

专题: 华南师范大学建校暨物理学科建立 90 周年

机器学习回归不确定性揭示自驱动 活性粒子的群集相变*

郭唯琛 艾保全[†] 贺亮[‡]

(华南师范大学物理学院, 理论物理研究所, 广州 510006)

(2023 年 5 月 30 日收到; 2023 年 7 月 16 日收到修改稿)

本文发展了一种利用逆统计问题中的回归不确定性来自动探索物质相的新方法. 以自驱动活性粒子的群集相变为例, 展示了对于这一类涉及非平衡、非晶格、一阶相变等复杂要素的多体系统, 在训练人工神经网络处理其中的逆统计问题回归任务, 成功重构出系统的噪声强度这一参数之后, 回归结果的不确定性关于实际噪声强度的分布具有非平庸的规律性, 可用于揭示该系统中的群集相变, 并自动提取相变的临界噪声强度. 本文还与两种基于神经网络分类能力的常见方法进行直接对比, 讨论了它们的异同和各自特点. 结果表明, 本文发展的新方法不仅具有使用效率较高和所需预设的物理知识较少等实用优势, 而且更有在理论层面较为自然地与传统物理概念建立联系的可能性, 对于跨领域的不同物理系统都有良好的通用性和有效性.

关键词: 机器学习, 相变, 非平衡多体系统, 逆统计问题**PACS:** 07.05.Mh, 05.70.Fh, 05.70.Ln, 02.30.Zz**DOI:** 10.7498/aps.72.20230896

1 引言

近年来, 基于人工神经网络 (artificial neural network, ANN) 的机器学习技术已越来越多地为凝聚态和统计物理领域的研究提供帮助. 尤其是 2017 年 Melko 和 Carrasquilla^[1] 以及 van Nieuwenburg 等^[2] 分别报道了自动探索物质相的“留白法” (learning with blanking)^[1,2] 和“混淆法” (learning by confusion)^[2] 之后, 这两种利用 ANN 处理分类任务的强大能力的方法以及它们的一些衍生方法^[1–11], 已成功地许多不同物质相的存在性提供了数据驱动的新证据, 并为对应相变点的参数临界值提供了数据驱动的新估计. 这样的成功案例遍及凝聚态和统计物理领域的各种物理系统, 也包括涉及非平

衡^[3,4]、拓扑缺陷^[5,6]、强关联费米子^[7,8] 等复杂要素的情况. 对于经典系统和量子系统, 这一类机器学习方法不仅能处理由数值模拟产生的原始数据, 还可以协助分析由实验观测得到的原始数据^[9–11]. 然而, 由于 ANN 的底层工作机制至今仍未得到足够清晰的解释, ANN 通过直接拟合原始数据而给出的分“类”结果与被研究的物理系统中的物质“相”的理论联系往往难以捉摸^[12–14].

面对这一困难, 值得注意的是 ANN 除了具有强大的处理分类任务的能力, 还同样具有强大的处理回归任务的能力, 而回归任务的结果通常具有明确的物理意义. 例如, 相应于研究一个物理系统时的正向思维“给定系统参数的取值, 求系统的可能状态”, 所谓的逆统计问题 (inverse statistical problem, ISP)^[15] 指的是“给定一个具体的系统状态,

* 国家自然科学基金 (批准号: 12275089, 12075090)、广东省自然科学基金 (批准号: 2023A1515012800, 2022A1515010449) 和科技部重点研发计划 (批准号: 2022YFA1405304) 资助的课题.

[†] 通信作者. E-mail: aibq@scau.edu.cn

[‡] 通信作者. E-mail: liang.he@scau.edu.cn

求它可能对应的系统参数值”,这就是一种典型的回归任务. 如果用 ANN 处理 ISP 的回归, ANN 的输出值就不再是处理分类任务时难以捉摸的“类”,而是被重构的系统参数值本身. 事实上, ANN 处理回归任务的能力及其与传统物理概念的直接联系,已经开始被物理学家关注. 尤其是 Tegmark 等^[16–19]探索了其自动构建物理理论的可能性,发现 ANN 可以在一些相关的回归任务中提取出系统的运动方程^[16]、对称性^[17]、守恒律^[18]等等,甚至还用 ANN 重建了《费曼物理学讲义》中涉及的 100 公式^[19]. 这些有趣的研究成果体现了 ANN 的回归结果比 ANN 的分类结果更具有联系物理的可能性.

由此,基于 ANN 回归的自动探索物质相的机器学习方法也逐渐开始出现,例如最近刚刚出现的利用 ISP 中的回归不确定性的方法 (learning from regression uncertainty, LFRU 方法)^[20]. 这一方法自动探索物质相的能力及其与传统物理概念的直接联系在 Ising 模型和 Clock 模型中已得到了初步验证^[20],为机器学习在相变研究中的应用提供了新的视角. 然而,该方法的通用性仍需要进一步检验,尤其是面对非平衡、非晶格的系统中的一阶相变, LFRU 方法仍然有效吗? 我们知道,与连续相变不同,在一阶相变的临界点处,系统的关联长度不发散,这带来了更丰富的临界物理现象,但也使得它们的跨尺度普适性质难以被重整化群^[21]等强大的物理学传统研究方法刻画. 另一方面,与平衡系统不同,非平衡多体系统中细致平衡的缺失,同样带来了更丰富的临界物理现象,例如湍流的无规律行为^[22],但这也同样使得相关领域缺乏较为通用的研究方法^[23]. 考虑到 ANN 的数据处理和信息挖掘能力本身是足够普适的,这类情景正是基于 ANN 的机器学习技术的用武之地. 如果基于 ANN 回归的 LFRU 方法 (以及基于 ANN 分类的“留白法”和“混淆法”) 能在不额外增加针对非平衡、非晶格的系统中的一阶相变的特殊设计的情况下,有效处理这类复杂多体系统,实现自动探索其中的物质相,那么这将为非平衡多体系统中的相变研究提供一个具有较强通用性的工具箱,有助于更好地揭示这类系统中的丰富的临界物理现象.

本文在一个由 Vicsek 模型描述的自驱动活性粒子系统^[24–27]中具体研究 LFRU 方法的通用性. 这是一个具有外部噪声的随机动力学模型,最初用于模拟鸟类在较低能见度的恶劣天气下的集群飞

行,也是统计物理领域关于自驱动活性粒子系统的基础模型之一,具有丰富的集体动力学行为和自组织现象^[24–27]. 这种非平衡多体系统的噪声强度的改变会引发一个从低噪声的群集相 (flocking phase, 所有粒子的运动方向大致相同) 到高噪声的无序相 (disordered phase, 系统保有旋转对称性) 的一阶相变^[24–27] (如图 1 所示). 研究发现,即使这涉及非平衡、非晶格、一阶相变等复杂要素, ANN 仍可以通用地被直接训练用于处理该系统中的 ISP 回归任务,成功重构出该系统中的噪声强度,如图 2(b) 所示. 进一步考察 ANN 在这个任务中的回归不确定性,发现它关于被重构的实际噪声强度的分布具有规律性,其曲线呈现 M 字形,如图 3(a) 所示. 最重要的是,研究发现由 ANN 自主得到的 M 字形曲线,可以用于自动探索物质相. M 字形曲线不仅揭示群集相变的存在,而且其中间的极小值所在的位置,正对应于该相变的临界噪声强度. 在先前研究^[20]的基础上,上述新发现清晰展现了 LFRU 方法对于跨领域的不同物理系统具有良好的通用性. 我们也检验了“混淆法”和“留白法”这两种基于 ANN 分类的典型方法,用于研究自驱动活性粒子的群集相变的效果,在方法的使用效率、所需预设的物理知识、联系物理的可能性等方面,对比讨论了 LFRU 方法与它们的异同和各自特点.

2 自驱动活性粒子的群集相变和逆统计问题

2.1 物理系统

本文研究的多体物理系统由 N 个在二维 $L \times L$ 空间中运动的自驱动粒子组成,这些粒子的集体行为由一组随机动力学方程描述^[24–27]:

$$\mathbf{v}_i(t + \Delta t) = v_0 \vartheta \left(\eta \mathcal{N}_i \boldsymbol{\xi}_i + \sum_{j \in A_i} \mathbf{v}_j(t) \right), \quad (1)$$

其中, Δt 是离散的时间间隔, v_0 是粒子的速率 (设为常数), ϑ 是矢量的归一化算符, $\vartheta(\mathbf{w}) = \mathbf{w}/|\mathbf{w}|$, A_i 是以 i 粒子所在位置为圆心的半径为 r 的区域, \mathcal{N}_i 是位于 A_i 区域内的粒子数 (包括 i 粒子自身), $\boldsymbol{\xi}_i$ 是随机指向的单位矢量噪声, η 为表征环境扰动影响程度的噪声强度系数. 这是一个标准的由外部噪声影响的 Vicsek 模型. 在系统密度 $\rho = N/L^2$ 取定的情况下,噪声强度的改变会引发一个从低噪

声的群集相到高噪声的无序相的一阶相变^[24–27]. 通过计算该系统的群速度 $\bar{v} = \left| \sum_{j=1}^N \mathbf{v}_j / (Nv_0) \right|$ 作为序参量, 可看到 \bar{v} 在相变点处发生突变. 为不失一般性, 以 $N = 2048$, $\rho = 2$, $L = 32$, $v_0 = 0.5$, $r = 1$ 的系统为例. 在该参数取值下, 如图 2(a) 所示, \bar{v} 的突变发生于 $\eta_c \approx 0.626$. 本文的具体目标是利用基于 ANN 的机器学习技术, 从图 1 所示的原始数据中自动提取出这一临界噪声强度 η_c .

2.2 机器学习

要利用基于 ANN 的机器学习技术研究该物理系统, 无论是让 ANN 处理分类任务还是回归任务, 都首先需要将数据整理为适合 ANN 进行分析的形式. 这当然有很多不同的做法. 我们选择类似于 ANN 在人脸识别等图像处理领域的用法, 将数

据整理为图像的形式 (如图 1 所示). 每个数据样本图像中的每个圆形标记表示二维空间中的一个自驱动粒子, 其空间分布表示自驱动粒子的瞬时空间分布, 其颜色分布表示自驱动粒子的运动方向的瞬时角度分布. 由于本文将直接使用一个工业界成熟的深度残差网络架构 (residual neural network, ResNet)^[28], 其默认的输入尺寸是 $3 \times 224 \times 224$, 其中的 3 对应于彩色图像的 RGB 三通道, 因而本文数据样本图像的尺寸也为 224×224 像素. 这些样本被分配为 3 组, 构成 3 个不同的数据集: 训练集、验证集、测试集.

所谓 ANN 的训练, 指的是若干次遍历训练集的样本, 每次遍历时, ANN 作为一个 $3 \times 224 \times 224 \rightarrow 1$ (用于 ISP 回归任务) 或 $3 \times 224 \times 224 \rightarrow 2$ (用于二元分类任务) 的映射, 对每个样本都给出

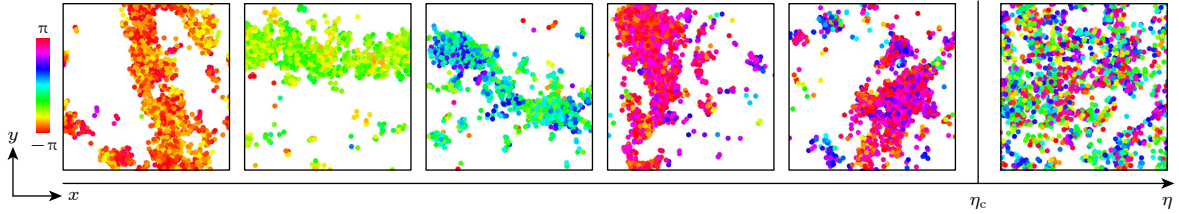


图 1 数值模拟生成的对应于不同噪声强度 η 的典型样本. 样本中的每个圆形标记表示二维空间中的一个自驱动粒子, 其空间分布表示自驱动粒子的瞬时空间分布, 其颜色分布表示自驱动粒子的运动方向的瞬时角度分布. 此处作为示例的样本中, 左边的 5 个样本处于群集相, 最右边的样本处于无序相

Fig. 1. Typical samples corresponding to different noise levels that are generated by numerical simulations. In every sample, each of the circular markers represents a single self-propelled particle in the two-dimensional space, with their spatial distribution representing the instantaneous spatial distribution of self-propelled particles, and their color distribution representing the instantaneous angular distribution of directions of motion of these self-propelled particles. Among the samples shown here for instance, the five samples in the left are in the flocking phase, and the rightmost one is in the disordered phase.

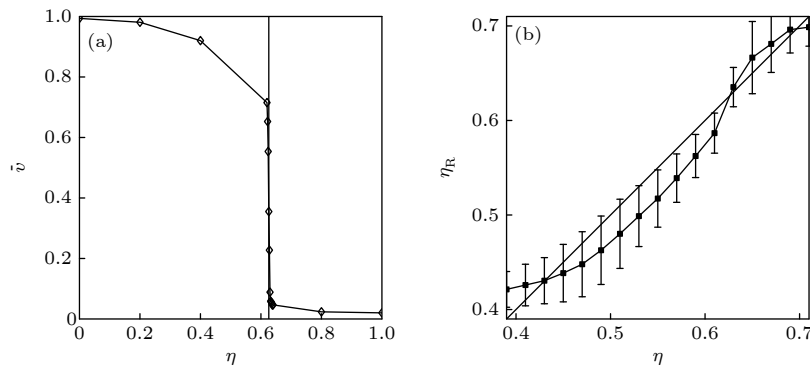


图 2 自驱动活性粒子系统中的 ISP (a) 系统的群速度 \bar{v} 关于噪声强度 η 的依赖关系, \bar{v} 在 $\eta_c = 0.626 \pm 0.006$ 的突变表明系统在该噪声强度处发生一阶相变; (b) 训练完成的 ANN 给出的重构噪声强度 η_R 关于实际噪声强度 η 的依赖关系, 误差棒表示回归不确定性 $U(\eta)$, 对角线表示理想的回归结果 $\eta_R = \eta$

Fig. 2. Inverse statistical problem in a self-propelled active particle system: (a) Noise level dependence of the system's global group velocity, whose jump at $\eta_c = 0.626 \pm 0.006$ characterizes the first-order flocking phase transition; (b) noise level dependence of the reconstructed noise level predicted by the well-trained ANN. The error bars represent the regression uncertainty $U(\eta)$, and the diagonal line represent the ideal regression result $\eta_R = \eta$.

1 个 (ISP 回归) 或 2 个 (二元分类) 相应的值作为输出结果. 基于 ANN 的输出结果, 计算一个损失函数, 并按照反向传播规则, 以梯度下降等方式优化 ANN 中的大量可训练系数的取值, 从而最小化损失函数, 这就是 ANN 的训练. 对于 ISP 回归任务, 损失函数可以是输出结果与标签之间的均方误差, 其中标签指的是我们为每个样本标注的参考答案, 即实际的噪声强度值 η . 对于二元分类任务, 损失函数可以是输出结果与标签之间的交叉熵函数, 这里标签则是甲类或乙类 (关于二元分类, 详见 3.2 节和 3.3 节). 为了提高 ANN 的泛化能力, 训练后最终采用的可训练系数的取值并不是在训练集实现损失函数最小的那一组, 而是在验证集实现损失函数最小的那一组. 带着这组最终取定的可训练系数, 训练完成的 ANN 将被应用于测试集, 以评估其实际应用效果. 下文首先讨论基于 ANN 的 ISP 回归.

2.3 逆统计问题

LFRU 方法利用的是 ISP 中的回归不确定性. 要使用这一方法研究自驱动活性粒子的群集相变, 首先需要构建一个 ISP 回归任务让 ANN 尝试处理. 在上述的非平衡多体系统中, 相应于正向思维的给定噪声强度 η , 求处于稳态的可能的系统状态 (位置 x, y 分布与速度 v 分布, 如图 1 所示), 一个比较自然的 ISP 是: 给定一个具体的系统状态, 求它可能对应的噪声强度 η . 这是一个统计推断问题, 推断得到的重构噪声强度记为 η_R . 由于原始数据是由随机动力学方程演化得到的, 不可避免在不同的噪声强度下出现极其类似的样本. 这意味着对于在同一噪声强度 η 下生成的不同样本, ANN (或其他方法) 给出的重构噪声强度值 η_R 不会完全一样. 这就带来了回归不确定性 $U(\eta)$. 我们可以用 ANN 回归结果的标准差来刻画这一不确定性, 即

$$U(\eta) = \sqrt{\langle (\eta_R - \langle \eta_R \rangle)^2 \rangle}, \quad (2)$$

其中, $\langle \cdot \rangle$ 表示对测试集所有属于同一噪声强度 η 的样本取平均. 对于 ISP 本身, 回归任务的核心目标之一是尽量减少这个不确定性, 但其存在是系统性的, 因而不可能被真正减少到零. 再考虑到这是一个涉及非平衡、非晶格、一阶相变等复杂要素的情况, 如何有效地实现 Vicsek 模型的 ISP, 本身就

是一个非平庸的问题. 传统方法研究 ISP 主要集中于 Ising 模型等简单情况^[15], 还通常要使用平均场^[29]或最大似然估计^[30]等稍具针对性的方法. 这里直接使用 ANN 进行 ISP 回归.

3 群集相变临界噪声强度的自动提取

3.1 LFRU 方法: 回归不确定性中的相变信号

本文使用的数据集涉及 $\eta \in [0.39, 0.71]$ 范围内以 $\Delta\eta = 0.02$ 为间隔的 17 个不同的噪声强度值, 对于每个噪声强度值, 有 2000 个样本用于训练, 500 个样本用于验证, 2500 个样本用于测试. 17 个 η 的总共 34000 个训练集样本, 在训练过程中被遍历 20 次, 并作相应的验证, 最终得到一个训练完成的 ANN, 在测试集评估其回归结果. 如图 2(b) 所示的回归结果取自 20 个独立训练、独立验证的 ResNet 的平均测试结果. 可以看到, ANN 给出的重构噪声强度 η_R 虽然不能完美贴合于实际噪声强度 η , 但也差得不远. 由于 ISP 仅仅是利用其中的回归不确定性自动探索物质相的一个中间过程, 目标不在于 ISP 本身, 因此这里不讨论图 2(b) 所示的回归结果与传统方法得到的回归结果的对比, 也不评判各种研究 ISP 的方法的优劣. 我们关注的是: 对于这样的具有一阶相变的复杂系统, ANN 可以克服诸如亚稳态等等的对于 ISP 回归的潜在干扰, “学会了”该系统中的噪声强度 η 这一参数. 这意味着其输出值确实可以视为与噪声强度 η 具有直接的物理联系, 使得进一步得到的自动探索物质相的结果也有了联系物理的可能性.

确认了 ANN 可以实现 ISP 回归之后, 考察 ANN 在这个任务中的回归不确定性 $U(\eta)$, 也就是图 2(b) 的误差棒. 这在图 2(b) 中并不明显, 图 3(a) 所示为 $U(\eta)$ 关于噪声强度 η 的依赖关系, 可以清晰地看到 $U(\eta)$ 的分布具有规律性, 其曲线呈现 M 字形, 并且中间的极小值所在的位置 $\tilde{\eta}_c = 0.63 \pm 0.01$, 并不是位于整个参数区域 $\eta \in [0.39, 0.71]$ 的正中间附近, 而是恰好对应于系统的临界噪声强度 $\eta_c \approx 0.626$ (图 3 中的竖线表示由 \bar{v} 的突变位置给出的临界噪声强度 η_c , 即传统方法得到的相变点). 这说明 LFRU 方法能够成功地从如图 1 所示的原始数据中自动提取出自驱动活性粒子的群集相变临界噪声强度 η_c .

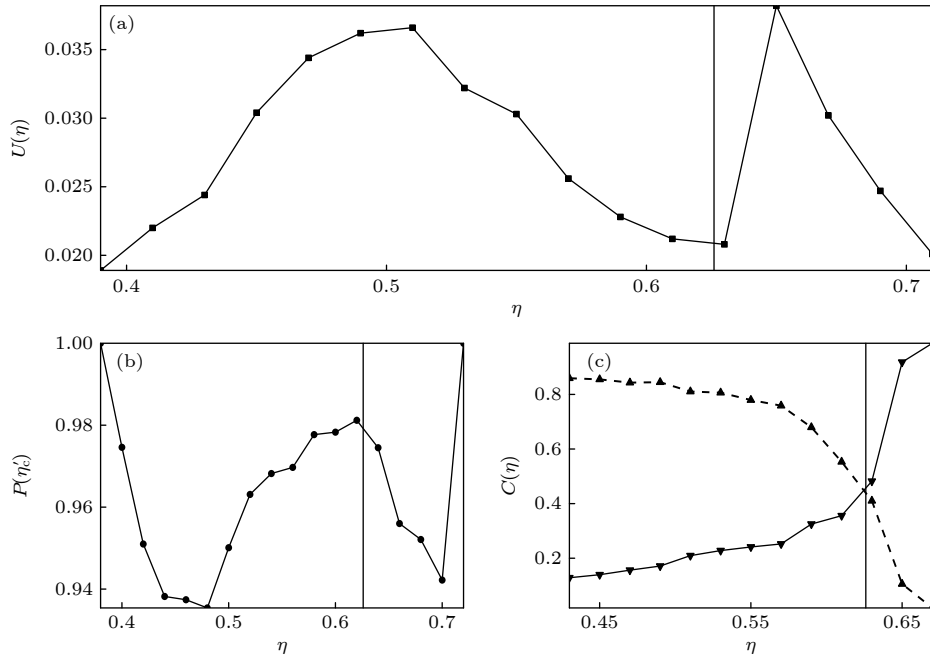


图 3 三种机器学习方法揭示自驱动活性粒子的群集相变 (a) 基于回归不确定性的 LFRU 方法; (b) “混淆法”; (c) “留白法”
Fig. 3. Revealing the flocking phase transition of self-propelled active particles via applying three different machine learning approaches: (a) The LFRU approach; (b) the “learning by confusion” approach; (c) the “learning with blanking” approach.

这与我们关于 LFRU 方法的研究^[20]中在 Ising 模型和 Clock 模型中发现的情况类似, 说明这一方法对于跨领域的不同物理系统具有良好的通用性. 利用 ANN 处理回归任务的强大能力及其与传统物理概念的直接联系, 研究者只需提供每个样本的实际参数值, 训练 ANN 处理 ISP 回归任务, 训练完成后的回归不确定性就可以用于自动探索物理相. 若 ISP 的参数区间内只有一个相, 回归不确定性的曲线只会呈现一个平庸的单峰^[20]. 当曲线呈现 M 字形, 这就揭示了该参数区间存在相变, 相变临界点可以直接从中间的极小值处提取.

3.2 “混淆法”

作为直接的对比, 使用两个典型的利用了 ANN 处理分类任务的强大能力的机器学习方法, 研究同样的非平衡多体系统中的群集相变. 要将 ANN 训练用于分类任务 (具体来说, 是二元分类任务), 需要将损失函数改换为交叉熵函数, 并且此时 ANN 对每个样本的输出应有 2 个值 (C_1, C_2), 它们具有概率的性质. 甲类和乙类对应的样本标签分别为 (1, 0) 和 (0, 1), 因而这 2 个输出值可以分别理解为 ANN 将一个样本识别为甲类或乙类的信心. 例如, 输出 (0.6, 0.4) 意味着 ANN 有六成的把握认为该样本属于甲类, 有四成的把握认为该样本属于

乙类. 很自然, 当 $C_1 > C_2$, ANN 对于该样本的分类结果即为甲类, $C_1 < C_2$ 则为乙类.

首先检验所谓的“混淆法”^[2]. 这个构思巧妙的方法, 利用的是 ANN 在面对不同程度上背离物理事实的混淆标签时的不同表现. 首先需要假定一个任意的噪声强度值 η'_c , 人为规定满足 $\eta < \eta'_c$ 的样本为甲类, 满足 $\eta > \eta'_c$ 的样本为乙类. 由于 η'_c 是任意假定的, 其对应的二元分类任务 (区分甲、乙两类的样本) 与这个系统中实际的物理相 (群集相、无序相) 不具有明确的理论联系. 训练完成后, 在测试集评估 ANN 针对这一任意假定的二元分类任务的表现. 在测试集的全部 m 个生成于不同噪声强度 η 的样本中, 若 ANN 成功识别了 m' 个, 计算出对应于 η'_c 的识别成功率 $P(\eta'_c) = m'/m$. 然后, 假定一系列不同的 η'_c , 分别重复上述的过程, 就可以得到 $P(\eta'_c)$ 关于 η'_c 取值的依赖关系.

如图 3(b) 所示, $P(\eta'_c)$ 的曲线呈现 W 字形. 考虑到对于任意 η'_c , 只要 η'_c 不符合物理上实际的临界噪声强度 η_c , 总会存在一些令 ANN 感到“混淆”的标签. 以 $\eta'_c > \eta_c$ 的情况为例 ($\eta'_c < \eta_c$ 的情况也类似), 那些满足 $\eta_c < \eta < \eta'_c$ 的样本与满足 $\eta > \eta'_c$ 的样本同处于无序相, 却在这一任意假定的二元分类任务中被贴上了不同标签. ANN 要如何理解二者之间纯属虚构的“不同”呢? 与此同时, 它们与

$\eta < \eta_c$ 的处于有序相的样本在物理上有显著的区别, 却被贴上了相同标签. ANN 又要如何理解二者的“相同”呢? 这些背离物理事实的标签限制了 ANN 的分类表现, 使得 $P(\eta'_c)$ 不会很高. 而对于同样的数据集, 当任意假定的 η'_c 越接近实际的临界噪声强度 η_c , 混淆标签就会越少, $P(\eta'_c)$ 也就越有机会取得更高的值. 这意味着当 $\eta'_c = \eta_c$ 时, $P(\eta'_c)$ 应具有非平庸的极大值. 因此, “混淆法”将 $P(\eta'_c)$ 的 W 形曲线的峰值对应的 η'_c 视为 ANN 给出的临界噪声强度预测值 $\tilde{\eta}_c$. 在图 3(b) 的计算中, 20 个独立训练、独立验证的 ResNet 的平均测试结果给出的预测值为 $\tilde{\eta}_c = 0.62 \pm 0.01$, 与 $\eta_c \approx 0.626$ 基本一致. 即, 该方法也能够成功地从图 1 所示的原始数据中, 自动提取出自驱动活性粒子的群集相变临界噪声强度 η_c .

3.3 “留白法”

现在我们检验所谓的“留白法”^[1,2]. 该方法直接利用 ANN 识别不同物质相的能力. 当所有样本都被贴上符合物理事实的标签 (也就是 3.2 节中提到的 $\eta'_c = \eta_c$ 的情况), 即使仅将极低和极高噪声的样本用于训练, 而将中间噪声“留白”^[1,2], ANN 仍然可以轻松完成相应的二元分类任务. 这里仅将 $\eta = 0.39, 0.41$ (甲类) 和 $\eta = 0.69, 0.71$ (乙类) 的样本用于训练、验证. 训练完成后, 在 $\eta \in [0.43, 0.67]$ 的测试集评估 ANN 识别甲、乙两类的样本的信心.

图 3(c) 的结果取自 20 个独立训练、独立验证的 ResNet 的平均测试结果, 其中虚线、实线分别表示 ANN 将样本识别为甲类、乙类的平均信心 $C(\eta)$ (同一 η 的所有测试样本的平均) 关于样本对应 η 的依赖关系. 两条线在 $\tilde{\eta}_c \approx 0.625$ 交叉. 由于在相变点处, 一个系统的瞬时状态既可能看起来像是处于群集相, 也可能看起来像是处于无序相, 因而“留白法”将 ANN 的平均分类信心取得 $C_1(\eta) = C_2(\eta)$ 的交叉点对应的 η 视为 ANN 给出的临界噪声强度预测值 $\tilde{\eta}_c$. 这个预测值也与传统方法得到的相变点 $\eta_c \approx 0.626$ 基本一致, 说明该方法同样能够从图 1 所示的原始数据中提取 η_c .

4 三种方法对比讨论

图 3 展示了在不额外增加针对非平衡、非晶格的系统中的一阶相变的特殊设计的情况下, 基于 ANN

回归的 LFRU 方法和基于 ANN 分类的“混淆法”“留白法”都能很方便地直接应用于这类复杂多体系统, 提取其中的相变临界值, 这为相关研究提供了一种具有较强通用性的工具箱. 现在进一步讨论三种方法各自的特点.

4.1 使用效率

效率是任何一个实用方法的基本要求. 用于回归和分类任务的 ANN, 其在网络架构上的区别仅在于输出值的个数略有不同 (ISP 回归为 1 个, 二元分类为 2 个), 这导致它们将同样的数据集遍历一次的计算复杂度是几乎相等的. 其训练过程用到的损失函数的计算复杂度也差不多, 且二者的收敛速度接近^[20]. 因此在应用 LFRU 方法和“混淆法”的过程中, 训练每个 ANN 的用时基本相同. 然而“混淆法”训练一个 ANN 只能得到对应于一个 η'_c 取值的识别成功率 $P(\eta'_c)$, LFRU 方法训练一个 ANN 却可以直接得到完整的回归不确定性 $U(\eta)$ 曲线, 这使得前者自动探索物质相的总体用时多于 LFRU 方法. 而“留白法”用于训练、验证的数据集可以远小于另外两种方法, 因此其总体用时是三者中最短的. 但这当然是有代价的, 它需要一些预设的物理知识, 并不能真正实现自动探索物质相.

4.2 预设的物理知识

要想自动探索物质相, 本文机器学习方法不应需要关于物质相和相变的预设的物理知识. 这涉及机器学习的“监督”概念. 在机器学习术语中, “监督学习算法”指的是涉及作为参考答案的标签的机器学习算法. 在这个意义下, 三者作为机器学习算法而言都是有监督的. 然而, 在 2.3 节和 3.1 节可看到, 在将 LFRU 方法应用于揭示自驱动活性粒子的群集相变时, 标签是噪声强度值, 而 LFRU 方法的真正目标是提取临界噪声强度 η_c . 这些标签仅提供关于 ISP 的预设的物理知识, 却完全不涉及物质相和相变. 因此, 对于机器学习在相变研究中的应用而言, LFRU 方法可以被视为一种无监督的方法. 在同样的意义下, 应用“混淆法”时的标签也不是关于临界噪声强度 η_c 的参考答案, 因此该方法通常也被视为一种无监督的方法^[2]. 但值得注意的是, 其二元分类暗含了“系统中最多可能存在两个相”的预先判断, 这使得它需要经过一定的改造之后才可以用于处理具有明显中间相的复杂多体系统^[6].

而“留白法”则将两条线的交叉点视为相变点,这预设了“系统中有且仅有一个相变”,使其在不经改造的情况下,不仅难以处理具有中间相的系统,甚至也无法排除相共存 (phase coexistence) 或平缓过渡 (crossover) 情况的干扰. 无论是相变、相共存、平缓过渡,都会让 ANN 的二元分类信心的曲线相交^[1]. 总之,三种方法之中,“留白法”需要预设的物理知识最多,“混淆法”次之, LFRU 方法则最少.

4.3 联系物理的可能性

自动探索物质相的另一个对机器学习方法的要求是具有可解释性. 由于 ANN 的底层工作机制至今仍未得到足够清晰的解释^[12–14], 这里不考虑机器学习技术本身的可解释性, 三种方法都把 ANN 视为一个黑箱映射. 但在这种情况下, 我们仍希望这些机器学习方法给出的结果能与传统物理概念建立联系. 基于 ANN 分类的“混淆法”在一定程度上就缺乏这样的联系. 该方法提取临界噪声强度 η_c 的最后一步是由研究者而非由 ANN 完成的, 即把甲类与乙类的分界点直接视为群集相与无序相的相变临界点. 这相当于研究者事后向 ANN 补充注入关于该群集相变的物理知识, 事实上削弱了该方法的无监督性. 对于真正待研究系统中的未知相变, 这样的做法缺乏足够的理由. 与之不同的是, 由于 ANN 能够“学会”该系统中的噪声强度 η , 基于 ANN 回归的 LFRU 方法可以较为自然地与传统物理概念建立联系. ANN 处理 ISP 回归时的输出值就是这个被重构的系统参数值 η 本身, 而这些输出值的统计性质 (例如回归不确定性) 则是系统本身的统计性质的体现. 当 ANN 的输出值的统计性质出现特殊的行为, 例如当回归不确定性 $U(\eta)$ 在某个 $\eta = \tilde{\eta}_c$ 处出现非平庸的极小值, 有理由相信系统的统计性质在此处也具有特殊的行为, 这就带来了将该极小值对应的噪声强度 $\tilde{\eta}_c$ 视为群集相变临界噪声强度 η_c 的合理性. 此外, 在关于 LFRU 方法的研究^[20]中发现 Ising 模型和 Clock 模型的回归不确定性与系统的响应性质具有理论上的联系, 预期在 Vicsek 模型中也存在类似的理论联系. 这种在数值上和理论上联系物理的可能性, 是基于 ANN 分类的方法不容易提供的.

5 结 论

在训练 ANN 处理由 Vicsek 模型描述的自驱

动活性粒子系统中的 ISP 回归任务之后, 发现 ANN 的回归不确定性其实隐藏着关于这个非平衡多体系统的群集相变的物理信息. 回归不确定性的 M 字形曲线印证了该一阶相变的存在, 并给出了临界噪声强度值的数据驱动的新估计 $\tilde{\eta}_c = 0.63 \pm 0.01$, 与传统方法得到的相变点 $\eta_c \approx 0.626$ 相符. 这展现了本文发展的利用 ISP 中的回归不确定性的 LFRU 方法用于自动探索物质相的有效性、高效性、对于跨领域的不同物理系统的良好通用性. 该方法与“混淆法”和“留白法”可以相辅相成, 共同构成一个具有较强通用性的工具箱. 对于那些给传统研究方法带来较大挑战的复杂系统, 特别是涉及非平衡、非晶格、一阶相变等复杂要素的情况, 本文讨论的机器学习方法提供了数据驱动实现自动探索物质相的新的视角. ANN 处理 ISP 回归任务的强大能力及其与传统物理概念的直接联系, 这使得我们有机会在接下来的系统性研究中构建回归不确定性与自驱动活性粒子系统的统计性质特别是响应性质的理论联系, 以期在更复杂的相变研究中进一步发挥 LFRU 方法在物理可解释性方面的潜在优势.

参考文献

- [1] Melko R G, Carrasquilla J 2017 *Nat. Phys.* **13** 431
- [2] van Nieuwenburg E P L, Liu Y H, Huber S D 2017 *Nat. Phys.* **13** 435
- [3] Guo W C, Ai B Q, He L 2021 *EPL* **136** 48002
- [4] Venderley J, Khemani V, Kim E A 2018 *Phys. Rev. Lett.* **120** 257204
- [5] Beach M J S, Golubeva A, Melko R G 2018 *Phys. Rev. B* **97** 045207
- [6] Lee S S, Kim B J 2019 *Phys. Rev. E* **99** 043308
- [7] Ch'ng K, Carrasquilla J, Melko R G, Khatami E 2017 *Phys. Rev. X* **7** 031038
- [8] Broecker P, Carrasquilla J, Melko R G, Trebst S 2017 *Sci. Rep.* **7** 1
- [9] Carrasquilla J, 2020 *Adv. Phys.* **X** 5 1797528
- [10] Yu L W, Zhang S Y, Shen P X, Deng D L 2023 *Fundamental Research* (In Press)
- [11] Rem B S, Käming N, Tarnowski M, Asteria L, Fläschner N, Becker C, Sengstock K, Weitenberg C 2019 *Nat. Phys.* **15** 917
- [12] Gökmen D E, Ringel Z, Huber S D, Koch-Janusz M 2021 *Phys. Rev. Lett.* **127** 240603
- [13] Gökmen D E, Ringel Z, Huber S D, Koch-Janusz M 2021 *Phys. Rev. E* **104** 064106
- [14] Miles C, Bohrdt A, Wu R, Chiu C, Xu M, Ji G, Greiner M, Weinberger K Q, Demler E, Kim E A 2021 *Nat. Commun.* **12** 3905
- [15] Nguyen H C, Zecchina R, Berg J 2017 *Adv. Phys.* **66** 197
- [16] Udrescu S M, Tegmark M 2021 *Phys. Rev. E* **103** 043307
- [17] Liu Z, Tegmark M 2022 *Phys. Rev. Lett.* **128** 180201
- [18] Liu Z, Tegmark M 2021 *Phys. Rev. Lett.* **126** 180604

- [19] Udrescu S M, Tegmark M 2020 *Sci. Adv.* **6** eaay2631
- [20] Guo W C, He L 2023 *New J. Phys.* **25** 083037
- [21] Binder K 1987 *Rep. Prog. Phys.* **50** 783
- [22] Falkovich G, Gawędzki K, Vergassola M 2001 *Rev. Mod. Phys.* **73** 913
- [23] Jarzynski C 2015 *Nat. Phys.* **11** 105
- [24] Vicsek T, Czirók A, Ben-Jacob E, Cohen I, Shochet O 1995 *Phys. Rev. Lett.* **75** 1226
- [25] Toner J, Tu Y, Ramaswamy S 2005 *Ann. Phys.* **318** 170
- [26] Grégoire G, Chaté H 2004 *Phys. Rev. Lett.* **92** 025702
- [27] Chaté H, Ginelli F, Grégoire G, Raynaud F 2008 *Phys. Rev. E* **77** 046113
- [28] He K, Zhang X, Ren S, Sun J 2016 *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* Las Vegas, USA, June 27–30, 2016 p770
- [29] Nguyen H C, Berg J 2012 *Phys. Rev. Lett.* **109** 050602
- [30] Jo J, Hoang D T, Periwat V 2020 *Phys. Rev. E* **101** 032107

SPECIAL TOPIC — The 90th Anniversary of South China Normal University and Physics Discipline

Reveal flocking phase transition of self-propelled active particles by machine learning regression uncertainty^{*}

Guo Wei-Chen Ai Bao-Quan[†] He Liang[‡]

(*Institute of Theory Physics, School of Physics, South China Normal University, Guangzhou 510006, China*)

(Received 30 May 2023; revised manuscript received 16 July 2023)

Abstract

We develop the neural network based “learning from regression uncertainty” approach for the automatic detection of phases of matter in nonequilibrium active systems. Taking the flocking phase transition of self-propelled active particles described by the Vicsek model for example, we find that after training a neural network for solving the inverse statistical problem, i.e. for performing the regression task of reconstructing the noise level from given samples of such a nonequilibrium many-body complex system’s steady state configurations, the uncertainty of regression results obtained by the well-trained network can actually be utilized to reveal possible phase transitions in the system under study. The noise level dependence of regression uncertainty is assumed to be in a non-trivial M-shape, and its valley appears at the critical point of the flocking phase transition. By directly comparing this regression-based approach with the widely-used classification-based “learning by confusion” and “learning with blanking” approaches, we show that our approach has practical effectiveness, efficiency, good generality for various physical systems across interdisciplinary fields, and a greater possibility of being interpretable via conventional notions of physics. These approaches can complement each other to serve as a promising generic toolbox for investigating rich critical phenomena and providing data-driven evidence on the existence of various phase transitions, especially for those complex scenarios associated with first-order phase transitions or nonequilibrium active systems where traditional research methods in physics could face difficulties.

Keywords: machine learning, phase transition, nonequilibrium many-body system, inverse statistical problem

PACS: 07.05.Mh, 05.70.Fh, 05.70.Ln, 02.30.Zz

DOI: 10.7498/aps.72.20230896

^{*} Project supported by the National Science Foundation of China (Grant Nos. 12275089, 12075090), the Basic and Applied Research Foundation of Guangdong Province, China (Grant Nos. 2023A1515012800, 2022A1515010449), and the National Key Research and Development Program of China (Grant No. 2022YFA1405304).

[†] Corresponding author. E-mail: aibq@scnu.edu.cn

[‡] Corresponding author. E-mail: liang.he@scnu.edu.cn

机器学习回归不确定性揭示自驱动活性粒子的群集相变

郭唯琛 艾保全 贺亮

Reveal flocking phase transition of self-propelled active particles by machine learning regression uncertainty

Guo Wei-Chen Ai Bao-Quan He Liang

引用信息 Citation: *Acta Physica Sinica*, 72, 200701 (2023) DOI: 10.7498/aps.72.20230896

在线阅读 View online: <https://doi.org/10.7498/aps.72.20230896>

当期内容 View table of contents: <http://wulixb.iphy.ac.cn>

您可能感兴趣的其他文章

Articles you may be interested in

基于机器学习 J_1 - J_2 反铁磁海森伯自旋链相变点的识别方法

Identifying phase transition point of J_1 - J_2 antiferromagnetic Heisenberg spin chain by machine learning

物理学报. 2021, 70(23): 230701 <https://doi.org/10.7498/aps.70.20210711>

基于波动与扩散物理系统的机器学习

Machine learning based on wave and diffusion physical systems

物理学报. 2021, 70(14): 144204 <https://doi.org/10.7498/aps.70.20210879>

通过机器学习实现基于摩擦纳米发电机的自驱动智能传感及其应用

Self-powered sensing based on triboelectric nanogenerator through machine learning and its application

物理学报. 2022, 71(7): 078702 <https://doi.org/10.7498/aps.71.20211632>

铅基钙钛矿铁电晶体高临界转变温度的机器学习研究

High critical transition temperature of lead-based perovskite ferroelectric crystals: A machine learning study

物理学报. 2019, 68(21): 210502 <https://doi.org/10.7498/aps.68.20190942>

机器学习辅助绝热量子算法设计

Machine learning assisted quantum adiabatic algorithm design

物理学报. 2021, 70(14): 140306 <https://doi.org/10.7498/aps.70.20210831>

基于机器学习的非线性局部Lyapunov向量集合预报订正

Machine learning based method of correcting nonlinear local Lyapunov vectors ensemble forecasting

物理学报. 2022, 71(8): 080503 <https://doi.org/10.7498/aps.71.20212260>