

# 关于机器学习的五篇论文

## 文献综述

### 摘要

本综述旨在梳理机器学习发展历程中三种核心范式的确立与演进。通过深入剖析《Learning representations by back-propagating errors》、《Support-Vector Networks》、《Statistical Modeling: The Two Cultures》与《Random Forests》这四篇里程碑式论文，以及一篇关于机器学习解释性的论文《机器学习的可解释性》。本文系统阐述了连接主义中的误差反向传播机制、统计学习理论下的最大间隔原则、算法建模文化的哲学思想以及集成学习中的双重随机性策略。综述揭示了机器学习从追求生物可解释性到强调统计泛化性，再到构建复杂、实用算法模型的内在发展逻辑，并对未来发展趋势进行了展望。

**关键词：**机器学习；反向传播；支持向量机；随机森林；解释性

### 1 引言

机器学习作为人工智能的核心驱动力，其发展并非一蹴而就，而是在关键的理论突破和哲学思辨中逐步成形的。上世纪 80 年代至本世纪初，是机器学习从理论萌芽走向实践繁荣的关键时期。在此期间，几篇开创性的论文为后续研究奠定了坚实的基石，并清晰地勾勒出不同的技术路径。Rumelhart 等人提出的反向传播算法，解决了训练多层网络的核心难题，开启了连接主义的复兴之路；Cortes 与 Vapnik 提出的支持向量机，将统计学习理论付诸实践，展现了最大化间隔原则的强大泛化能力；Breiman 在《两种文化》中的深刻思辨，从哲学层面划分了数据建模与算法建模的界限，解放了机器学习研究者的思想；而他本人提出的随机森林算法，则完美体现了算法文化的精髓，成为集成学习的典范。

本文旨在通过对这五篇论文的核心思想、历史贡献与内在关联进行系统性回顾，梳理机器学习经典范式的演进脉络，从而为理解当代机器学习，特别是深度学习的由来与发展，提供一个清晰的历史视角和理论基础。

### 2 主题内容

#### 2.1 机器学习的可解释性

初学者在学习机器学习时首先就要面对可解释性与预测准确性的两难抉择，上述的《Statistical Modeling: The Two Cultures》这篇论文虽然有力地支持了全面倒向追求预测准确性的选择，但并不是说明可解释性就是可以被舍弃的要素。

随着机器学习尤其是深度学习的快速发展，模型复杂度不断增加，导致可解释性急

剧下降。黑盒模型严重阻碍了机器学习在医疗、金融等高风险领域的应用，比如欧盟GDPR等法规明确要求算法决策必须具备可解释性。这很容易理解，机器学习虽然可以保证“结果”，但是过程的“黑盒”在一些领域是无法容忍的。这使我今后不得不在面对某些实际问题时还需要使用一些方法去理解模型的结构。

这篇论文对可解释性进行定义，区分了Interpretation、Explanation和Understanding三个概念，从模型解释、预测结果解释和模仿者模型三个角度进行了数学形式化。

### 三大技术框架

#### 1. 模型解释技术

- 基于规则的解释：使用决策树等可解释模型提取规则
- 激活值最大化：通过优化输入寻找最大激活模式
- 隐层神经元分析：通过可视化和重构分析隐藏层特征
- 分离式表征：评估隐藏层与语义概念的契合度
- 注意力机制：解释输入与输出之间的对齐关系

#### 2. 预测结果解释技术

- 敏感度分析：研究输入变化对输出的影响
- 泰勒分解：通过泰勒展开分解函数值
- 相关度传播：反向传播相关度得分（如LRP方法）

#### 3. 模仿者模型技术

- 线性分类器拟合：如LIME方法
- 模型压缩：训练浅层网络模拟深度网络
- 知识蒸馏：从大模型蒸馏出小模型
- 其他方法：如GBTmimic模型等

提出了四个核心评估指标：

- 解释一致性：相似输入应产生相似解释
- 解释选择性：通过特征移除评估解释质量
- 解释准确性：解释模型自身的预测精度
- 解释保真度：解释模型模仿原模型的准确程度

在神经网络大行其道的当下，模型的可解释性显示受到很大的挑战，当下的复杂模型往往以牺牲可解释性为代价。该论文的主要价值是为研究者提供了全面的技术路线图，对推动可解释机器学习在实际应用中的落地具有重要意义，特别是在医疗、金融等高风险决策领域。

## 2.2 连接主义的复兴：反向传播与表示学习

1986年，Rumelhart, Hinton和Williams在《Learning representations by back-propagating errors》一文中，清晰地阐述了反向传播算法。该算法的核心价值在于，它利用链式法则将输出层的误差梯度反向传播至网络的每一层，从而为调整隐藏层的连接权重提供了可计算的依据。

这一突破的意义是革命性的。首先，它解决了多层感知机的训练难题，使得构建和使用深度神经网络成为可能。其次，也是更深远的一点，是它实现了“表示学习”。论文表明，神经网络无需人工设计特征，其隐藏层能通过数据驱动的方式，自动学习到从低级到高级的层次化特征表示。这为后来深度学习在计算机视觉、自然语言处理等领域的成功奠定了最核心的训练基础，尽管其巨大潜力因当时算力和数据的限制而迟来了近二十年。

## 2.3 统计学习理论的辉煌：支持向量机与最大间隔原则

作为统计学习理论的杰出代表，Cortes和Vapnik在1995年的《Support-Vector Networks》中提出了成熟的支持向量机模型。SVM的核心思想是结构风险最小化，其目标是找到一个不仅能划分数据，而且能最大化两类数据间隔的分类超平面。

该论文的关键创新在于引入了“软间隔”概念，通过允许部分样本点误分类，极大地增强了模型对噪声和非线性可分数据的处理能力，提升了其实用性。结合“核技巧”，SVM能够隐式地将数据映射到高维特征空间，从而高效地解决非线性分类问题。SVM以其坚实的理论根基、在小样本数据集上的卓越性能以及优美的数学模型，在随后的近二十年里成为了机器学习领域的主流算法之一，展示了统计理论指导下的模型所能达到的泛化性能高度。

## 2.4 算法建模文化的宣言：Breiman 的哲学思辨

2001年，Leo Breiman在《Statistical Modeling: The Two Cultures》中进行的并非技术革新，而是一次深刻的哲学思辨。他犀利地指出统计学界存在两种文化：一是基于预设数据生成模型的“数据模型文化”，二是专注于预测准确性的“算法模型文化”。

这篇论文解释了在使用机器学习算法之前，统计学界总是试图在一个数据集上找到一个符合预期的data model，但真实的情况并不是那么理想，单一的data model总有瓶颈，某个data model可能与另一个完全不同的data model得到完全相同的误差结果，这令人困惑，但这就是现实。因为data model过度关注模型的可解释性，那么模型的结构就必须简单，这样就无法处理复杂的特征信息，从而使预测准确性很差。Breiman（论文的作者）解释道：每一个data model都是从某一个角度理解这些特征，而这些角度不尽相同，所以会出现误差结果相似，但模型差异较大的情况。Breiman还在文中使用“罗生门”解释这种现象。

在做特征工程时为了可以让测试集适用data model，就必须做牺牲，而这常常使一些重要特征被忽略，使得一个data model只能从一个角度片面的理解实际问题，而机器学习算法如SVN，随机森林等不要求模型具有良好的可解释性，只要追求预测准确，所以要尽可能多的提取特征并使用复杂的模型，尽量使有用的特征都囊括进来，同时避免过拟合，由此训练出来的模型具有比data model更好的预测准确性。这样一来模型就成了“黑匣子”，Breiman解释说我们并不需要为此担心，因为我们的目标是让模型具有更好的预测准确性。

这篇论文虽然并不像一篇技术论文，但带给我的启发是巨大的，它启发我在以后的

学习中不再需要为理解模型中复杂的结构而担忧，避免了在未来的某一天完成一个很好的试验之后反过来为理解模型的复杂性而绞尽脑汁

## 2.5 集成学习的典范：随机森林与双重随机性

同样是Breiman，在提出哲学思辨的同一年，他用《Random Forests》论文提供了一个算法建模文化的完美范例。随机森林通过构建大量决策树并进行集成，其核心创新在于引入了“双重随机性”：在构建每棵树时，不仅使用Bootstrap抽样对数据样本进行随机采样，还在每个节点分裂时随机选择部分特征。

这种设计巧妙地降低了森林中所有树之间的相关性。根据Breiman的分析，集成模型的误差取决于单棵树的强度与树间相关性的平衡。随机森林通过双重随机性，在仅轻微牺牲单棵树强度的情况下，极大地降低了树间相关性，从而通过模型平均效应显著提升了整体的泛化能力，并有效防止了过拟合。此外，其内置的袋外估计方法为模型评估和特征重要性分析提供了无偏、高效的工具。随机森林以其开箱即用的高性能、强鲁棒性成为机器学习实践中最可靠、最常用的算法之一。

## 3 总结与展望

从反向传播开启的连接主义道路，到SVM代表的统计学习理论的辉煌，再到《两种文化》为复杂算法模型正名，最终由随机森林这样的实用化集成模型将算法文化的理念推向高峰。它们共同构成了现代机器学习多元而统一的基石，另一方面随着“黑匣子”的模型大行其道，不能简单的抛弃可解释性，在一些领域，可解释性同样重要。

## 4 参考文献

- [1] Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533-536.
- [2] Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine learning*, 20(3), 273-297.
- [3] Breiman, L. (2001). Statistical modeling: The two cultures. *Statistical science*, 16(3), 199-231.
- [4] Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
- [5] 陈珂锐,孟小峰. 机器学习的可解释性[J]. 计算机研究与发展,2020,57(9): 1971-1986.
- [6] 李渭子,侯磊. 知识图谱研究综述[J]. 山西大学学报(自然科学版),2017,(3): 454-459.
- [7] 马帅,刘建伟,左信. 图神经网络综述[J]. 计算机研究与发展,2022,59(1): 47-80.
- [8] 周培诚,程堃,姚西文,等. 高分辨率遥感影像解译中的机器学习范式[J]. 遥感学报,2021,25(1): 182-197.
- [9] 刘辰屹,徐明伟,耿男,等. 基于机器学习的智能路由算法综述[J]. 计算机研究与发展,2020,57(4): 671-687.
- [10] 李国良,周煊赫,孙佶,等. 基于机器学习的数据库技术综述[J]. 计算机学报,2020,43: 2019-2049.
- [11] 米晓希,汤爱涛,朱雨晨,等. 机器学习技术在材料科学领域中的应用进展[J]. 材料导报,2021,35(15): 15115-15124.
- [12] 梁英凯,商枫楠,陈桥,等. 基于机器视觉与机器学习的火龙果重量估计[J]. 食品与机械,2023,39(7): 99-103.
- [13] 钟宇,徐燕,刘德祥,等. 基于计算机视觉和机器学习的真伪卷烟包装鉴别[J]. 烟草科技,2020,53(5): 83-92.
- [14] 周若彤,谭凯,杨建儒,等. 结合SAM 视觉分割模型与随机森林机器学习的无人机影像盐沼植被

- “精灵圈”提取[J]. 海洋学报,2024,46(5): 116-126.
- [15] 董成烨,李东方,冯槐区,等. 基于机器视觉和机器学习技术的浙贝母外观品质等级区分[J]. 浙江大学学报(农业与生命科学版),2023,49(6): 881-892.
- [16] 付永民,范磊,李长进,等. 基于计算机视觉与机器学习的烟丝杂质图像级联检测方法[J]. 轻工学报,2023,38(4): 113-121.
- [17] 顾文君,李强,周易,等. 基于机器视觉与机器学习的苹果体积质量测量系统研究[J]. 工业控制计算机,2025,38(9): 76-78.
- [18] 焦李成,杨淑媛,刘芳,等. 神经网络七十年: 回顾与展望[J]. 计算机学报,2016,39(8): 1697-1716.
- [19] 蓝金辉,王迪,申小盼. 卷积神经网络在视觉图像检测的研究进展[J]. 仪器仪表学报,2020,41: 167-182.
- [20] 关胜晓. 机器视觉及其应用发展[J]. 自动化博览,2005,(3): 88-92.