

# Machine Learning and Computational Mathematics

Weinan E<sup>i</sup>

Princeton University and Beijing Institute of Big Data Research

*In memory of Professor Feng Kang (1920-1993)*

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Machine learning-based algorithms for problems in computational science</b>	<b>3</b>
2.1	Nonlinear multi-grid method and protein folding . . . . .	4
2.2	Molecular dynamics with <i>ab initio</i> accuracy . . . . .	6
<b>3</b>	<b>Machine learning-based algorithms for high dimensional problems in scientific computing</b>	<b>9</b>
3.1	Stochastic control . . . . .	9
3.2	Nonlinear parabolic PDEs . . . . .	11
3.3	Moment closure for kinetic equations modeling gas dynamics . . . . .	13
<b>4</b>	<b>Mathematical theory of machine learning</b>	<b>15</b>
4.1	An introduction to neural network-based supervised learning . . . . .	15
4.2	Approximation theory . . . . .	17
4.3	Estimation error . . . . .	19
4.4	A priori estimates for regularized models . . . . .	21
<b>5</b>	<b>Machine learning from a continuous viewpoint</b>	<b>22</b>
5.1	Representation of functions . . . . .	22
5.2	The stochastic optimization problem . . . . .	23
5.3	Optimization: Gradient flows . . . . .	24
5.4	Discretizing the gradient flows . . . . .	24
5.5	The optimal control problem for flow-based representation . . . . .	25
<b>6</b>	<b>Concluding remarks</b>	<b>28</b>

## Abstract

Neural network-based machine learning is capable of approximating functions in very high dimension with unprecedented efficiency and accuracy. This has opened up many exciting new possibilities, not just in traditional areas of artificial intelligence, but also in scientific computing and computational science. At the same time, machine learning has also acquired the reputation of being a set of “black box”

---

<sup>i</sup>weinan@math.princeton.edu

type of tricks, without fundamental principles. This has been a real obstacle for making further progress in machine learning.

In this article, we try to address the following two very important questions: (1) How machine learning has already impacted and will further impact computational mathematics, scientific computing and computational science? (2) How computational mathematics, particularly numerical analysis, can impact machine learning? We describe some of the most important progress that has been made on these issues. Our hope is to put things into a perspective that will help to integrate machine learning with computational mathematics.

## 1 Introduction

Neural network-based machine learning (ML) has shown very impressive success on a variety of tasks in traditional artificial intelligence. This includes classifying images, generating new images such as (fake) human faces and playing sophisticated games such as Go. A common feature of all these tasks is that they involve objects in very high dimension. Indeed when formulated in mathematical terms, the image classification problem is a problem of approximating a high dimensional function, defined on the space of images, to the discrete set of values corresponding to the category of each image. The dimensionality of the input space is typically 3 times the number of pixels in the image, where 3 is the dimensionality of the color space. The image generation problem is a problem of generating samples from an unknown high dimensional distribution, given a set of samples from that distribution. The Go game problem is about solving a Bellman-like equation in dynamic programming, since the optimal strategy satisfies such an equation. For sophisticated games such as Go, this Bellman-like equation is formulated on a huge space.

All these are made possible by the ability to accurately approximate high dimensional functions, using modern machine learning techniques. This opens up new possibilities for attacking problems that suffer from the “curse of dimensionality” (CoD): As dimensionality grows, computational cost grows exponentially fast. This CoD problem has been an essential obstacle for the scientific community for a very long time.

Take, for example, the problem of solving partial differential equations (PDEs) numerically. With traditional numerical methods such as finite difference, finite element and spectral methods, we can now routinely solve PDEs in three spatial dimensions plus the temporal dimension. Most of the PDEs currently studied in computational mathematics belong to this category. Well known examples include the Poisson equation, the Maxwell equation, the Euler equation, the Navier-Stokes equations, and the PDEs for linear elasticity. Sparse grids can increase our ability to handling PDEs to, say 8 to 10 dimensions. This allows us to try solving problems such as the Boltzmann equation for simple molecules. But we are totally lost when faced with PDEs, say in 100 dimension. This makes it essentially impossible to solve Fokker-Planck or Boltzmann equations for complex molecules, many-body Schrödinger, or the Hamilton-Jacobi-Bellman equations for realistic control problems.

This is exactly where machine learning can help. Indeed, starting with the work in [20, 10, 21], machine learning-based numerical algorithm for solving high dimensional

PDEs and control problems has been one of the most exciting new developments in recent years in scientific computing, and this has opened up a host of new possibilities for computational mathematics. We refer to [17] for a review of this exciting development.

Solving PDEs is just the tip of the iceberg. There are many other problems for which CoD is the main obstacle, including:

- classical many-body problem, e.g. protein folding
- turbulence. Even though turbulence can be modeled by the three dimensional Navier-Stokes equation, it has so many active degrees of freedom that an effective model for turbulence should involve many variables.
- solid mechanics. In solid mechanics, we do not even have the analog of the Navier-Stokes equation. Why is this the case? Well, the real reason is that the behavior of solids is essentially a multi-scale problem that involves scales from atomistic all the way to macroscopic.
- multi-scale modeling. In fact most multi-scale problems for which there is no separation of scales belong to this category. An immediate example is the dynamics of polymer fluids or polymer melts.

Can machine learning help for these problems? More generally, *can we extend the success of machine learning beyond traditional AI?* We will try to convince the reader that this is indeed the case for many problems.

Besides being extremely powerful, neural network-based machine learning has also got the reputation of being a set of tricks instead of a set of systematic scientific principles. Its performance depends sensitively on the value of the hyper-parameters, such as the network widths and depths, the initialization, the learning rates, etc. Indeed just a few years ago, parameter tuning was considered to be very much of an art. Even now, This is still the case for some tasks. Therefore a natural question is: Can we understand these subtleties and propose better machine learning models whose performance is more robust?

In this article, we review what has been learned on these two issues. We discuss the impact that machine learning has already made or will make on computational mathematics, and how the ideas from computational mathematics, particularly numerical analysis, can be used to help understanding and better formulating machine learning models. On the former, we will mainly discuss the new problems that can now be addressed using ML-based algorithms. Even though machine learning also suggests new ways to solve some traditional problems in computational mathematics, we will not say much on this front.

## 2 Machine learning-based algorithms for problems in computational science

In this and the next section, we will discuss how neural network models can be used to develop new algorithms. For readers who are not familiar with neural networks, just

think of them as being some replacement of polynomials. We will discuss neural networks afterwards.

## 2.1 Nonlinear multi-grid method and protein folding

In traditional multi-grid method [5], say for solving the linear systems of equation that arise from some finite difference or finite element discretization of a linear elliptic equation, our objective is to minimize a quadratic function like

$$I_h(\mathbf{u}_h) = \frac{1}{2} \mathbf{u}_h^T L_h \mathbf{u}_h - \mathbf{f}_h^T \mathbf{u}_h$$

Here  $h$  is the grid size of the discretization. The basic idea of the multi-grid method is to iterate between solving this problem and a reduced problem on a coarser grid with grid size  $H$ . In order to do this, we need the following

- a projection operator:  $P : \mathbf{u}_h \rightarrow \mathbf{u}_H$ , that maps functions defined on the fine grid to functions defined on the coarse grid.
- the effective operator at scale  $H$ :  $L_H = P^T L_h P$ . This defines the objective function on the coarse grid:

$$I_H(\mathbf{u}_H) = \frac{1}{2} \mathbf{u}_H^T L_H \mathbf{u}_H - \mathbf{f}_H^T \mathbf{u}_H$$

- a prolongation operator  $Q : \mathbf{u}_H \rightarrow \mathbf{u}_h$ , that maps functions defined on the coarse grid to functions defined on the fine grid. Usually one can take  $Q$  to be  $P^T$ .

The key idea here is coarse graining, and iterating between the fine scale and the coarse-grained problem. The main components in coarse graining is a set of coarse-grained variables and the effective coarse-grained problem. Formulated this way, these are obviously general ideas that can be relevant for a wide variety of problems. In practice, however, the difficulty lies in how to obtain the effective coarse-grained problem, a step that is trivial for linear problems, and this is where machine learning can help.

We are going to use the protein folding problem as an example to illustrate the general idea for nonlinear problems.

Let  $\{\mathbf{x}_j\}$  be the positions of the atoms in a protein and the surrounding solvent, and  $U = U(\{\mathbf{x}_j\})$  be the potential energy of the combined protein-solvent system. The potential energy consists of the energies due to chemical bonding, Van der Waals interaction, electro-static interaction, etc. The protein folding problem is to find the “ground state” of the energy  $U$ :

$$\text{“Minimize” } U.$$

Here we have added quotation marks since really what we want to do is to sample the distribution

$$\rho_\beta = \frac{1}{Z} e^{-\beta U}, \beta = (k_B T)^{-1}$$

To define the coarse-grained problem, we assume that we are given a set of collective variables:  $\mathbf{s} = (s_1, \dots, s_n)$ ,  $s_j = s_j(\mathbf{x}_1, \dots, \mathbf{x}_N)$ , ( $n < N$ ). One possibility is to use the

dihedral angles as the coarse-grained variables. In principle, one may also use machine learning methods to learn the “best” set of coarse-grained variables but this direction will not be pursued here.

Having defined the coarse-grained variables, the effective coarse-grained problem is simply the free energy associated with this set of coarse-grained variables:

$$A(\mathbf{s}) = -\frac{1}{\beta} \ln p(\mathbf{s}), \quad p_{\beta}(\mathbf{s}) = \frac{1}{Z} \int e^{-\beta U(\mathbf{x})} \delta(\mathbf{s}(\mathbf{x}) - \mathbf{s}) d\mathbf{x},$$

Unlike the case for linear problems, for which the effective coarse-grained model is readily available, in the current situation, we have to find the function  $A$  first.

The idea is to approximate  $A$  by neural networks. The issue here is how to obtain the training data.

Contrary to most standard machine learning problems where the training data is collected beforehand, in applications to computational science and scientific computing, the training data is collected “on-the-fly” as learning proceeds. This is referred to as the “concurrent learning” protocol [11]. In this regard, the standard machine learning problems for which the training data is collected beforehand are examples of “sequential learning”. The key issue for concurrent learning is an efficient algorithm for generating the data in the best way. The training dataset should on one hand be representative enough and on the other hand be as small as possible.

A general procedure for generating such datasets is suggested in [11]. It is called the EELT (exploration- examination-labeling-training) algorithm and it consists of the following steps:

- exploration: exploring the  $\mathbf{s}$  space. This can be done by sampling  $\frac{1}{Z} e^{-\beta A(\mathbf{s})}$  with the current approximation of  $A$ .
- examination: for each state explored, decide whether that state should be labeled. One way to do this is to use an *a posteriori error estimator*. One possible such a posteriori error estimator is the variance of the predictions of an ensemble of machine learning models, see [41].
- labeling: compute the mean force (say using restrained molecular dynamics)

$$\mathbf{F}(\mathbf{s}) = -\nabla_{\mathbf{s}} A(\mathbf{s}).$$

from which the free energy  $A$  can be computed using standard thermodynamic integration.

- training: train the appropriate neural network model. To come up with a good neural network model, one has to take into account the symmetries in the problem. For example, if we coarse grain a full atom representation of a collection of water molecules by eliminating the positions of the hydrogen atoms, then the free energy function for the resulting system should have permutation symmetry and this should be taken into account when designing the neural network model (see the next subsection).

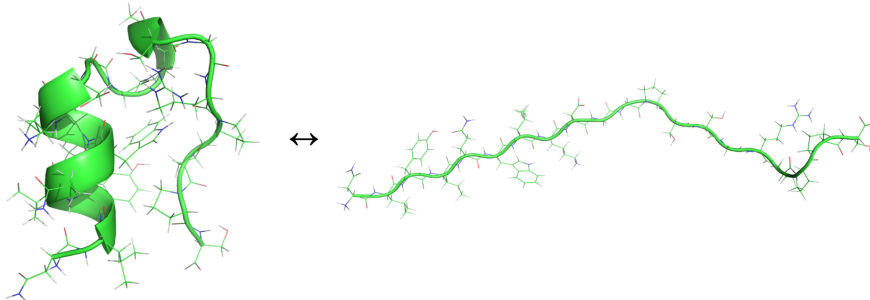


Figure 1: The folded and extended states of Trp-cage, reproduced with permission from [37]

This can also be viewed as a nonlinear multi-grid algorithm in the sense that it iterates between sampling  $p_\beta$  on the space of the coarse-grained variables and the (constrained) Gibbs distribution  $\rho_\beta$  for the full atom description.

This is a general procedure that should work for a large class of nonlinear “multi-grid” problems.

Shown in Figure 1 is the extended and folded structure of Trp-cage. This is a small protein with 20 amino acids. We have chosen the 38 dihedral angles as the collective variables. The full result is presented in [37].

## 2.2 Molecular dynamics with *ab initio* accuracy

Molecular dynamics is a way of studying the behavior of molecular and material systems by tracking the trajectories of all the nuclei in the system. The dynamics of the nuclei is assumed to obey Newton’s law, with some potential energy function (typically called potential energy surface or PES)  $V$  that models the effective interaction between the nuclei:

$$m_i \frac{d^2 \mathbf{x}_i}{dt^2} = -\nabla_{\mathbf{x}_i} V, \quad V = V(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N),$$

How can we get the function  $V$ ? Traditionally, there have been two rather different approaches. The first is to compute the inter-atomic forces ( $-\nabla V$ ) on the fly using quantum mechanics models, the most popular one being the density functional theory (DFT). This is known as the Car-Parrinello molecular dynamics or *ab initio* molecular dynamics [6, 8]. This approach is quite accurate but is also very expensive, limiting the size of the system that one can handle to about 1000 atoms, even with high performance supercomputers. The other approach is to come up with empirical potentials. Basically one guesses a functional form of  $V$  with a small set of fitting parameters which are then determined by a small set of data. This approach is very efficient but unreliable. This dilemma between accuracy and efficiency has been an essential road block for molecular dynamics for a long time.

With machine learning, the new paradigm is to use DFT to generate the data, and then use machine learning to generate an approximation to  $V$ . This approach has the

Systems			Al		Mg		Al-Mg alloy	
Type	Lattice	#atom	#Confs	#Data	#Confs	#Data	#Confs	#Data
Bulk	FCC	32	15,174,000	1,326	15,174,000	860	39,266,460	7,313
	HCP	16	15,174,000	908	15,174,000	760	18,999,900	2,461
	Diamond	16	5,058,000	1,026	5,058,000	543	5,451,300	2,607
	SC	8	5,058,000	713	5,058,000	234	2,543,940	667
Surface	FCC (100)	12	3,270,960	728	3,270,960	251	62,203,680	1,131
	FCC (110)	16 <sup>a</sup> ,20 <sup>b</sup>	3,270,960	838	3,270,960	353	10,744,2720	2,435
	FCC (111)	12	3,270,960	544	3,270,960	230	62,203,680	1,160
	HCP (0001)	12	3,270,960	39	3,270,960	109	62,203,680	176
	HCP (1010)	12	3,270,960	74	3,270,960	167	62,203,680	203
	HCP (1120)	16 <sup>a</sup> ,20 <sup>b</sup>	3,270,960	293	3,270,960	182	107,442,720	501
sum			60,089,760	6,489	60,089,760	3,689	529,961,760	18,654

<sup>a</sup>Pure Al

<sup>b</sup>Mg and Al-Mg alloy

Figure 2: The results of EELT algorithm: Number of configurations explored vs. the number of data points labeled. Only a very small percentage of the configurations are labeled. Reproduced with permission from Linfeng Zhang. See also [40].

potential to produce an approximation to  $V$  that is as accurate as the DFT model and as efficient as the empirical potentials.

To achieve this goal, we have to address two issues. The first is the generation of data. The second is coming up with the appropriate neural network model. These two issues are the common features for all the problems that we discuss here.

The issue of adaptive data generation is very much the same as before. The EELT procedure can still be used. The details of how each step is implemented is a little different. We refer to [40] for details.

Figure 2 shows the effect of using the EELT algorithm. As one can see, a very small percentage of the configurations explored are actually labeled. For the Al-Mg example, only  $\sim 0.005\%$  configurations explored by are selected for labeling.

For the second issue, the design of appropriate neural networks, the most important considerations are:

1. Extensiveness, the neural network should be extensive in the sense that if we want to extend the system, we just have to extend the neural network accordingly. One way of achieving this is suggested by Behler and Parrinello [4].
2. Preserving the symmetry. Besides the usual translational and rotational symmetry, one also has the permutational symmetry: If we relabel a system of copper atoms, its potential energy should not change. It makes a big difference in terms of the accuracy of the neural network model whether one takes these symmetries into account (see [23] and Figure 3).

One very nice and general way of addressing the symmetry problem is to design the neural network model as the composition of two networks: An embedding network followed by a fitting network. The task for the embedding network is to represent enough symmetry-preserving functions to be fed into the fitting network [39].

With these issues properly addressed, one can come up with very satisfactory neural network-based representation of  $V$  (see Figure 4). This representation is named Deep Potential [23, 39] and the Deep Potential-based molecular dynamics is named DeePMD [38]. As has been demonstrated recently in [30, 25], DeePMD, combined with state of the



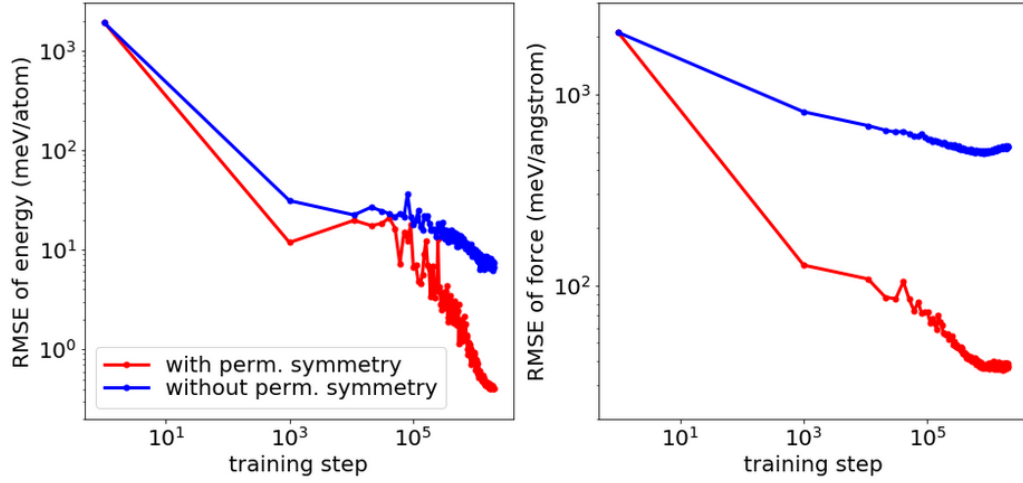


Figure 3: The effect of symmetry preservation on testing accuracy. Shown in red are the results of poor man’s way of imposing symmetry (see main text for explanation). One can see that testing accuracy is drastically improved. Reproduced with permission from Linfeng Zhang.

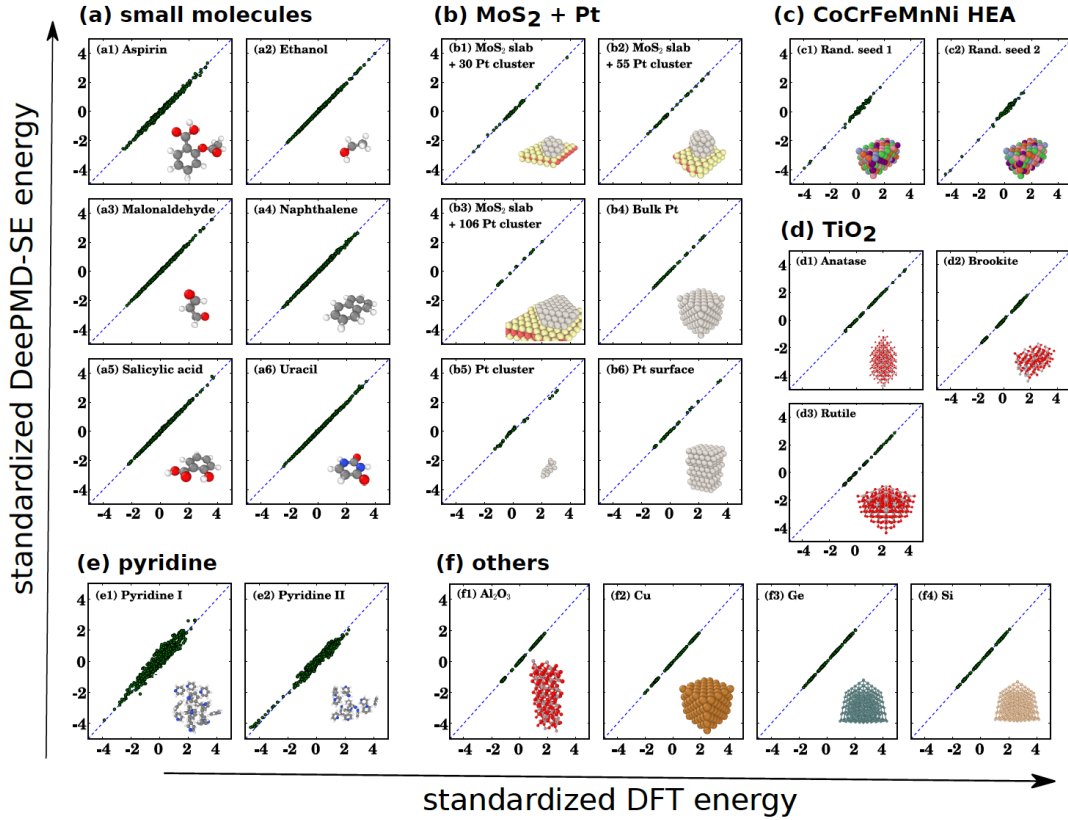


Figure 4: The test accuracy of the Deep Potential for a wide variety of systems. Reproduced with permission from Linfeng Zhang. See also [39].



art high performance supercomputers, can help to increase the size of the system that one can model with *ab initio* accuracy by 5 orders of magnitude.

### 3 Machine learning-based algorithms for high dimensional problems in scientific computing

#### 3.1 Stochastic control

The first application of machine learning for solving high dimensional problems in scientific computing was presented in [20]. Stochastic control was chosen as the first example due to its close analogy with machine learning. Consider the discrete stochastic dynamical system:

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

Here  $s_t$  and  $a_t$  are respectively the state and control at time  $t$ ,  $\xi_t$  is the noise at time  $t$ . Our objective is to solve:

$$\min_{\{a_t\}_{t=0}^{T-1}} \mathbb{E}_{\{\xi_t\}} \left\{ \sum_{t=0}^{T-1} c_t(s_t, a_t) + c_T(s_T) \right\} \quad (2)$$

within the set of feedback controls:

$$a_t = A_t(s_t). \quad (3)$$

We approximate the functions  $A_t$  by neural network models:

$$A_t(s) \approx \tilde{A}_t(s|\theta_t), \quad t = 0, \dots, T-1 \quad (4)$$

The optimization problem (2) then becomes:

$$\min_{\{\theta_t\}_{t=0}^{T-1}} \mathbb{E}_{\{\xi_t\}} \left\{ \sum_{t=0}^{T-1} c_t(s_t, \tilde{A}_t(s_t|\theta_t)) + c_T(s_T) \right\}. \quad (5)$$

Unlike the situation in standard supervised learning, here we have  $T$  set of neural networks to be trained simultaneously. The network architecture is shown in Figure 5

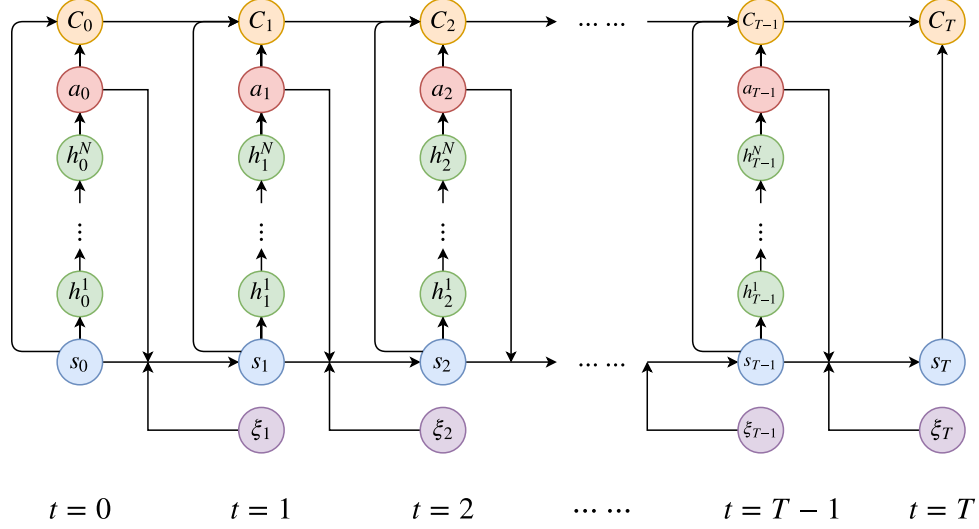


Figure 5: Network architecture for solving stochastic control in discrete time. The whole network has  $(N + 1)T$  layers in total that involve free parameters to be optimized simultaneously. Each column (except  $\xi_t$ ) corresponds to a sub-network at  $t$ . Reproduced with permission from Jiequn Han. See also [20].

Compared with the standard setting for machine learning, one can see a clear analogy in which (1) plays the role for the residual networks and the noise  $\{\xi_t\}$  plays the role of data. Indeed, stochastic gradient descent (SGD) can be readily used to solve the optimization problem (5).

An example of the application of this algorithm is shown in Figure 6 for the problem of energy storage with multiple devices. Here  $n$  is the number of devices. For more details, we refer to [20].

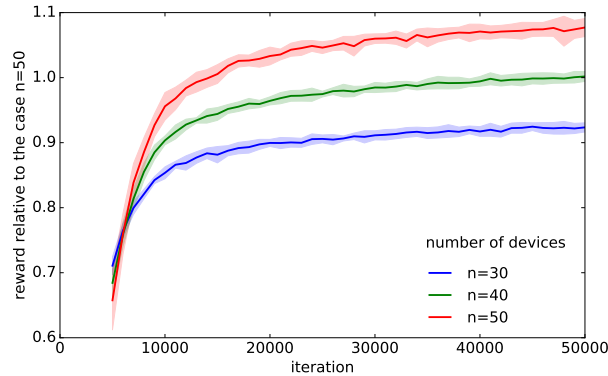


Figure 6: Relative reward for the energy storage problem. The space of control function is  $\mathbb{R}^{n+2} \rightarrow \mathbb{R}^{3n}$  for  $n = 30, 40, 50$ , with multiple equality and inequality constraints. Reproduced with permission from Jiequn Han. See also [20].

### 3.2 Nonlinear parabolic PDEs

Consider parabolic PDEs of the form:

$$\frac{\partial u}{\partial t} + \frac{1}{2} \sigma \sigma^T : \nabla_x^2 u + \mu \cdot \nabla u + f(\sigma^T \nabla u) = 0, \quad u(T, x) = g(x)$$

We study a terminal-value problem instead of initial-value problem since one of the main applications we have in mind is in finance. To develop machine learning-based algorithms, we would like to first reformulate this as a stochastic optimization problem. This can be done using backward stochastic differential equations (BSDE) [33].

$$\inf_{Y_0, \{Z_t\}_{0 \leq t \leq T}} \mathbb{E} |g(X_T) - Y_T|^2, \quad (6)$$

$$s.t. \quad X_t = \xi + \int_0^t \mu(s, X_s) ds + \int_0^t \Sigma(s, X_s) dW_s, \quad (7)$$

$$Y_t = Y_0 - \int_0^t h(s, X_s, Y_s, Z_s) ds + \int_0^t (Z_s)^T dW_s. \quad (8)$$

It can be shown that the unique minimizer of this problem is the solution to the PDE with:

$$Y_t = u(t, X_t) \quad \text{and} \quad Z_t = \sigma^T(t, X_t) \nabla u(t, X_t). \quad (9)$$

With this formulation, one can develop a machine learning-based algorithm along the following lines, adopting the ideas for the stochastic control problems discussed earlier [10, 21]:

- After time discretization, approximate the unknown functions

$$X_0 \mapsto u(0, X_0) \quad \text{and} \quad X_{t_j} \mapsto \sigma^T(t_j, X_{t_j}) \nabla u(t_j, X_{t_j})$$

by feedforward neural networks  $\psi$  and  $\phi$ .

- Using the BSDE, one constructs an approximation  $\hat{u}$  that takes the paths  $\{X_{t_n}\}_{0 \leq n \leq N}$  and  $\{W_{t_n}\}_{0 \leq n \leq N}$  as the input data and gives the final output, denoted by  $\hat{u}(\{X_{t_n}\}_{0 \leq n \leq N}, \{W_{t_n}\}_{0 \leq n \leq N})$  as an approximation to  $u(t_N, X_{t_N})$ .
- The error in the matching between  $\hat{u}$  and the given terminal condition defines the expected loss function

$$l(\theta) = \mathbb{E} \left[ \left| g(X_{t_N}) - \hat{u}(\{X_{t_n}\}_{0 \leq n \leq N}, \{W_{t_n}\}_{0 \leq n \leq N}) \right|^2 \right].$$

This algorithm is called the Deep BSDE method.

As applications, let us first study a stochastic control problem, but we now solve this problem using the Hamilton-Jacobi-Bellman (HJB) equation. Consider the well-known LQG (linear quadratic Gaussian) problem at dimension  $d = 100$ :

$$dX_t = 2\sqrt{\lambda} m_t dt + \sqrt{2} dW_t, \quad (10)$$

with the cost functional:  $J(\{m_t\}_{0 \leq t \leq T}) = \mathbb{E}[\int_0^T \|m_t\|_2^2 dt + g(X_T)]$ . The corresponding HJB equation is given by

$$\frac{\partial u}{\partial t} + \Delta u - \lambda \|\nabla u\|_2^2 = 0 \quad (11)$$

Using the Hopf-Cole transform, we can express the solution in the form:

$$u(t, x) = -\frac{1}{\lambda} \ln \left( \mathbb{E} \left[ \exp \left( -\lambda g(x + \sqrt{2} W_{T-t}) \right) \right] \right). \quad (12)$$

This can be used to calibrate the accuracy of the Deep BSDE method.

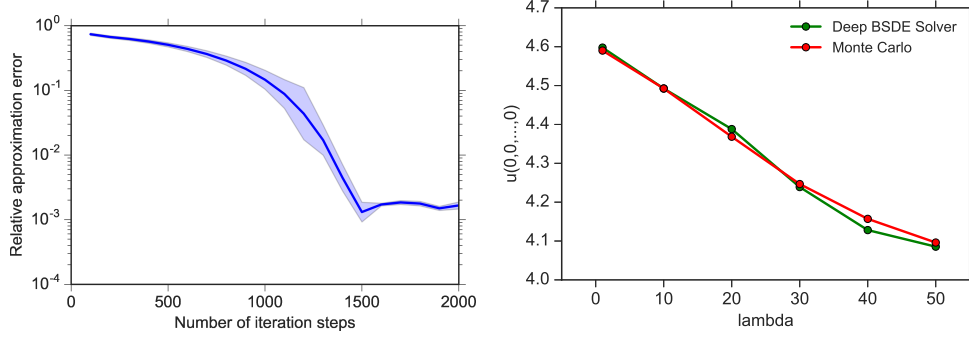


Figure 7: Left: Relative error of the deep BSDE method for  $u(t=0, x=(0, \dots, 0))$  when  $\lambda = 1$ , which achieves 0.17% in a runtime of 330 seconds. Right: Optimal cost  $u(t=0, x=(0, \dots, 0))$  against different  $\lambda$ . Reproduced with permission from Jiequn Han. See also [21].

As a second example, we study the Black-Scholes equation with default risk:

$$\frac{\partial u}{\partial t} + \Delta u - (1 - \delta) Q(u(t, x)) u(t, x) - R u(t, x) = 0$$

where  $Q$  is some nonlinear function. This form of modeling the default risk was proposed and/or used in the literature for low dimensional situation ( $d = 5$ , see [21] for references). The Deep BSDE method was used for this problem with  $d = 100$  [21].

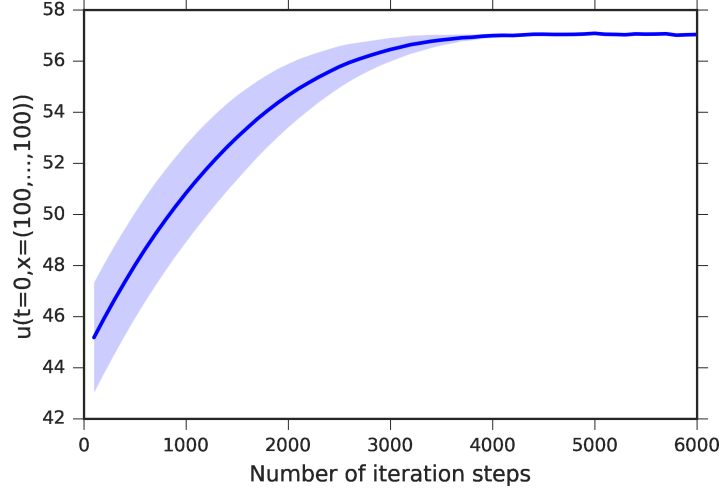


Figure 8: The solution of the Black-Scholes equation with default risk at  $d = 100$ . The Deep BSDE method achieves a relative error of size 0.46% in a runtime of 617 seconds. Reproduced with permission from Jiequn Han. See also [21].

The Deep BSDE method has been applied to pricing basket options, interest rate-dependent options, Libor market model, Bermudan Swaption, barrier option (see [17] for references).

### 3.3 Moment closure for kinetic equations modeling gas dynamics

The dynamics of gas can be modeled very accurately by the well-known Boltzmann Equation:

$$\partial_t f + \mathbf{v} \cdot \nabla_{\mathbf{x}} f = \frac{1}{\varepsilon} Q(f), \quad \mathbf{v} \in \mathbb{R}^3, \quad \mathbf{x} \in \Omega \subset \mathbb{R}^3, \quad (13)$$

where  $f$  is the phase space density function,  $\varepsilon$  is the Knudsen number:

$$\varepsilon = \frac{\text{mean free path}}{\text{macroscopic length}},$$

$Q$  is the collision operator that models the collision process between gas particles. When  $\varepsilon \ll 1$ , this can be approximated by the Euler equation:

$$\partial_t \mathbf{U} + \nabla_{\mathbf{x}} \cdot \mathbf{F}(\mathbf{U}) = 0, \quad (14)$$

where

$$\mathbf{U} = (\rho, \rho \mathbf{u}, E)^T, \quad \rho = \int f \, d\mathbf{v}, \quad \mathbf{u} = \frac{1}{\rho} \int f \mathbf{v} \, d\mathbf{v}.$$

and

$$\mathbf{F}(\mathbf{U}) = (\rho \mathbf{u}, \rho \mathbf{u} \otimes \mathbf{u} + pI, (E + p)\mathbf{u})^T$$

Euler's equation can be obtained by projecting the Boltzmann equation on to the low order moments involved, and making use of the ansatz that the distribution function  $f$  is close to be local Maxwellian.

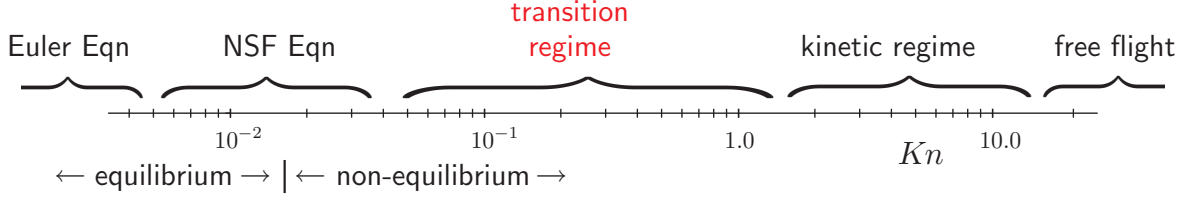


Figure 9: The different regimes of gas dynamics. Reproduced with permission from Jiequn Han. See also [22].

What happens when  $\varepsilon$  is not small? In this case, a natural idea is to seek some generalization of Euler’s equation using more moments. This program was initiated by Harold Grad who constructed the well-known 13-moment system using the moments of  $\{1, \mathbf{v}, (\mathbf{v} - \mathbf{u}) \otimes (\mathbf{v} - \mathbf{u}), |\mathbf{v} - \mathbf{u}|^2(\mathbf{v} - \mathbf{u})\}$ . This line of work has encountered several difficulties. First, there is no guarantee that the equations obtained are well-posed. Secondly there is always the “closure problem”: When projecting the Boltzmann equation on a set of moments, there are always terms which involve moments outside the set of moments considered. In order to obtain a closed system, one needs to model these terms in some way. For Euler’s equation, this is done using the local Maxwellian approximation. This is accurate when  $\varepsilon$  is small, but is no longer so when  $\varepsilon$  is not small. It is highly unclear what should be used as the replacement.

In [22], Han et al developed a machine learning-based moment method. The overall objective is to construct an uniformly accurate (generalized) moment model. The methodology consists of two steps:

1: Learn a set of optimal generalized moments through an auto-encoder. Here by optimality we mean that the set of generalized moments retains a maximum amount of information about the original distribution function and can be used to recover the distribution function with a minimum loss of accuracy. This can be done as follows: Find an encoder  $\Psi$  and a decoder  $\Phi$  that recovers the original  $f$  from  $\mathbf{U}, \mathbf{W}$

$$\mathbf{W} = \Psi(f) = \int \mathbf{w} f d\mathbf{v}, \quad \Phi(\mathbf{U}, \mathbf{W})(\mathbf{v}) = \mathbf{h}(\mathbf{v}; \mathbf{U}, \mathbf{W}).$$

$$\text{Minimize}_{\mathbf{w}, \mathbf{h}} \mathbb{E}_{f \sim \mathcal{D}} \|f - \Phi(\Psi(f))\|^2.$$

$\mathbf{U}$  and  $\mathbf{W}$  form the set of generalized hydrodynamic variables that we will use to model the gas flow.

2: Learn the fluxes and source terms in the PDE for the projected PDE. The effective PDE for  $\mathbf{U}$  and  $\mathbf{W}$  can be obtained by formally projecting the Boltzmann equation on this set of (generalized) moments. This gives us a set of PDEs of the form:

$$\begin{cases} \partial_t \mathbf{U} + \nabla_{\mathbf{x}} \cdot \mathbf{F}(\mathbf{U}, \mathbf{W}; \varepsilon) = 0, \\ \partial_t \mathbf{W} + \nabla_{\mathbf{x}} \cdot \mathbf{G}(\mathbf{U}, \mathbf{W}; \varepsilon) = \mathbf{R}(\mathbf{U}, \mathbf{W}; \varepsilon). \end{cases} \quad (15)$$

where  $\mathbf{F}(\mathbf{U}, \mathbf{W}; \varepsilon) = \int \mathbf{v} \mathbf{U} f d\mathbf{v}$ ,  $\mathbf{G}(\mathbf{U}, \mathbf{W}; \varepsilon) = \int \mathbf{v} \mathbf{W} f d\mathbf{v}$ ,  $\mathbf{R}(\mathbf{U}, \mathbf{W}; \varepsilon) = \varepsilon^{-1} \int \mathbf{W} Q(f) d\mathbf{v}$ . Our task now is to learn  $\mathbf{F}, \mathbf{G}, \mathbf{R}$  from the original kinetic equation.

Again the important issues are (1) get an optimal dataset, and (2) enforce the physical constraints. Here two notable physical constraints are (1) conservation laws and (2) symmetries. Conservation laws are automatically respected in this approach. Regarding symmetries, besides the usual static symmetry, there is now a new dynamic symmetry: the Galilean invariance. These issues are all discussed in [22]. We also refer to [22] for numerical results for the models obtained this way.

## 4 Mathematical theory of machine learning

While neural network-based machine learning has demonstrated a wide range of very impressive successes, it has also acquired a reputation of being a “black magic” rather than a solid scientific technique. This is due to the fact that (1) we lack a basic understanding of the fundamental reasons behind its success; (2) the performance of these models and algorithms is quite sensitive to the choice of the hyper-parameters such as the architecture of the network and the learning rate; and (3) some techniques, such as batch normalization, does appear to be a black magic.

To change this situation, we need to (1) improve our understanding of the reasons behind the success and the fragility of neural network-based models and algorithms and (2) find ways to formulate more robust models and design more robust algorithms. In this section we address the first issue. The next section will be devoted to the second issue.

Here is a list of the most basic questions that we need to address:

- Why does it work in such high dimension?
- Why simple gradient descent works for training neural network models?
- Is over-parametrization good or bad?
- Why does neural network modeling require such extensive parameter tuning?

At this point, we do not yet have clear answers to all these questions. But some coherent picture is starting to emerge. We will focus on the problem of supervised learning, namely approximating a target function using a finite dataset. For simplicity, we will limit ourselves to the case when the physical domain of interest is  $X = [0, 1]^d$ .

### 4.1 An introduction to neural network-based supervised learning

The basic problem in supervised learning is as follows: Given a natural number  $n \in \mathbb{N}$  and a sequence  $\{(\mathbf{x}_j, y_j) = (\mathbf{x}_j, f^*(\mathbf{x}_j)), j \in \{1, 2, \dots, n\}\}$ , of pairs of input-output data, we want to recover the target function  $f^*$  as accurately as possible. We will assume that the input data  $\{\mathbf{x}_j, j \in \{1, 2, \dots, n\}\}$ , is sampled from the probability distribution  $\mu$  on  $\mathbb{R}^d$ .

**Step 1. Choose a hypothesis space.** This is a set of trial functions  $\mathcal{H}_m$  where  $m \in \mathbb{N}$  is the dimensionality of  $\mathcal{H}_m$ . One might choose piecewise polynomials or wavelets.



In modern machine learning the most popular choice is neural network functions. Two-layer neural network functions (one input layer, one output layer which usually does not count, and one hidden layer) take the form

$$f_m(\mathbf{x}, \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle) \quad (16)$$

where  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  is a fixed scalar nonlinear function and where  $\theta = \{(a_j, w_j)_{j \in \{1, 2, \dots, m\}}\}$  are the parameters to be optimized (or trained). A popular example for the nonlinear function  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  is the ReLU (rectified linear unit) activation function:  $\sigma(z) = \max\{z, 0\}$ , for all  $z \in \mathbb{R}$ . We will restrict our attention to this activation function. Roughly speaking, deep neural network (DNN) functions are obtained if one composes two-layer neural network functions several times. One important class of DNN models are residual neural networks (ResNet). They closely resemble discretized ordinary differential equations and take the form

$$\mathbf{z}_{l+1} = \mathbf{z}_l + \sum_{j=1}^M \mathbf{a}_{j,l} \sigma(\langle \mathbf{z}_l, \mathbf{w}_{j,l} \rangle), \quad \mathbf{z}_0 = \mathbf{V} \tilde{\mathbf{x}}, \quad f_L(\mathbf{x}, \theta) = \langle \alpha, \mathbf{z}_L \rangle \quad (17)$$

for  $l \in \{0, 1, \dots, L-1\}$  where  $L, M \in \mathbb{N}$ . Here the parameters are  $\theta = (\alpha, \mathbf{V}, (\mathbf{a}_{j,l})_{j,l}, (\mathbf{w}_{j,l})_{j,l})$ . ResNets are the model of choice for truly deep neural network models.

**Step 2. Choose a loss function.** The primary consideration for the choice of the loss function is to fit the data. Therefore the most obvious choice is the  $L^2$  loss:

$$\hat{\mathcal{R}}_n(f) = \frac{1}{n} \sum_{j=1}^n |f(\mathbf{x}_j) - y_j|^2 = \frac{1}{n} \sum_{j=1}^n |f(\mathbf{x}_j) - f^*(\mathbf{x}_j)|^2. \quad (18)$$

This is also called the “empirical risk”. Sometimes we also add regularization terms.

**Step 3. Choose an optimization algorithm.** The most popular optimization algorithms in machine learning are different versions of the gradient descent (GD) algorithm, or its stochastic analog, the stochastic gradient descent (SGD) algorithm. Assume that the objective function we aim to minimize is of the form

$$F(\theta) = \mathbb{E}_{\xi \sim \nu} [l(\theta, \xi)]. \quad (19)$$

The simplest form of SGD iteration takes the form

$$\theta_{k+1} = \theta_k - \eta \nabla l(\theta_k, \xi_k), \quad (20)$$

for  $k \in \mathbb{N}_0$  where  $\{\xi_k, k \in \mathbb{N}_0 = \{0, 1, 2, \dots\}\}$  is a sequence of i.i.d. random variables sampled from the distribution  $\nu$  and  $\eta$  is the learning rate which might also change during the course of the iteration. In contrast, GD takes the form

$$\theta_{k+1} = \theta_k - \eta \nabla \mathbb{E}_{\xi \sim \nu} [l(\theta_k, \xi)]. \quad (21)$$

Obviously this form of SGD can be adapted to loss functions of the form (18) which can be regarded as an expectation with  $\nu$  being the empirical distribution on the training dataset. This DNN-SGD paradigm is the heart of modern machine learning.

## 4.2 Approximation theory

The simplest way of approximating functions is to use polynomials. For polynomial approximation, there are two kinds of theorems. The first is the Weierstrass' Theorem which asserts that continuous functions can be uniformly approximated by polynomials on compact domains. The second is Taylor's Theorem which tells us that the rate of convergence depends on the smoothness of the target function.

Using the terminology in neural network theory, Weierstrass' Theorem is the "Universal Approximation Theorem" (UAT). It is a useful fact. But Taylor's Theorem is more useful since it tells us something about the rate of convergence. The form of Taylor's Theorem used in approximation theory are the direct and inverse approximation theorems which assert that a given function can be approximated by polynomials with a particular rate if and only if certain norms of that function are finite. This particular norm, which measures the regularity of the function, is the key quantity that characterizes this approximation scheme. For piecewise polynomials, these norms are some Besov space norms [36]. For  $L^2$ , a typical result looks like follows:

$$\inf_{f \in \mathcal{H}_m} \|f - f_m\|_{L^2(X)} \leq C_0 h^\alpha \|f\|_{H^\alpha(X)} \quad (22)$$

Here  $H^\alpha$  stands for the  $\alpha$ -th order Sobolev norm [36],  $m$  is the number of degrees of freedom. On a regular grid, the grid size is given by

$$h \sim m^{-1/d} \quad (23)$$

An important thing to notice is that the convergence rate in (22) suffers from CoD: If we want to reduce the error by a factor of  $\epsilon$ , we need to increase  $m$  by a factor  $m \sim \epsilon^{-d}$  if  $\alpha = 1$ . For  $d = 100$  which is not very high dimension by the standards of machine learning, this means that we have to increase  $m$  by a factor of  $\epsilon^{-100}$ . This is why polynomials and piecewise polynomials are not useful in high dimensions.

Another way to appreciate this is as follows. The number of monomials of degree  $p$  in dimension  $d$  is  $C_{p+d}^d$ . This grows very fast for large values of  $d$  and  $p$ .

What should we expect in high dimension? One example that we can learn from is Monte Carlo methods for integration. Consider the problem of approximating

$$I(g) = \mathbb{E}_{\mathbf{x} \sim \mu} g(\mathbf{x})$$

using

$$I_m(g) = \frac{1}{m} \sum_j g(\mathbf{x}_j)$$

where  $\{\mathbf{x}_j, j \in [m]\}$  is a set of i.i.d samples of the probability distribution  $\mu$ . A direct computation gives

$$\mathbb{E}(I(g) - I_m(g))^2 = \frac{\text{var}(g)}{m}, \quad \text{var}(g) = \mathbb{E}_{\mathbf{x} \sim \mu} g^2(\mathbf{x}) - (\mathbb{E}_{\mathbf{x} \sim \mu} g(\mathbf{x}))^2$$

This exact relation tells us two things. (1) The convergence rate of Monte Carlo integration is independent of dimension. (2) The error constant is given by the variance of the integrand. Therefore to reduce error, one has to do variance reduction.

Had we used grid-based quadrature rules, the accuracy would have also suffered from CoD.

It is possible to improve the Monte Carlo rate by more sophisticated ways of choosing  $\{\mathbf{x}_j, j \in [m]\}$ , say using number-theoretic-based quadrature rules. But these typically give rise to an  $O(1/d)$  improvement for the convergence rate and it diminishes as  $d \rightarrow \infty$ .

Based on these considerations, we aim to find function approximations that satisfy:

$$\inf_{f \in \mathcal{H}_m} \mathcal{R}(f) = \inf_{f \in \mathcal{H}_m} \|f - f^*\|_{L^2(d\mu)}^2 \lesssim \frac{\|f^*\|_*^2}{m}$$

The natural questions are then:

- How can we achieve this? That is, what kind of hypothesis space should we choose?
- What should be the “norm”  $\|\cdot\|_*$  (associated with the choice of  $\mathcal{H}_m$ )? Here we put norm in quotation marks since it does not have to be a real norm. All we need is that it controls the approximation error.

Regarding the first question, let us look an illustrative example.

Consider the following Fourier representation of the function  $f$  and its approximation  $f_m$  (say FFT-based):

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

Here  $\{\boldsymbol{\omega}_j\}$  is a fixed grid, e.g. uniform grid. For this approximation, we have

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

which suffers from CoD.

Now consider the alternative representation

$$f(\mathbf{x}) = \int_{\mathbb{R}^d} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})} \pi(d\boldsymbol{\omega}) = \mathbb{E}_{\boldsymbol{\omega} \sim \pi} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})} \quad (24)$$

where  $\pi$  is a probability distribution. Now to approximate  $f$ , it is natural to use Monte Carlo. Let  $\{\boldsymbol{\omega}_j\}$  be an i.i.d. sample of  $\pi$ ,  $f_m(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a(\boldsymbol{\omega}_j) e^{i(\boldsymbol{\omega}_j, \mathbf{x})}$ , then we have

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

This approximation does not suffer from CoD. Notice that  $f_m(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\boldsymbol{\omega}_j^T \mathbf{x})$  is nothing but a two-layer neural network with activation function  $\sigma(z) = e^{iz}$  (here  $a_j = a(\boldsymbol{\omega}_j)$ ).

We believe that this simple argument is really at the heart of why neural network models do so well in high dimension.

Now let us turn to a concrete example of the kind of approximation theory for neural network models. We will consider two-layer neural networks.

$$\mathcal{H}_m = \{f_m(\mathbf{x}) = \frac{1}{m} \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x})\}, \quad \theta = \{(a_j, \mathbf{w}_j), j \in [m]\}$$

Consider function  $f : X = [0, 1]^d \mapsto \mathbb{R}$  of the following form

$$f(\mathbf{x}) = \int_{\Omega} a\sigma(\mathbf{w}^T \mathbf{x}) \rho(da, d\mathbf{w}) = \mathbb{E}_{(a, \mathbf{w}) \sim \rho}[a\sigma(\mathbf{w}^T \mathbf{x})], \quad \mathbf{x} \in X$$

where  $\Omega = \mathbb{R}^1 \times \mathbb{R}^{d+1}$ ,  $\rho$  is a probability distribution on  $\Omega$ . Define:

$$\|f\|_{\mathcal{B}} = \inf_{\rho \in P_f} (\mathbb{E}_{\rho}[a^2 \|\mathbf{w}\|_1^2])^{1/2}$$

where  $P_f := \{\rho : f(\mathbf{x}) = \mathbb{E}_{\rho}[a\sigma(\mathbf{w}^T \mathbf{x})]\}$ . This is called the Barron norm [2, 12, 13]. The space

$$\mathcal{B} = \{f \in C^0 : \|f\|_{\mathcal{B}} < \infty\}$$

is called the Barron space [2, 12, 13] (see also [3, 26, 14]).

In analogy with classical approximation theory, we can also prove some direct and inverse approximation theorem [13].

**Theorem 1** (Direct Approximation Theorem). *If  $\|f\|_{\mathcal{B}} < \infty$ , then for any integer  $m > 0$ , there exists a two-layer neural network function  $f_m$  such that*

$$\|f - f_m\|_{L^2(X)} \lesssim \frac{\|f\|_{\mathcal{B}}}{\sqrt{m}}$$

**Theorem 2** (Inverse Approximation Theorem). *Let*

$$\mathcal{N}_C \stackrel{\text{def}}{=} \left\{ \frac{1}{m} \sum_{k=1}^m a_k \sigma(\mathbf{w}_k^T \mathbf{x}) : \frac{1}{m} \sum_{k=1}^m |a_k|^2 \|\mathbf{w}_k\|_1^2 \leq C^2, m \in \mathbb{N}^+ \right\}.$$

*Let  $f^*$  be a continuous function. Assume there exists a constant  $C$  and a sequence of functions  $f_m \in \mathcal{N}_C$  such that*

$$f_m(\mathbf{x}) \rightarrow f^*(\mathbf{x})$$

*for all  $\mathbf{x} \in X$ . Then there exists a probability distribution  $\rho^*$  on  $\Omega$ , such that*

$$f^*(\mathbf{x}) = \int a\sigma(\mathbf{w}^T \mathbf{x}) \rho^*(da, d\mathbf{w}),$$

*for all  $\mathbf{x} \in X$  and  $\|f^*\|_{\mathcal{B}} \leq C$ .*

### 4.3 Estimation error

Another issue we have to worry about is the performance of the machine learning model outside the training dataset. This issue also shows up in classical approximation theory. Illustrated in Figure 10 is the classical Runge phenomenon for polynomial interpolation on a uniform grid. One can see that away from the grid points, the error of the interpolant can be very large. This is a situation that we would like to avoid.

What we do in practice is to minimize the training error:

$$\hat{\mathcal{R}}_n(\theta) = \frac{1}{n} \sum_j (f(\mathbf{x}_j, \theta) - f^*(\mathbf{x}_j))^2$$

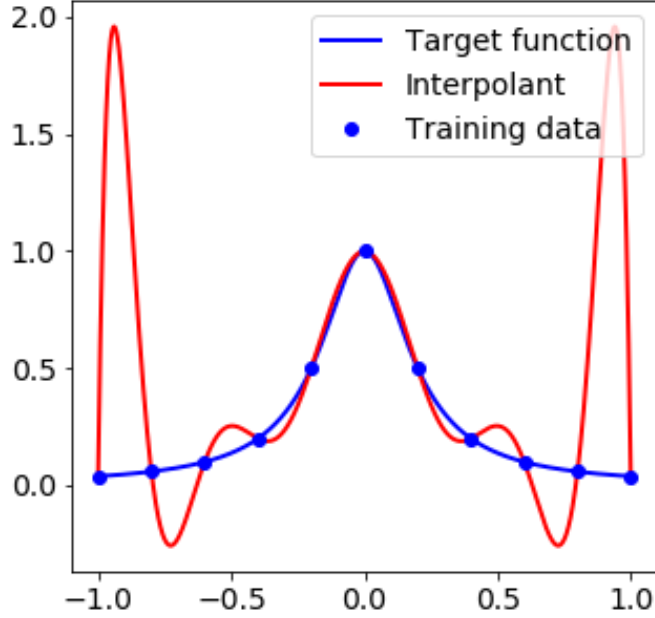


Figure 10: The Runge phenomenon:  $f^*(x) = \frac{1}{1+25x^2}$ . Reproduced with permission from Chao Ma.

but we are interested in the testing error, which is a sampled version of the population risk:

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

The question is how we can control the difference between these two errors.

One way of doing this is to use the notion of Rademacher complexity. The important fact for us here is that the Rademacher complexity controls the difference between training and testing errors (also called the “generalization gap”). Indeed, let  $\mathcal{H}$  be a set of functions, and  $S = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  be a dataset. Then, up to logarithmic terms, we have

$$\sup_{h \in \mathcal{H}} \left| \mathbb{E}_{\mathbf{x}} [h(\mathbf{x})] - \frac{1}{n} \sum_{i=1}^n h(\mathbf{x}_i) \right| \sim \text{Rad}_S(\mathcal{H})$$

where the **Rademacher complexity** of  $\mathcal{H}$  with respect to  $S$  is defined as

$$\text{Rad}_S(\mathcal{H}) = \frac{1}{n} \mathbb{E}_{\xi} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^n \xi_i h(\mathbf{x}_i) \right], \quad (25)$$

where  $\{\xi_i\}_{i=1}^n$  are i.i.d. random variables taking values  $\pm 1$  with equal probability.

The question then becomes to bound the Rademacher complexity of a hypothesis space. For the Barron space, we have [2]:

**Theorem 3.** Let  $\mathcal{F}_Q = \{f \in \mathcal{B}, \|f\|_{\mathcal{B}} \leq Q\}$ . Then we have

$$\text{Rad}_S(\mathcal{F}_Q) \leq 2Q \sqrt{\frac{2 \ln(2d)}{n}}$$

where  $n = |S|$ , the size of the dataset  $S$ .

#### 4.4 A priori estimates for regularized models

Consider the regularized model

$$\mathcal{L}_n(\theta) = \hat{\mathcal{R}}_n(\theta) + \lambda \sqrt{\frac{\log(2d)}{n}} \|\theta\|_{\mathcal{P}}, \quad \hat{\theta}_n = \argmin \mathcal{L}_n(\theta) \quad (26)$$

where the path norm is defined by:

$$\|\theta\|_{\mathcal{P}} = \left( \frac{1}{m} \sum_{k=1}^m |a_k|^2 \|\mathbf{w}_k\|_1^2 \right)^{1/2}$$

The following result was proved in [12]:

**Theorem 4.** Assume  $f^* : X \mapsto [0, 1] \in \mathcal{B}$ . There exist constants  $C_0$ , such that for any  $\delta > 0$ , if  $\lambda \geq C_0$ , then with probability at least  $1 - \delta$  over the choice of the training dataset, we have

$$\mathcal{R}(\hat{\theta}_n) \lesssim \frac{\|f^*\|_{\mathcal{B}}^2}{m} + \lambda \|f^*\|_{\mathcal{B}} \sqrt{\frac{\log(2d)}{n}} + \sqrt{\frac{\log(1/\delta) + \log(n)}{n}}.$$

Similar approximation theory and a priori error estimates have been proved for other machine learning models. Here is a brief summary of these results.

- Random feature model: The corresponding function space is the reproducing kernel Hilbert space (RKHS).
- Residual networks (ResNets): The corresponding function space is the so-called flow-induced space introduced in [13].
- Multi-layer neural networks: A candidate for the corresponding function space is the multi-layer space introduced in [15].

What is really important is the “norms” that control the approximation error and the generalization gap. These quantities are defined for functions in the corresponding spaces. After the approximation theorems and Rademacher complexity estimates are in place, one can readily prove a theorem of the following type for regularized models: Up to logarithmic terms, the minimizers of the regularized models satisfy:

$$\mathcal{R}(\hat{f}) \lesssim \frac{\Gamma(f^*)}{m} + \frac{\gamma(f^*)}{\sqrt{n}}$$

where  $m$  is the number of free parameters,  $n$  is the size of the training dataset. Note that for the multilayer spaces, the results proved in [15] are not as sharp.

We only discussed the analysis of the hypothesis space. There are many more other questions. We refer to [19] for more discussion on the current understanding of neural network-based machine learning.

## 5 Machine learning from a continuous viewpoint

Now we turn to alternative formulations of machine learning. Motivated by the situation for PDEs, we would like to first formulate machine learning in a continuous setting and then discretize to get concrete models and algorithms. The key here is that continuous problems that we come up with should be nice mathematical problems. For PDEs, this is accomplished by requiring them to be “well-posed”. For problems in calculus of variations, we require the problem to be “convex” in some sense and lower semi-continuous. The point of these requirements is to make sure that the problem has a unique solution. Intuitively, for machine learning problems, being “nice” means that the variational problem should have a simple landscape. How to formulate this precisely is an important research problem for the future.

As was pointed out in [16], the key ingredients for the continuous formulation are as follows:

- representation of functions (as expectations)
- formulating the variational problem (as expectations)
- optimization, e.g. gradient flows

### 5.1 Representation of functions

Two kinds of representations are considered in [16]: Integral transform-based representation and flow-based representation. The simplest integral-transform based representation is a generalization of (24):

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

Here  $\theta$  denotes the parameters in the model:  $\theta$  can be  $a(\cdot)$  or the prob distributions  $\pi$  or  $\rho$ .

This representation corresponds to two-layer neural networks. A generalization to multi-layer neural networks is presented in [15].

Next we turn to flow-based representation:

$$\frac{d\mathbf{z}}{d\tau} = \mathbb{E}_{\mathbf{w} \sim \pi_\tau} \mathbf{a}(\mathbf{w}, \tau) \sigma(\mathbf{w}^T \mathbf{z}) \quad (27)$$

$$= \mathbb{E}_{(a, \mathbf{w}) \sim \rho_\tau} \mathbf{a} \sigma(\mathbf{w}^T \mathbf{z}) \quad (28)$$

$$= \mathbb{E}_{\mathbf{u} \sim \rho_\tau} \phi(\mathbf{z}, \mathbf{u}), \quad \mathbf{z}(0, \mathbf{x}) = \mathbf{x} \quad (29)$$

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

In this representation, the parameter  $\theta$  can be either  $\{a_\tau(\cdot)\}$  or  $\{\pi_\tau\}$  or  $\{\rho_\tau\}$



## 5.2 The stochastic optimization problem

Stochastic optimization problems are of the type:

$$\min_{\theta} \mathbb{E}_{\mathbf{w} \sim \nu} g(\theta, \mathbf{w})$$

These kinds of problems can readily be approached by stochastic algorithms, which is a key component in machine learning. For example, instead of the gradient descent algorithm:

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathbb{E}_{\mathbf{w} \sim \nu} g(\theta, \mathbf{w}) = \theta_k - \eta \nabla_{\theta} \sum_{j=1}^n g(\theta, \mathbf{w}_j)$$

one can use the stochastic gradient descent:

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

where  $\{\mathbf{w}_k\}$  is a set of random variables sampled from  $\nu$ .

The following are some examples of the stochastic optimization problems that arise in modern machine learning:

- Supervised learning: In this case, the minimization of the population risk becomes

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

- Eigenvalue problems for quantum many-body Hamiltonian:

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

Here  $H$  is the Hamiltonian of the quantum system.

- Stochastic control problems:

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

Substituting the representations discussed earlier to these expressions for the stochastic optimization problems, we obtain the final variational problem that we need to solve.

One can either discretize these variational problems directly and then solve the discretized problem using some optimization algorithms, or one can write down continuous forms of some optimization algorithms, typically gradient flow dynamics, and then discretize these continuous flows. We are going to discuss the second approach.

### 5.3 Optimization: Gradient flows

To write continuous form of the gradient flows, we draw some inspiration from statistical physics. Take the supervised learning as an example. We regard the population risk as being the “free energy”, and following Halperin and Hohenberg [24], we divide the parameters into two kinds, conserved and non-conserved. For example,  $a$  is a non-conserved parameter and  $\pi$  is conserved since its total integral has to be 1.

For non-conserved parameter, as was suggested in [24], one can use the “model A” dynamics:

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

which is simply the usual  $L^2$  gradient flow.

For conserved parameters such as  $\pi$ , one should use the “model B” dynamics which works as follows: First define the “chemical potential”

$$V = \frac{\delta \mathcal{R}}{\delta \pi}.$$

From the chemical potential, one obtains the velocity field  $\mathbf{v}$  and the current  $\mathbf{J}$ :

$$\mathbf{J} = \pi \mathbf{v}, \quad \mathbf{v} = -\nabla V$$

The continuity equation then gives us the gradient flow dynamics:

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

This is also the gradient flow under the Wasserstein metric.

### 5.4 Discretizing the gradient flows

To obtain practical models, one needs to discretize these continuous problems. The first step is to replace the population risk by the empirical risk using the training data.

The more non-trivial issue is how to discretize the gradient flows in the parameter space. The parameter space has the following characteristics: (1) It has a simple geometry – unlike the real space which may have a complicated geometry. (2) It is also usually high dimensional. For these reasons, the most natural numerical methods for the discretization in the parameter space is the particle method which is the dynamic version of Monte Carlo. Smoothed particle method might be helpful to improve the performance. In relatively low dimensions, one might also consider the spectral method, particularly some sparse version of the spectral method, due to the relative simple geometry of the parameter space.

Take for example the discretization of the conservative flow for the integral transform-based representation. With the representation:  $f(\mathbf{x}; \theta) = \mathbb{E}_{(a, \mathbf{w}) \sim \rho} a \sigma(\mathbf{w}^T \mathbf{x})$ , the gradient flow equation becomes:

$$\partial_t \rho = \nabla(\rho \nabla V), \quad V = \frac{\delta \mathcal{R}}{\delta \rho} \tag{30}$$

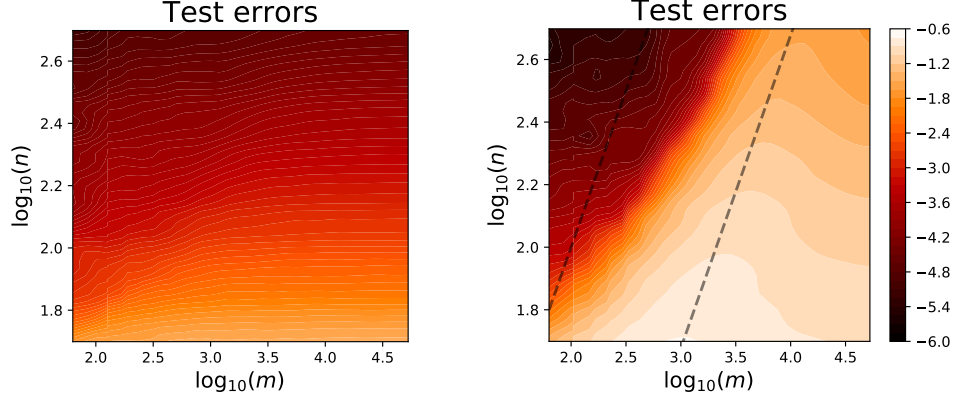


Figure 11: **(Left)** continuous viewpoint; **(Right)** conventional NN models. Target function is a single neuron. Reproduced with permission from Lei Wu.

The particle method discretization is based on:

$$\rho(a, \mathbf{w}, t) \sim \frac{1}{m} \sum_j \delta_{(a_j(t), \mathbf{w}_j(t))} = \frac{1}{m} \sum_j \delta_{\mathbf{u}_j(t)}$$

where  $\mathbf{u}_j(t) = (a_j(t), \mathbf{w}(t))$ . One can show that in this case, (30) reduces to

$$\frac{d\mathbf{u}_j}{dt} = -\nabla_{\mathbf{u}_j} I(\mathbf{u}_1, \dots, \mathbf{u}_m)$$

where

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

This is exactly gradient descent for “scaled” two-layer neural networks. In this case, the continuous formulation also coincides with the “mean-field” formulation for two-layer neural networks [7, 32, 34, 35].

The scaling factor  $1/m$  in front of the “scaled” two-layer neural networks is actually quite important and makes a big difference for the test performance of the network. Shown in Figure 11 is the heat map of the test error for two-layer neural network models with and without this scaling factor. The target function is the simple single neuron function:  $f^*(\mathbf{x}) = \sigma(x_1)$ . The important observation is that in the absence of this scaling factor, the test error shows a “phase transition” between a phase where the neural network model performs like an associated random feature model and another phase where it shows much better performance than the random feature model. Such a phase transition is one of the reasons that choosing the right set of hyper-parameters, here the network width  $m$ , is so important. However, if one uses the scaled form, this phase transition phenomenon is avoided and the performance is more robust [31].

## 5.5 The optimal control problem for flow-based representation

The flow-based representation naturally leads to a control problem. This viewpoint has been used explicitly or implicitly in machine learning for quite some time (see for example

[27]). The back propagation algorithm, for example, is an example of control-theory based algorithm. Another more recent example is the development of maximum principle-based algorithm, first introduced in [28]. Despite these successes, we feel that there is still a lot of room for using the control theory viewpoint to develop new algorithms.

We consider the flow-based representation in a slightly more general form

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

where  $\mathbf{z}$  is the state,  $\rho_\tau$  is the control at time  $\tau$ . Our objective is to 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 as before

$$f(\mathbf{x}) = \mathbf{1}^T \mathbf{z}(1, \mathbf{x}) \quad (31)$$

One most important result for this control problem is Pontryagin's maximum principle (PMP). To state this result, let us define the Hamiltonian  $H : \mathbb{R}^d \times \mathbb{R}^d \times \mathcal{P}_2(\Omega) \mapsto \mathbb{R}$  as

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

Pontryagin's maximum principle asserts that the solutions of the control problem must satisfy:

$$\rho_\tau = \operatorname{argmax}_\rho \mathbb{E}_{\mathbf{x}} [H(\mathbf{z}_\tau^{t,\mathbf{x}}, \mathbf{p}_\tau^{t,\mathbf{x}}, \rho)], \quad \forall \tau \in [0, 1], \quad (32)$$

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

$$\begin{aligned} \frac{d\mathbf{z}_\tau^{t,\mathbf{x}}}{d\tau} &= \nabla_{\mathbf{p}} H = \mathbb{E}_{\mathbf{u} \sim \rho_\tau(\cdot; t)} [\boldsymbol{\phi}(\mathbf{z}_\tau^{t,\mathbf{x}}, \mathbf{u})] \\ \frac{d\mathbf{p}_\tau^{t,\mathbf{x}}}{d\tau} &= -\nabla_{\mathbf{z}} H = \mathbb{E}_{\mathbf{u} \sim \rho_\tau(\cdot; t)} [\nabla_{\mathbf{z}}^T \boldsymbol{\phi}(\mathbf{z}_\tau^{t,\mathbf{x}}, \mathbf{u}) \mathbf{p}_\tau^{t,\mathbf{x}}], \end{aligned} \quad (33)$$

with the boundary conditions:

$$\mathbf{z}_0^{t,\mathbf{x}} = \mathbf{x} \quad (34)$$

$$\mathbf{p}_1^{t,\mathbf{x}} = -2(f(\mathbf{x}; \rho(\cdot; t)) - f^*(\mathbf{x}))\mathbf{1}. \quad (35)$$

Pontryagin's maximum principle is slightly stronger than the KKT condition for the stationary points in that (32) is a statement of optimality rather than criticality. In fact (32) also holds when the parameters are discrete and this has been used in [29] to develop efficient numerical algorithms for this case.

With the help of PMP, it is also easy to write down the gradient descent flow for the optimization problem. Formally, one can simply write down the gradient descent flow for (32) for each  $\tau$ :

$$\partial_t \rho_\tau(\mathbf{u}, t) = \nabla \cdot (\rho_\tau(\mathbf{u}, t) \nabla V(\mathbf{u}; \rho)), \quad \forall \tau \in [0, 1], \quad (36)$$

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  $\{(\mathbf{z}_\tau^{t,\mathbf{x}}, \mathbf{p}_\tau^{t,\mathbf{x}})\}$  are defined as before by the forward/backward equations.

To discretize the gradient flow, we can simply use:

- forward Euler for the flow in  $\tau$  variable, with step size  $1/L$  ( $L$  is the number of grid points in the  $\tau$  variable);
- particle method for the gradient descent dynamics, with  $M$  samples in each  $\tau$ -grid point.

This gives us

$$\mathbf{z}_{l+1}^{t,x} = \mathbf{z}_l^{t,x} + \frac{1}{LM} \sum_{j=1}^M \phi(\mathbf{z}_l^{t,x}, \mathbf{u}_l^j(t)), \quad l = 0, \dots, L-1 \quad (37)$$

$$\mathbf{p}_l^{t,x} = \mathbf{p}_{l+1}^{t,x} + \frac{1}{LM} \sum_{j=1}^M \nabla_{\mathbf{z}} \phi(\mathbf{z}_{l+1}^{t,x}, \mathbf{u}_{l+1}^j(t)) \mathbf{p}_{l+1}^{t,x}, \quad l = 0, \dots, L-1 \quad (38)$$

$$\frac{d\mathbf{u}_l^j(t)}{dt} = -\mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{w}}^T \phi(\mathbf{z}_l^{t,x}, \mathbf{u}_l^j(t)) \mathbf{p}_l^{t,x}]. \quad (39)$$

This recovers the GD algorithm (with back-propagation) for the (scaled) ResNet:

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

We call this “scaled” ResNet because of the presence of the factor  $1/(LM)$ .

In a similar spirit, one can also obtain an algorithm using PMP [28]. Adopting the terminology in control theory, this kind of algorithms are called “method of successive approximation” (MSA). The basic MSA is as follows:

- Initialize:  $\theta^0 \in \mathcal{U}$
- For  $k = 0, 1, 2, \dots$ :

– Solve

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

– 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  $\tau \in [0, 1]$

In practice, this basic version does not perform as well as the “extended MSA” which works in the same way as the MSA except that the Hamiltonian is replaced by an extended Hamiltonian [28]:

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

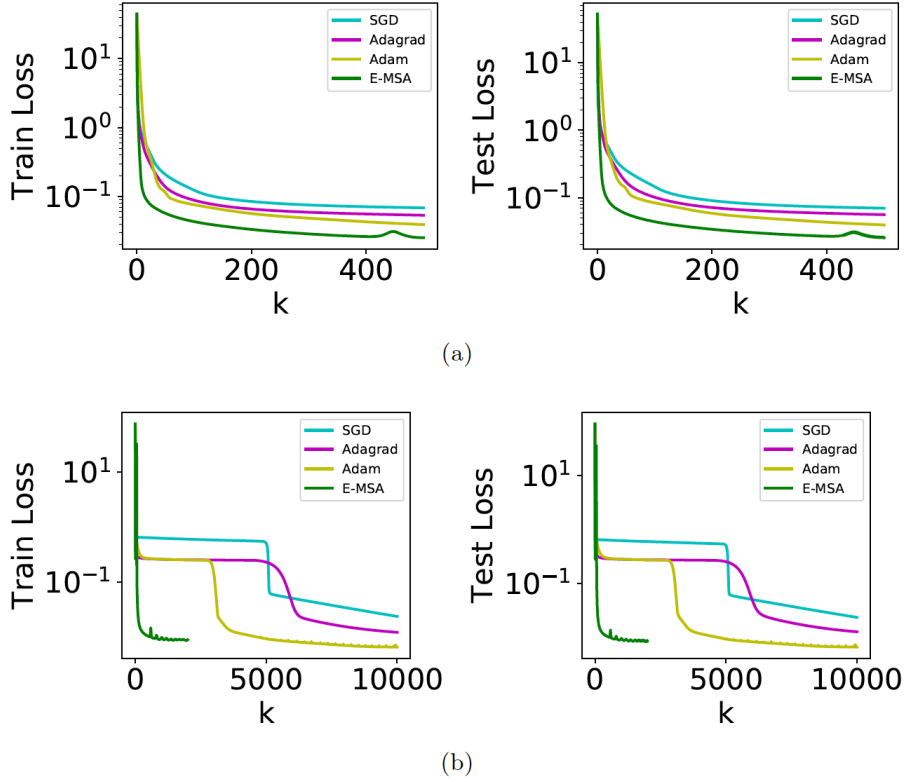


Figure 12: Comparison of the extended MSA with different versions of stochastic gradient descent algorithms. The top figures show results for small initialization. The bottom figures show results for bigger initialization. Reproduced with permission from Qianxiao Li. See also [28].

Figure 12 shows the results of the extended MSA compared with different versions of SGD for two kinds of initialization. One can see that in terms of the number of iterations, extended MSA outperforms all the SGDs. In terms of wall clock time, the advantage of the extended MSA is diminished significantly. This is possibly due to the inefficiencies in the implementation of the optimization algorithm (here the BFGS) used for solving (32). We refer to [28] for more details. In any case, it is clear that there is a lot of room for improvement.

## 6 Concluding remarks

In conclusion, we have discussed a wide range of problems for which machine learning-based algorithms have made and/or will make a significant difference. These problems are relatively new to computational mathematics. We believe strongly that machine learning-based algorithms will also significantly impact the way we solve more traditional problems in computational mathematics. However, research in this direction is still at a very early stage.

Another important area that machine learning can be of great help is multi-scale modeling. The moment-closure problem discussed above is an example in this direction.

There are many more possible applications, see [8]. Machine learning seems to be able to provide the missing link in making advanced multi-scale modeling techniques really practical. For example in the heterogeneous multi-scale method (HMM) [1, 9], one important component is to extract the relevant macro-scale information from micro-scale simulation data. This step has always been a major obstacle in HMM. It seems quite clear that machine learning techniques can of great help here.

We also discussed how the viewpoint of numerical analysis can help to improve the mathematical foundation of machine learning as well as propose new and possibly more robust formulations. In particular, we have given a taste of how high dimensional approximation theory should look like. We also demonstrated that commonly used machine learning models and training algorithms can be recovered from some particular discretization of continuous models, in a scaled form. From this discussion, one can see that neural network models are quite natural and rather inevitable.

What have we really learned from machine learning? Well, it seems that the most important new insight from machine learning is the representation of functions as expectations. We reproduce them here for convenience:

- integral-transform based:

$$f(\mathbf{x}) = \mathbb{E}_{(\mathbf{a}, \mathbf{w}) \sim \rho} \mathbf{a} \sigma(\mathbf{w}^T \mathbf{x})$$

$$f(\mathbf{x}) = \mathbb{E}_{\theta_L \sim \pi_L} a_{\theta_L}^{(L)} \sigma(\mathbb{E}_{\theta_{L-1} \sim \pi_{L-1}} \dots \sigma(\mathbb{E}_{\theta_1 \sim \pi_1} a_{\theta_2, \theta_1}^1 \sigma(a_{\theta_1}^0 \cdot \mathbf{x})) \dots)$$

- flow-based:

$$\frac{d\mathbf{z}}{d\tau} = \mathbb{E}_{(\mathbf{a}, \mathbf{w}) \sim \rho_\tau} \mathbf{a} \sigma(\mathbf{w}^T \mathbf{z}), \quad \mathbf{z}(0, \mathbf{x}) = \mathbf{x} \quad (40)$$

$$f(\mathbf{x}, \theta) = \mathbf{1}^T \mathbf{z}(1, \mathbf{x}) \quad (41)$$

From the viewpoint of computational mathematics, this suggests that the central issue will move from specific discretization schemes to more effective representations of functions.

This review is rather sketchy. Interested reader can consult the three review articles [17, 18, 19] for more details.

**Acknowledgement:** I am very grateful to my collaborators for their contribution to the work described here. In particular, I would like to express my sincere gratitude to Roberto Car, Jiequn Han, Arnulf Jentzen, Qianxiao Li, Chao Ma, Han Wang, Stephan Wojtowytsch, and Lei Wu for the many discussions that we have had on the issues discussed here. This work is supported in part by a gift to the Princeton University from iFlytek as well as the ONR grant N00014-13-1-0338.

## References

- [1] Assyr Abdulle, Weinan E, Bjorn Engquist and Eric Vanden-Eijnden, The heterogeneous multiscale methods, *Acta Numerica*, vol. 21, pp. 1-87, 2012.



- [2] Francis Bach, “Breaking the curse of dimensionality with convex neural networks”, *Journal of Machine Learning Research*, 18(19):1-53, 2017.
- [3] Andrew R. Barron, “Universal approximation bounds for superpositions of a sigmoidal function”, *IEEE Transactions on Information theory*, 39(3):930-945, 1993.
- [4] Jörg Behler and Michele Parrinello, “Generalized neural-network representation of high-dimensional potential-energy surfaces”, *Physical review letters*, 98(14):146401, 2007.
- [5] Achi Brandt, “Multiscale scientific computation: review 2001”. In Barth, T.J., Chan, T.F. and Haimes, R. (eds.): *Multiscale and Multiresolution Methods: Theory and Applications*, Springer Verlag, Heidelberg, 2001, pp. 196.
- [6] Roberto Car and Michele Parrinello, “Unified approach for molecular dynamics and density-functional theory”, *Physical Review Letters*, 55(22):2471, 1985.
- [7] Lenaïc Chizat and Francis Bach, “On the global convergence of gradient descent for over-parameterized models using optimal transport”, In *Advances in neural information processing systems*, pages 3036-3046, 2018.
- [8] Weinan E, *Principles of Multiscale Modeling*, Cambridge University Press, 2011.
- [9] Weinan E and Bjorn Engquist, *The heterogeneous multiscale methods*, *Comm. Math. Sci.*, vol. 1, no. 1, pp. 87-132, 2003.
- [10] Weinan E, Jiequn Han and Arnulf Jentzen, “Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations”, *Communications in Mathematics and Statistics* 5, 4 (2017), 349-380.
- [11] Weinan E, Jiequn Han, and Linfeng Zhang, “Integrating Machine Learning with Physics-Based Modeling”, <https://arxiv.org/pdf/2006.02619.pdf>, 2020.
- [12] Weinan E, Chao Ma and Lei Wu, “A priori estimates of the population risk for two-layer neural networks”, *Communications in Mathematical Sciences*, 17(5):1407-1425, 2019; arXiv:1810.06397, 2018.
- [13] Weinan E, Chao Ma and Lei Wu, “Barron spaces and the flow-induced function spaces for neural network models”, arXiv:1906.08039, 2019.
- [14] Weinan E and Stephan Wojtowytsch, “Representation formulas and pointwise properties for Barron functions”, arXiv:2006.05982, 2020.
- [15] Weinan E and Stephan Wojtowytsch, “On the Banach spaces associated with multi-layer ReLU networks: Function representation, approximation theory and gradient descent dynamics”, <https://arxiv.org/pdf/2007.15623.pdf>, 2020.
- [16] Weinan E, Chao Ma, and Lei Wu, “Machine learning from a continuous viewpoint”, arXiv:1912.12777, 2019.

- [17] Weinan E, Jiequn Han and Arnulf Jentzen, “Algorithms for solving high dimensional PDEs: From nonlinear Monte Carlo to machine learning”, <https://arxiv.org/pdf/2008.13333.pdf>, 2020.
- [18] Weinan E, Jiequn Han and Linfeng Zhang, “Integrating machine learning with physics-based modeling”, <https://arxiv.org/pdf/2006.02619.pdf>, 2020.
- [19] Weinan E, Chao Ma, Stephan Wojtowytsch and Lei Wu, “Towards a mathematical understanding of machine learning: What is known and what is not”, <https://arxiv.org/pdf/2009.10713.pdf>.
- [20] Jiequn Han and Weinan E, “Deep learning approximation for stochastic control problems”, Deep Reinforcement Learning Workshop, NIPS (2016), <https://arxiv.org/pdf/1611.07422.pdf>.
- [21] Jiequn Han, Arnulf Jentzen and Weinan E, “Solving high-dimensional partial differential equations using deep learning”, *Proceedings of the National Academy of Sciences*, 115, 34 (2018), 8505-8510.
- [22] Jiequn Han, Chao Ma, Zheng Ma, Weinan E, “Uniformly Accurate Machine Learning Based Hydrodynamic Models for Kinetic Equations”, *Proceedings of the National Academy of Sciences*, 116 (44) 21983-21991; DOI: 10.1073/pnas.1909854116, 2019.
- [23] Jiequn Han, Linfeng Zhang, Roberto Car, and Weinan E, “Deep potential: a general representation of a many-body potential energy surface”, *Communications in Computational Physics*, 23(3):629-639, 2018.
- [24] Pierre C Hohenberg and Bertrand I Halperin, “Theory of dynamic critical phenomena”, *Reviews of Modern Physics*, 49(3):435, 1977.
- [25] Weile Jia, Han Wang, Mohan Chen, Denghui Lu, Jiduan Liu, Lin Lin, Roberto Car, Weinan E, and Linfeng Zhang, “Pushing the limit of molecular dynamics with ab initio accuracy to 100 million atoms with machine learning”, *arXiv:2005.00223*, 2020.
- [26] Jason M Klusowski and Andrew R Barron, “Risk bounds for high-dimensional ridge function combinations including neural networks”, *arXiv:1607.01434*, 2016.
- [27] Yann LeCun, “A theoretical framework for back propagation”, In: Touretzky, D., Hinton, G., Sejnowski, T. (eds.) *Proceedings of the 1988 connectionist models summer school*, Carnegie-Mellon University, Morgan Kaufmann, 1989.
- [28] Qianxiao Li, Long Chen, Cheng Tai and Weinan E, “Maximum Principle Based Algorithms for Deep Learning”, *Journal of Machine Learning Research*, vol.18, no.165, pp.1-29, 2018, <https://arxiv.org/pdf/1710.09513.pdf>.
- [29] Qianxiao Li and Shuji Hao, “An Optimal Control Approach to Deep Learning and Applications to Discrete-Weight Neural Networks”, *Proceedings of the 35th International Conference on Machine Learning*, 2018.

- [30] Denghui Lu, Han Wang, Mohan Chen, Jiduan Liu, Lin Lin, Roberto Car, Weinan E, Weile Jia, and Linfeng Zhang, “86 pflops deep potential molecular dynamics simulation of 100 million atoms with ab initio accuracy”, arXiv:2004.11658, 2020.
- [31] Chao Ma, Lei Wu, and Weinan E, “The quenching-activation behavior of the gradient descent dynamics for two-layer neural network models”, arXiv:2006.14450, 2020.
- [32] Song Mei, Andrea Montanari, and Phan-Minh Nguyen, “A mean field view of the landscape of two-layer neural networks”, *Proceedings of the National Academy of Sciences*, 115(33):E7665-E7671, 2018.
- [33] Etienne Pardoux and Shige Peng, “Backward stochastic differential equations and quasilinear parabolic partial differential equations”, in *Stochastic partial differential equations and their applications* (Charlotte, NC, 1991), vol. 176 of *Lecture Notes in Control and Inform. Sci.*, Springer, Berlin, 1992, pp. 200-217.
- [34] Grant Rotskoff and Eric Vanden-Eijnden, “Parameters as interacting particles: long time convergence and asymptotic error scaling of neural networks”, In *Advances in neural information processing systems*, pages 7146-7155, 2018.
- [35] Justin Sirignano and Konstantinos Spiliopoulos, “Mean field analysis of neural networks: A central limit theorem”, arXiv:1808.09372, 2018.
- [36] Hans Triebel, *Theory of Function Spaces*, Birkhäuser, 1983.
- [37] Han Wang, Linfeng Zhang and Weinan E, unpublished.
- [38] Linfeng Zhang, Jiequn Han, Han Wang, Roberto Car, and Weinan E, “Deep potential molecular dynamics: A scalable model with the accuracy of quantum mechanics”, *Physical Review Letters*, 120:143001, Apr 2018.
- [39] Linfeng Zhang, Jiequn Han, Han Wang, Wissam A Saidi, Roberto Car, and Weinan E, “End-to-end symmetry preserving inter-atomic potential energy model for finite and extended systems”, In *Advances of the Neural Information Processing Systems (NIPS)*, 2018.
- [40] Linfeng Zhang, De-Ye Lin, Han Wang, Roberto Car, and Weinan E, “Active learning of uniformly accurate interatomic potentials for materials simulation”, *Physical Review Materials*, 3(2):023804, 2019.
- [41] Linfeng Zhang, Han Wang and Weinan E, “Reinforced dynamics for the enhanced sampling in large atomic and molecular systems. I. Basic Methodology”, *J. Chem. Phys.*, vol. 148, pp.124113, 2018.