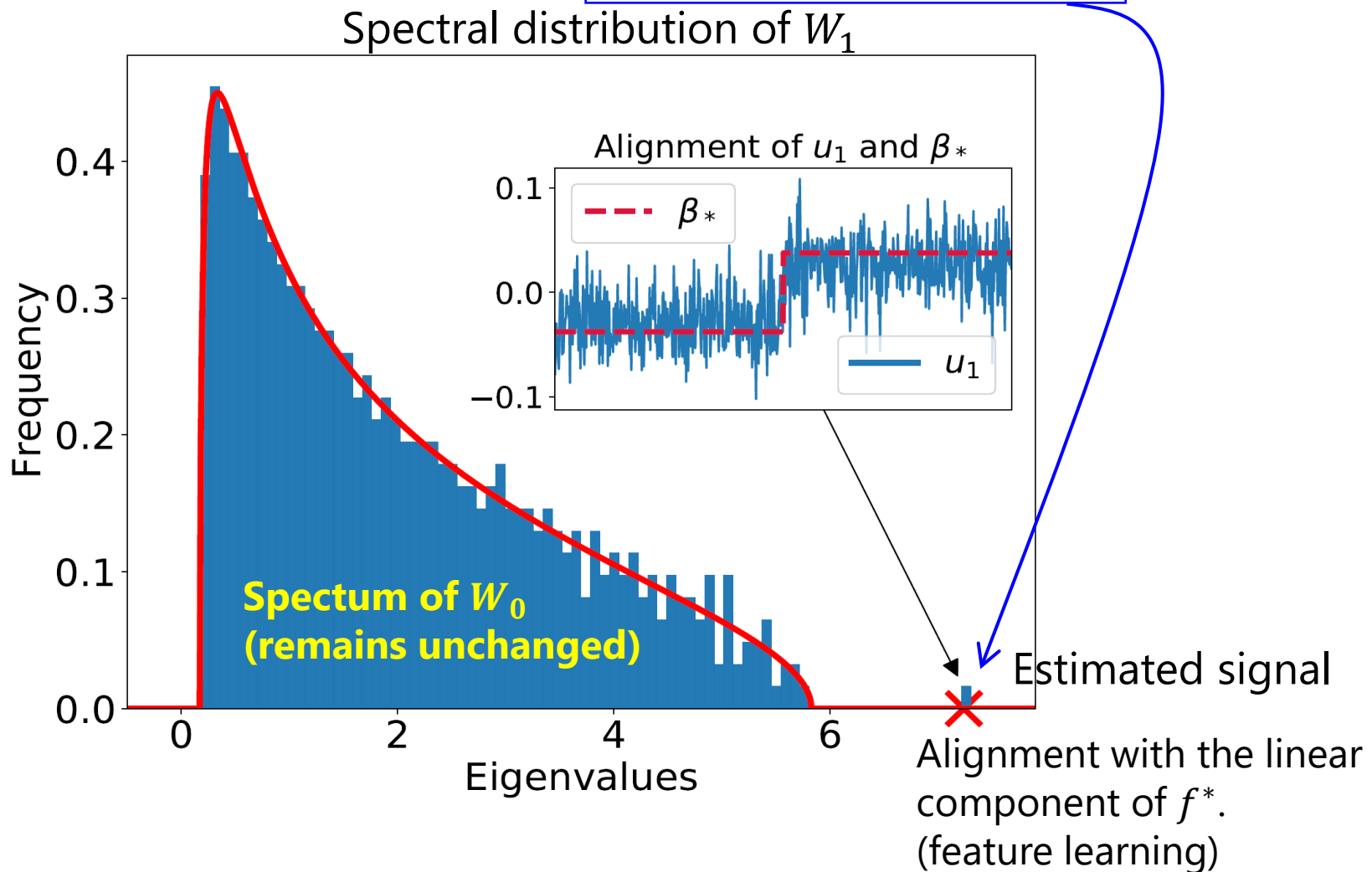
Spike

Effect of large step-size update

Reference

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

Theorem: almost rank 1



- (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), ...] $\mathcal{R}_{XX}(f) = \mathbb{E}[(f^*(X) - \hat{f}_{XX}(X))^2]$

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

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

Nonlinear part cannot be trained!

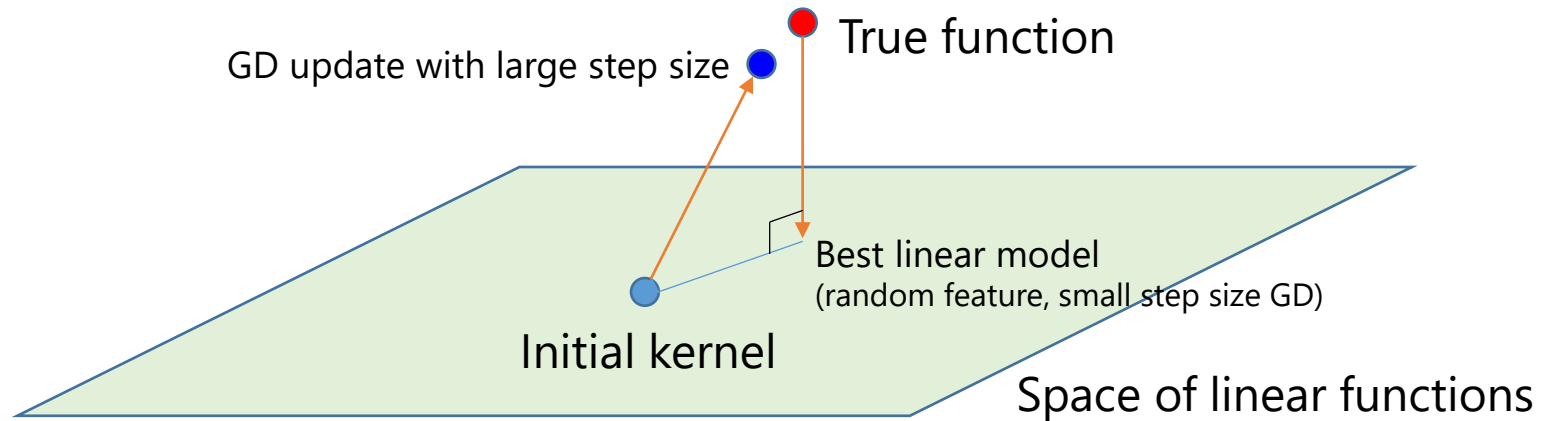
$$P_{>1}f^* := (I - P_{\leq 1})f^*$$

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

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



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

Known as maximal update parameterization (μP) [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 $\sigma = \sigma^* = \text{tanh}$.

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

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

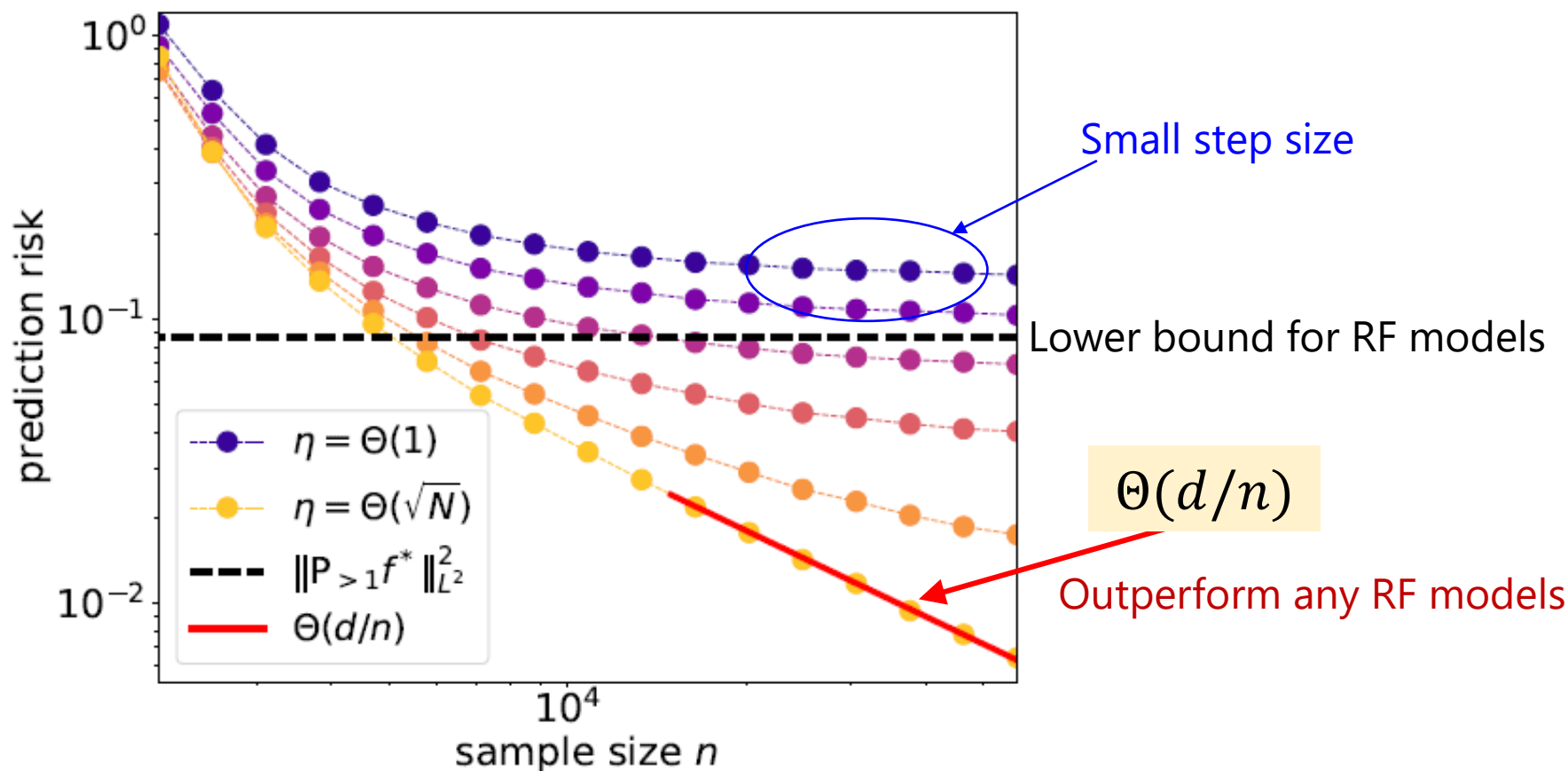
Large learning rate yields feature learning and can be better than the small step size regime if $\tau^* \ll \|P_{>1} f^*\|^2$.

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 $\sigma_\epsilon = 0.1$.