

# 机器学习和计算数学

渭南市<sup>1</sup>

普林斯顿大学和北京大数据研究所

纪念冯康教授 (1920–1993 年)

## 内容

1	导言	2
2	基于机器学习的计算 SCI 问题算法 恩斯	3
2.1	非线性多网格法和蛋白质折叠.....	4
2.2	具有从头计算精度的分子动力学.....	6
3	基于机器学习的高维问题算法 科学计算	9
3.1	随机控制.....	9
3.2	非线性抛物型 PDE.....	11
3.3	动力学方程模拟气体动力学的动量闭包.....	13
4	机器的数学理论	

## 摘要

基于神经网络的机器学习能够以前所未有的效率和精度在非常高的维度上逼近函数。这开辟了许多令人兴奋的新可能性,不仅在传统的人工智能领域,而且在科学计算和计算科学领域。同时,机器学习也获得了作为一套“黑匣子”类型的技巧的声誉,没有基本的原则。这是在机器学习方面取得进一步进展的真正障碍。

<sup>1</sup>渭南@math.princeton.edu

在本文中，我们试图解决以下两个非常重要的问题：（1）机器学习已经如何影响并将进一步影响计算数学、科学计算和计算科学？（2）计算数学，特别是数值分析如何影响机器学习？我们介绍在这些问题上取得的一些最重要进展。我们的希望是把事物放在一个有助于将机器学习与计算数学结合起来的角度的角度。

## 1 引言

基于神经网络的机器学习(ML)在传统人工智能中的各种任务上取得了令人印象深刻的成功。这包括对图像进行分类，生成新的图像，如（假）人脸，并玩复杂的游戏，如围棋。所有这些任务的一个共同特点是它们涉及到非常高维的对象。事实上，当用数学术语表示时，图像分类问题是一个将在图像空间上定义的高维函数逼近每个图像类别对应的离散值集的问题。输入空间的维数通常是图像中像素数的3倍，其中3是颜色空间的维数。图像生成问题是一个从未知的高维分布中生成样本的问题，给定一组来自该分布的样本。围棋博弈问题是在动态规划中求解一个类贝尔曼方程，因为最优策略满足这样一个方程。对于像Go这样复杂的游戏，这个类似于Bellman的方程是在一个巨大的空间上制定的。

所有这些都是通过使用现代机器学习技术精确地逼近高维函数的能力来实现的。这为解决遭受“维度诅咒”(CoD)的问题开辟了新的可能性：随着维度的增长，计算成本成倍增长。长期以来，这个CoD问题一直是科学界的一个基本障碍。

以数值求解偏微分方程(PDE)为例。利用传统的数值方法，如有限差分、有限元和谱方法，我们现在可以常规地求解三个空间维数加时间维数的PDE。目前在计算数学中研究的大多数PDE都属于这一类。众所周知的例子包括泊松方程、麦克斯韦方程、欧拉方程、Navier-Stokes方程和线性弹性的PDE。稀疏网格可以增加我们处理PDE的能力，例如8到10维。这使我们能够尝试解决问题，如简单分子的Boltzmann方程。但是当我们面对PDE时，我们完全迷失了，比如说在100维。这使得基本上不可能求解复杂分子的Fokker-Planck或Boltzmann方程、多体薛定谔方程或实际控制问题的Hamilton-Jacobi-Bellman方程。

这正是机器学习可以帮助的地方。事实上，从[20, 10, 21]的工作开始，基于机器学习的求解高维PDE和控制问题的数值算法一直是近年来科学计算中最令人兴奋的新发展之一，这为计算数学开辟了一系列新的可能性。我们请[17]回顾这一令人兴奋的发展。

解决PDE只是冰山一角。还有许多其他问题，其中CoD是主要障碍，包括：

- 经典的多体问题，例如。蛋白质折叠
- 湍流。尽管湍流可以用三维Navier-Stokes方程来建模，但它具有如此多的主动自由度，一个有效的湍流模型应该涉及许多变量。
- 固体力学。在固体力学中，我们甚至没有Navier-Stokes方程的模拟。为什么是这样？那么，真正的原因是固体的行为本质上是一个多尺度的问题，它涉及从原子化到宏观的尺度。
- 多尺度建模。事实上，没有区分尺度的大多数多尺度问题都属于这一类。一个直接的例子是聚合物流体或聚合物熔体的动力学。

*机器学习能帮助解决这些问题吗？更广泛地说，我们能否将机器学习成功扩展到传统AI*

之外？我们将设法使读者相信，许多问题确实如此。

基于神经网络的机器学习除了具有强大的功能外，还获得了一套技巧而不是一套系统的科学原理的声誉。它的性能敏感地取决于超参数的值，如网络宽度和深度、初始化、学习速率等。事实上，就在几年前，参数调整被认为是一门艺术。即使现在，一些任务仍然是这样。因此，一个自然的问题是：我们能否理解这些微妙之处，并提出更好的机器学习模型，其性能更健壮？

在本文中，我们回顾了在这两个问题上学到的东西。我们讨论了机器学习已经或将对计算数学产生的影响，以及如何利用计算数学的思想，特别是数值分析，帮助理解和更好地制定机器学习模型。在此基础上，我们将主要讨论基于 ML 算法的新问题。尽管机器学习也提出了解决计算数学中一些传统问题的新方法，但在这方面不会说太多。

## 2 基于机器学习的计算科学问题算法

在本节和下一节中，我们将讨论如何使用神经网络模型来开发新的算法。对于不熟悉神经网络的读者来说，只要把它们看作是多项式的一些替换。之后我们将讨论神经网络。

### 2.1 非线性多网格法和蛋白质折叠

在传统的多网格方法[5]中，对于求解由线性椭圆方程的某种有限差分或有限元离散而产生的方程的线性系统，我们的目标是使二次函数最小化

$$I_h(u_h) = \frac{1}{2} L_h u_h - F_h u_h$$

这里  $h$  是离散化的网格大小。多网格方法的基本思想是在求解这个问题和网格大小为  $H$  的粗网格上的一个约简问题之间迭代。为了做到这一点，我们需要以下几点

- 投影算子  $P: U_h \rightarrow U_H$ ，它将在精细网格上定义的函数映射到在粗网格上定义的函数。
- 标度  $H$ :  $L$  的有效算子  $L_H = P^T L P$ 。这定义了粗网格上的目标函数：

$$I_H(u_H) = \frac{1}{2} L_H u_H - F_H u_H$$

- 延长算子  $Q: U_H \rightarrow U_h$ ，它将在粗网格上定义的函数映射到在细网格上定义的函数。通常一个人可以把  $Q$  取为  $P^T$ 。

这里的关键思想是粗粒化，在细尺度和粗粒度问题之间迭代。粗粒度中的主要成分是一组粗粒度变量和有效的粗粒度问题。通过这种方式，这些显然是一般的想法，可以与各种各样的问题相关。然而，在实践中，困难在于如何获得有效的粗粒度问题，对于线性问题来说，这是一个微不足道的步骤，而这是机器学习可以帮助的地方。

我们将以蛋白质折叠问题为例来说明非线性问题的一般思想。

设  $\{x_j\}$  是蛋白质和周围溶剂中原子的位置， $U=U(\{x_j\})$  是组合蛋白质-溶剂体系的势能。势能由化学键、范德华相互作用、静电相互作用等引起的能量组成。蛋白质折叠问题是寻找能量  $U$  的“基态”：

“最小化”  $U$ 。

在这里，我们增加了引号，因为我们真正想做的是抽样分发

所=：“即顶所=(如 t

为了定义粗粒度问题，我们假设我们得到了一组集合变量： $s=(s_i, \dots, s_n)$ ,  $S_j=S_j(x_i, \dots, x_n)$ , ( $n < N$ )。一种可能是使用

二面角作为粗粒度变量。原则上，人们也可以使用机器学习方法来学习粗粒度变量的“最佳”集合，但这里不会追求这个方向。

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

定义了粗粒度变量后，有效的粗粒度问题只是与这组粗粒度变量相关的自由能：

与线性问题的情况不同，在目前的情况下，我们必须首先找到函数  $A$ 。

其思想是用神经网络近似  $A$ 。这里的问题是如何获取培训数据。

与大多数预先收集训练数据的标准机器学习问题相反，在计算科学和科学计算的应用中，随着学习的进行，训练数据被“即时”收集。这被称为“并发学习”协议[11]。在这方面，预先收集训练数据的标准机器学习问题是“顺序学习”的例子”。并行学习的关键问题是以最佳方式生成数据的有效算法。培训数据集一方面应具有足够的代表性，另一方面应尽可能小。

[11]提出了生成这类数据集的一般程序。它被称为 EELT（探索检查-标记-训练）算法，它包括以下步骤：

- 探索：探索  $S$  空间。这可以通过采样  $e^{-\beta A(\mathbf{s})}$  来完成用  $A$  的电流近似。
- 检查：对于每个探索的状态，决定是否应该标记该状态。这样做的一种方法是使用后验误差估计器。这种后验误差估计的一个可能是机器学习模型集合预测的方差，见[41]。
- 标记：计算平均力（例如使用约束分子动力学）

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

其中自由能  $A$  可以用标准热力学积分计算。

训练：训练适当的神经网络模型。要想建立一个好的神经网络模型，就必须考虑到问题中的对称性。例如，如果我们通过消除氢原子的位置来粗化水分子集合的全原子表示，那么所得到的系统的自由能函数应该具有排列对称性，在设计神经网络模型时应该考虑到这一点（见下一小节）。

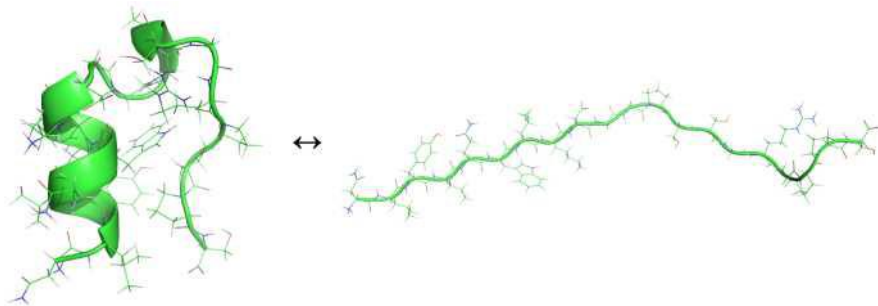


图 1: Trp 笼的折叠和扩展状态, 经[37]许可复制

这也可以看作是一种非线性的多网格算法, 因为它迭代粗粒度变量空间上的采样和全原子描述的 (受约束的) 吉布斯分布伽。

这是一个一般的过程, 应该适用于一大类非线性 “多网格” 问题。

图 1 显示的是 Trp 笼的扩展和折叠结构。这是一种含有 20 个氨基酸的小蛋白质。我们选择了 38 个二面角作为集体变量。 [37]给出了全部结果。

分子动力学是通过跟踪系统中所有核的轨迹来研究分子和材料系统行为的一种方法。假设原子核的动力学服从牛顿定律, 一些势能函数 (通常称为势能面或 PES)  $V$  来模拟原子核之间的有效相互作用:

## 2.2 具有从头计算精度的分子动力学

$$d^2X/dt^2 = -\nabla V(\mathbf{X}), \quad \mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N),$$

怎样才能得到函数  $V$ ? 传统上, 有两种相当不同的方法。第一种是利用量子力学模型计算飞行中的原子间力 ( $-W$ ), 最流行的是密度泛函理论 (DFT)。这被称为 Car-Parrinello 分子动力学或从头算分子动力学 [6, 8]。这种方法是相当准确的, 但也是非常昂贵的, 限制了系统的大小, 一个人可以处理大约 1000 个原子, 即使高性能超级计算机。另一种方法是提出经验潜力。基本上, 一个人猜测一个函数形式的  $V$  与一组小的拟合参数, 然后由一组小的数据确定。这种方法非常有效, 但不可靠。这种准确性和效率之间的困境长期以来一直是分子动力学的一个重要障碍。

对于机器学习, 新的范式是使用 DFT 生成数据, 然后使用机器学习生成对  $V$  的近似。这种方法有

系统		Al		镁		铝镁合金		
类型	格	#原子	奥菲	#数据	#C6ins	#数	奥菲	#数据
	联邦通信委	3	15, 174, 000	1, 32	15, 174, 000	860	39, 266, 46	7, 313
	HCP	1	15, 174, 000	908	15, 174, 000	760	18, 999, 90	2, 461
	钻石	1	5, 058, 000	1, 02	5, 058, 000	543	5, 451, 30	2, 607
	sc	8	5, 058, 000	713	5, 058, 000	234	2, 543, 94	667

地面	FCC (100)	1	3,270,960	728	3,270,960	251	62,203,68	1,131
	FCC (110)	16 <sup>a</sup> , 20 <sup>b</sup>	3,270,960	838	3,270,960	353	10,744,272	2,435
	FCC (111)	1	3,270,960	544	3,270,960	230	62,203,68	1,160
	HCP (0001)	1	3,270,960	3	3,270,960	109	62,203,68	176
	HCP (1010)	1	3,270,960	7	3,270,960	167	62,203,68	203
	HCP (1120)	16 <sup>a</sup> , 20 <sup>b</sup>	3,270,960	293	3,270,960	182	107,442,72	501
和			60,089,760	6,489	60,089,760	3,689	529,961,76	18,654

<sup>a</sup> 纯铝

<sup>b</sup> 镁和铝镁合金

图 2: EELT 算法的结果: 探索的配置数量与标记的数据点的数量。只有很小比例的配置被标记。经张林峰许可转载。 另见[40]。

产生与  $V$  的近似的潜力, 它与 DFT 模型一样精确, 并且与经验电位一样有效。

为了实现这一目标, 我们必须解决两个问题。一是数据的生成。二是提出了合适的神经网络模型。这两个问题是我们在这里讨论的所有问题的共同特征。

自适应数据生成的问题与以前非常相同。仍然可以使用 EELT 程序。每个步骤如何实现的细节有点不同。详情请参阅[40]。

图 2 显示了使用 EELT 算法的效果。正如人们所看到的, 探索的配置中很小一部分实际上被标记。对于 Al-Mg 示例, 只选择 $\sim 0.005\%$ 的配置进行标记。

对于第二个问题, 设计合适的神经网络, 最重要的考虑因素是:

1. 扩展, 神经网络应该是广泛的, 因为如果我们想扩展系统, 我们只需要相应地扩展神经网络。Behler 和 Parrinello[4]提出了实现这一目标的一种方法。
2. 保持对称性。除了通常的平移对称性和旋转对称性外, 还具有置换对称性: 如果我们重新标记一个铜原子系统, 它的势能不应该改变。它在神经网络模型的准确性方面有很大的不同, 是否考虑到这些对称性 (见[23]和图 3)。

解决对称性问题的一种非常好的方法是将神经网络模型设计为两个网络的组成: 嵌入网络和拟合网络。嵌入网络的任务是表示足够多的保持对称性的函数, 以便输入到拟合网络[39]中。

通过适当地解决这些问题, 我们可以提出非常满意的基于神经网络的  $V$  表示 (见图 4)。这种表示被命名为深度势[23, 39], 基于深度势的分子动力学被命名为 DeePMD[38]。正如最近在 [30, 25]中所证明的那样, DeePMD 结合了状态

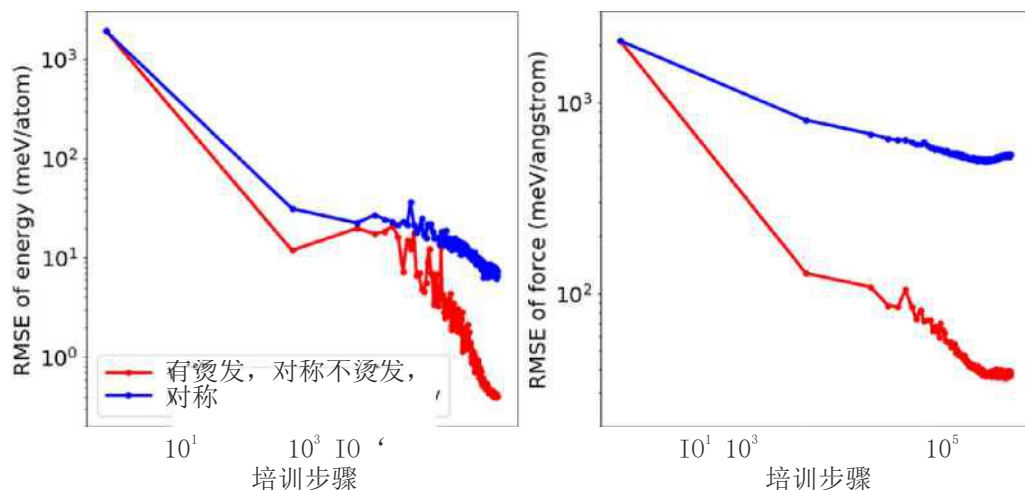


图 3: 对称保存对测试精度的影响。红色显示是穷人强加对称的结果（解释见正文）。可以看出，测试精度大大提高。 经张林峰许可转载。

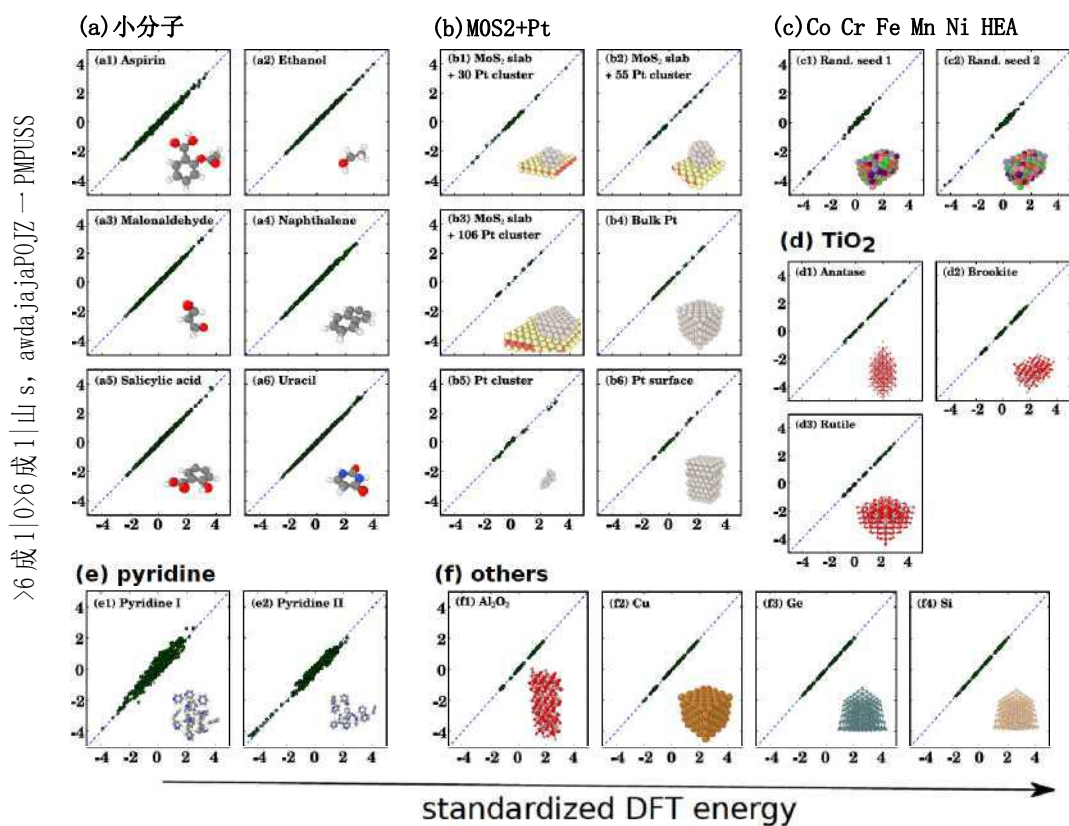


图 4: 各种系统的深势测试精度。 经张林峰许可转载。 另见[39]。





艺术高性能超级计算机，可以帮助增加系统的大小，一个人可以建模的从头精度 5 个数量级。

### 3 基于机器学习的科学计算高维问题算法

#### 3.1 随机控制

[20]介绍了机器学习在科学计算中解决高维问题的第一个应用。随机控制由于其与机器学习的相似性而被选择为第一个例子。考虑离散随机动力系统：

$$s_{t+1} = s_t + b(s_t, a_t) + \epsilon_t$$

这里  $s_t$  和  $a_t$  分别是时间  $t$  的状态和控制， $\epsilon_t$  是时间  $t$  的噪声。我们的目标是解决：

$$\min_{a_t} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t (c(s_t, a_t) + V(s_{t+1}) - V(s_t)) \right] \quad (2)$$

在一组反馈控件中：

$$a_t = \pi(s_t) \quad (3)$$

我们用神经网络模型逼近函数：

$$\pi(s_t) \approx \pi(s_t; \theta), \quad t=0, \dots, T-1 \quad (4)$$

优化问题 (2) 则变为：

$$\min_{\theta} \mathbb{E} \left[ \sum_{t=0}^{T-1} \gamma^t (c(s_t, \pi(s_t; \theta)) + V(s_{t+1}) - V(s_t)) \right] \quad (5)$$

与标准监督学习的情况不同，这里我们有  $T$  组神经网络要同时训练。网络架构如图 5 所示

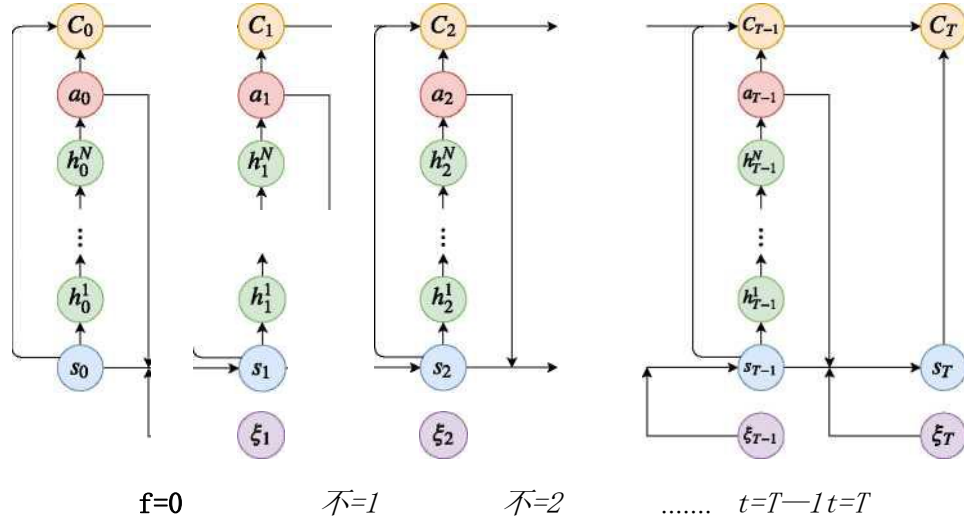


图 5: 求解离散时间随机控制的网络体系结构。整个网络总共有  $(N+1)T$  层, 涉及自由参数, 需要同时优化。每一列 (除  $\xi$ ) 对应于一个子网络在  $T$ 。转载得到了杰群汉的许可。另见[20]。

与机器学习的标准设置相比, 可以看到一个明确的类比, 其中 (1) 对残差网络起作用, 噪声  $\{\xi_t\}$  起着数据的作用。事实上, 随机梯度下降 (SGD) 可以很容易地用于解决优化问题 (5)。

该算法的应用示例如图 6 所示, 用于多个设备的储能问题。这里  $n$  是设备的数量。更多细节, 我们参考[20]。

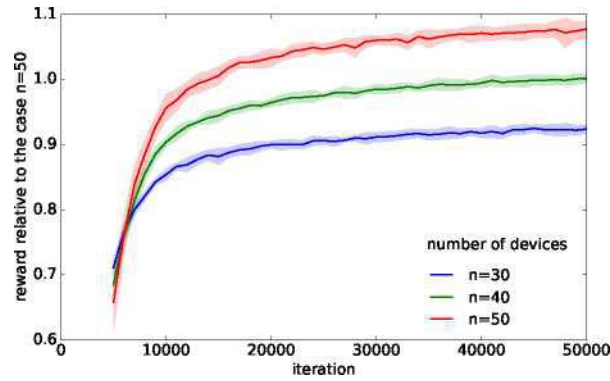


图 6: 储能问题的相对奖励。控制函数的空间为  $\mathbb{R}^{n+2} \times \mathbb{R}^{3n}$  对于  $n=30, 40, 50$ , 具有多重平等和不平等约束。经韩洁群许可转载。另见[20]。

### 3.2 非线性抛物线 PDE

考虑抛物型 PDE 的形式:

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

我们研究的是一个终端价值问题，而不是初始值问题，因为我们考虑的主要应用之一是金融。为了开发基于机器学习的算法，我们首先将其重新定义为一个随机优化问题。这可以使用反向随机微分方程 (BSDE) [33] 来完成。

$$\inf_{\pi} \mathbb{E} [ |g(X_T) - Y_T|^2 ], \quad (6)$$

$$s. \text{。} \quad X = e^{\int_0^t \mu_s ds + \int_0^t \sigma_s dW_s}, \quad \mathbb{E}(s, X_s) dW_s = 0 \quad (7)$$

$$Y_t = \int_t^T \rho_s ds + g(X_T), \quad \mathbb{E}(Z | \mathcal{F}_t) = 0 \quad (8)$$

可以看出，这个问题的唯一极小值是 PDE 的解有:

$$u = u(t, X_t) \quad \text{还有} \quad Z_t = \sigma(t, X_t) \nabla u(t, X_t). \quad (9)$$

利用这种公式，可以沿着以下路线开发一种基于机器学习的算法，采用前面 [10, 21] 讨论的随机控制问题的思想:

时间离散后，近似未知函数

$$u(0, x_0) \quad \text{还有} \quad X_t = e^{\int_0^t \mu_s ds + \int_0^t \sigma_s dW_s} \quad (10)$$

通过前馈神经网络明和

- 使用 BSDE，构造一个逼近  $u$ ，它接受路径  $\{X_t\}_{0 \leq t \leq T}$  和  $\{W_t\}_{0 \leq t \leq T}$  作为输入数据，并给出最终输出，用  $u(\{X_t\})$  表示  $u(t, X_t)$  的近似值。
- 在  $u$  与给定终端条件之间匹配的误差定义了预期损失函数

$$J(u) = \mathbb{E} [ |g(X_T) - u(\{X_t\})|^2 ]$$

这种算法被称为深度 BSDE 方法。

作为应用，让我们首先研究一个随机控制问题，但我们现在使用 Hamilton-Jacobi-Bellman (HJB) 方程来解决这个问题。考虑维数  $d=100$  的著名 LQ-G (线性二次高斯) 问题:

$$dX_t = (A_t X_t + B_t u_t) dt + \sigma_t dW_t, \quad (10)$$

成本函数:  $J(\{m_i\}_{0 \leq i \leq T}) = \mathbb{E}[\int_0^T \sigma^2(X_t) dt + g(X_T)]$ 。 相应的

给出了 HJB 方程

$$-\partial_t u - \frac{1}{2} \sigma^2(x) \partial_x^2 u + \lambda(u - \mathbb{E}[u]) = 0 \tag{11}$$

利用 Hopf-Cole 变换, 可以用形式表示解:

$$u(t, x) = -\lambda \ln \left( \mathbb{E}[\exp(-\lambda g(X_T)) | \mathcal{F}_t] \right). \tag{12}$$

这可以用来校准准确性 的 深 B SDE 方法。

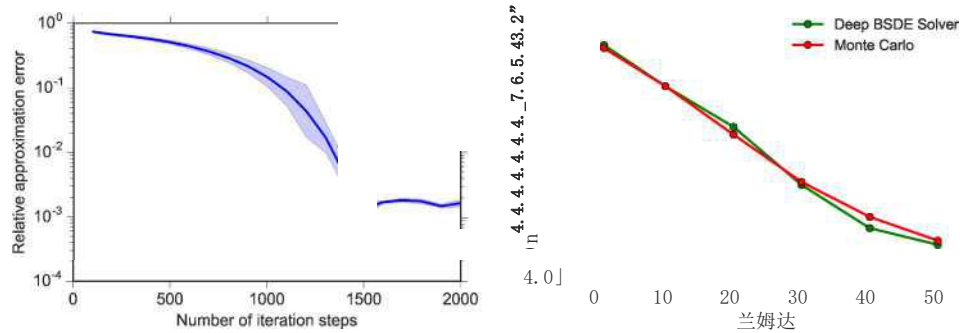
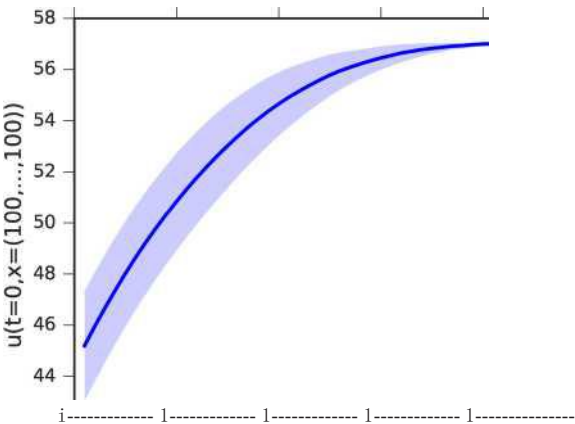


图 7: 左: 当  $\lambda=1$  时,  $u(t=0, x=(0, \dots, 0))$  的深 B SDE 方法的相对误差, 在 330 秒的运行时达到 0.17。右: 针对不同  $\lambda$  的最优成本  $u(t=0, x=(0, \dots, 0))$ 。经韩洁群许可转载。 另见[21]。

作为第二个例子, 我们研究了默认风险的 Black-Scholes 方程:

$$\partial_t u + \frac{1}{2} \sigma^2(x) \partial_x^2 u + \lambda(u - \mathbb{E}[u]) - r u(t, x) = 0$$

其中  $Q$  是一些非线性函数。 这种建模默认风险的形式是针对低维情况提出和/或在文献中使用的 ( $d=5$ , 参考见[21])。 在  $d=100$ [21]的情况下, 对这一问题采用了深 B SDE 方法。



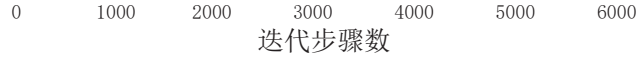


图 8: 默认风险在  $d=100$  的 Black-Scholes 方程的解。深度 BSDE 方法在运行时 617 秒内实现了 0.46% 大小的相对误差。 经韩洁群许可转载。 另见[21]。

深 B SDE 方法已应用于定价篮子选项、利息评级依赖选项、Libor 市场模型、百慕大交换、障碍选项（参考见[17]）。

### 3.3 气体动力学动力学方程的动量闭包

气体的动力学可以用著名的 Boltzmann 方程非常精确地建模：

$$D_t f + v \cdot \nabla_x f = -Q(f), \quad v \in \mathbb{R}^3, \quad x \in \mathbb{Q} \subset \mathbb{R}^3, \quad (13)$$

其中  $f$  是相位空间密度函数， $g$  是 Knudsen 数：

意思是自由的道路

宏观长度

$q$  是模拟气体粒子碰撞过程的碰撞算子。 当  $g \ll 1$  时，这可以用欧拉方程近似：

$$d_t U + \nabla_x \cdot F(U) = 0, \quad (14)$$

在哪里

$$U = \left( \int_{\mathbb{R}^3} f dv, \int_{\mathbb{R}^3} p f dv, \int_{\mathbb{R}^3} p^2 f dv \right)^T, \quad p = |v|,$$

还有

$$F(U) = \left( \int_{\mathbb{R}^3} p^2 f dv, \int_{\mathbb{R}^3} p^3 f dv, \int_{\mathbb{R}^3} p^4 f dv \right)^T, \quad (E + p)u = \int_{\mathbb{R}^3} p^4 f dv.$$

通过将玻尔兹曼方程投影到所涉及的低阶矩上, 并利用分布函数  $f$  接近局部 Maxwellian 的 ansatz, 可以得到欧拉方程。



图 9: 气体动力学的不同状态。 经韩洁群许可转载。 另见[22]。

当 $\epsilon$ 不小的时候会发生什么？ 在这种情况下，一个自然的想法是寻求一些推广欧拉方程使用更多的矩。这个程序是由哈罗德·格拉德发起的，他利用 $\{1, v, (v-u) \otimes (v-u), |v-u|\}$ 的时刻构造了著名的 13 矩系统 $^{(v-u)}$ 。这一行工作遇到了几个困难。首先，不能保证得到的方程是合适的。其次，总是有“闭包问题”：当将玻尔兹曼方程投影到一组矩上时，总是有一些术语涉及到所考虑的矩集之外的矩。为了获得一个封闭的系统，人们需要以某种方式对这些术语进行建模。对于欧拉方程，这是使用局部麦克斯韦近似来完成的。当 $\epsilon$ 不小时，这是准确的，但当 $\epsilon$ 不小时，就不再是这样了。目前还不清楚应该用什么来代替。

在[22]，Han 等人开发了一种基于机器学习的矩量法。总体目标是构造一个均匀精确（广义）弯矩模型。方法包括两个步骤：

1: 通过自动编码器学习一组最优广义矩。在这里，通过最优性，我们意味着广义矩集合保留了关于原始分布函数的最大信息量，并且可以用于以最小的精度损失恢复分布函数。这可以做如下：找到编码器中和解码器小，从  $U, W$  恢复原始  $f$

$$\hat{W} = \hat{f}(f) = \int w f dv, \quad \forall U, W \quad (v) = h(v; U, W).$$

$$\text{尽量减少 } \|W, h_e\|_f^2 \text{ 以中 } (f) \text{ 中 } \|\cdot\|^2.$$

$U$  和  $W$  构成了一组广义水动力变量，我们将用来模拟气体流动。

2: 学习 PDE 中投影 PDE 的通量和源项。通过在这组（广义）矩上正式投影玻尔兹曼方程，可以得到  $U$  和  $W$  的有效 PDE。这给了我们一组形式的 PDE：

$$\begin{aligned} \epsilon \frac{DT}{Dt} U + V_x - F(U, W; \epsilon) &= 0, \\ \text{我知道 } W + V_x - G(U, W; \epsilon) &= R(U, W; \epsilon). \end{aligned} \quad (15)$$

其中  $F(U, W; \epsilon) = \int v U f dv$ ,  $G(U, W; \epsilon) = \int v W f dv$ ,  $R(U, W; \epsilon) = \epsilon^{-1} \int f W Q(f) dv$ 。我们现在的任务是从原始动力学方程中学习  $F, G, R$ 。

同样重要的问题是（1）获得最优数据集，（2）强制物理约束。这里两个值得注意的物理约束是（1）守恒定律和（2）对称性。在这种方法中，自然会尊重保护法律。关于对称性，除了通常的静态对称性外，现在还有一种新的动态对称性：伽利略不变性。这些问题都在[22]讨论。我们还参考[22]的数值结果的模型，以获得这种方式。

## 4 机器学习的数学理论

虽然基于神经网络的机器学习已经展示了广泛的非常令人印象深刻的成功，但它也获得了“黑魔法”的声誉，而不是坚实的科学技术。这是由于（1）我们对其成功背后的基本原因缺乏基本的理解；（2）这些模型和算法的性能对超参数的选择非常敏感，例如网络的体系结构和学习速率；（3）

一些技术，如批量归一化，似乎是一种黑魔法。

为了改变这种情况，我们需要（1）提高我们对基于神经网络的模型和算法成功背后的原因和脆弱性的理解；（2）寻找方法来制定更稳健的模型和设计更稳健的算法。在本节中，我们讨论第一个问题。下一节将专门讨论第二个问题。

下面列出我们需要解决的最基本问题：

- 为什么它在如此高的维度上工作？
- 为什么简单梯度下降用于训练神经网络模型？
- 过度标准化是好是坏？
- 为什么神经网络建模需要如此广泛的参数调优？

在这一点上，我们尚未对所有这些问题有明确的答案。但一些连贯的画面开始浮现。我们将重点讨论监督学习的问题，即使用有限数据集逼近目标函数。为了简单起见，我们将把自己限制在物理感兴趣的领域是  $X=[0, 1]$  的情况下<sup>4</sup>。

## 4.1 介绍基于神经网络的监督学习

监督学习的基本问题如下：给定输入输出数据对的自然数  $n \in \mathbb{N}$  和序列  $\{(X_j) = (x_j, f^*(X_j)), j \in \{1, 2, \dots, n\}$ ，我们希望尽可能准确地恢复目标函数  $f$ 。我们将假设输入数据  $\{x_j, j \in \{1, 2, \dots, n\}$ ，是从  $\mathbb{R}$  上的概率分布中采样的<sup>4</sup>。

第一步。选择一个假设空间。这是一组试验函数，其中  $m, n$  是  $H$  的维数<sup>4</sup>。人们可以选择分段多项式或小波。



在现代机器学习中，最流行的选择是神经网络函数。 双层神经网络函数（一个输入层，一个通常不计算的输出层，一个隐藏层）采取这种形式

$$f(\mathbf{x}) = \sum_{j=1}^M a_j \phi_j(\mathbf{x}) \quad (16)$$

其中  $a_j: \mathbb{R} \rightarrow \mathbb{R}$  是一个固定的标量非线性函数，其中  $0 = \{(\mathbf{A}_j, \mathbf{W}_j)\}$  是要优化的参数（或训练的参数）。非线性函数  $a_j: \mathbb{R} \rightarrow \mathbb{R}$  的一个流行例子是 ReLU（校正线性单元）激活函数： $a(z) = \max\{z, 0\}$ ，对于所有  $z \in \mathbb{R}$ 。我们将把注意力限制在这个激活函数上。粗略地说，如果一个人多次组合两层神经网络函数，则得到深度神经网络 (DNN) 函数。 其中一类重要的 DNN 模型是残差神经网络 (ResNet)。 它

$$f(\mathbf{x}) = \sum_{j=1}^M a_j(\mathbf{x}) \phi_j(\mathbf{x}), \quad \phi_j(\mathbf{x}) = \max\{0, \mathbf{W}_j^T \mathbf{x} + b_j\} \quad (17)$$

它与离散常微分方程非常相似，并采取这种形式

对于  $l \in \{0, 1, \dots, L-1\}$ ，其中  $L, M, N$ 。这里的参数是  $0 = (\mathbf{a}, \mathbf{V}, \{\phi_j\})$ ， $(\mathbf{W}_j, \mathbf{b}_j)$ 。 ResNet 是真正深层神经网络模型的首选模型。

第二步。 选择一个损失函数。 损失函数的选择的首要考虑是拟合数据。 因此，最明显的选择是  $L^2$  损失：

$$L(\mathbf{f}) = \frac{1}{n} \sum_{i=1}^n |f(\mathbf{x}_i) - y_i|^2 = \frac{1}{n} \sum_{i=1}^n |f(\mathbf{x}_i) - y_i|^2 \quad (18)$$

这也被称为“经验风险”。 有时我们还添加正则化项。

第三步。 选择一个优化算法。 机器学习中最流行的优化算法是不同版本的梯度下降 (GD) 算法，或其随机模拟，随机梯度下降 (SGD) 算法。 假设我们目标最小化的目标函数是形式的

$$F(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; \mathbf{x}_i, y_i) \quad (19)$$

最简单的 SGD 迭代形式

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

对于  $k \in \mathbb{N}$ ，其中  $\{\eta_k\}_{k \in \mathbb{N}}$  是一个 i.i.d 序列。 从分布  $\nu$  和  $n$  中采样的随机变量是在迭代过程中也可能发生变化的学习速率。 相比之下，GD 采取的形式

$$\theta_{k+1} = \theta_k - \eta \nabla F(\theta_k) \quad (21)$$

显然，这种形式的 SGD 可以适应形式 (18) 的损失函数，它可以看作是一个期望， $\nu$  是训练数据集上的经验分布。 这种 DNN-SGD 范式是现代机器学习的核心。

## 4.2 近似理论

逼近函数最简单的方法是使用多项式。 对于多项式逼近，有两种定理。 第一个是 Weierstrass 定理，它断言连续函数可以用紧域上的多项式统一逼近。 其次是泰勒定理，它告诉我们收敛速度取决于目标函数的光滑性。

利用神经网络理论中的术语，Weierstrass 定理是“普遍逼近定理”(UAT)。这是一个有用的事实。但泰勒定理更有用，因为它告诉我们一些关于收敛速度的东西。近似理论中使用的泰勒定理的形式是直接和逆近似定理，它断言给定的函数可以用具有特定速率的多项式逼近，当且仅当该函数的某些范数是有限的。这个特殊的范数，它测量函数的规律性，是表征这个近似方案的关键量。对于分段多项式，这些范数是一些贝索夫空间规范[36]。为了  $L^2$  一个典型的结果如下：

$$\|f - p_m\|_{L^2(X)} \leq C \|f\|_{H^s} m^{-s/d} \quad (22)$$

给你  $s$  代表第一阶 Sobolev 范数[36]， $m$  是自由度的数目。在规则网格上，网格大小由

$$m \propto N^d \quad (23)$$

需要注意的一件重要的事情是，(22) 中的收敛速度受到 CoD 的影响：如果我们想将误差减少一个  $e$  因数，我们需要将  $m$  增加一个因数  $m \leftarrow m e^{d/e}$  如果  $e=1$ 。对于  $d=100$ ，根据机器学习标准，它不是很高的维数，这意味着我们必须将  $m$  增加一个  $e$  的因子  $^{100}$ 。这就是为什么多项式和分段多项式在高维中不有用的原因。

另一种理解这一点的方法如下。度  $p$  在维数  $d$  中的单项式个数为  $C: \frac{d!}{p!(d-p)!}$ 。对于  $d$  和  $p$  的大值来说，这增长非常快。

在高维中我们应该期望什么？我们可以借鉴的一个例子是用于集成的蒙特卡罗方法。考虑近似问题

$$I(g) = \int_{\mathcal{X}} g(x) p(x) dx$$

使用

$$\hat{I}_m(g) = \frac{1}{m} \sum_{j=1}^m g(x_j)$$

其中  $\{x_j\}_{j=1}^m$  是直接计算给出的概率分布  $p$  的一组 i. i. d 样本。

$$E[(\hat{I}_m(g) - I(g))^2] = \frac{1}{m} \text{var}(g) = \frac{1}{m} (E[g^2(x)] - (E[g(x)])^2)$$

这个确切的关系告诉我们两件事。(1) 蒙特卡罗积分的收敛速度与维数无关。(2) 误差常数由被积方差给出。因此，为了减少误差，必须做方差减少。

如果我们使用基于网格的正交规则，精度也会受到 CoD 的影响。

通过更复杂的选择  $\{x_j, j \in [m]\}$  的方法，例如使用基于数论的求积规则，可以提高蒙特卡罗率。但这些通常会导致收敛速度的  $O(1/d)$  提高，并且随着  $dT_x$  的减小而减小。

基于这些考虑，我们的目标是找到满足：的函数近似：

$$\|f - f_m\|_{L^2(X)} \leq \frac{1}{\sqrt{m}} \sum_{j=1}^m |a(\omega_j)| e^{i(\omega_j, x)}$$

自然的问题是：

我们如何做到这一点？即我们应该选择什么样的假设空间？

什么应该是“规范”  $\|\cdot\|$  (与  $H$  的选择有关)？在这里，我们把规范放在引号中，因为它不必是真正的规范。我们只需要它控制近似误差。

关于第一个问题，让我们举一个说明性的例子。

$$f(x) = \int_{\mathbb{R}^d} a(\omega) e^{i(\omega, x)} d\mu(\omega)$$

考虑函数  $f$  及其近似  $f_m$  的傅里叶表示 (例如基于 FFT 的)：

这里  $\{\omega_j\}$  是一个固定的网格，例如，均匀的网格。对于这个近似，我们有  $\|f - f_m\|_{L^2(X)} \leq C \frac{1}{\sqrt{m}} \sum_{j=1}^m |a(\omega_j)| e^{i(\omega_j, x)}$  患有 CoD。

现在考虑替代代表

$$f(x) = \int_{\mathbb{R}^d} a(\omega) e^{i(\omega, x)} d\mu(\omega) = \int_{\mathbb{R}^d} a(\omega) e^{i(\omega, x)} d\nu(\omega) \quad (24)$$

其中  $\nu$  是概率分布。现在要近似  $f$ ，使用蒙特卡罗是很自然的。设  $\{\omega_j\}$  为 i. i. d。样本  $n$ ， $FM(X) = \frac{1}{n} \sum_{j=1}^n a(\omega_j) e^{i(\omega_j, x)}$  那我们就有了

$$E \|f(x) - FM(X)\|^2 =$$

这种近似不受 CoD 的影响。请注意， $FM(X)$  和工业 1AJA( $\hat{J}X$ ) 只不过是一个具有激活函数  $a(Z) = \max(0, Z)$  的两层神经网络 (这里是  $AJ = \max(0, \cdot)$ )。

我们认为，这个简单的论点确实是为什么神经网络模型在高维中表现如此出色的核心。

现在让我们来举一个神经网络模型近似理论的具体例子。我们将考虑两层神经网络。

$$f_m(x) = \sum_{j=1}^m a(w_j \cdot x) \phi_j \quad \{(a_j, w_j), j \in [m]\}$$

考虑函数  $f: X=[0, 1]^d$  以下表格的 TR

$$f(x) = \int_Q AA(Wx)p(da, dw) = E_{(a,w) \sim p}[aa(w^T x)], \quad x$$

其中  $Q=\mathbb{R}^1 \times \mathbb{R}^{d-1}$ ,  $p$  是  $Q$  上的概率分布。定义:

$$\|f\|_B = \left( \int_Q p(dw) \|f(w)\|^2 \right)^{1/2}$$

其中  $Pf = \{p: f(x) = E_p[aa(Wx)]\}$ 。这被称为巴伦规范[2, 12, 13]。空间

$$B = \{f \in C^0 \mid \|f\|_B < \infty\}$$

被称为 Barron 空间[2, 12, 13] (另见[3, 26, 14])。

类比经典逼近理论, 还可以证明一些直接逼近定理和逆近似定理[13]。

**定理 1 (直接逼近定理)** 如果  $\|f\|_B < \infty$ , 那么对于任意整数  $m > 0$ , 存在两层神经网络函数  $f_m$  就

$$f_m(x) = \int_Q p(dw) \sum_{i=1}^m a_i(w^T x)$$

定理 2 (逆逼近定理) 设

$$NC = \{f \in C^0 \mid \|f\|_B < \infty\}$$

设  $f^*$  为连续函数。假设存在一个常数  $C$  和一个函数序列  $f_m \in NC$ , 使得

对于所有的  $x \in X$ 。然后在  $Q$  上存在概率分布  $p^*$ , 这样

$$f^*(x) = \int_Q aa(w^T x) p^*(da, dw),$$

对于所有  $x \in X$  和  $\|f\|_B < \infty$ ,  $m, C$ 。

## 4.3 估计误差

我们必须担心的另一个问题是训练数据集之外的机器学习模型的性能。这一问题也出现在经典近似理论中。图 10 说明了在均匀网格上多项式插值的经典 Runge 现象。可以看出, 远离网格点, 插值的误差可能很大。这是我们希望避免的情况。

我们在实践中所做的是尽量减少训练误差:

$$R_n(f) = E(f(x_j) - f^*(x_j))^2$$

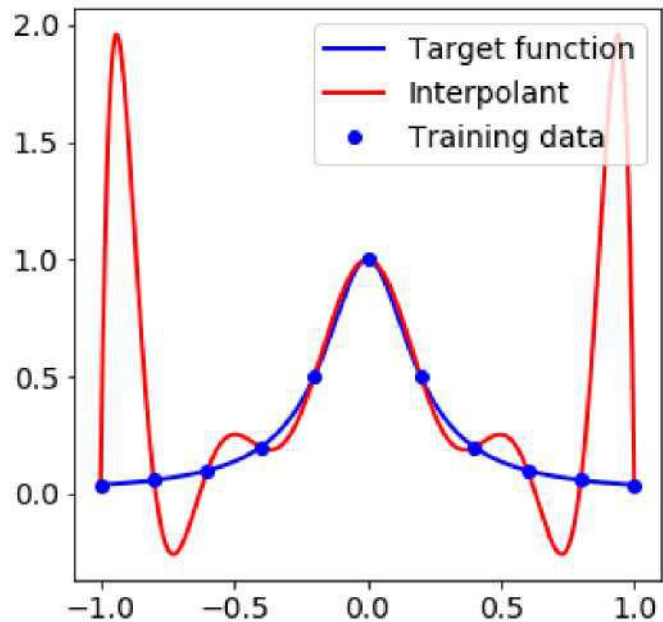


图 10: 龙格现象:  $f^*(X)$  搭训 Reproduced 得到潮马的许可。

但我们对测试错误感兴趣，这是人口风险的抽样版本：

$$M(f) = E \int_{\mathcal{X}} f^*(x)^2 dx$$

问题是我们如何控制这两个错误之间的差异。

这样做的一个方法是使用拉德马赫复杂性的概念。对我们来说，这里的重要事实是，拉德马赫的复杂性控制了训练和测试错误之间的差异（也称为“泛化差距”）。事实上，让  $H$  是一组

$$H = \{h(x) = \sum_{i=1}^n \alpha_i \phi_i(x) \mid \alpha_i \in \{-1, 1\}\}$$

函数， $S = (x_1, x_2, \dots, x_n)$  成为一个数据集。然后，直到对数项，我们有

$$R(H) \leq E \sum_{i=1}^n \sum_{j=1}^n \mathbb{E} [h(x_i) h(x_j)] \quad (25)$$

其中， $H$  对  $S$  的拉德马赫复杂性被定义为

在那里  $\{x_i\}_{i=1}^n$  是 i. i. d 的。随机变量取值  $\pm 1$ ，概率相等。

然后问题变成了约束一个假设空间的拉德马赫复杂性。对于巴伦空间，我们有[2]：

定理 3。 设  $f \in q = \{f \in B, \|f\|/b < q\}$ 。 我们就有了

雷达 (FQ)  $< 2Q^n H^n$  其中  $n = |S|$ , 数据集  $S$  的大小。

## 4.4 正规化模型的先验估计

考虑正则化模型

$$L_n(\theta) = r_n(\theta) + \lambda \sum_{i=1}^n |w_i|, R = \arg\min L_n(\theta)$$

其中路径范数由：定义：

$$\|w\|_1 = \sum_{i=1}^n |w_i| \quad (26)$$

[12]证明了以下结果：

定理 4。 假设  $f^* \in XT[0, 1] \in B$ 。 存在常数  $C_0$ , 因此对于任何常数  $\delta > 0$ , 如果  $A > C_0$ , 那么概率至少为  $1 - \delta$  超过训练数据集的选择, 我们有

$$L_n(f^*) \leq L_n(\hat{f}) + O\left(\frac{\log(2d)}{n}\right) + O\left(\frac{\log(n)}{n}\right).$$

类似的近似理论和先验误差估计已经被证明适用于其他机器学习模型。 以下是对这些结果的简要总结。

随机特征模型：相应的函数空间是再现核 Hilbert 空间 (RKHS)。

剩余网络 (Res 网)：相应的函数空间是在 [13] 中引入的所谓流动诱导空间。

多层神经网络：相应的函数空间的候选是 [15] 中引入的多层空间。

真正重要的是控制近似误差和泛化间隙的“规范”。这些量是为相应空间中的函数定义的。在逼近定理和 Rademach 复杂度估计到位后，人们可以很容易地证明正则化模型的以下类型的定理：直到对数项，正则化模型的最小化子满足：

$$RF \leq \frac{1}{\sqrt{mn}} + \frac{1}{\sqrt{n}}$$

其中  $m$  是自由参数的数目， $n$  是训练数据集的大小。 请注意，对于多层空间，[15] 中证明的结果不是那么尖锐。

我们只讨论了假设空间的分析。 还有许多其他问题。 我们参考 [19] 来更多地讨论目前对基于神经网络的机器学习的理解。

## 从连续的角度 5 机器学习

现在我们转向机器学习的替代配方。 由于 PDE 的情况，我们希望首先在连续设置中制定机器学习，然后离散化，以获得具体的模型和算法。 这里的关键是，我们提出的连续问题应该是好的数学问题。 对于 PDE，这是通过要求它们“适口”来实现的”。 对于变分问题，我们要求问题在某种意义上是“凸”的，并且是较低的半连续的。 这些要求的要点是确保问题有唯一的解决

方案。直观地说，对于机器学习问题，“好”意味着变分问题应该有一个简单的景观。如何准确地表述这一点，是未来研究的一个重要问题。

正如[16]所指出的，连续配方的关键成分如下：

- 功能的表示（作为期望）
- 制定变分问题（作为期望）
- 优化，例如。 梯度流动

## 5.1 职能的表述

[16]中考虑了两种表示：基于积分变换的表示和基于流的表示。 最简单的基于积分变换的表示是（24）的推广）：

$$f(x; \theta) = \int a(w) a(w^T x) n(dw)$$

$$a(w) a(w^T x)$$

$$= \int p_{aa}(w^T x)$$

$$p^{(+)}, u)$$

这里  $\theta$  表示模型中的参数： $\theta$  可以是  $a(\cdot)$  或 Prob 分布  $n$  或  $p$ 。

这种表示对应于两层神经网络。 [15]对多层神经网络进行了推广。

接下来我们转向基于流的表示：

$$\frac{d}{dt} a(w, t) = -a(w, t) a(w^T z) \quad (27)$$

$$= E(a, w) \text{ 阵 } A A(w^T z) \quad (28)$$

$$= \text{阵} \odot (z, u), z(0, x) = x \quad (29)$$

$$\frac{d}{dt} (x, 0) = -1(1, x)$$

在此表示中，参数  $\theta$  可以是  $\{a_i\}$  或  $\{n_i\}$  或  $\{p_i\}$

## 5.2 随机优化问题

随机优化问题类型有：

随机算法是机器学习中的关键组成部分，可以很容易地解决这些问题。例如，代替梯度下降算法：

$$w_{k+1} = w_k - \frac{1}{n} \sum_{j=1}^n g(w_k, w_j)$$

人们可以使用随机梯度下降：

$$w_{k+1} = w_k - g(w_k, w_j)$$

其中  $\{w_k\}$  是从  $v$  采样的一组随机变量。

以下是现代机器学习中的随机优化问题的一些例子：

监督学习：在这种情况下，人口风险最小化

$$R(f) = \mathbb{E} \int (f(x) - f^*(x))^2 dx$$

- 量子多体哈密顿量的特征值问题：

$$H \psi = E \psi, \quad H = \sum_{i,j} t_{ij} a_i^\dagger a_j + \sum_i \epsilon_i a_i^\dagger a_i + \sum_{i,j,k} V_{ijk} a_i^\dagger a_j^\dagger a_k$$

这里  $H$  是量子系统的哈密顿量。

$$J_t = \sum_{i=1}^n \lambda_i \phi_i(x_t) \neq \mathbb{E} \left\{ \sum_{t=0}^{T-1} \gamma^t J_t \right\}$$

- 随机控制问题：

将前面讨论的表示替换为随机优化问题的这些表达式，我们得到了我们需要解决的最终变分问题。

人们可以直接对这些变分问题进行离散化，然后使用一些优化算法求解离散化问题，也可以写下一些优化算法的连续形式，通常是梯度流动动力学，然后对这些连续流动进行离散化。我们将讨论第二种方法。

## 5.3 优化：梯度流动

为了写出梯度流动的连续形式，我们从统计物理中得到了一些启示。以监督学习为例。我们认为人口风险是“自由能”，根据 Halperin 和 Hohenberg 的[24]，我们将参数分为保守和非保守两种。例如， $a$  是一个非保守参数， $n$  是保守的，因为它的总积分必须是 1。



对于非保守参数，如[24]所建议的，可以使用“模型 A”动力学：

$$\frac{da}{DT} = \frac{b}{\bar{E}} R$$

这只是通常的  $L^2$  梯度流动。

对于  $n$  这样的保守参数，应该使用“模型 B”动力学，其工作原理如下：首先定义“化学势”

$$\frac{b}{\bar{E}} R$$

由化学势得速度场  $v$  和电流  $J$ ：

$$J = NV, \quad v = -\nabla W$$

连续性方程给出了梯度流动动力学：

这也是 Wasserstein 度量下的梯度流。

## 5.4 分散梯度流动

为了获得实际的模型，需要对这些连续问题进行离散化。 第一步，利用训练数据用经验风险代替人口风险。

更重要的问题是如何离散参数空间中的梯度流。 参数空间具有以下特点：（1）它具有简单的几何形状-与实际空间不同，它可能具有复杂的几何形状。（2）它通常也是高维的。 由于这些原因，参数空间离散最自然的数值方法是粒子法，这是蒙特卡罗的动态版本。 雾化粒子法可能有助于提高性能。 在相对较低的维数中，由于参数空间的几何形状相对简单，人们也可以考虑谱方法，特别是谱方法的一些稀疏版本。

例如，基于积分变换的表示的保守流的离散化。 用表示法：  $f(x; 0) = E_{(\theta, \mu)} \mathbb{P}_x$ ，梯度流动方程变为：

$$DTP = \nabla \cdot (\hat{\nabla} \Psi), \quad \Psi = bR \quad (30)$$

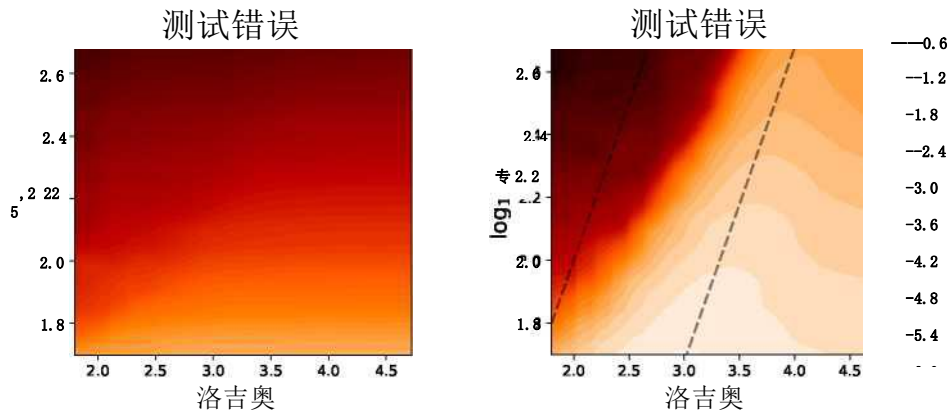


图 11：（左）连续视点；（右）常规 NN 模型。 目标函数是单个神经元。 经雷武许可转载。

粒子法离散化是基于：

$$p_{am} E_8(a_j(t), w(t)) = m E_{\text{知}}(*)$$

其中  $U_j(T) = (A_j(T), w(T))$ 。 可以证明在这种情况下，（30）减少到

$$\frac{1}{m} (u_1, \dots, u_m)$$

在哪里

$$I(u_1, \dots, u_m) = m, \quad u_j = j, \quad w_j, \quad m^{(x)} = \sum_j a_j b_j^{(w, x)}$$

这正是“缩放”两层神经网络的梯度下降。 在这种情况下，连续公式也与两层神经网络 [7, 32, 34, 35] 的“平均场”公式相吻合。

在“缩放”两层神经网络前面的缩放因子  $1/m$  实际上是相当重要的，对网络的测试性能有很大的影响。 图 11 显示了具有和不具有此缩放因子的两层神经网络模型测试误差的热图。 目标函数是简单的单神经元函数：  $f^*(X) = a(X_i)$ 。 重要的观察是，在没有这种缩放因子的情况下，测试误差在神经网络模型表现为相关随机特征模型的相位和显示比随机特征模型更好的性能的另一个相位之间显示出“相变”。 这样的相变是选择正确的超参数集的原因之一，这里的网络宽度  $m$  是如此重要。 然而，如果使用缩放形式，则避免了这种相变现象，[31] 性能更健壮。

## 5.5 基于流量的表示的最佳控制问题

基于流的表示自然会导致控制问题。 这种观点在机器学习中已经显式或隐式地使用了相当长的一段时间(例如，参见

[27])。例如，反向传播算法是基于控制理论的算法的一个例子。另一个最近的例子是最大基于原理的算法的发展，第一次介绍在[28]。尽管取得了这些成功，但我们认为仍然有很大的空间使用控制理论的观点来开发新的算法。

我们认为基于流的表示形式稍微更一般

$$\begin{aligned} & \text{新西兰} \\ & = \text{博士, 臥 } z, u), \quad z(0, x) = x \end{aligned}$$

其中  $z$  是状态,  $P_t$  是时间  $t$  的控制。我们的目标是在  $\{p\}$  上最小化  $R$

$$R(\{p\}) = \int_0^T (\dot{z}^T - \dot{z}^T * \dot{z})^2 dt = \int_0^T (\dot{z}^T - \dot{z}^T * \dot{z})^2 dt \quad \text{如}$$

就像以前一样

$$\dot{z}^T(x) = 1 \quad \text{兹}(1, x) \quad (31)$$

这个控制问题的一个最重要的结果是 Pontryagin 的最大原理 (PMP)。为了说明这一结果, 让我们将哈密顿量  $H: x \mapsto P_2(Q): TR$  定义为  $H(z, p, u) = Eu = [p^T \dot{z}, u]$ 。

Pontryagin 的最大原理断言控制问题的解必须满足:

$$P_t = \arg \max_{p \in \mathbb{R}^n} \{ p^T \dot{z}(t) - H(z(t), p, u) \}, \quad \forall t \in [0, 1], \quad (32)$$

每个  $x$ ,  $\{z^*(t)\}$  由前向/后向方程定义:  $dz^T, x$

$$\dot{z} = VPH = Eu_{z, t} [0(z^T, u)] dp^T, x$$

具有边界条件的 DR:

$$P_T = -V_z H = Eu^* (z^T, u) p^T, \quad (33)$$

$$\frac{dz}{dt} = x$$

$$p^T = -\dot{z}^T / \dot{z}^T; P(\cdot; t)) / \dot{z}^T, \quad (34)$$

$$(35)$$

对于平稳点, Pontryagin 的最大原则略强于 KKT 条件, 因为 (32) 是最优性而不是临界性的陈述。事实上, (32) 也适用于参数是离散的, 这在[29]中已经被用来为这种情况开发有效的数值算法。

在 PMP 的帮助下, 也很容易写出优化问题的梯度下降流。形式上, 可以简单地写下每个  $t$ : (32) 的梯度下降流:

$$\dot{P}_t = \nabla_p H(z(t), p, u), \quad \forall t \in [0, 1], \quad (36)$$

在哪里

$$\dot{z}^T(u, p) = \dot{z}^T[\hat{p}(z^T, u), p, u] \quad \text{片 } h$$

和  $\{z^*(t)\}$  由前向/后向方程定义。

为了离散梯度流, 我们可以简单地使用:

- 前向欧拉表示  $t$  变量中的流，步长为  $1/L$  ( $L$  是  $t$  变量中网格点的数目)；
- 梯度下降动力学的粒子法，每个  $t$  格点有  $M$  个样本。

这给了我们

$$z_T + \text{我} = \text{中兴通讯, } x + \text{巾伝件}, \quad U_j(t), \quad 1 = 0, \dots, \quad 1 - 1 \quad (37)$$

$$j=i$$

$$P \text{笋} = P; \quad +1 + LMEV_{Z^0}(W+1, \quad U_j + \text{---}^{(t)}) p; \quad +i, \quad 1 = 0, \dots, \quad 1-1 \quad (38)$$

$$j=1$$

$$\text{衅} = -\text{Ex}[\text{大众臥才}, U_j(T) p \text{門}]. \quad (39)$$

这将恢复（缩放）ResNet 的 GD 算法（反向传播）：

$$z \text{我} + \text{我} = z + \text{巾}(2)^{uz}.$$

由于因子  $1/(LM)$  的存在，我们称之为“缩放”ResNet)。

本着类似的精神，人们也可以使用 PMP[28] 获得算法。采用控制理论中的术语，这种算法被称为“逐次逼近法” (MSA)。基本 MSA 如下：

初始化：  $0^\circ \text{ cU}$

关于  $k=0, 1, 2, \dots$ ：

解决  $d$

解决

$$\begin{aligned} DP_0 \quad -V_Z H(\text{球 } pk, \text{ 郎}, pi) &= -2 \left( \int (x; \text{假}) - \int (x) \right) 1 \\ d T \end{aligned}$$

$$\sim_P^{H(Z)} \quad t, p_0, \quad z_0 = x$$

-Set

为每个人  $[0, 1]$

$$^o r^{*i} = \arg\max \text{心}^{H(z)} \text{不}, p_0, \text{ 幻}$$

在实践中，这个基本版本的性能不如“扩展 MSA”，它的工作方式与 MSA 相同，只是哈密顿被扩展的哈密顿[28]所取代：

$$^{H(z)} , p^{0, v}, q) : = \frac{H(z, \quad, 0) - \quad}{2 \text{入皿-}/} (z, 0) \quad \text{『-}^1 \text{入恆}^+ z \quad pW) \text{。}$$

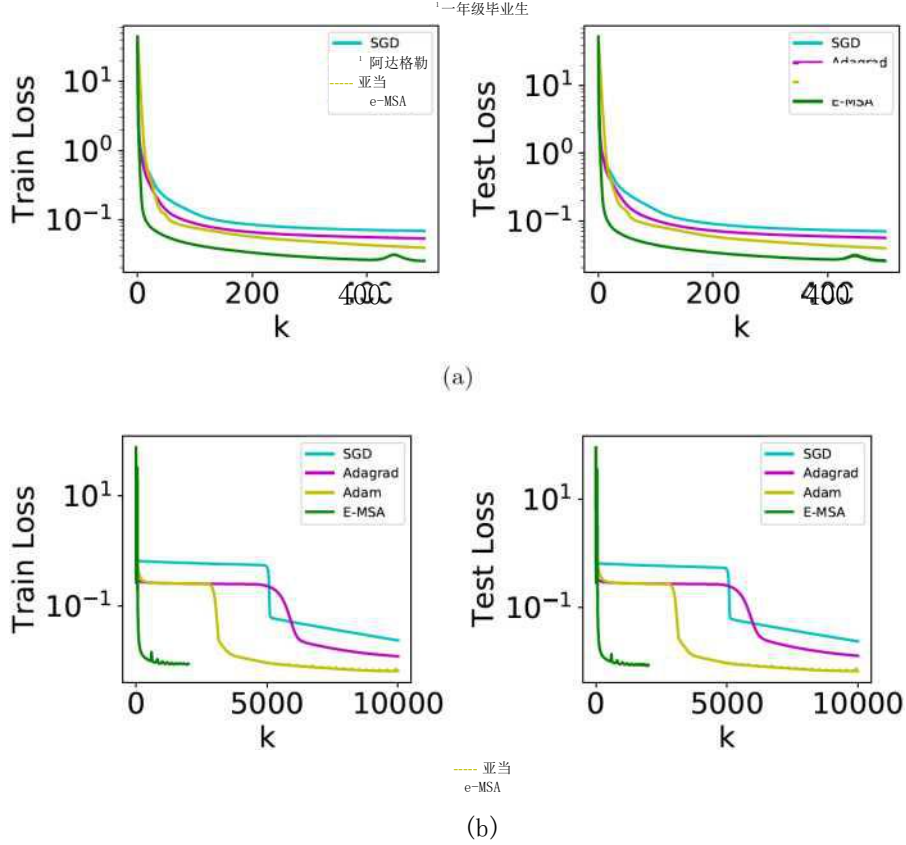


图 12: 不同版本随机梯度下降算法的扩展 MSA 比较。顶部的数字显示了小初始化的结果。底部数字显示了更大初始化的结果。经李千晓许可转载。另见[28]。

图 12 显示了扩展 MSA 与不同版本的 SGD 在两种初始化中的结果。可以看到, 在迭代次数方面, 扩展 MSA 的性能优于所有 SGD。在挂钟时间方面, 扩展 MSA 的优势明显减弱。这可能是由于用于求解 (32) 的优化算法(这里是 BFGS)的实现效率低下。详情请参阅[28]。无论如何, 显然还有很大的改进空间。

## 6 结束语

最后, 我们讨论了基于机器学习的算法已经和/或将产生显著差异的广泛问题。这些问题对于计算数学来说是比较新的。我们坚信, 基于机器学习的算法也会显著影响我们解决计算数学中更传统问题的方式。然而, 这方面的研究仍处于非常早期的阶段。

机器学习的另一个重要领域是多尺度建模。上面讨论的闭矩问题就是这个方向的一个例子。

有更多可能的应用程序, 见[8]。机器学习似乎能够提供缺失的环节, 使先进的多尺度建模技术真正实用。例如, 在异构多尺度方法(HMM)[1, 9]中, 一个重要的组成部分是从微观尺度模拟数据中提取相关的宏观尺度信息。这一步骤一直是 HMM 的主要障碍。显然, 机器学习技术在这里有很大的帮助。

我们还讨论了数值分析的观点如何有助于改善机器学习的数学基础, 并提出了新的和可能更稳健的公式。特别是, 我们已经尝到了高维近似理论应该是什么样的。我们还证明, 常用的机

机器学习模型和训练算法可以从连续模型的某些特定离散化中恢复，以缩放的形式。从这一讨论中可以看出，神经网络模型是相当自然和相当不可避免的。

我们从机器学习中真正学到了什么？嗯，似乎机器学习最重要的新见解是将函数表示为期望。我们在这里复制它们是为了方便：

- 基于积分变换：

$$f(x) = E_{\mathbf{a}}[AA(W+x)]$$

$$f^{(x)} = \text{电} \quad a(\text{區姑知 } \mathbf{a}^{\circ, -x}) \dots)$$

- 基于流程：

$$\begin{aligned} & \text{新西兰} \\ & dT = E(\mathbf{a}, \mathbf{w}) \text{阵} aa(\mathbf{w}^T \mathbf{z}), \quad z(0, x) = x \\ & f(\mathbf{w}) = 1 \&^{(1, x)} \end{aligned} \quad \begin{aligned} & (40) \\ & (41) \end{aligned}$$

从计算数学的角度来看，这表明中心问题将从特定的离散化方案转向更有效的函数表示。

这篇评论相当粗略。有兴趣的读者可以参考三篇评论文章[17, 18, 19]了解更多细节。

致谢：我非常感谢我的合作者对本文所述工作的贡献。特别是，我要衷心感谢罗伯托·卡尔、杰群·汉、阿努夫·詹岑、钱晓丽、潮马、韩王、斯蒂芬·沃伊托维奇和雷武对这里讨论的问题进行了多次讨论。这项工作的一部分支持来自 iFlytek 给普林斯顿大学的礼物以及 ONR 赠款 N00014-13-1-0338。

## 参考资料

- [1] Assyr Abdulle, 渭南 E, Bjorn Engquist 和 Eric Vanden-Eijnden, 异源多尺度方法, Acta Numerica, vol. 21, pp. 1-87, 2012.
- [2] 弗朗西斯·巴赫, “用凸神经网络打破维数的诅咒”, “机器学习研究杂志”, 18(19): 1-53, 2017 年。
- [3] 安德鲁·R·巴伦, “乙状结肠函数叠加的通用近似界”, IEEE 信息论交易, 39(3): 930-945, 1993 年。
- [4] 乔格·贝勒和米歇尔·帕里内洛, “高维势能表面的广义神经网络表示”, 物理评论信件, 98(14): 146401, 2007 年。
- [5] Achi Brandt, “多尺度科学计算: 回顾 2001 年”。在 Barth, T. J., Chan, T. F. 和 Haimes, R. (编辑.): 多尺度多分辨率方法: 理论与应用, Springer Verlag, Heidelberg, 2001, pp. 196.
- [6] 罗伯托·卡 and 米歇尔·帕里内洛, “分子动力学和密度功能理论的统一方法”, “物理评论信函”, 55(22): 2471, 1985 年。

- [7] 莱纳奇·奇扎特和弗朗西斯·巴赫, “关于使用最优传输的过度参数化模型的梯度下降的全局收敛性”, “神经信息处理系统的进展”, 第 3036-3046 页, 2018 年。
- [8] 渭南 E, 多尺度建模原理, 剑桥大学出版社, 2011。
- [9] 渭南 E 和 BjornEngquist, 异构多尺度方法, Comm. 数学. SCI, 第二卷. 1, 不. 1, pp. 87-132, 2003.
- [10] 渭南 E, 韩杰群和阿努夫·金岑, “基于深度学习的高维抛物偏微分方程和后向随机微分方程的数值方法”, 数学与统计通信 5, 4 (2017), 349-380。
- [11] 渭南 E, 韩洁群, 张林峰, “将机器学习与基于物理的建模相结合”, <https://arxiv.org/pdf/2006.02619.pdf>, 2020 年。
- [12] 渭南 E, 马超和雷武, “两层神经网络种群风险的先验估计”, “数学科学中的通信”, 17 (5): 1407-1425, 2019 年; ArXiv: 1810.06397, 2018 年。
- [13] 渭南 E, 马超和吴磊, “Barron 空间和神经网络模型的流诱导函数空间”, arXiv: 1906.08039, 2019。
- [14] 渭南 E 和 StephanWojtowytsch, “Barron 函数的表示公式和点性质”, arXiv: 2006.05982, 2020。
- [15] 渭南 E 和 StephanWojtowytsch, “关于与多层 ReLU 网络相关的 Banach 空间: 函数表示、近似理论和梯度下降动力学”, <https://arxiv.org/pdf/2007.15623.pdf>, 2020 年。
- [16] 渭南 E, 马超, 雷武, “机器学习从一个连续的观点”, ar 希夫: 1912.12777, 2019 年。
- [17] 渭南 E, 韩洁群和 ArnulfJentzen, “求解高维 PDE 的算法: 从非线性蒙特卡罗到机器学习”, <https://arxiv.org/pdf/2008.13333.pdf>, 2020 年。
- [18] 渭南 E, 韩洁群和张林峰, “将机器学习与基于物理的建模相结合”, <https://arxiv.org/pdf/2006.02619.pdf>, 2020 年。
- [19] 渭南 E, 马超, 斯蒂芬·沃伊托维奇和雷武, “对机器学习的数学理解: 什么是已知的, 什么是不知道的”, <https://arxiv.org/pdf/2009.10713.pdf>.
- [20] 韩洁群和渭南 E, “随机控制问题的深度学习近似”, 深度强化学习讲习班, NIPS (2016), <https://arxiv.org/pdf/1611.07422.pdf>.
- [21] 韩杰群, 阿努夫·金岑和渭南 E, “用深度学习求解高维偏微分方程”, “国家科学院学报”, 115, 34 (2018), 8505-8510。
- [22] 韩杰群, 马超, 郑马, 渭南 E, “均匀精确的基于机器学习的动力学方程流体模型”, 国家科学院学报, 116 (44) 21983-21991; DOI: 10.1073/PNAS. 1909854116, 2019。
- [23] 韩洁群, 张林峰, 罗伯托·卡尔和渭南 E, “深势: 多体势能面的一般表示”, 计算物理中

- 的通信, 23 (3): 629-639, 2018。
- [24] 皮埃尔·霍恩伯格和伯特兰·I·哈尔佩林, “动态临界现象理论”, “现代物理学评论”, 49 (3): 435, 1977。
- [25] 贾伟乐, 王汉, 陈墨汉, 陆登辉, 刘吉端, 林林, 罗伯托·卡尔, 渭南 E, 张林峰, “用机器学习将分子动力学的从头精度极限推到一亿个原子”, arXiv: 2005.00223, 2020。
- [26] Jason M Klusowski 和 Andrew R Barron, “包括神经网络在内的高维脊函数组合的风险界限”, ar Xiv: 1607.01434, 2016。
- [27] 延勒村, “反向传播的理论框架”, 载于: 图雷茨基, D., Hinton, G., Sejnowski, T. (编辑。) 1988 年连接主义模型暑期学校论文集, 卡内基-梅隆大学, 摩根考夫曼, 1989 年。
- [28] 李千晓, 龙晨, 程泰和渭南 E, “基于最大原理的深度学习算法”, “机器学习研究杂志”, 第二卷。18, 不。165 页。2018 年 1 月 1 日至 29 日, <https://arxiv.org/pdf/1710.09513.pdf>。
- [29] 李千晓和郝树基, “一种最优的深度学习控制方法和应用于离散视觉神经网络”, 第 35 届国际机器学习会议记录, 2018 年。
- [30] 鲁登辉, 王汉, 陈墨汉, 刘吉端, 林林, 罗伯托卡, 渭南 E, 贾伟乐, 张林峰, “86Pflops 深势分子动力学模拟 1 亿原子的从头计算精度”, arXiv: 2004.11658, 2020。
- [31] 马超, 雷武, 渭南 E, “两层神经网络模型梯度下降动力学的猝灭激活行为”, arXiv: 2006.14450, 2020。
- [32] 宋梅, Andrea Montanari 和 Phan-Minh Nguyen, “两层神经网络景观的平均视野”, “中国科学院学报”, 115 (33): E7665-E7671, 2018。
- [33] Etienne Pardoux 和 Shige Peng, “反向随机微分方程和拟线性抛物偏微分方程”, 在随机偏微分方程及其应用(夏洛特, NC, 1991), 第二卷。176 在控制和通知的讲座笔记。Sci., Springer, Berlin, 1992, pp. 200-217。
- [34] 格兰特·罗茨科夫和埃里克·范登-埃因登, “作为相互作用粒子的参数: 神经网络的长时间收敛和渐近误差缩放”, “神经信息处理系统的进展”, 第 7146-7155 页, 2018 年。
- [35] 贾斯汀·西里尼亚诺和康斯坦蒂诺斯·斯皮里奥普洛斯, “神经网络的平均场分析: 中心极限定理”, arXiv: 1808.09372, 2018 年。
- [36] 汉斯·崔贝尔, 函数空间理论, Birkhauser, 1983 年。
- [37] 韩旺, 张林峰和渭南 E, 未发表。
- [38] 张林峰, 韩洁群, 王汉, 罗伯托·卡, 渭南 E, “深势分子动力学: 具有量子力学准确性的可扩展模型”, “物理评论信函”, 120: 143001, 2018 年 4 月。
- [39] 张林峰, 韩杰群, 王汉旺, 维萨姆·A·赛迪, 罗伯托·卡, 渭南 E, “有限系统和扩展系



统原子间势能模型的端到端对称性保持”，“神经信息处理系统的进展”，2018 年。

- [40] 张林峰, 林德业, 王汉, 罗伯托, 渭南 E, “主动学习均匀精确的原子间势用于材料模拟”, “物理评论材料”, 3 (2): 023804, 2019.
- [41] 张林峰, 王汉和渭南 E, “增强动力学在大原子和分子系统中的增强采样。i. 基本方法”, J. 化学。 菲斯, 第二卷。 148, pp. 124113, 2018.