Neural Networks as Interacting Particle Systems



Eric Vanden-Eijnden Courant Institute

Advances in Computational Statistical Physics, CIRM 2018 Joint work with Grant Rotskoff Ref: arXiv:1805.00915

The unreasonable effectiveness of machine learning

- (Deep) neural networks have led to extraordinary progress in speech and image recognition, language processing and translation, object detection, etc.
- Problems assumed to be intractable a decade ago are now routine. Are image / speech recognition inherently low dimensional?
- Alternatively, can neural networks accurately represent high-dimensional data / functions?

 $f: \{0,1\}^{N^2} \to \{0,1,\ldots,9,\text{NaN}\}$

- Standard representations, like Galerkin truncations or finite element decompositions, are linear and cannot be scaled to high dimensional problem *curse of dimensionality.*
- In contrast, neural networks are nonlinear in their adjusting parameters their effective dimensionality is much higher than the number of these parameters.

Basic roadmap of neural network training

Approximate a target function *f* : Ω → ℝ defined on Ω ⊆ ℝ^d by a neural network representation, e.g. a single-layer network with sigmoid nonlinearity

$$f_n(x) = \sum_{i=1}^n c_i h(a_i \cdot x + b_i), \quad h(z) = 1/(1 + e^{-z})$$

non-linear approximation

• Measure the approximation error via the *loss function*

$$\ell(f, f_n) = \frac{1}{2} \int_{\Omega} |f(x) - f_n(x)|^2 d\mu(x) = \frac{1}{2} \mathbb{E}_{\text{data}} |f - f_n|^2$$

• In practice, estimate $\ell(f, f_n)$ via the *empirical loss function*

$$\ell_P(f, f_n) = rac{1}{P} \sum_{p=1}^{P} |f(x_p) - f_n(x_p)|^2, \quad \{x_p\}_{p=1}^{P} = \text{iid drawn from } \mu.$$

 $z(t + \Delta t) = z(t) - \Delta t \nabla_z L_P(z(t)),$

• Train the network via stochastic gradient descent (SGD) to minimize the loss over the parameters

non-convex optimization problem

where z denotes the network parameters (e.g. $z = (a_1, b_1, c_1, \dots, a_n, b_n, c_n)$) and $L_P(z) = \ell_P(f, f_n)$ the empirical loss viewed as function of z.

A test case: Spherical 3-spin model

$$f(\boldsymbol{x}) = \frac{1}{d} \sum_{p,q,r=1}^{d} a_{p,q,r} x_p x_q x_r, \qquad \boldsymbol{x} \in S^{d-1}\left(\sqrt{d}\right) \subset \mathbb{R}^d$$

 $a_{p,q,r} \sim N(0,1)$

Number of critical points exponential in *d* (Ben Arous, etc.)



 $x_i(\theta) = \sqrt{d}\cos(\theta), \qquad x_j(\theta) = \sqrt{d}\sin(\theta), \qquad x_k(\theta) = 0 \quad \forall k \neq i, j.$

Two grid points per dimension = $2^{25} = 33,554,432$ grid points

Representation power of neural networks

Universal Approximation Theorems (Barron, Cybenko, Park, others)

- Say that neural network representations are dense in the space of square-integrable target functions. *There is a neural network approximation arbitrarily close to any such function.*
- The theorems *do not* answer: How do we construct the representation?
 - 1. Can the network be trained (i.e. how should we get the parameters)?
 - 2. Do the typical machine learning algorithms converge?
 - 3. How does the error scale with the network size?

Neural networks as particle systems

• Using $f_n(x) = \frac{1}{n} \sum_{i=1}^n c_i \varphi(x, y_i)$, the loss function $\ell(f, f_n) = \frac{1}{2} \int_{\Omega} |f(x) - f_n(x)|^2 d\mu(x)$ can be expanded as

$$\ell(f, f_n) = C_f - \frac{1}{n} \sum_{i=1}^n c_i F(\boldsymbol{y}_i) + \frac{1}{2n^2} \sum_{i,j=1}^n c_i c_j K(\boldsymbol{y}_i, \boldsymbol{y}_j), \qquad C_f = \frac{1}{2} \int_{\Omega} |f(\boldsymbol{x})|^2 d\mu(\boldsymbol{x})$$

where

$$F(\boldsymbol{y}) = \int_{\Omega} f(\boldsymbol{x}) \varphi(\boldsymbol{x}, \boldsymbol{y}) d\mu(\boldsymbol{x}), \qquad K(\boldsymbol{y}, \boldsymbol{z}) = \int_{\Omega} \varphi(\boldsymbol{x}, \boldsymbol{y}) \varphi(\boldsymbol{x}, \boldsymbol{z}) d\mu(\boldsymbol{x}) \equiv K(\boldsymbol{z}, \boldsymbol{y}).$$

• Minimization of the loss over $\{c_i, y_i\}_{i=1}^n$ is a complex, presumably non-convex optimization problem.

Neural networks as particle systems

• Using $f_n(x) = \frac{1}{n} \sum_{i=1}^n c_i \varphi(x, y_i)$, the loss function $\ell(f, f_n) = \frac{1}{2} \int_{\Omega} |f(x) - f_n(x)|^2 d\mu(x)$ can be expanded as

$$\ell(f, f_n) = C_f - \frac{1}{n} \sum_{i=1}^n c_i F(\boldsymbol{y}_i) + \frac{1}{2n^2} \sum_{i,j=1}^n c_i c_j K(\boldsymbol{y}_i, \boldsymbol{y}_j), \qquad C_f = \frac{1}{2} \int_{\Omega} |f(\boldsymbol{x})|^2 d\mu(\boldsymbol{x})$$

interaction potential

where

$$F(y) = \int_{\Omega} f(x)\varphi(x,y)d\mu(x), \qquad K(y,z) = \int_{\Omega} \varphi(x,y)\varphi(x,z)d\mu(x) \equiv K(z,y).$$
one-body interaction two-body interaction

- Minimization of the loss over $\{c_i, y_i\}_{i=1}^n$ is a complex, presumably non-convex optimization problem.
- Parameters = particles; loss function = interaction potential

Non-equilibrium dynamics of training / optimization

• Gradient descent (GD) dynamics over loss function: $f_n(t, x)$

on:
$$f_n(t, \boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^n C_i(t) \varphi(\boldsymbol{x}, \boldsymbol{Y}_i(t))$$
 with

$$\frac{dY_i}{dt} = C_i \nabla F(Y_i) - \frac{1}{n} \sum_{j=1}^n C_i C_j \nabla K(Y_i, Y_j),$$

$$\frac{dC_i}{dt} = F(Y_i) - \frac{1}{n} \sum_{j=1}^n C_j K(Y_i, Y_j)$$
 (to be generalized to SGD)

• Empirical distribution:

$$ho_n(t, \boldsymbol{y}, c) = rac{1}{n} \sum_{i=1}^n \delta(c - C_i(t)) \delta(\boldsymbol{y} - \boldsymbol{Y}_i(t)), \qquad f_n(t, \boldsymbol{x}) = \int_{D imes R} c \varphi(\boldsymbol{x}, \boldsymbol{y})
ho_n(t, \boldsymbol{y}, c) d\boldsymbol{y} dc$$

• Dynamics = McKean-Vlasov equation

$$\partial_t \rho_n = \nabla \cdot \left(-c \nabla F \rho_n + \int_{D \times \mathbb{R}} cc' \nabla K(\boldsymbol{y}, \boldsymbol{y}') \rho'_n \rho_n d\boldsymbol{y}' dc' \right) \\ + \partial_c \left(-F \rho_n + \int_{D \times \mathbb{R}} c' K(\boldsymbol{y}, \boldsymbol{y}') \rho'_n \rho_n d\boldsymbol{y}' dc' \right)$$

Asymptotic convexity in the mean field limit

- McKean-Vlasov equation = gradient descent flow in Wasserstein metric (Otto, Villani, Serfaty, etc.)
- Mean-field limit: $\rho_n(t) \rightharpoonup \rho_0(t)$ as $n \to \infty$, where $\rho_0(t)$ descents on quadratic energy:

$$\partial_t \rho_0 = \nabla \cdot \left(\rho_0 \nabla \frac{\delta \mathcal{E}}{\delta \rho_0} \right) + \partial_c \left(\rho_0 \partial_c \frac{\delta \mathcal{E}}{\delta \rho_0} \right)$$
 Propagation of chaos

$$\begin{split} \mathcal{E}[\rho_0] &= C_f - \int_{D \times \mathbb{R}} cF \rho_0 dy dc + \frac{1}{2} \int_{(D \times \mathbb{R})^2} cc' K(y, y') \rho_0 \rho'_0 dy \\ &= \frac{1}{2} \int_{\Omega} \left(f(x) - \frac{1}{2} \int_{D \times \mathbb{R}} c\varphi(x, y) \rho_0 dy dc \right)^2 d\mu(x) \end{split}$$

 Ruggedness of the microscopic landscape viewed by particles / parameters disappears at the level of their empirical distribution.

Asymptotic convexity

Not the full story: particular shape of interaction potential from network architecture matters to avoid stationary points

Similar results appeared recently in:

Mei, Montanari, & Nguyen arXiv:1804.06561; Sirigano & Spiliopoulos arXiv:1805.01053

Law of Large Numbers (LLN)

• $\rho_n(t) \rightharpoonup \rho_0(t)$ as $n \to \infty$, where $f_0(t, x) = \int_{D \times \mathbb{R}} c\varphi(x, y) \rho_0(t, y, c) dy dc$ satisfies gradient flow over loss function:

$$\partial_t f_0(t, x) = -\int_{\Omega} M([\rho_0(t)], x, x') \left(f_0(t, x') - f(x') \right) d\mu(x') = -\int_{\Omega} M([\rho_0(t)], x, x') D_{f_0(t, x')} \ell(f, f_0(t)) d\mu(x')$$
(1)

where $D_{f(x)}$ denotes the gradient with respect to f(x) in the $L^2(\Omega, \mu)$ -norm and

$$M([\rho], \boldsymbol{x}, \boldsymbol{x}') = \int_{D \times \mathbb{R}} \left(c^2 \nabla_{\boldsymbol{y}} \varphi(\boldsymbol{x}, \boldsymbol{y}) \nabla_{\boldsymbol{y}} \varphi(\boldsymbol{x}', \boldsymbol{y}) + \varphi(\boldsymbol{x}, \boldsymbol{y}) \varphi(\boldsymbol{x}', \boldsymbol{y}) \right) \rho(\boldsymbol{y}, c) d\boldsymbol{y} dc.$$

Prop 1 (LLN) Let
$$f_n(t) = \frac{1}{n} \sum_{i=1}^n C_i(t) \varphi(\cdot, Y_i(t))$$
. Then

$$\lim_{n \to \infty} f_n(t) = f_0(t) \quad \text{almost surely}$$
where $f_0(t)$ solves (1) and satisfies

$$\lim_{t \to \infty} f_0(t) = f \quad \text{a.e. in } \Omega \quad \text{if} \quad \lim_{t \to \infty} \int_{\mathbb{R}} \rho_0(\cdot, c) dc > 0 \quad \text{a.e. in } D$$

• Dynamical variant of Universal Representation Theorem which also indicates how to realize it in practice.

• Condition $\int_{\mathbb{R}} \rho_0(\cdot, c) dc > 0$ a.e. *D* generic under small perturbations.

Error scaling

parameters drawn independently at t=0

Prop 2 (CLT) Let $f_n(t) = \frac{1}{n} \sum_{i=1}^n C_i(t) \varphi(\cdot, Y_i(t))$ with well-prepared initial conditions. Then for any $\overline{\xi} < 1$ and any $a_n > 0$ such that $a_n / \log n \to \infty$ as $n \to \infty$, we have $\lim_{n \to \infty} n^{\overline{\xi}} (f_n(a_n) - f) = 0 \quad \text{almost surely} \quad \text{if} \quad \lim_{t \to \infty} \int_{\mathbb{R}} \rho_0(\cdot, c) dc > 0 \quad \text{a.e. in } D$

Healing of errors by training:

- $\xi = \frac{1}{2}$ initially CLT scaling of *iid* initial conditions;
- $\xi = 1$ after optimization.

Error scaling — Central Limit Theorem (CLT)

$$d\mathbf{Y}_{i} = C_{i}\nabla F(\mathbf{Y}_{i})dt - \frac{1}{n}\sum_{j=1}^{n}C_{i}C_{j}\nabla K(\mathbf{Y}_{i},\mathbf{Y}_{j})dt + (\beta n)^{-1}\nabla \log m(\mathbf{Y}_{i},C_{i})dt + \sqrt{2(\beta n)^{-1}}d\mathbf{W}_{i},$$

$$dC_{i} = F(\mathbf{Y}_{i})dt - \frac{1}{n}\sum_{j=1}^{n}C_{j}K(\mathbf{Y}_{i},\mathbf{Y}_{j})dt + (\beta n)^{-1}\partial_{c}\log m(\mathbf{Y}_{i},C_{i})dt + \sqrt{2(\beta n)^{-1}}dW_{i}'$$

$$low-temperature regime$$

 $\mathcal{E}_{n}[\rho] = \frac{1}{2} \int_{\Omega} \left(f(\boldsymbol{x}) - \frac{1}{2} \int_{D \times \mathbb{R}} c\varphi(\boldsymbol{x}, \boldsymbol{y}) \rho_{0} d\boldsymbol{y} dc \right)^{2} d\mu(\boldsymbol{x}) + (\beta n)^{-1} \int_{D \times \mathbb{R}} \rho \log(\rho/m) d\boldsymbol{y} dc$

entropic correction

Prop 3 (CLT at finite temperature) Let $f_n(t) = \frac{1}{n} \sum_{i=1}^n C_i(t) \varphi(\cdot, Y_i(t))$ with well-prepared initial condition at t = T. Then

$$\lim_{n \to \infty} \lim_{T \to -\infty} n \left(f_n(t) - f \right) = f_1(t) \quad \text{in law}$$

where $f_1(t)$ is a Gaussian process such that: for any $\chi \in L^2(\Omega, \mu)$,

$$\mathbb{E} \int_{\Omega} \chi(x) f_1(t,x) d\mu(x) = eta^{-1} \int_{\Omega} \chi(x) \epsilon^*(x) d\mu(x) \ \mathbb{E} \left(\int_{\Omega} \chi(x) \left(f_1(t,x) - eta^{-1} \epsilon^*(x)
ight) d\mu(x)
ight)^2 = eta^{-1} \int_{\Omega} |\chi(x)|^2 d\mu(x)$$

Discrete training set and stochastic gradient descent

• Loss function approximated by *empirical loss*

$$\ell_P(f, f_n) = rac{1}{P} \sum_{p=1}^{P} |f(x_p) - f_n(x_p)|^2, \quad \{x_p\}_{p=1}^{P} = \text{iid drawn from } \mu.$$

• Stochastic gradient descent (SGD):

$$Z(t + \Delta t) = Z(t) - \nabla_z L_P(Z(t)) \Delta t$$

where $\Delta t > 0$ is some time-step and $L_P(z) = n\ell_P(f, f_n)$ is the empirical loss viewed as a function of the parameters $z = (c_1, y_1, \dots, c_n, y_n)$,

• If training set $\{x_p\}_{p=1}^P$ is redrawn from μ independently at every time step t, leads to an *effective SDE*:

$$dZ = -\nabla_z L(Z)dt + \sqrt{\theta} dB$$

where $L(z) = n\ell(f, f_n(z))$, and the quadaric variation of the noise is the covariance of the empirical loss

$$\langle d\boldsymbol{B}, d\boldsymbol{B} \rangle = \mathbb{E} \left(\nabla_{\boldsymbol{z}} (L_{P=1}(\boldsymbol{z}) - n\ell(f, f_n(\boldsymbol{z}))) \right)^{\otimes 2} dt$$

and

$$\theta = \Delta t / P =$$
time step / batch size

LLN and CLT for SGD

• Setting $\theta = \Delta t/P = an^{-2\alpha}$ with a > 0, $\alpha > 0$ (e.g. $P = O(n^{2\alpha})$), the empirical distribution satisfies

$$\partial_{t}\rho_{n} = \nabla \cdot \left(-c\nabla F\rho_{n} + \int_{D\times\mathbb{R}} cc'\nabla K(\boldsymbol{y},\boldsymbol{y}')\rho_{n}'\rho_{n}d\boldsymbol{y}'dc' \right) + \partial_{c} \left(-F\rho_{n} + \int_{D\times\mathbb{R}} c'K(\boldsymbol{y},\boldsymbol{y}')\rho_{n}'\rho_{n}d\boldsymbol{y}'dc' \right) \\ + \frac{1}{2}an^{-2\alpha}\nabla\nabla : \left(\rho_{n}c^{2}A_{2}([f_{n}(t) - f], \boldsymbol{y}, \boldsymbol{y}) \right) + \frac{1}{2}an^{-2\alpha}\partial_{c}^{2} \left(\rho_{n}A_{0}([f_{n}(t) - f], \boldsymbol{y}, \boldsymbol{y}) \right) \\ + an^{-2\alpha}\partial_{c}\nabla \cdot \left(\rho_{n}cA_{1}([f_{n}(t) - f], \boldsymbol{y}, \boldsymbol{y}) \right) \\ + \sqrt{an^{-\alpha}}\dot{\eta}_{n}(t, \boldsymbol{y}, c)$$
 additional terms higher order; noise term dominates those.

Prop 7 (LLN & CLT for SGD) Let $f_n(t) = \frac{1}{n} \sum_{i=1}^n C_i(t) \varphi(\cdot, Y_i(t))$ with $\{Y_i(t), C_i(t)\}_{i=1}^n$ solution to limiting SDE with $\theta = an^{-2\alpha}$, $a > 0 \ \alpha \in (0, 1)$. Then for any $a_n > 0$ such that $a_n / \log n \to \infty$ as $n \to \infty$, we have $\lim_{n \to \infty} n^{\alpha} (f_n(a_n) - f_0(a_n)) = 0 \quad \text{almost surely}$ where $f_0(t) \to f$ as $t \to \infty$

- if $\alpha < 1$, training set error dominates
- if $\alpha = 1$, training set error and discretization error are in the same scale;
- if $\alpha > 1$, no gain (and convergence in time may be impeded)

3-spin model on the high-dimensional sphere

• Spherical 3-spin model: $f : S^{d-1}(\sqrt{d}) \to \mathbb{R}$ given by

$$f(\boldsymbol{x}) = \frac{1}{d} \sum_{p,q,r=1}^{d} a_{p,q,r} x_p x_q x_r, \qquad \boldsymbol{x} \in S^{d-1}(\sqrt{d}) \subset \mathbb{R}^d$$

where the coefficients $\{a_{p,q,r}\}_{p,q,r=1}^d$ are independent Gaussian random variables with mean zero and variance one.

- Complex function with a number of critical points that grows exponentially with the dimensionality d.
- Previous works drew a parallel between the glassy 3-spin function and generic loss functions.
- In contrast, we use the 3-spin function as a difficult target for approximation by neural networks, that is:

▷ we train networks to learn f with a particular realization of $a_{p,q,r}$, and;

 \triangleright we study the accuracy of that representation as a function of the number of particles n.



At time 2E6, the batch size is increased to initiate an additional quench.

Error scaling for single layer neural network with sigmoid nonlinearities. Dashed line are $\propto n^{-1}$

Conclusions

- Neural networks are a potentially powerful tool for computational physics and applied mathematics. They can massively reduce the cost of representing functions in high dimensional spaces.
- Viewing the parameters as interacting particles, we can demonstrate that the loss landscape is asymptotically convex and stochastic gradient descent converges to an energy minimizer dynamical generalization of Universal Representation Theorems.
- The approximation error can be identified up to a constant, which shows that its scaling is universal.
- This offers exciting possibilities in scientific computing (free energy methods, quantum variational energy calculations, PDE solving, etc.) that are only beginning to be explored.



QUAND ON NE SAIT PAS DU L'ON VA, IL FAUT Y ALLER !!... ... ET LE PLUS VITE POSSIBLE. "That's another thing we've learned from your Nation," said Mein Herr, "map-making. But we've carried it much further than you. What do you consider the largest map that would be really useful?"

"About six inches to the mile."

"Only six inches!" exclaimed Mein Herr. "We very soon got to six yards to the mile. Then we tried a hundred yards to the mile. And then came the grandest idea of all! We actually made a map of the country, on the scale of a mile to the mile!"

"Have you used it much?" I enquired.

"It has never been spread out, yet," said Mein Herr: "the farmers objected: they said it would cover the whole country, and shut out the sunlight! So we now use the country itself, as its own map, and I assure you it does nearly as well.

Lewis Carroll, Sylvie and Bruno, 1889 - 1893.

Functional formulation in the limit of large n

• Assume that the set $\{(c_i, y_i)\}_{i=1}^n$ is such that for some signed density G on D we have

$$f_n = \frac{1}{n} \sum_{i=1}^n c_i \varphi(\cdot, \boldsymbol{y}_i) \to \tilde{f} = \int_D \varphi(\cdot, \boldsymbol{y}) G(\boldsymbol{y}) d\boldsymbol{y} \quad \text{as} \quad n \to \infty$$

• Make the loss function quadratic in G: $\ell(f, \tilde{f}) = C_f - \int_D F(y)G(y)dy + \frac{1}{2}\int_{D \times D} K(y, z)G(y)G(z)dydz$

Universal Representation Theorem (Barron, Cybenko, Park,...) *If the kernel* φ *is discriminating, then given any* $f \in L^2(\Omega, \mu)$ *and* $\epsilon > 0$ *:*

$$\exists f^* \text{ such that } \|f - f^*\|_{L^2(\Omega,\mu)} \leq \epsilon$$

and f^* can be represented as

$$\int_D arphi(\cdot,oldsymbol{y}) G^*(oldsymbol{y}) doldsymbol{y} = f^*$$
 a.e. in Ω

where G^* solves

$$F^*(\boldsymbol{y}) = \int_D K(\boldsymbol{y}, \boldsymbol{z}) G^*(\boldsymbol{z}) d\boldsymbol{z}$$
 with $F^*(\boldsymbol{y}) = \int_\Omega f^*(\boldsymbol{x}) \varphi(\boldsymbol{x}, \boldsymbol{y}) d\mu(\boldsymbol{x})$

The function f^* can also be realized as $f^* = \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} c_j \varphi(\cdot, y_j)$ for some choice of $\{y_i, c_i\}_{i \in \mathbb{N}}$.

Stationary points of McKean-Vlasov equation

The stationary points of the McKean-Vlasov equation satify

$$-F(\boldsymbol{y}) + \int_{D_0} K(\boldsymbol{y}, \boldsymbol{y}') d\gamma_0(\boldsymbol{y}') = 0$$
 for ν_0 -almost all $\boldsymbol{y} \in D_0 = \operatorname{supp} \nu_0$

where

$$\lim_{t\to\infty}\int_{D\times\mathbb{R}}\chi(\boldsymbol{y})\rho_0(t,\cdot,c)d\boldsymbol{y}dc = \int_D\chi(\boldsymbol{y})d\nu_0(\boldsymbol{y}) \qquad \lim_{t\to\infty}\int_{D\times\mathbb{R}}\chi(\boldsymbol{y})c\rho_0(t,\cdot,c)d\boldsymbol{y}dc = \int_D\chi(\boldsymbol{y})d\gamma_0(\boldsymbol{y})$$

• D_0 may be a (singular) subset of D; however, if $\varphi(\cdot, y)$ is discriminating in D_0 , i.e. if

$$\int_{\Omega} g(x)\varphi(x,\cdot)d\mu(x) = 0 \quad \text{a.e. in } D_0 \quad \Rightarrow \quad g = 0 \quad \text{a.e. in } \Omega$$
(2)

then

Universal Representation Thm

$$\int_{D_0} \varphi(\cdot, \boldsymbol{y}) d\gamma_0(\boldsymbol{y}) = f \quad \text{a.e. in } \Omega.$$
(3)

Even if φ(·, y) is not discriminating in D₀, we can continously deform D₀ into a set such that (2) and hence (3) hold. As a result these equations always hold in the presence of noise (as in SGD).