

Machine Learning from a Continuous Viewpoint

Weinan E

Princeton University

Joint work with:

Chao Ma, Lei Wu

<https://arxiv.org/pdf/1912.12777.pdf>

Standard ML problems for which we are given the dataset:

- ① **Supervised learning:** Given $S = \{(\mathbf{x}_j, y_j = f^*(\mathbf{x}_j)), j \in [n]\}$, learn f^* .
Objective: Minimize population risk over the “hypothesis space”

$$\mathcal{R}(f) = \mathbb{E}_{\mathbf{x} \sim \mu} (f(\mathbf{x}) - f^*(\mathbf{x}))^2$$

- ② **Dimension reduction:** Given $S = \{\mathbf{x}_j, j \in [n]\} \subset \mathbb{R}^D$ sampled from μ , find mapping: $\Phi : \mathbb{R}^D \rightarrow \mathbb{R}^d, (d \ll D)$ that best preserves all important features of μ .

“Auto-encoder”: Minimize reconstruction error $\mathbf{x} - \tilde{\mathbf{x}}$: $\tilde{\mathbf{x}} = \Psi(\mathbf{z}), \mathbf{z} = \Phi(\mathbf{x})$

$$\mathcal{R}(f) = \mathbb{E}_{\mathbf{x} \sim \mu} (\mathbf{x} - \Psi(\mathbf{z}))^2 = \mathbb{E}_{\mathbf{x} \sim \mu} (\mathbf{x} - \Psi(\Phi(\mathbf{x})))^2$$

Non-standard ML problem, no dataset given beforehand:

① **Ground state of quantum many-body problem:**

Let $H = -\frac{\hbar^2}{2m}\Delta + V$ be the Hamiltonian operator of the system

$$I(\phi) = \frac{(\phi, H\phi)}{(\phi, \phi)} = \mathbb{E}_{\mathbf{x} \sim \mu_\phi} \frac{\phi(\mathbf{x})H\phi(\mathbf{x})}{\phi(\mathbf{x})^2}, \quad \mu_\phi(d\mathbf{x}) = \frac{1}{(\phi, \phi)} \phi^2(\mathbf{x})d\mathbf{x}$$

subject to the constraint imposed by Pauli exclusion principle.

② **Stochastic control problems:**

$$s_{t+1} = s_t + b_t(s_t, a_t) + \xi_{t+1},$$

s_t = state at time t , a_t = control at time t , ξ_t = i.i.d. noise

$$L(\{a_t\}_{t=0}^{T-1}) = \mathbb{E}_{\{\xi_t\}} \left\{ \sum_{t=0}^{T-1} c_t(s_t, a_t(s_t)) + c_T(s_T) \right\},$$

Look for feedback control: $a_t = F(t, s_t), t = 0, 1, \dots, T - 1$.

Remark: High dimensionality

Benchmark: High dimensional integration

$$I(g) = \int_{X=[0,1]^d} g(\mathbf{x}) d\mu, \quad I_m(g) = \frac{1}{m} \sum_j g(\mathbf{x}_j)$$

Grid-based quadrature rules:

$$I(g) - I_m(g) \sim \frac{C(g)}{m^{\alpha/d}}$$

Appearance of $1/d$ in the exponent of m : **Curse of dimensionality (CoD)!**

If we want $m^{-\alpha/d} = 0.1$, then $m = 10^{d/\alpha} = 10^d$, if $\alpha = 1$.

Monte Carlo: $\{\mathbf{x}_j, j \in [m]\}$ is uniformly distributed in X .

$$\mathbb{E}(I(g) - I_m(g))^2 = \frac{\text{var}(g)}{m}, \quad \text{var}(g) = \int_X g^2(\mathbf{x}) d\mathbf{x} - \left(\int_X g(\mathbf{x}) d\mathbf{x} \right)^2$$

However, $\text{var}(g)$ can be very large in high dimension. Variance reduction!

Overall strategy:

Formulate a “nice” continuous problem, then discretize to get concrete models/algorithms.

- For PDEs, “nice” = well-posed.
- For calculus of variation problems, “nice” = “convex”, lower semi-continuous.

For machine learning, “nice” = variational problem has simple landscape.

How do we represent a function? An illustrative example

Traditional approach for Fourier transform:

$$f(\mathbf{x}) = \int_{\mathbb{R}^d} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})} d\boldsymbol{\omega}, \quad f_m(\mathbf{x}) = \frac{1}{m} \sum_j a(\boldsymbol{\omega}_j) e^{i(\boldsymbol{\omega}_j, \mathbf{x})}$$

$\{\boldsymbol{\omega}_j\}$ is a fixed grid, e.g. uniform.

$$\|f - f_m\|_{L^2(X)} \leq C_0 m^{-\alpha/d} \|f\|_{H^\alpha(X)}$$

“New” approach: Let π be a probability distribution and

$$f(\mathbf{x}) = \int_{\mathbb{R}^d} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})} \pi(d\boldsymbol{\omega}) = \mathbb{E}_{\boldsymbol{\omega} \sim \pi} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})}$$

Let $\{\boldsymbol{\omega}_j\}$ be an i.i.d. sample of π , $f_m(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a(\boldsymbol{\omega}_j) e^{i(\boldsymbol{\omega}_j, \mathbf{x})}$,

$$\mathbb{E} |f(\mathbf{x}) - f_m(\mathbf{x})|^2 = m^{-1} \text{var}(f)$$

$f_m(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\boldsymbol{\omega}_j^T \mathbf{x}) =$ two-layer neural network with activation function $\sigma(z) = e^{iz}$.

Integral transform-based representation

Let σ be a scalar nonlinear function (activation function), e.g. $\sigma = \text{ReLU}$

Consider functions represented in the form:

$$\begin{aligned} f(\mathbf{x}; \theta) &= \int_{\mathbb{R}^d} a(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x}) \pi(d\mathbf{w}) \\ &= \mathbb{E}_{\mathbf{w} \sim \pi} a(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x}) \\ &= \mathbb{E}_{(a, \mathbf{w}) \sim \rho} a \sigma(\mathbf{w}^T \mathbf{x}) \end{aligned}$$

θ = parameter to be optimized:

- $\theta = a(\cdot)$ corresponds to a feature-based model
- $\theta = \rho$ corresponds to a two-layer neural network-like model.
- $\theta = (a(\cdot), \pi(\cdot))$, a new model

- Fourier method: $\pi \sim \frac{1}{N} \sum_j \delta_{\omega_j}$ where $\{\omega_j\}$ lives on a uniform lattice. Optimize $a(\cdot)$:

$$f(\mathbf{x}; \theta) \sim f_m(\mathbf{x}) = \frac{1}{m} \sum_j a(\mathbf{w}_j) \sigma(\mathbf{w}_j^T \mathbf{x})$$

- Neural network-based method: $\rho \sim \frac{1}{N} \sum_j \delta_{(a_j, \omega_j)}$ ($\{\omega_j\}$ is also optimized):

$$f(\mathbf{x}; \theta) \sim f_m(\mathbf{x}) = \frac{1}{m} \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x})$$

then optimize (say, using L-BFGS) — this is more in line with traditional numerical analysis (e.g. nonlinear finite element or meshless methods).

For truly large datasets, we need to use stochastic algorithms

- objective function are all expressed as expectations:

$$\mathcal{R}(\theta) = \mathbb{E}_{\mathbf{x} \sim \mu} (f(\mathbf{x}; \theta) - f^*(\mathbf{x}))^2$$

$$\mathcal{R}(\theta_1, \theta_2) = \mathbb{E}_{\mathbf{x} \sim \mu} (\mathbf{x} - \Psi(\Phi(\mathbf{x}; \theta_1); \theta_2))^2$$

$$I(\theta) = \mathbb{E}_{\mathbf{x} \sim \mu_\theta} \frac{\phi(\mathbf{x}; \theta) H \phi(\mathbf{x}; \theta)}{\phi(\mathbf{x}; \theta)^2}$$

- Gradient descent (GD) can be readily converted to stochastic gradient descent (SGD).
Let $F = F(\theta) = \mathbb{E}_{\mathbf{x} \sim \mu} g(\theta, \mathbf{x})$ be the objective function:

$$\text{GD : } \theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathbb{E}_{\mathbf{x} \sim \mu} g(\theta_k, \mathbf{x})$$

$$\text{SGD : } \theta_{k+1} = \theta_k - \eta \nabla_{\theta} g(\theta_k, \mathbf{x}_k)$$

where $\{\mathbf{x}_k\}$ are i.i.d. random samples.

Optimization: Defining gradient flows

“Free energy” = $\mathcal{R}(f) = \mathbb{E}_{\mathbf{x} \sim \mu} (f(\mathbf{x}) - f^*(\mathbf{x}))^2$

$$f(\mathbf{x}) = \int a(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x}) \pi(d\mathbf{w}) = \mathbb{E}_{\mathbf{w} \sim \pi} a(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x})$$

Follow Halperin and Hohenberg (1977):

- a = non-conserved, use “model A” dynamics (Allen-Cahn):

$$\frac{\partial a}{\partial t} = -\frac{\delta \mathcal{R}}{\delta a}$$

- π = conserved (probability density), use “model B” (Cahn-Hilliard):

$$\frac{\partial \pi}{\partial t} + \nabla \cdot \mathbf{J} = 0$$

$$\mathbf{J} = \pi \mathbf{v}, \mathbf{v} = -\nabla V, V = \frac{\delta \mathcal{R}}{\delta \pi}.$$

Gradient flow for the feature-based model

Fix π , optimize a .

$$\partial_t a(\mathbf{w}, t) = -\frac{\delta \mathcal{R}}{\delta a}(\mathbf{w}, t) = -\int a(\tilde{\mathbf{w}}, t) K(\mathbf{w}, \tilde{\mathbf{w}}) \pi(d\tilde{\mathbf{w}}) + \tilde{f}(\mathbf{w})$$

$$K(\mathbf{w}, \tilde{\mathbf{w}}) = \mathbb{E}_{\mathbf{x}}[\sigma(\mathbf{w}^T \mathbf{x}) \sigma(\tilde{\mathbf{w}}^T \mathbf{x})], \quad \tilde{f}(\mathbf{w}) = \mathbb{E}_{\mathbf{x}}[f^*(\mathbf{x}) \sigma(\mathbf{w}^T \mathbf{x})]$$

This is an integral equation with a symmetric positive definite kernel.

Decay estimates due to convexity: Let $f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{w} \sim \pi} a^*(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x})$,

$$I(t) = \frac{1}{2} \|a(\cdot, t) - a^*(\cdot)\|^2 + t(\mathcal{R}(a(t)) - \mathcal{R}(a^*))$$

Then we have

$$\frac{dI}{dt} \leq 0, \quad \mathcal{R}(a(t)) \leq \frac{C_0}{t}$$

Conservative gradient flow

Optimize ρ :

$$f(\mathbf{x}) = \mathbb{E}_{\mathbf{u} \sim \rho} \phi(\mathbf{x}, \mathbf{u})$$

Example: $\mathbf{u} = (a, \mathbf{w})$, $\phi(\mathbf{x}, \mathbf{u}) = a\sigma(\mathbf{w}^T \mathbf{x})$

$$\partial_t \rho = \nabla(\rho \nabla V)$$

$$V(\mathbf{u}) = \frac{\delta \mathcal{R}}{\delta \rho}(\mathbf{u}) = \mathbb{E}_{\mathbf{x}}[(f(\mathbf{x}) - f^*(\mathbf{x}))\phi(\mathbf{x}, \mathbf{u})] = \int K(\mathbf{u}, \tilde{\mathbf{u}}) \rho(d\tilde{\mathbf{u}}) - \tilde{f}(\mathbf{u})$$

- This is the mean-field equation derived by Chizat and Bach (2018), Mei, Montanari and Nguyen (2018), Rotskoff and Vanden-Eijnden (2018), Sirignano and Spiliopoulos (2018), by studying the continuum limit of two-layer neural networks.
- Does not satisfy displacement convexity.

Optimize (a, π) (a = non-conservative, π = conservative):

$$\partial_t a(\mathbf{w}, t) = -\frac{\delta \mathcal{R}}{\delta a}(\mathbf{w}, t) = -\int a(\tilde{\mathbf{w}}, t) K(\mathbf{w}, \tilde{\mathbf{w}}) \pi(d\tilde{\mathbf{w}}, t) + \tilde{f}(\mathbf{w})$$

$$\partial_t \pi = \nabla(\pi \nabla \tilde{V}), \quad V(\mathbf{w}) = \frac{\delta \mathcal{R}}{\delta \pi}(\mathbf{w})$$

Discretizing the gradient flows

- Discretizing the population risk (into the empirical risk) using data
- Discretizing the gradient flow
 - particle method – the dynamic version of Monte Carlo
 - smoothed particle method – analog of vortex blob method
 - spectral method – very effective in low dimensions

We will see that gradient descent algorithm (GD) for random feature and neural network models are simply the particle method discretization of the gradient flows discussed before.

Particle method for the feature-based model

$$\partial_t a(\mathbf{w}, t) = -\frac{\delta \mathcal{R}}{\delta a}(\mathbf{w}) = -\int a(\tilde{\mathbf{w}}, t) K(\mathbf{w}, \tilde{\mathbf{w}}) \pi(d\tilde{\mathbf{w}}) + \tilde{f}(\mathbf{w})$$

$$\pi(d\mathbf{w}) \sim \frac{1}{m} \sum_j \delta_{\mathbf{w}_j}, a(\mathbf{w}_j, t) \sim a_j(t)$$

Discretized version:

$$\frac{d}{dt} a_j(t) = -\frac{1}{m} \sum_k K(\mathbf{w}_j, \mathbf{w}_k) a_k(t) + \tilde{f}(\mathbf{w}_j)$$

This is exactly the GD for the random feature model.

$$f(\mathbf{x}) \sim f_m(\mathbf{x}) = \frac{1}{m} \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x})$$

Particle method for the conservative flow

$$\partial_t \rho = \nabla(\rho \nabla V) \quad (1)$$

Particle method discretization:

$$\rho(t, d\mathbf{u}) \sim \frac{1}{m} \sum_j \delta_{\mathbf{u}_j(t)}$$

Define the loss function

$$I(\mathbf{u}_1, \dots, \mathbf{u}_m) = \mathcal{R}(f_m), \quad f_m(\mathbf{x}) = \frac{1}{m} \sum_j \phi(\mathbf{x}, \mathbf{u}_j)$$

Lemma: Given a set of initial data $\{\mathbf{u}_j^0, j \in [m]\}$. The solution of (1) with initial data $\rho(0) = \frac{1}{m} \sum_{j=1}^m \delta_{\mathbf{u}_j^0}$ is given by

$$\rho(t) = \frac{1}{m} \sum_{j=1}^m \delta_{\mathbf{u}_j(t)}$$

where the particles $\{\mathbf{u}_j(\cdot), j \in [m]\}$ solves:

$$\frac{d\mathbf{u}_j}{dt} = -\nabla_{\mathbf{u}_j} I(\mathbf{u}_1, \dots, \mathbf{u}_m), \quad \mathbf{u}_j(0) = \mathbf{u}_j^0, \quad j \in [m]$$

This is exactly the GD dynamics for two-layer neural networks.

Comparison with conventional NN models

Continuous viewpoint (in this case same as mean-field): $f_m(\mathbf{x}) = \frac{1}{m} \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x})$

Conventional NN models: $f_m(\mathbf{x}) = \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x})$

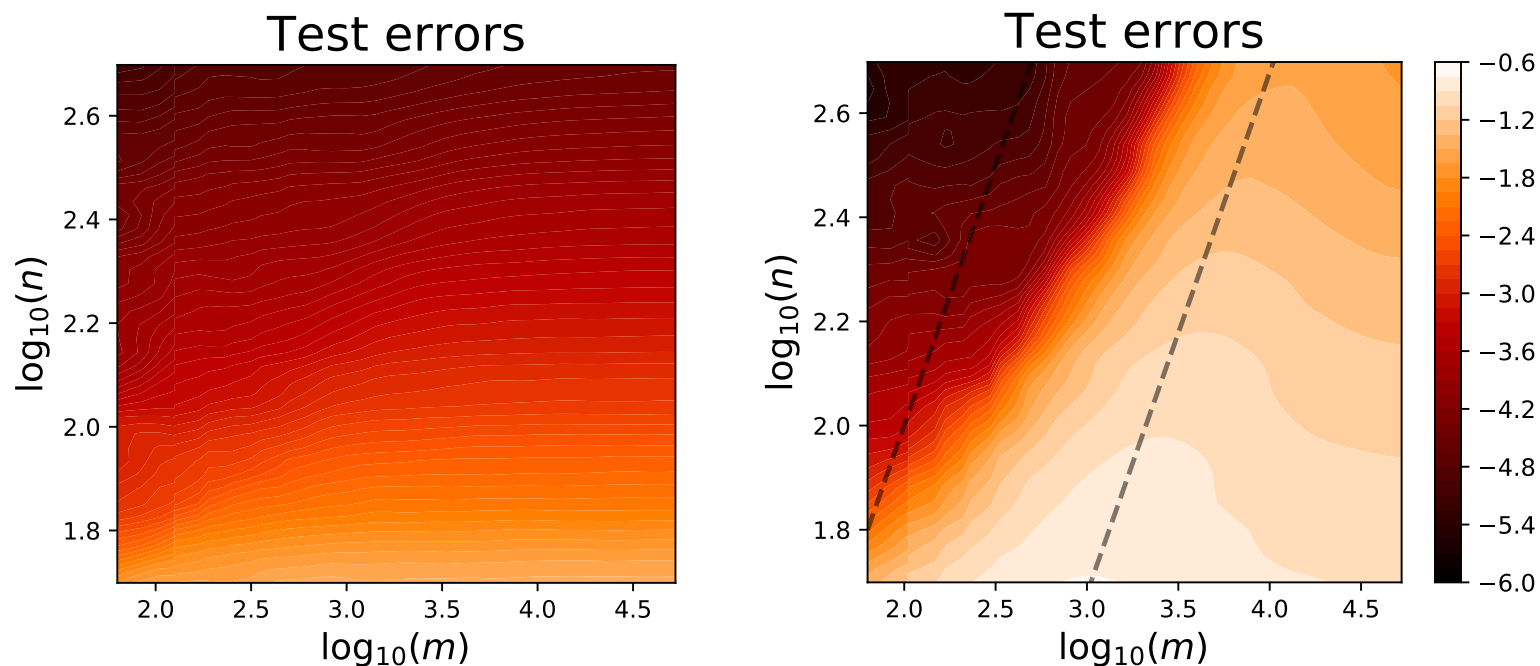


Figure: (Left) continuous viewpoint; (Right) conventional NN models. Target function is a single neuron $f^*(\mathbf{x}) = \sigma(\mathbf{e}_1^T \mathbf{x})$.

Flow-based representation

Continuous dynamical system viewpoint (E (2017), Haber and Ruthotto (2017), “Neural ODEs” (Chen et al, 2018))

$$\frac{d\mathbf{z}}{d\tau} = \mathbf{g}(\tau, \mathbf{z}), \quad \mathbf{z}(0) = \mathbf{x}$$

The flow-map at time 1: $\mathbf{x} \rightarrow \mathbf{z}(\mathbf{x}, 1)$.

Trial functions:

$$f = \alpha^T \mathbf{z}(1)$$

Will take $\alpha = \mathbf{1}$ for simplicity.

How do we choose the form of g ?

The correct form of g is given by (E, Ma and Wu, 2019):

$$\mathbf{g}(\tau, \mathbf{z}) = \mathbb{E}_{\mathbf{w} \sim \pi_\tau} \mathbf{a}(\mathbf{w}, \tau) \sigma(\mathbf{w}^T \mathbf{z})$$

where $\{\pi_\tau\}$ is a family of probability distributions.

$$\begin{aligned} \frac{d\mathbf{z}}{d\tau} &= \mathbb{E}_{\mathbf{w} \sim \pi_\tau} \mathbf{a}(\mathbf{w}, \tau) \sigma(\mathbf{w}^T \mathbf{z}) \\ &= \mathbb{E}_{(\mathbf{a}, \mathbf{w}) \sim \rho_\tau} \mathbf{a} \sigma(\mathbf{w}^T \mathbf{z}) \end{aligned}$$

Discretize: We obtain the “residual neural network” model:

$$\mathbf{z}_{l+1} = \mathbf{z}_l + \frac{1}{LM} \sum_{j=1}^M \mathbf{a}_{j,l} \sigma(\mathbf{z}_l^T \mathbf{w}_{j,l}), l = 1, 2, \dots, L - 1, \quad \mathbf{z}_0 = V \tilde{\mathbf{x}}$$
$$f_L(\mathbf{x}) = \mathbf{1}^T \mathbf{z}_L$$

Compositional law of large numbers

Consider the following compositional scheme:

$$\begin{aligned}z_{0,L}(\mathbf{x}) &= \mathbf{x}, \\z_{l+1,L}(\mathbf{x}) &= z_{l,L}(\mathbf{x}) + \frac{1}{LM} \sum_{k=1}^M \mathbf{a}_{l,k} \sigma(\mathbf{w}_{l,k}^T z_{l,L}(\mathbf{x})),\end{aligned}$$

$(\mathbf{a}_{l,k}, \mathbf{w}_{l,k})$ are pairs of vectors i.i.d. sampled from a distribution ρ .

Theorem (E, Ma and Wu 2019)

Assume that

$$\mathbb{E}_{\rho} \|\|\mathbf{a}\|\|\mathbf{w}^T\|\|_F^2 < \infty$$

where for a matrix or vector \mathbf{A} , $|\mathbf{A}|$ means taking element-wise absolute value for \mathbf{A} . Define \mathbf{z} by

$$\begin{aligned}z(\mathbf{x}, 0) &= \mathbf{x}, \\ \frac{d}{d\tau} z(\mathbf{x}, t) &= \mathbb{E}_{(\mathbf{a}, \mathbf{w}) \sim \rho} \mathbf{a} \sigma(\mathbf{w}^T z(\mathbf{x}, \tau)).\end{aligned}$$

Then we have

$$z_{L,L}(\mathbf{x}) \rightarrow z(\mathbf{x}, 1)$$

almost surely as $L \rightarrow +\infty$.

The optimal control problem

In a slightly more general form

$$\frac{d\mathbf{z}}{d\tau} = \mathbb{E}_{\mathbf{u} \sim \rho_\tau} \boldsymbol{\phi}(\mathbf{z}, \mathbf{u}), \quad \mathbf{z}(\mathbf{x}, 0) = \mathbf{x}$$

\mathbf{z} = state, ρ_τ = control at time τ .

The objective : Minimize \mathcal{R} over $\{\rho_\tau\}$

$$\mathcal{R}(\{\rho_\tau\}) = \mathbb{E}_{\mathbf{x} \sim \mu} (f(\mathbf{x}) - f^*(\mathbf{x}))^2 = \int_{\mathbb{R}^d} (f(\mathbf{x}) - f^*(\mathbf{x}))^2 d\mu$$

where

$$f(\mathbf{x}) = \mathbf{1}^T \mathbf{z}(\mathbf{x}, 1)$$

Pontryagin's maximum principle

Define the Hamiltonian $H : \mathbb{R}^d \times \mathbb{R}^d \times \mathcal{P}_2(\Omega) \rightarrow \mathbb{R}$ as

$$H(\mathbf{z}, \mathbf{p}, \mu) = \mathbb{E}_{\mathbf{u} \sim \mu}[\mathbf{p}^T \phi(\mathbf{z}, \mathbf{u})].$$

The solutions of the control problem must satisfy:

$$\rho_\tau = \operatorname{argmax}_\rho \mathbb{E}_{\mathbf{x}}[H(\mathbf{z}_\tau^{t,\mathbf{x}}, \mathbf{p}_\tau^{t,\mathbf{x}}, \rho)], \quad \forall \tau \in [0, 1],$$

and for each \mathbf{x} , $(\mathbf{z}_\tau^{t,\mathbf{x}}, \mathbf{p}_\tau^{t,\mathbf{x}})$ are defined by the forward/backward equations:

$$\begin{aligned} \frac{d\mathbf{z}_\tau^{t,\mathbf{x}}}{d\tau} &= \nabla_{\mathbf{p}} H = \mathbb{E}_{\mathbf{u} \sim \rho_\tau(\cdot; t)}[\phi(\mathbf{z}_\tau^{t,\mathbf{x}}, \mathbf{u})] \\ \frac{d\mathbf{p}_\tau^{t,\mathbf{x}}}{d\tau} &= -\nabla_{\mathbf{z}} H = \mathbb{E}_{\mathbf{u} \sim \rho_\tau(\cdot; t)}[\nabla_{\mathbf{z}}^T \phi(\mathbf{z}_\tau^{t,\mathbf{x}}, \mathbf{u}) \mathbf{p}_\tau^{t,\mathbf{x}}]. \end{aligned}$$

$$f(\mathbf{x}) = \mathbf{1}^T \mathbf{z}(\mathbf{x}, 1)$$

with the boundary conditions:

$$\begin{aligned} \mathbf{z}_0^{t,\mathbf{x}} &= \mathbf{x} \\ \mathbf{p}_1^{t,\mathbf{x}} &= 2(f(\mathbf{x}; \rho(\cdot; t)) - f^*(\mathbf{x}))\mathbf{1}. \end{aligned}$$

Gradient flow for flow-based models

Define the Hamiltonian $H : \mathbb{R}^d \times \mathbb{R}^d \times \mathcal{P}_2(\Omega) \mapsto \mathbb{R}$ as

$$H(\mathbf{z}, \mathbf{p}, \mu) = \mathbb{E}_{\mathbf{u} \sim \mu}[\mathbf{p}^T \phi(\mathbf{z}, \mathbf{u})].$$

The gradient flow for $\{\rho_\tau\}$ is given by

$$\partial_t \rho_\tau(\mathbf{u}, t) = \nabla \cdot (\rho_\tau(\mathbf{u}, t) \nabla V(\mathbf{u}; \rho)), \quad \forall \tau \in [0, 1],$$

where

$$V(\mathbf{u}; \rho) = \mathbb{E}_{\mathbf{x}} \left[\frac{\delta H}{\delta \rho} (\mathbf{z}_\tau^{t, \mathbf{x}}, \mathbf{p}_\tau^{t, \mathbf{x}}, \rho_\tau(\cdot; t)) \right],$$

and for each \mathbf{x} , $(\mathbf{z}_\tau^{t, \mathbf{x}}, \mathbf{p}_\tau^{t, \mathbf{x}})$ are defined by the forward/backward equations:

$$\begin{aligned} \frac{d\mathbf{z}_\tau^{t, \mathbf{x}}}{d\tau} &= \nabla_{\mathbf{p}} H = \mathbb{E}_{\mathbf{u} \sim \rho_\tau(\cdot; t)} [\phi(\mathbf{z}_\tau^{t, \mathbf{x}}, \mathbf{u})] \\ \frac{d\mathbf{p}_\tau^{t, \mathbf{x}}}{d\tau} &= -\nabla_{\mathbf{z}} H = \mathbb{E}_{\mathbf{u} \sim \rho_\tau(\cdot; t)} [\nabla_{\mathbf{z}}^T \phi(\mathbf{z}_\tau^{t, \mathbf{x}}, \mathbf{u}) \mathbf{p}_\tau^{t, \mathbf{x}}]. \end{aligned}$$

with the boundary conditions:

$$\begin{aligned} \mathbf{z}_0^{t, \mathbf{x}} &= \mathbf{x} \\ \mathbf{p}_1^{t, \mathbf{x}} &= 2(f(\mathbf{x}; \rho(\cdot; t)) - f^*(\mathbf{x})) \mathbf{1}. \end{aligned}$$

Discretize the gradient flow

- forward Euler for the flow in τ variable, step size $1/L$.
- particle method for the GD dynamics, M samples in each layer

$$\mathbf{z}_{l+1}^{t,\mathbf{x}} = \mathbf{z}_l^{t,\mathbf{x}} + \frac{1}{LM} \sum_{j=1}^M \phi(\mathbf{z}_l^{t,\mathbf{x}}, \mathbf{u}_l^j(t)), \quad l = 0, \dots, L-1$$

$$\mathbf{p}_l^{t,\mathbf{x}} = \mathbf{p}_{l+1}^{t,\mathbf{x}} + \frac{1}{LM} \sum_{j=1}^M \nabla_{\mathbf{z}} \phi(\mathbf{z}_{l+1}^{t,\mathbf{x}}, \mathbf{u}_{l+1}^j(t)) \mathbf{p}_{l+1}^{t,\mathbf{x}}, \quad l = 0, \dots, L-1$$

$$\frac{d\mathbf{u}_l^j(t)}{dt} = -\mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{w}}^T \phi(\mathbf{z}_l^{t,\mathbf{x}}, \mathbf{u}_l^j(t)) \mathbf{p}_l^{t,\mathbf{x}}].$$

This recovers the gradient descent algorithm (with back-propagation) for the ResNet:

$$\mathbf{z}_{l+1} = \mathbf{z}_l + \frac{1}{LM} \sum_{j=1}^M \phi(\mathbf{z}_l, \mathbf{u}_l).$$

Max principle-based training method

Qianxiao Li, Long Chen, Cheng Tai and Weinan E (2017):

Basic “method of successive approximation” (MSA):

Initialize: $\theta^0 \in \mathcal{U}$

For $k = 0, 1, 2, \dots$:

- Solve

$$\frac{d\mathbf{z}_\tau^k}{d\tau} = \nabla_{\mathbf{p}} H(\mathbf{z}_\tau^k, \mathbf{p}_\tau^k, \theta_\tau^k), \quad \mathbf{z}_0^k = V \mathbf{x}$$

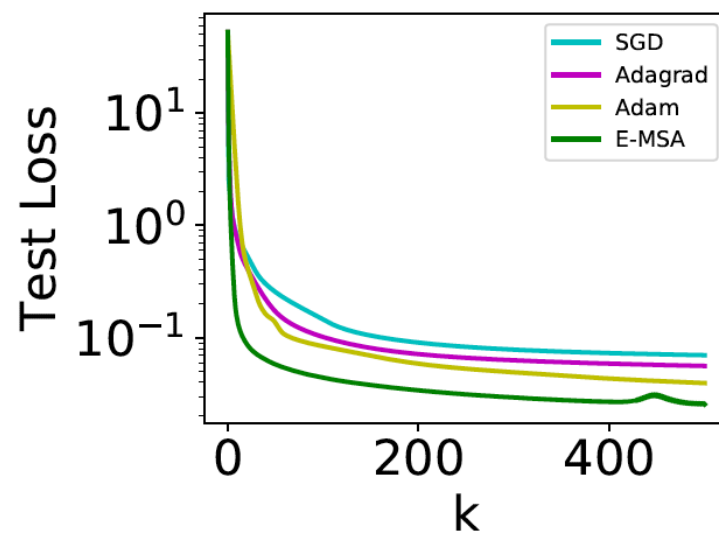
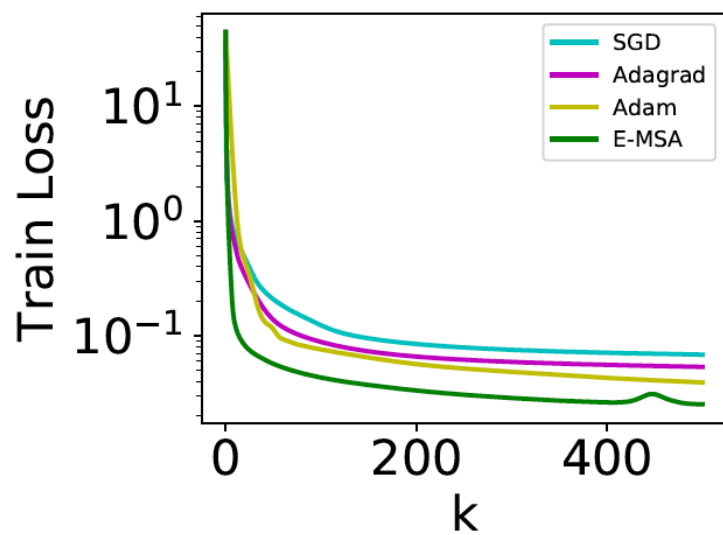
- Solve

$$\frac{d\mathbf{p}_\tau^k}{d\tau} = -\nabla_{\mathbf{z}} H(\mathbf{z}_\tau^k, \mathbf{p}_\tau^k, \theta_\tau^k), \quad \mathbf{p}_1^k = 2(f(\mathbf{x}; \theta^k) - f^*(\mathbf{x})) \mathbf{1}$$

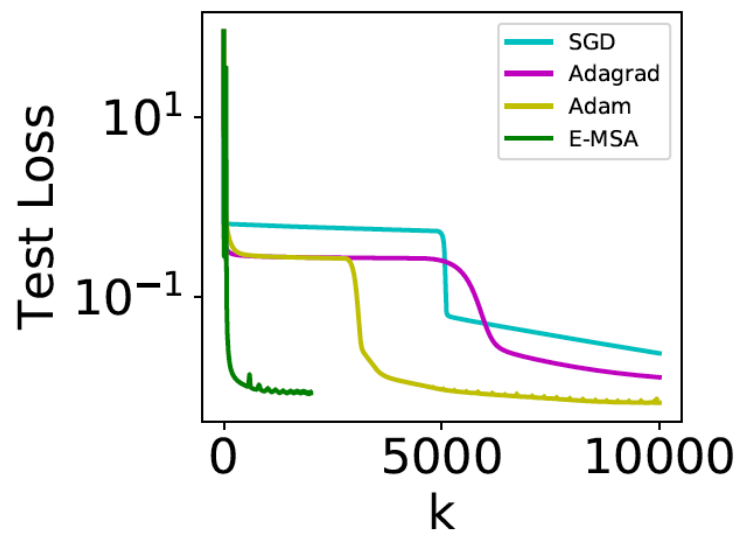
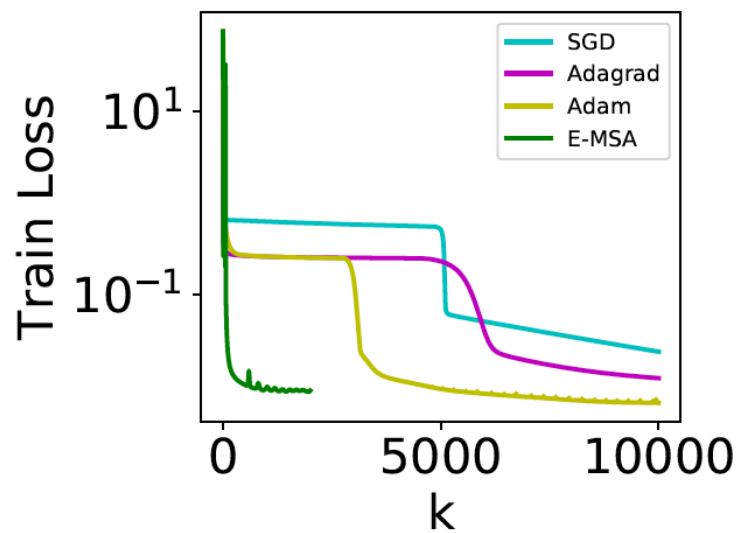
- Set $\theta_\tau^{k+1} = \operatorname{argmax}_{\theta \in \Theta} H(\mathbf{z}_\tau^k, \mathbf{p}_\tau^k, \theta)$, for each $\tau \in [0, 1]$

Extended MSA:

$$\tilde{H}(\mathbf{z}, \mathbf{p}, \theta, \mathbf{v}, \mathbf{q}) := H(\mathbf{z}, \mathbf{p}, \theta) - \frac{1}{2}\rho \|\mathbf{v} - f(\mathbf{z}, \theta)\|^2 - \frac{1}{2}\rho \|\mathbf{q} + \nabla_{\mathbf{z}} H(\mathbf{z}, \mathbf{p}, \theta)\|^2.$$



(a)



(b)

Comparison between GD and maximum principle

Maximum principle:

$$\rho_\tau = \operatorname{argmax}_\rho \mathbb{E}_x[H(\mathbf{z}_\tau^{t,x}, \mathbf{p}_\tau^{t,x}, \rho)], \quad \forall \tau \in [0, 1],$$

GD:

$$\partial_t \rho_\tau(\mathbf{u}, t) = \nabla \cdot \left(\rho_\tau(\mathbf{u}, t) \nabla \mathbb{E}_x \left[\frac{\delta H}{\delta \rho} (\mathbf{z}_\tau^{t,x}, \mathbf{p}_\tau^{t,x}, \rho_\tau(\mathbf{u}; t)) \right] \right), \quad \forall \tau \in [0, 1],$$

Hybrid:

Introducing a different time scale for optimization step: One time forward/backward propagation every k steps of optimization.

- $k = 1$, usual GD or SGD
- $k = \infty$, maximum principle

“Mean-field” vs. “continuous”

- They sometimes give rise to the same models (e.g. two-layer NNs).

They are DIFFERENT viewpoints.

- mean field: discrete \rightarrow continuous by taking the limit (more like interacting particles in stat phys)
- continuous formulation: continuous \rightarrow discrete by discretization (more like the usual numerical analysis situation)

“Continuous” formulation tries to formulate the “first principles” of ML.

It allows us to think “outside the box” about ML.

- alternative discretization (spectral)
- alternative model

Flow-based random feature model

$$\frac{dz}{d\tau} = \mathbb{E}_{\mathbf{w} \sim \pi_\tau} \mathbf{a}(\mathbf{w}, \tau) \sigma(\mathbf{w}^T \mathbf{z})$$

Random feature model: Fix $\{\pi_\tau\}$, optimize over $\{\mathbf{a}(\cdot, \tau)\}$.

- Continuous formulation tries to formulate the “first principles” of ML.
- It gives rise to new variational and PDE-like problems.
- The performance of these models/algorithms are more stable to the choice of hyper-parameters (no “phase transition”).
- It gives rise to new models and new algorithms, more suited to the specific application.

It allows us to think “outside the box” about ML.

The crucial point is the representation of functions as some form of expectations.

- integral transform-based
- flow-based

Other representations?

In a way, neural networks are very natural.