

A Mean-Field Theory of Lazy Training
in Two-Layer Neural Nets:
Entropic Regularization and Controlled
McKean–Vlasov Dynamics

Maxim Raginsky

(joint work with **Belinda Tzen**)

University of Illinois at Urbana-Champaign

April 2020

Motivation

Infinitely wide neural nets

- ▶ mean-field view: abstract away finite-size effects, focus on universal phenomena
- ▶ training dynamics: PDE on the space of probability measures (distributional or nonlinear dynamics)
- ▶ finite-neuron setting can be obtained by sampling/discretization

Effects of massive overparametrization and hyperparameter choices

- ▶ *lazy training* (Chizat–Oyallon–Bach, 2019): weights barely move from random Gaussian initialization, while risk decreases
- ▶ *Neural Tangent Kernel* (Jacot–Gabriel–Hongler, 2018): during training, neural net stays close to its linearization around initialization

In this talk: an alternative mean-field theory of lazy training via *entropic regularization*.

Related Work

Related Work

- ▶ Mean-field limit:

Nitanda–Suzuki, 2017; Chizat and Bach, 2018;

Mei–Montanari–Nguyen, 2018; Rotskoff and Vanden–Eijnden, 2018;

Mei–Misiakiewicz–Montanari, 2019; Sirignano and Spiliopoulos, 2020

Related Work

- ▶ Mean-field limit:
Nitanda–Suzuki, 2017; Chizat and Bach, 2018;
Mei–Montanari–Nguyen, 2018; Rotskoff and Vanden–Eijnden, 2018;
Mei–Misiakiewicz–Montanari, 2019; Sirignano and Spiliopoulos, 2020
- ▶ Lazy training and the Neural Tangent Kernel:
Jacot–Gabriel–Hongler, 2018; Allen–Zhu–Li–Liang, 2019;
Chizat–Oyallon–Bach, 2019

Related Work

- ▶ Mean-field limit:
Nitanda–Suzuki, 2017; Chizat and Bach, 2018;
Mei–Montanari–Nguyen, 2018; Rotskoff and Vanden–Eijnden, 2018;
Mei–Misiakiewicz–Montanari, 2019; Sirignano and Spiliopoulos, 2020
- ▶ Lazy training and the Neural Tangent Kernel:
Jacot–Gabriel–Hongler, 2018; Allen–Zhu–Li–Liang, 2019;
Chizat–Oyallon–Bach, 2019
- ▶ Universal approximation by transporting weights:
Ji–Telgarsky–Xian, 2020

Function Approximation by Two-Layer Neural Nets

- ▶ Two-layer net with N hidden units:

$$\hat{f}_N(x; \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; w^i), \quad w^i \in \mathbb{R}^d$$

Function Approximation by Two-Layer Neural Nets

- ▶ Two-layer net with N hidden units:

$$\hat{f}_N(x; \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; w^i), \quad w^i \in \mathbb{R}^d$$

- ▶ Approximation risk for a target function $f : \mathcal{X} \rightarrow \mathbb{R}$:

$$R_N(\mathbf{w}) := \|f - \hat{f}_N(\cdot; \mathbf{w})\|_{L^2(\pi)}^2 = \int_{\mathcal{X}} \pi(dx) (f(x) - \hat{f}_N(x; \mathbf{w}))^2$$

Function Approximation by Two-Layer Neural Nets

- ▶ Two-layer net with N hidden units:

$$\hat{f}_N(x; \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; w^i), \quad w^i \in \mathbb{R}^d$$

- ▶ Approximation risk for a target function $f : \mathcal{X} \rightarrow \mathbb{R}$:

$$R_N(\mathbf{w}) := \|f - \hat{f}_N(\cdot; \mathbf{w})\|_{L^2(\pi)}^2 = \int_{\mathcal{X}} \pi(dx) (f(x) - \hat{f}_N(x; \mathbf{w}))^2$$

- ▶ The risk depends only on the empirical distribution of \mathbf{w} :

$$R_N(\mathbf{w}) = R_0 + \frac{2}{N} \sum_{i=1}^N \tilde{f}(w^i) + \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N K(w^i, w^j)$$

where $R_0 := \mathbf{E}_{\pi}[f^2(X)]$,

$$\tilde{f}(w) := -\mathbf{E}_{\pi}[f(X)\sigma(X; w)],$$

$$K(w, \tilde{w}) := \mathbf{E}_{\pi}[\sigma(X; w)\sigma(X; \tilde{w})]$$

Mean-Field Limit: A Continuum of Neurons

- ▶ From finite population of neurons to a continual ensemble:

$$\hat{\mu}_{\mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \delta_{w^i} \quad \longrightarrow \quad \mu \in \mathcal{P}(\mathbb{R}^d)$$

$$\hat{f}_N(x; \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; w^i) \quad \longrightarrow \quad \hat{f}(x; \mu) = \int_{\mathbb{R}^d} \mu(dw) \sigma(x; w)$$

Mean-Field Limit: A Continuum of Neurons

- ▶ From finite population of neurons to a continual ensemble:

$$\hat{\mu}_{\mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \delta_{w^i} \longrightarrow \mu \in \mathcal{P}(\mathbb{R}^d)$$

$$\hat{f}_N(x; \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; w^i) \longrightarrow \hat{f}(x; \mu) = \int_{\mathbb{R}^d} \mu(dw) \sigma(x; w)$$

- ▶ L^2 risk:

$$R(\mu) = \|f - \hat{f}(\cdot; \mu)\|_{L^2(\pi)}^2 = R_0 + 2 \int \tilde{f} d\mu + \int K d(\mu \otimes \mu)$$

— *convex quadratic* in the measure μ

Mean-Field Limit: A Continuum of Neurons

- ▶ From finite population of neurons to a continual ensemble:

$$\hat{\mu}_{\mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \delta_{w^i} \quad \longrightarrow \quad \mu \in \mathcal{P}(\mathbb{R}^d)$$

$$\hat{f}_N(x; \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; w^i) \quad \longrightarrow \quad \hat{f}(x; \mu) = \int_{\mathbb{R}^d} \mu(dw) \sigma(x; w)$$

- ▶ L^2 risk:

$$R(\mu) = \|f - \hat{f}(\cdot; \mu)\|_{L^2(\pi)}^2 = R_0 + 2 \int \tilde{f} d\mu + \int K d(\mu \otimes \mu)$$

— *convex quadratic* in the measure μ

- ▶ Back to finite N : $R_N(\mathbf{w}) = R(\hat{\mu}_{\mathbf{w}})$

Entropy-Regularized Risk in the Mean-Field Limit

Free energy:

$$F_{\beta, \tau}(\mu) := \frac{1}{2}R(\mu) + \frac{\tau}{\beta}D_{\text{KL}}(\mu \parallel \gamma_{\tau})$$

where $\gamma_{\tau} = \mathcal{N}(0, \tau I_d)$ is a Gaussian “prior” on the weights

Entropy-Regularized Risk in the Mean-Field Limit

Free energy:

$$F_{\beta, \tau}(\mu) := \frac{1}{2}R(\mu) + \frac{\tau}{\beta}D_{\text{KL}}(\mu \parallel \gamma_{\tau})$$

where $\gamma_{\tau} = \mathcal{N}(0, \tau I_d)$ is a Gaussian “prior” on the weights

Some motivation:

Entropy-Regularized Risk in the Mean-Field Limit

Free energy:

$$F_{\beta, \tau}(\mu) := \frac{1}{2}R(\mu) + \frac{\tau}{\beta}D_{\text{KL}}(\mu \parallel \gamma_{\tau})$$

where $\gamma_{\tau} = \mathcal{N}(0, \tau I_d)$ is a Gaussian “prior” on the weights

Some motivation:

► If $\mu = \mathcal{N}(\theta, \tau I_d)$, then

$$F_{\beta, \tau}(\mu) = \frac{1}{2}R(\mu) + \frac{1}{2\beta}\|\theta\|_2^2$$

— a form of ridge regression

Entropy-Regularized Risk in the Mean-Field Limit

Free energy:

$$F_{\beta, \tau}(\mu) := \frac{1}{2}R(\mu) + \frac{\tau}{\beta}D_{\text{KL}}(\mu \parallel \gamma_{\tau})$$

where $\gamma_{\tau} = \mathcal{N}(0, \tau I_d)$ is a Gaussian “prior” on the weights

Some motivation:

- ▶ If $\mu = \mathcal{N}(\theta, \tau I_d)$, then

$$F_{\beta, \tau}(\mu) = \frac{1}{2}R(\mu) + \frac{1}{2\beta}\|\theta\|_2^2$$

— a form of ridge regression

- ▶ Talagrand’s entropy-transport inequality

$$D_{\text{KL}}(\mu \parallel \gamma_{\tau}) \geq \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \gamma_{\tau}),$$

where $\mathcal{W}_2(\mu, \nu) := \min_{X \sim \mu, Y \sim \nu} \sqrt{\mathbf{E}[\|X - Y\|_2^2]}$, gives

$$F_{\beta, \tau}(\mu) \geq \frac{1}{2}R(\mu) + \frac{1}{2\beta} \mathcal{W}_2^2(\mu, \gamma_{\tau})$$

Entropy-Regularized Risk: Static Formulation

Assumptions:

(i) f, σ bounded; (ii) $\nabla_w \sigma(X; w)$ subgaussian ($X \sim \pi$); (iii) \tilde{f}, K Lipschitz with Lipschitz gradients

Theorem (Tzen–Raginsky, 2020) The free energy $F_{\beta, \tau}(\cdot)$ admits a *unique* minimizer $\mu^* = \mu_{\beta, \tau}^*$, such that:

1. The Boltzmann fixed-point condition holds:

$$\frac{d\mu^*}{d\gamma_\tau}(w) \propto \exp\left(-\frac{\beta}{\tau}\Psi(w; \mu^*)\right)$$

where $\Psi(w; \mu) := \tilde{f}(w) + \int K(w, \tilde{w})\mu(d\tilde{w})$

2. The optimal weight distribution is nearly Gaussian:

$$D_{\text{KL}}(\mu^* \parallel \gamma_\tau) \leq \frac{\kappa\beta^2}{\tau}, \quad \mathcal{W}_2^2(\mu^*, \gamma_\tau) \leq \kappa\beta^2$$

3. If $f = \hat{f}(\cdot; \mu^\circ)$, then $R(\mu^*) \leq \frac{2\tau}{\beta} D_{\text{KL}}(\mu^\circ \parallel \gamma_\tau)$.

Some Consequences

- ▶ Weak regularization: $\beta \rightarrow \infty$ while $\tau = o(\beta)$ gives

$$\liminf_{\beta \uparrow \infty} \inf_{\mu} F_{\beta, \tau}(\mu) = \frac{1}{2} \inf_{\mu} R(\mu)$$

Some Consequences

- ▶ Weak regularization: $\beta \rightarrow \infty$ while $\tau = o(\beta)$ gives

$$\liminf_{\beta \uparrow \infty} \inf_{\mu} F_{\beta, \tau}(\mu) = \frac{1}{2} \inf_{\mu} R(\mu)$$

- ▶ Strong regularization: both β and τ are small but $\tau \ll \beta$:

$$\tau = \varepsilon^2/d, \quad \beta = \sqrt{\tau d} = \varepsilon$$

$$D_{\text{KL}}(\mu^* \|\gamma_{\tau}) \leq \kappa d, \quad \mathcal{W}_2^2(\mu^*, \gamma_{\tau}) \leq \kappa \varepsilon$$

Some Consequences

- ▶ Weak regularization: $\beta \rightarrow \infty$ while $\tau = o(\beta)$ gives

$$\liminf_{\beta \uparrow \infty} \inf_{\mu} F_{\beta, \tau}(\mu) = \frac{1}{2} \inf_{\mu} R(\mu)$$

- ▶ Strong regularization: both β and τ are small but $\tau \ll \beta$:

$$\tau = \varepsilon^2/d, \quad \beta = \sqrt{\tau d} = \varepsilon$$

$$D_{\text{KL}}(\mu^* \parallel \gamma_{\tau}) \leq \kappa d, \quad \mathcal{W}_2^2(\mu^*, \gamma_{\tau}) \leq \kappa \varepsilon$$

- ▶ Realizable case: if $f = \hat{f}(\cdot; \mu^{\circ})$ with $D_{\text{KL}}(\mu^{\circ} \parallel \gamma_{\tau}) = O(d)$, then

$$R(\mu^*) \leq \frac{2\tau}{\beta} D_{\text{KL}}(\mu^{\circ} \parallel \gamma_{\tau}) = \frac{2\kappa\varepsilon^2/d}{\varepsilon} \cdot O(d) = O(\varepsilon)$$

Lazy Training?

- ▶ Consider random Gaussian initialization:

$$\mathbf{W} = (W^1, \dots, W^N) \sim (\gamma_\tau)^{\otimes N} \quad \longrightarrow \quad \hat{f}_N(x; \mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; W^i)$$

Lazy Training?

- ▶ Consider random Gaussian initialization:

$$\mathbf{W} = (W^1, \dots, W^N) \sim (\gamma_\tau)^{\otimes N} \quad \longrightarrow \quad \hat{f}_N(x; \mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; W^i)$$

- ▶ With high probability, there exist good approximations to “low-complexity” f near $\hat{f}_N(\cdot; \mathbf{W})$ (Allen-Zhu-Li-Liang, 2019)

Lazy Training?

- ▶ Consider random Gaussian initialization:

$$\mathbf{W} = (W^1, \dots, W^N) \sim (\gamma_\tau)^{\otimes N} \quad \longrightarrow \quad \hat{f}_N(x; \mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; W^i)$$

- ▶ With high probability, there exist good approximations to “low-complexity” f near $\hat{f}_N(\cdot; \mathbf{W})$ (Allen-Zhu-Li-Liang, 2019)
- ▶ These good approximations can be found by gradient descent, and the weight updates do not drift too far from \mathbf{W}

Lazy Training?

- ▶ Consider random Gaussian initialization:

$$\mathbf{W} = (W^1, \dots, W^N) \sim (\gamma_\tau)^{\otimes N} \quad \longrightarrow \quad \hat{f}_N(x; \mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \sigma(x; W^i)$$

- ▶ With high probability, there exist good approximations to “low-complexity” f near $\hat{f}_N(\cdot; \mathbf{W})$ (Allen-Zhu-Li-Liang, 2019)
- ▶ These good approximations can be found by gradient descent, and the weight updates do not drift too far from \mathbf{W}
- ▶ Transport map interpretation: $W^i \mapsto T(W^i)$ for some $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ (Ji-Telgarsky-Xian, 2020):

$$\hat{f}_N(\cdot; \mathbf{W}) \longmapsto \hat{f}_N(\cdot; T(\mathbf{W}))$$

such that $\|T(W^i) - W^i\|_2$ is small

Corollary (Tzen–Raginsky, 2020) For β sufficiently small, there exists a Lipschitz transport map $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$, such that all of the following holds with probability at least $1 - \delta$ for $\mathbf{W} = (W^1, \dots, W^N) \sim (\gamma_\tau)^{\otimes N}$:

1. The transported weights $T(W^i)$ are uniformly close to the i.i.d. Gaussian weights W^i :

$$\max_{i \in [N]} \|T(W^i) - W^i\|_2 \leq \kappa\beta + \kappa\sqrt{\tau(\log N + \log(1/\delta))};$$

2. The transported neural net $\hat{f}_N(\cdot; T(\mathbf{W}))$ satisfies

$$\|f - \hat{f}_N(\cdot; T(\mathbf{W}))\|_{L^2(\pi)} \leq \|f - \hat{f}(\cdot; \mu_{\beta, \tau}^*)\|_{L^2(\pi)} + \kappa\sqrt{\frac{\log(1/\delta)}{N}};$$

3. The transported neural net $\hat{f}_N(\cdot; T(\mathbf{W}))$ and the random Gaussian neural net $\hat{f}_N(\cdot; \mathbf{W})$ are close in $L^2(\pi)$ norm:

$$\left\| \hat{f}_N(\cdot; T(\mathbf{W})) - \hat{f}_N(\cdot; \mathbf{W}) \right\|_{L^2(\pi)}^2 \leq \kappa\beta + \kappa\sqrt{\frac{\tau \log(1/\delta)}{N}}.$$

The Proof: Some Intuition

- ▶ Optimal μ^* (Boltzmann fixed-point condition):

$$\mu^*(dw) \propto \exp\left(-\frac{1}{2\tau}\|w\|_2^2 - \frac{\beta}{\tau}\Psi(w; \mu^*)\right) dw$$

The Proof: Some Intuition

- ▶ Optimal μ^* (Boltzmann fixed-point condition):

$$\mu^*(dw) \propto \exp\left(-\frac{1}{2\tau}\|w\|_2^2 - \frac{\beta}{\tau}\Psi(w; \mu^*)\right) dw$$

- ▶ Since $\nabla\Psi$ is Lipschitz, $w \mapsto \frac{1}{2}\|w\|_2^2 + \beta\Psi(w; \mu^*)$ is **strongly convex** for all $\beta < \beta_0$; thus, the optimal (Brenier–McCann) transport map $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that achieves $\mathcal{W}_2(\mu^*, \gamma_\tau)$ is Lipschitz by **Caffarelli's regularity theorem**.

The Proof: Some Intuition

- ▶ Optimal μ^* (Boltzmann fixed-point condition):

$$\mu^*(dw) \propto \exp\left(-\frac{1}{2\tau}\|w\|_2^2 - \frac{\beta}{\tau}\Psi(w; \mu^*)\right) dw$$

- ▶ Since $\nabla\Psi$ is Lipschitz, $w \mapsto \frac{1}{2}\|w\|_2^2 + \beta\Psi(w; \mu^*)$ is **strongly convex** for all $\beta < \beta_0$; thus, the optimal (Brenier–McCann) transport map $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that achieves $\mathcal{W}_2(\mu^*, \gamma_\tau)$ is Lipschitz by **Caffarelli's regularity theorem**.
- ▶ Key consequences:

The Proof: Some Intuition

- ▶ Optimal μ^* (Boltzmann fixed-point condition):

$$\mu^*(dw) \propto \exp\left(-\frac{1}{2\tau}\|w\|_2^2 - \frac{\beta}{\tau}\Psi(w; \mu^*)\right) dw$$

- ▶ Since $\nabla\Psi$ is Lipschitz, $w \mapsto \frac{1}{2}\|w\|_2^2 + \beta\Psi(w; \mu^*)$ is **strongly convex** for all $\beta < \beta_0$; thus, the optimal (Brenier–McCann) transport map $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that achieves $\mathcal{W}_2(\mu^*, \gamma_\tau)$ is Lipschitz by **Caffarelli’s regularity theorem**.

- ▶ Key consequences:

- ▶ can sample $W^1, \dots, W^N \stackrel{\text{i.i.d.}}{\sim} \gamma_\tau$, then transport to μ^* :

$$\mathbf{W} \mapsto T(\mathbf{W}) = (T(W^1), \dots, T(W^N)) \sim (\mu^*)^{\otimes N};$$

high-prob. guarantee for L^2 risk via refined Maurey’s empirical method (Ji–Telgarsky–Xian, 2020)

The Proof: Some Intuition

- ▶ Optimal μ^* (Boltzmann fixed-point condition):

$$\mu^*(dw) \propto \exp\left(-\frac{1}{2\tau}\|w\|_2^2 - \frac{\beta}{\tau}\Psi(w; \mu^*)\right) dw$$

- ▶ Since $\nabla\Psi$ is Lipschitz, $w \mapsto \frac{1}{2}\|w\|_2^2 + \beta\Psi(w; \mu^*)$ is **strongly convex** for all $\beta < \beta_0$; thus, the optimal (Brenier–McCann) transport map $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$ that achieves $\mathcal{W}_2(\mu^*, \gamma_\tau)$ is Lipschitz by **Caffarelli’s regularity theorem**.

- ▶ Key consequences:

- ▶ can sample $W^1, \dots, W^N \stackrel{\text{i.i.d.}}{\sim} \gamma_\tau$, then transport to μ^* :

$$\mathbf{W} \mapsto T(\mathbf{W}) = (T(W^1), \dots, T(W^N)) \sim (\mu^*)^{\otimes N};$$

high-prob. guarantee for L^2 risk via refined Maurey’s empirical method (**Ji–Telgarsky–Xian, 2020**)

- ▶ the functions

$$w \mapsto \|T(w) - w\|_2 \quad \text{and} \quad \mathbf{w} \mapsto \|\hat{f}_N(\cdot; \mathbf{w}) - \hat{f}_N(\cdot; T(\mathbf{w}))\|_2^2$$

are Lipschitz, so we can use Gaussian concentration inequalities

Entropy-Regularized Risk: Dynamic Formulation

In a nutshell: we will construct a flow of measures $\boldsymbol{\mu}^* = (\mu_t^*)_{0 \leq t \leq 1}$ with densities ρ_t^* , such that:

1. $\mu_0^* = \delta_0$ and $\mu_1^* = \mu^*$
2. The evolution of $t \mapsto \mu_t^*$ is governed by two coupled PDEs:

$$\partial_t \rho_t^*(w) = \nabla_w \cdot (\rho_t^*(w) \nabla_w V^*(w, t)) + \frac{\tau}{2} \Delta_w \rho_t^*(w) \quad (1a)$$

$$\partial_t V^*(w, t) = -\frac{\tau}{2} \Delta_w V^*(w, t) + \frac{1}{2} \|\nabla_w V^*(w, t)\|_2^2 \quad (1b)$$

on $\mathbb{R}^d \times [0, 1]$, where (1a) is the **forward** (Fokker–Planck) equation with initial condition $\rho_0^*(\cdot) = \delta(\cdot)$ and (1b) is the **backward** (Hamilton–Jacobi–Bellman) equation with terminal condition

$$V^*(w, 1) = \beta \left(R_0 + \int \tilde{f} d\mu_1^* + \Psi(w; \mu_1^*) \right).$$

3. $\boldsymbol{\mu}^*$ arises as a solution of a stochastic control problem on the space of measures.

The McKean–Vlasov Control Problem

- ▶ Controlled Itô SDE:

$$dW_t^{\mathbf{u}} = u_t dt + \sqrt{\tau} dB_t, \quad W_0 \equiv 0; t \in [0, 1]$$

where the control $\mathbf{u} = (u_t)_{t \in [0,1]}$ is a progressively measurable drift

The McKean–Vlasov Control Problem

- ▶ Controlled Itô SDE:

$$dW_t^{\mathbf{u}} = u_t dt + \sqrt{\tau} dB_t, \quad W_0 \equiv 0; t \in [0, 1]$$

where the control $\mathbf{u} = (u_t)_{t \in [0,1]}$ is a progressively measurable drift

- ▶ Flow of measures $\boldsymbol{\mu}^{\mathbf{u}} = (\mu_t^{\mathbf{u}})_{t \in [0,1]}$, where $\mu_t^{\mathbf{u}} := \text{Law}(W_t^{\mathbf{u}})$

The McKean–Vlasov Control Problem

- ▶ Controlled Itô SDE:

$$dW_t^{\mathbf{u}} = u_t dt + \sqrt{\tau} dB_t, \quad W_0 \equiv 0; t \in [0, 1]$$

where the control $\mathbf{u} = (u_t)_{t \in [0,1]}$ is a progressively measurable drift

- ▶ Flow of measures $\boldsymbol{\mu}^{\mathbf{u}} = (\mu_t^{\mathbf{u}})_{t \in [0,1]}$, where $\mu_t^{\mathbf{u}} := \text{Law}(W_t^{\mathbf{u}})$

Optimal control problem: minimize the total cost

$$J_{\beta, \tau}(\mathbf{u}) := \mathbf{E} \left[\underbrace{\frac{1}{2} \int_0^1 \|u_t\|_2^2 dt}_{\text{control cost}} \right] + \underbrace{\frac{\beta}{2} R(\mu_1^{\mathbf{u}})}_{\text{terminal cost}}$$

The McKean–Vlasov Control Problem

- ▶ Controlled Itô SDE:

$$dW_t^{\mathbf{u}} = u_t dt + \sqrt{\tau} dB_t, \quad W_0 \equiv 0; t \in [0, 1]$$

where the control $\mathbf{u} = (u_t)_{t \in [0,1]}$ is a progressively measurable drift

- ▶ Flow of measures $\boldsymbol{\mu}^{\mathbf{u}} = (\mu_t^{\mathbf{u}})_{t \in [0,1]}$, where $\mu_t^{\mathbf{u}} := \text{Law}(W_t^{\mathbf{u}})$

Optimal control problem: minimize the total cost

$$J_{\beta, \tau}(\mathbf{u}) := \mathbf{E} \left[\underbrace{\frac{1}{2} \int_0^1 \|u_t\|_2^2 dt}_{\text{control cost}} \right] + \underbrace{\frac{\beta}{2} R(\boldsymbol{\mu}^{\mathbf{u}})}_{\text{terminal cost}}$$

- ▶ **Key point:** for a control \mathbf{u} , the terminal cost depends *nonlinearly* on the probability law of $W_1^{\mathbf{u}}$ — this is an instance of a *McKean–Vlasov control problem* (Carmona and Delarue)

Transport Map Interpretation

- ▶ Each control \mathbf{u} transports the Brownian path $B_{[0,1]} = (B_t)_{t \in [0,1]}$ to a random vector $W_1^{\mathbf{u}}$:

$$T^{\mathbf{u}} : B_{[0,1]} \longmapsto \int_0^1 u_t dt + \sqrt{\tau} B_1$$

Transport Map Interpretation

- ▶ Each control \mathbf{u} transports the Brownian path $B_{[0,1]} = (B_t)_{t \in [0,1]}$ to a random vector $W_1^{\mathbf{u}}$:

$$T^{\mathbf{u}} : B_{[0,1]} \mapsto \int_0^1 u_t dt + \sqrt{\tau} B_1$$

- ▶ Zero drift ($\mathbf{u} \equiv 0$): $T^0(B_{[0,1]}) = \sqrt{\tau} B_1 \sim \gamma_{\tau}$ (samples from the Gaussian prior)

Transport Map Interpretation

- ▶ Each control \mathbf{u} transports the Brownian path $B_{[0,1]} = (B_t)_{t \in [0,1]}$ to a random vector $W_1^{\mathbf{u}}$:

$$T^{\mathbf{u}} : B_{[0,1]} \mapsto \int_0^1 u_t dt + \sqrt{\tau} B_1$$

- ▶ Zero drift ($\mathbf{u} \equiv 0$): $T^0(B_{[0,1]}) = \sqrt{\tau} B_1 \sim \gamma_{\tau}$ (samples from the Gaussian prior)
- ▶ Any other \mathbf{u} affects the distribution $\mu_1^{\mathbf{u}}$, incurs the control cost

$$\mathbf{E} \left[\frac{1}{2} \int_0^1 \|u_t\|_2^2 dt \right] = \tau \cdot D_{\text{KL}}(\mathbf{P}^{\mathbf{u}} \| \mathbf{P}^0),$$

where $\mathbf{P}^{\mathbf{u}} := \text{Law}(W_{[0,1]}^{\mathbf{u}})$ and $\mathbf{P}^0 := \text{Law}(B_{[0,1]})$

Transport Map Interpretation

- ▶ Each control \mathbf{u} transports the Brownian path $B_{[0,1]} = (B_t)_{t \in [0,1]}$ to a random vector $W_1^{\mathbf{u}}$:

$$T^{\mathbf{u}} : B_{[0,1]} \mapsto \int_0^1 u_t dt + \sqrt{\tau} B_1$$

- ▶ Zero drift ($\mathbf{u} \equiv 0$): $T^0(B_{[0,1]}) = \sqrt{\tau} B_1 \sim \gamma_\tau$ (samples from the Gaussian prior)
- ▶ Any other \mathbf{u} affects the distribution $\mu_1^{\mathbf{u}}$, incurs the control cost

$$\mathbf{E} \left[\frac{1}{2} \int_0^1 \|u_t\|_2^2 dt \right] = \tau \cdot D_{\text{KL}}(\mathbf{P}^{\mathbf{u}} \parallel \mathbf{P}^0),$$

where $\mathbf{P}^{\mathbf{u}} := \text{Law}(W_{[0,1]}^{\mathbf{u}})$ and $\mathbf{P}^0 := \text{Law}(B_{[0,1]})$

- ▶ The *optimal control* \mathbf{u}^* will give $W_1^{\mathbf{u}^*} \sim \mu^*$, and

$$D_{\text{KL}}(\mathbf{P}^{\mathbf{u}^*} \parallel \mathbf{P}^0) = \min_{\mathbf{u}} D_{\text{KL}}(\mathbf{P}^{\mathbf{u}} \parallel \mathbf{P}^0) = D_{\text{KL}}(\mu^* \parallel \gamma_\tau)$$

— entropic optimal transport (or the Schrödinger bridge) problem!

Entropy-Regularized Risk: Dynamic Formulation

Theorem (Tzen–Raginsky, 2020) Let μ^\star be the (unique) minimizer of the free energy $F_{\beta,\tau}(\mu)$. Then the optimal controlled process solves the Itô SDE

$$dW_t = -\nabla_w V^\star(W_t, t) dt + \sqrt{\tau} dB_t, \quad t \in [0, 1]; \quad W_0 = 0,$$

$$\text{where } V^\star(w, t) := -\tau \log \mathbf{E} \left[\exp \left(-\frac{\beta}{\tau} \Psi(B_\tau; \mu^\star) \right) \middle| B_{\tau t} = w \right].$$

Under the optimal control \mathbf{u}^\star , W_1 is distributed according to μ^\star , and

$$\inf_{\mathbf{u}} J_{\beta,\tau}(\mathbf{u}) = \beta F_{\beta,\tau}(\mu^\star).$$

Entropy-Regularized Risk: Dynamic Formulation

Theorem (Tzen–Raginsky, 2020) Let μ^\star be the (unique) minimizer of the free energy $F_{\beta,\tau}(\mu)$. Then the optimal controlled process solves the Itô SDE

$$dW_t = -\nabla_w V^\star(W_t, t) dt + \sqrt{\tau} dB_t, \quad t \in [0, 1]; \quad W_0 = 0,$$

$$\text{where } V^\star(w, t) := -\tau \log \mathbf{E} \left[\exp \left(-\frac{\beta}{\tau} \Psi(B_\tau; \mu^\star) \right) \middle| B_{\tau t} = w \right].$$

Under the optimal control \mathbf{u}^\star , W_1 is distributed according to μ^\star , and

$$\inf_{\mathbf{u}} J_{\beta,\tau}(\mathbf{u}) = \beta F_{\beta,\tau}(\mu^\star).$$

The Schrödinger bridge problem:

- ▶ Let $\mathcal{U}(\mu^\star) := \{\mathbf{u} : W_1^{\mathbf{u}} \sim \mu^\star\}$
- ▶ Then the optimal drift $\mathbf{u}^\star \in \mathcal{U}(\mu^\star)$, and

$$\mathbf{E} \left[\frac{1}{2} \int_0^1 \|u_t^\star\|_2^2 dt \right] = \inf_{\mathbf{u} \in \mathcal{U}(\mu^\star)} \mathbf{E} \left[\frac{1}{2} \int_0^1 \|u_t\|_2^2 dt \right] = \tau D_{\text{KL}}(\mu^\star \| \gamma_\tau)$$

A Closer Look at the Optimal Control

- ▶ Optimal McKean–Vlasov dynamics:

$$dW_t = -\nabla_w V^*(W_t, t) dt + \sqrt{\tau} dB_t, \quad t \in [0, 1]; \quad W_0 = 0,$$

$$\text{where } V^*(w, t) := -\tau \log \mathbf{E} \left[\exp \left(-\frac{\beta}{\tau} \Psi(B_\tau; \mu^*) \right) \middle| B_{\tau t} = w \right].$$

A Closer Look at the Optimal Control

- ▶ Optimal McKean–Vlasov dynamics:

$$dW_t = -\nabla_w V^*(W_t, t) dt + \sqrt{\tau} dB_t, \quad t \in [0, 1]; \quad W_0 = 0,$$

$$\text{where } V^*(w, t) := -\tau \log \mathbf{E} \left[\exp \left(-\frac{\beta}{\tau} \Psi(B_\tau; \mu^*) \right) \middle| B_{\tau t} = w \right].$$

- ▶ Direct computation:

$$-\nabla_w V^*(W_t, t) = -\beta \int_{\mathbb{R}^d} \nabla \Psi(w; \mu^*) Q_{W_t, t}(dw)$$

where $Q_{w, t}(dv)$ is the Gibbs measure

$$Q_{w, t}(A) := \frac{\int_A \exp \left(-\frac{\|v-w\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(v; \mu^*) \right) dv}{\int_{\mathbb{R}^d} \exp \left(-\frac{\|v-w\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(v; \mu^*) \right) dv}$$

A Closer Look at the Optimal Control

- ▶ Optimal McKean–Vlasov dynamics:

$$dW_t = -\nabla_w V^*(W_t, t) dt + \sqrt{\tau} dB_t, \quad t \in [0, 1]; \quad W_0 = 0,$$

$$\text{where } V^*(w, t) := -\tau \log \mathbf{E} \left[\exp \left(-\frac{\beta}{\tau} \Psi(B_\tau; \mu^*) \right) \middle| B_{\tau t} = w \right].$$

- ▶ Direct computation:

$$-\nabla_w V^*(W_t, t) = -\beta \int_{\mathbb{R}^d} \nabla \Psi(w; \mu^*) Q_{W_t, t}(dw)$$

where $Q_{w, t}(dv)$ is the Gibbs measure

$$Q_{w, t}(A) := \frac{\int_A \exp \left(-\frac{\|v-w\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(v; \mu^*) \right) dv}{\int_{\mathbb{R}^d} \exp \left(-\frac{\|v-w\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(v; \mu^*) \right) dv}$$

- ▶ Gibbs flow $(Q_{w, t})_{t \in [0, 1]}$: interpolates between $Q_{0, 0} = \mu^*$ and $Q_{0, 1} = \delta_0$; becomes more concentrated as $t \rightarrow 1$:

$$\int_{\mathbb{R}^d} \|v - w\|^2 Q_{w, t}(dv) \leq \kappa(\beta^2 + \tau d)(1 - t).$$

Gradient Descent as a Greedy Heuristic

- ▶ Optimal McKean–Vlasov dynamics:

$$dW_t = -\beta \left(\int_{\mathbb{R}^d} \nabla \Psi(w; \mu^\star) Q_{W_t, t}(dw) \right) dt + \sqrt{\tau} dB_t$$

$$\text{where } Q_{W_t, t}(dw) = \frac{\exp \left(-\frac{\|w - W_t\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(w; \mu^\star) \right) dw}{\int_{\mathbb{R}^d} \exp \left(-\frac{\|w - W_t\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(W_t; \mu^\star) \right) dw}$$

is concentrated around W_t

Gradient Descent as a Greedy Heuristic

- ▶ Optimal McKean–Vlasov dynamics:

$$dW_t = -\beta \left(\int_{\mathbb{R}^d} \nabla \Psi(w; \mu^\star) Q_{W_t, t}(dw) \right) dt + \sqrt{\tau} dB_t$$

$$\text{where } Q_{W_t, t}(dw) = \frac{\exp \left(-\frac{\|w - W_t\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(w; \mu^\star) \right) dw}{\int_{\mathbb{R}^d} \exp \left(-\frac{\|w - W_t\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(W_t; \mu^\star) \right) dw}$$

is concentrated around W_t

- ▶ Greedy approximation:

$$\int_{\mathbb{R}^d} \nabla \Psi(w; \mu^\star) Q_{W_t, t}(dw) \approx \nabla \Psi(W_t; \mu^\star) \approx \nabla \Psi(W_t; \text{Law}(W_t))$$

Gradient Descent as a Greedy Heuristic

- ▶ Optimal McKean–Vlasov dynamics:

$$dW_t = -\beta \left(\int_{\mathbb{R}^d} \nabla \Psi(w; \mu^*) Q_{W_t, t}(dw) \right) dt + \sqrt{\tau} dB_t$$

$$\text{where } Q_{W_t, t}(dw) = \frac{\exp \left(-\frac{\|w - W_t\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(w; \mu^*) \right) dw}{\int_{\mathbb{R}^d} \exp \left(-\frac{\|w - W_t\|_2^2}{2\tau(1-t)} - \frac{\beta}{\tau} \Psi(W_t; \mu^*) \right) dw}$$

is concentrated around W_t

- ▶ Greedy approximation:

$$\int_{\mathbb{R}^d} \nabla \Psi(w; \mu^*) Q_{W_t, t}(dw) \approx \nabla \Psi(W_t; \mu^*) \approx \nabla \Psi(W_t; \text{Law}(W_t))$$

- ▶ Nonlinear dynamics:

$$d\hat{W}_t = -\beta \nabla \Psi(\hat{W}_t; \text{Law}(\hat{W}_t)) dt + \sqrt{\tau} dB_t, \quad \hat{W}_0 = 0; t \in [0, 1]$$

Question: can we show that \hat{W} tracks W with high probability?

SGD as a Greedy Heuristic

SGD update:

- ▶ Data: $X_1, X_2, \dots \stackrel{\text{i.i.d.}}{\sim} \pi$, $Y_k = f(X_k)$
- ▶ Generate iterates $\mathbf{W}_k = (W_k^1, \dots, W_k^N)$, $k = 0, 1, 2, \dots$:

$$W_{k+1}^i = W_k^i + \eta \beta (Y_{k+1} - \hat{f}_N(X_{k+1}, \mathbf{W}_k)) \nabla_w \sigma(X_{k+1}; W_k^i)$$

with $\mathbf{W}_0 \sim (\gamma_\tau)^{\otimes N}$ (Gaussian initialization)

SGD as a Greedy Heuristic

SGD update:

- ▶ Data: $X_1, X_2, \dots \stackrel{\text{i.i.d.}}{\sim} \pi, \quad Y_k = f(X_k)$
- ▶ Generate iterates $\mathbf{W}_k = (W_k^1, \dots, W_k^N), \quad k = 0, 1, 2, \dots:$

$$W_{k+1}^i = W_k^i + \eta \beta (Y_{k+1} - \hat{f}_N(X_{k+1}, \mathbf{W}_k)) \nabla_w \sigma(X_{k+1}; W_k^i)$$

with $\mathbf{W}_0 \sim (\gamma_\tau)^{\otimes N}$ (Gaussian initialization)

Note:

- ▶ no explicit regularization or additional noise
- ▶ Gaussian initialization, $W_0^i \stackrel{\text{i.i.d.}}{\sim} \gamma_\tau$
- ▶ we will be interested in the regime

$$\tau = \varepsilon^2/d, \quad \beta = \sqrt{\tau d} = \varepsilon$$

SGD Tracks the Optimal McKean–Vlasov Dynamics

Theorem (Tzen–Raginsky, 2020) Let $\boldsymbol{\mu}^* = \{\mu_t^*\}_{t \in [0,1]}$ be the flow of measures along the optimal McKean–Vlasov dynamics, with $\mu_0^* = \delta_0$ and $\mu_1^* = \mu^*$. Then, with probability at least $1 - \delta$,

$$\begin{aligned} \max_{0 \leq k \leq n} |R_N(\mathbf{W}_k) - R(\mu_{k\eta}^*)| &\leq \frac{\kappa}{N} + \kappa \sqrt{\frac{1}{N} \log \frac{N}{\delta}} + \kappa \beta \\ &+ \kappa(1 + \sqrt{\eta}\beta)(1 + \beta)e^{\kappa\beta} \left(\beta + \sqrt{\tau d \log \frac{Nd}{\delta\eta}} \right). \end{aligned}$$

SGD Tracks the Optimal McKean–Vlasov Dynamics

Theorem (Tzen–Raginsky, 2020) Let $\mu^\star = \{\mu_t^\star\}_{t \in [0,1]}$ be the flow of measures along the optimal McKean–Vlasov dynamics, with $\mu_0^\star = \delta_0$ and $\mu_1^\star = \mu^\star$. Then, with probability at least $1 - \delta$,

$$\begin{aligned} \max_{0 \leq k \leq n} |R_N(\mathbf{W}_k) - R(\mu_{k\eta}^\star)| &\leq \frac{\kappa}{N} + \kappa \sqrt{\frac{1}{N} \log \frac{N}{\delta}} + \kappa \beta \\ &+ \kappa(1 + \sqrt{\eta}\beta)(1 + \beta)e^{\kappa\beta} \left(\beta + \sqrt{\tau d \log \frac{Nd}{\delta\eta}} \right). \end{aligned}$$

Remarks:

- ▶ Similar to the bound of Mei–Misiakiewicz–Montanari (2019), except here $T = 1$
- ▶ If we take

$$\beta = \sqrt{\tau d} = \varepsilon \quad \text{and} \quad N \sim \frac{1}{\varepsilon^2},$$

then w.h.p. SGD will track optimal MKV dynamics to $O(\varepsilon)$ accuracy

Comparison with Nonlinear Dynamics

1. **This work:** Optimal McKean–Vlasov dynamics for $\mu_t^*(dw) = \rho_t^*(w) dw$

$$\partial_t \rho_t^*(w) = \nabla_w \cdot (\rho_t^*(w) \nabla_w V^*(w, t)) + \frac{\tau}{2} \Delta_w \rho_t^*(w)$$

$$\partial_t V^*(w, t) = -\frac{\tau}{2} \Delta_w V^*(w, t) + \frac{1}{2} \|\nabla_w V^*(w, t)\|_2^2$$

(with appropriate initial and terminal conditions)

- ▶ finite time horizon ($T = 1$)
- ▶ $\mu_1^* = \mu^*$ (exact optimality at $T = 1$)
- ▶ two coupled nonlinear PDEs, forward and backward in time

Comparison with Nonlinear Dynamics

1. **This work:** Optimal McKean–Vlasov dynamics for $\mu_t^*(dw) = \rho_t^*(w) dw$

$$\partial_t \rho_t^*(w) = \nabla_w \cdot (\rho_t^*(w) \nabla_w V^*(w, t)) + \frac{\tau}{2} \Delta_w \rho_t^*(w)$$

$$\partial_t V^*(w, t) = -\frac{\tau}{2} \Delta_w V^*(w, t) + \frac{1}{2} \|\nabla_w V^*(w, t)\|_2^2$$

(with appropriate initial and terminal conditions)

- ▶ finite time horizon ($T = 1$)
- ▶ $\mu_1^* = \mu^*$ (exact optimality at $T = 1$)
- ▶ two coupled nonlinear PDEs, forward and backward in time

2. **Existing work:** Nonlinear dynamics for $\mu_t(dw) = \rho_t(w) dw$

$$\partial_t \rho_t(w) = \beta \nabla_w \cdot (\rho_t(w) \nabla_w \Psi(w; \mu_t)) + \frac{\tau}{2} \Delta_w \rho_t(w)$$

- ▶ infinite time ($T \rightarrow \infty$)
- ▶ can guarantee convergence to μ^* as $t \rightarrow \infty$
- ▶ one PDE running forward in time

Comparison with Nonlinear Dynamics

1. **This work:** Optimal McKean–Vlasov dynamics for $\mu_t^*(dw) = \rho_t^*(w) dw$

$$\partial_t \rho_t^*(w) = \nabla_w \cdot (\rho_t^*(w) \nabla_w V^*(w, t)) + \frac{\tau}{2} \Delta_w \rho_t^*(w)$$

$$\partial_t V^*(w, t) = -\frac{\tau}{2} \Delta_w V^*(w, t) + \frac{1}{2} \|\nabla_w V^*(w, t)\|_2^2$$

(with appropriate initial and terminal conditions)

- ▶ finite time horizon ($T = 1$)
- ▶ $\mu_1^* = \mu^*$ (exact optimality at $T = 1$)
- ▶ two coupled nonlinear PDEs, forward and backward in time

2. **Existing work:** Nonlinear dynamics for $\mu_t(dw) = \rho_t(w) dw$

$$\partial_t \rho_t(w) = \beta \nabla_w \cdot (\rho_t(w) \nabla_w \Psi(w; \mu_t)) + \frac{\tau}{2} \Delta_w \rho_t(w)$$

- ▶ infinite time ($T \rightarrow \infty$)
- ▶ can guarantee convergence to μ^* as $t \rightarrow \infty$
- ▶ one PDE running forward in time

With proper scaling of τ and β , nonlinear dynamics can be seen as a greedy approximation to optimal MKV dynamics.

Preprint at <https://arxiv.org/abs/2002.01987>