# ● 医学信息研究 ●

# 机器学习方法在因果推断中混杂因素控制的应用\*

兰雨姗 郑 思 李 姣

(中国医学科学院/北京协和医学院医学信息研究所 北京 100020)

[摘要] 介绍医学研究中混杂因素对因果推断的影响及常见混杂因素识别方法,梳理机器学习方法在因果推断中控制混杂因素的应用,讨论应用机器学习方法在混杂因素控制中面临的机遇和挑战。

[关键词] 机器学习;因果推断;混杂因素控制

[中图分类号] R - 058 [文献标识码] A [DOI] 10. 3969/j. issn. 1673 - 6036. 2022. 11. 004

Machine Learning Methods for Confounding Factor Control in Causal Inference LAN Yushan, ZHENG Si, LI Jiao, Institute of Medical Information, Chinese Academy of Medical Sciences & Peking Union Medical College, Beijing 100020, China

[Abstract] The paper introduces the influence of confounding factors on causal inference in medical studies and the identification methods of common confounding factors, summarizes the application of machine learning methods for controlling confounding factors in causal inference, and discusses the opportunities and challenges of applying machine learning methods to control confounding factors.

(Keywords) machine learning; causal inference; confounding factor control

## 1 引言

因果关系是指某因素是结局发生的原因,因果推断反映一种在试验设计和分析过程中对混杂、偏移等的慎重考虑,从而在得出因果关系结论时排除各种干扰的影响,做出正确的结论<sup>[1]</sup>。混杂因素是

[修回日期] 2022-03-10

[作者简介] 兰雨姗,硕士研究生;郑思,副研究员;通 信作者:李姣,研究员,博士生导师。

[基金项目] 中国医学科学院医学与健康创新工程"医学知识管理与智能化知识服务关键技术研究"(项目编号: 2021-I2M-1-056);中国医学科学院医学与健康创新工程"医学人工智能算法评价标准库构建"(项目编号: 2018-I2M-AI-016)。

指某个既与暴露有关又与结局相关的因素,该因素可能会使暴露和结局之间的因果关系产生偏移。例如一项关于口服避孕药和心肌梗死之间关联的研究中,在未考虑混杂因素的情况下得到口服避孕药诱发心肌梗死的比值比(Odds Ratio,OR)为 2.20,而在控制年龄因素(将研究对象按照年龄是否小于40岁进行分组)的影响后,得到口服避孕药和心肌梗死之间的 OR 值为 2.79<sup>[2]</sup>。年龄这一混杂因素的存在减弱了口服避孕药和心肌梗死之间的关联。因此混杂因素的控制和识别是因果推断中的关键<sup>[3]</sup>。传统流行病研究中常采用限制、配对、随机化、工具变量等方法控制混杂因素,但随着医学大数据的不断积累,混杂因素的维度不断增加,传统方法难以较好地处理高维特征。因此越来越多的研究将机器学习算法引入混杂因素控制领域,希望借助机器

学习算法良好的分类和预测能力提升估计因果效应的能力。本文介绍在观察性研究数据中识别混杂因素的方法及如何利用机器学习方法对识别到的或潜在的混杂因素进行控制,见图1。

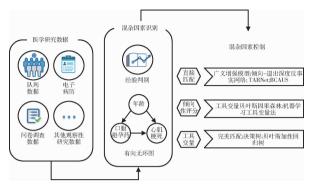


图 1 因果推断中混杂控制的流程

## 2 因果推断定义及其方法

## 2.1 定义

因果推断是研究变量间因果关系的学科,流行病学领域提出一些因果推断理论用于推断出暴露和结局间的因果关系。例如密尔氏法则,即求同法、求异法、共变法和排除法。1965 年流行病学病因研究中因果推断的9条准则被提出,即时间顺序、关联强度、剂量反应关系、可重复性、合理性、考虑可替代的解释、实验证据、关联的特异性、关联的一致性,该标准被简称为希尔准则(Hill's Criteria),仍广泛地用于人群研究中判断因果关系。

#### 2.2 方法

随着相关研究积累,产生了许多因果推断框架,其中结构性因果模型和潜在结果框架得到广泛应用。结构性因果模型由 Judea Pearl 教授于 1995年提出<sup>[4]</sup>,包括因果图和结构方程,可以用来描述一个系统的因果机制。因果图即有向无环图(Directed Acyclic Graph,DAG),一个有向无环图能唯一确定一个联合分布。结构方程<sup>[5]</sup>是一种建立、估计和检验因果关系模型的多元回归方法。从结构性因果模型中,Pearl 教授又提取出 Pearl 因果层次结构(Pearl Causal Hierarchy,PCH)<sup>[6]</sup>,该结构将因果推断分为3个层次,包括关联、干预和反事实推

论。反事实推论是因果推断中最高层次的问题,用于估计治疗的潜在效果,回答一个经典问题"如果采取不同的行动会怎么样"。当混杂因素存在时会对反事实问题的推导产生影响,例如上述研究中年龄的存在导致估计未服用避孕药的患者是否会患心肌梗死时,低估了口服避孕药对于患心肌梗死的影响。潜在结果框架由哈佛大学 Donald Rubin 教授提出,目的是估计不能被观察到的潜在结果,从而估计实际的干预效果,潜在结果框架又被称为鲁宾因果模型。干预的效果可以被定义为: ATE = E ( $Xi \mid Zi = 1$ ) -E ( $Xi \mid Zi = 0$ ) [7],其中,Z 代表研究所施加的干预措施,X 是研究的结局。干预的效果是干预措施对结局的因果效应。

## 3 混杂因素的识别

混杂因素是指与暴露对结局的影响相混淆的另一个风险或保护因素<sup>[8]</sup>。观察性研究中存在大量已知和未知的混杂因素,这些混杂因素存在时会歪曲暴露和结局之间真实的因果关联。例如在上述关于口服避孕药与心肌梗死的关联研究中,避孕药是要研究的暴露因素,心肌梗死是本研究的结局,而年龄是一个与口服避孕药对心肌梗死的影响相混淆的风险因素,即一个混杂因素。因此如何识别和控制混杂因素是进行因果推断时需要考虑的重要问题之一。

#### 3.1 利用经验识别

经验识别<sup>[9]</sup>是指根据已掌握的专业经验知识来 判断某个因素是否为混杂因素,当某个因素满足以 下条件时即可被判定为混杂因素:该因素与暴露因 素相关;该因素与结局相关;该因素不是暴露与结 局之间的中间变量。在许多疾病的病因因果推断 中,通常会考虑控制年龄、性别等常见的混杂因素 来得到合理的因果关系,但对于是否要控制某个外 部因素就需要慎重的考虑。需要参考专家意见,考 虑控制该外部因素后得到的结果与专业知识是否符 合,与同类研究和既往研究进行比较等,才能得出 最终结论。

### 3.2 利用有向无环图识别

有向无环图是关于暴露、结局以及其他相关变 量之间假设关系的图形表示, 当存在多个混杂因素 时,有向无环图可以将各变量间的关系更加直观地 表示出来,帮助研究人员更加全面地识别混杂因 素。有向无环图是由节点和箭头组成,每个箭头的 起点被称为父节点,而该箭头指向的节点被称为子 节点。当一个有向无环图中存在"后门"路径时, 提示研究中存在混杂因素, "后门"路径是指存在 一个混杂因素既指向研究因素又指向结局,从而导 致研究因素和结局之间表现出相关性。如前文图 1 所示,在有向无环图中存在年龄这一因素既指向口 服避孕药,又指向心肌梗死,该路径是一个"后 门"路径。"后门"路径表明年龄既与研究因素有 关又与结局有关,是研究中的一个混杂因素。控制 混杂因素的过程实际上是切断"后门"路径,从而 排除混杂因素的干扰。控制混杂因素的过程可以看 作是固定混杂因素的值, 当混杂因素的值给定后, 暴露和结局间的相关性与混杂因素无关,相关性能 够反映出因果性。

# 4 利用机器学习控制混杂因素的方法

混杂因素的存在有时会导致错误的因果关联,根据 Yule - Simpson 悖论,当忽略了第 3 个变量时,2 个变量间的相关性可能会从正相关变为负相关。这表明良好的混杂因素控制方法有助于推断出正确的因果关系。有学者<sup>[10]</sup>利用机器学习算法估计倾向性得分,从而达到控制混杂因素的目的。该研究利用模拟的队列数据比较 Logistic 回归和机器学习建立的倾向性评分模型,结果表明由机器学习方法构建的倾向性评分模型具有更小的估计误差,能够更好地控制混杂因素。传统流行病学中控制混杂因素的方法包括限制、匹配、随机化、分层分析、多元分析等,本研究介绍机器学习与匹配、倾向性评分及工具变量法相结合,提升控制混杂因素的能力。

## 4.1 基于机器学习的样本匹配方法

匹配是流行病学中常用的控制混杂因素的方·22·

法,可以确保某些变量的分布在暴露组和对照组之 间相同或尽可能相同,提高估计因果效应的效 率[11]。直接匹配的思想是通过计算处理样本和对照 样本间的距离, 距离越小的样本间差异越小, 将距 离暴露组样本最近的对照组样本进行匹配。根据数 据集的不同特点可以选择不同的距离函数,例如马 氏距离、欧式距离等。随着医学数据的不断积累, 协变量特征不断增加,基于机器学习算法的直接匹 配将会带来更好的匹配效果。树模型的原理是根据 变量特征对样本进行分组,利用回归树可以将样本 协变量特征与阈值依次进行比较, 从而将不同样本 纳入不同分组中,确保暴露组和对照组间协变量分 布相似或相同。有学者[12]提出交互树的概念,交互 树是利用随机森林直接识别影响因果效应异质性的 变量的重要程度。之后,有学者[13]将贝叶斯算法加 入到随机森林模型上,提出贝叶斯加性回归树 (Bayesian Additive Regression Trees, BART) o BART 模型可以自动识别变量之间的非线性关系,并且能 够估计异质性因果效应[[14],该方法在估计因果效 应方面得到广泛应用[15]。BART 模型对超参数规范 有显著的鲁棒性,并且适用于高维度情况[13]。研究 者陆续提出因果树[16]、因果森林[17]及广义随机森 林[18] 等方法,这些方法与传统的近邻匹配相比能够 更好地匹配样本,并且能够更加准确地估计异质性 因果效应。此外,深度学习算法也被用于样本匹 配。例如在最近邻匹配的思想上结合神经网络方法 提出完美匹配方法(Perfect Match, PM)<sup>[19]</sup>, 这是 一种通过训练神经网络来对样本进行直接匹配的方 法。PM 方法易于实现,可以应对不同场景和数据 集,并且不会增加任何超参数或计算复杂度。有学 者<sup>[20]</sup>利用搜索算法比较不同样本间的马氏距离,从 而提出一种多元匹配方法。该方法可以用于改善样 本匹配后协变量平衡的问题。其中,协变量平衡是 指暴露组和对照组的协变量具有相同的联合分布。 引入搜索算法后无需手工迭代检查不同样本间距 离,能够更加高效地进行匹配。另有学者[21]提出一 种用于估计个体因果效应的有限混合模型,该方法 能够用于分析潜在变量对于结局的因果效应,并且 可以对样本进行分类。

## 4.2 基于机器学习的倾向性评分方法

4.2.1 概述 倾向性评分 (Propensity Score, PS) 是一种仅使用协变量评分来衡量一个人接受治疗的 可能性的方法。倾向性评分的主要目标是实现协变 量的平衡,从而控制评估治疗或暴露的平均效果时 的混杂偏差,对标量倾向得分的调整足以消除由于 所有观测协变量造成的偏差[22]。倾向性评分技术可 以将高维度特征压缩为某一个复合特征, 可以直接 评价暴露组和对照组在背景特征方面是否相似[23]。 将倾向性评分用于因果推断过程可以达到控制混杂 效应的作用。合理地利用倾向性评分进行匹配、回 归、分层,可以在估计因果效应时减小选择偏倚的 影响,实现"事后随机化"。当一个暴露个体和一 个未暴露个体具有相同或相近的倾向性评分时,可 以认为这两个个体的治疗分配不受任何混杂因素的 影响,该暴露个体和未暴露个体之间的差异可以用 于回答反事实问题,从而推断暴露因素与结局之间 的因果关系。

4.2.2 基于树模型的倾向性评分方法 倾向性评 分匹配是最常见的利用倾向性评分的方法。利用决 策树可以对样本进行倾向性评分匹配,决策树是机 器学习中用于分类和回归的一种非监督学习方式, 可以帮助直接进行倾向性评分匹配。利用决策树对 一组个体进行分类时,决策树可以将个体划分为多 个叶子节点,在每个叶子节点内分类的所有数据点 都具有相似的分类概率,因此决策树可以直接将研 究对象分为暴露组和对照组[24]。此外,通用梯度回 归模型 (Generalized Boosted Regression Model, GBM)已应用于估计倾向性评分。在考虑大量预测 因素的情况下, GBM 也可以产生平衡不同组间协变 量分布的模型<sup>[25]</sup>。GBM 通过其迭代程序,找到使 处理组和对照组之间达到最佳平衡的倾向性评分模 型[26],从而实现研究对象的"随机化"。随机森 林、增强型分类和回归树等方法基于协变量对数据 进行递归分区来分配治疗组和对照组,从而合理评 估变量间的因果效应。其中, 随机森林的方法还可 以有效处理协变量中的缺失数据[27]。在流行病学研 究中,随机森林算法已经被广泛用于构建疾病预测 模型[28-29]。

4.2.3 基于深度学习的倾向性评分方法 将深度 学习算法引入倾向性评分匹配也是当前因果推断中 控制混杂的趋势之一。有学者[30]考虑利用生成式对 抗网络 (Generative Adversarial Networks, GAN) 生 成合成数据作为未被匹配的实验个体的对照,从而 推断出两组间的差异,判断因果关系。倾向-退出 的深度反事实网络 (Deep Counterfactual Network with Propensity - Dropout, DCN - PD) 是一种多任务 学习方法,可以用来减少暴露组和对照组之间的选 择偏移。有研究者使用一个深度多任务网络来模拟 被试的潜在结果,该网络包含一组事实和反事实结 果之间的共享层,以及一组特定结果层。通过倾向 性评分计算每个样本被排除的概率,并使用该概率 对网络进行正则化,最后利用正则化后的网络输出 暴露组和对照组,从而达到控制混杂的目的[31]。 4.2.4 基于神经网络倾向性评分方法 神经网络 算法也被用于倾向性评分匹配,一种可以用于估计 个体治疗效果的神经网络——TARNet 被提出[23]. 研究者利用一个双头多任务模型来估计二元治疗的 效果,其中每个样本被分配一个特定的权重,从而 平衡治疗组和对照组之间的混杂因素。有学者[32]在 TARNet 的基础上、提出 Dragonnet 方法、利用倾向 性评分对模型进行正则化修饰, 从而避免模型过拟 合。还有研究者提出 BCAUS<sup>[33]</sup> 方法,这是一种自 动因果推理方法,该模型通过对指定治疗的错误预 测以及逆概率加权变量之间的不平衡程度进行惩 罚,之后再与传统的基于倾向得分的方法结合使 用,估计变量间的因果效应。其中,逆概率加权 (Inverse Probability Weighting, IPW) 是根据倾向性 评分计算得到的,为研究对象接受其实际接受暴露 条件概率的倒数,即 $\frac{1}{1-PS}$ ,其中 PS 是该研究对象 的倾向性评分。逆概率加权与倾向性评分相同,用 干评估患者分配到暴露因素的概率,从而平衡暴露 组和对照组间的混杂因素。尽管利用倾向性评分可 以较好地控制研究中的混杂因素, 在利用倾向性评 分对研究对象进行匹配时, 容易将所有因素都当作 混杂因素,但并不是所有的因素都是混杂。因此利

用倾向性评分控制混杂因素容易导致数据间信息量的降低。为解决这一问题,有学者<sup>[34]</sup>提出一种数据驱动的变量分解(Data – Driven Variable Decomposition, D2VD)算法,该算法可以自动分离混杂因素和调整变量,同时估计治疗效果。

## 4.3 基于机器学习的工具变量法

工具变量(Instrumental Variable, IV)法是利用一个额外的工具变量识别变量间的因果关系。工具变量法早期主要应用于经济学和社会学领域,后来用于流行病学研究中混杂因素的控制。工具变量法的基本原理是利用一个只与暴露因素相关,而与其他混杂因素无关的变量来评估暴露和结局之间的因果关系。工具变量的质量是工具变量分析的核心问题<sup>[35]</sup>,因此越来越多的研究利用机器学习方法来构建工具变量,从而提高工具变量质量。有学者<sup>[36]</sup>将贝叶斯方法与工具变量法相结合,提出工具变量贝叶斯因果森林(Bayesian Causal Forest with Instrumental Variable,BCF - IV)方法,用于发现和估计

异质性因果效应。BCF-IV 是建立在贝叶斯加性回 归树上的一个半参数贝叶斯回归模型。有学者[37]提 出一种利用深度学习网络估计工具变量的方法,将 工具变量法分为两个阶段,通过调整损失函数来估 计因果效应。有学者<sup>[38]</sup>提出一种机器学习工具变量 (Machine Learning Instrument Variables, MLIV) 算 法,利用机器学习算法的非线性建模方式和正则化 方法来优化工具变量的构造过程, 能够在构造工具 变量的同时进行因果推断。此外,虽然工具变量法 可以帮助识别变量间的因果关系, 但当工具变量构 建缺乏精确性时会导致因果推断结果也缺乏一定的 精确性。常见的提高工具变量精确性的方法包括使 用多种工具或利用近似最优工具。有学者[39]提出一 种基于"多工具"的渐进方案,利用正则化的机器 学习方法构建工具变量,通过将套索回归的方法与 工具变量的估计相结合可以提高构建工具变量的收 益。机器学习方法能够提升混杂因素的控制效果, 但可能会存在过拟合等问题,见表1。

表 1 机器学习方法及其应用

表 1			
方法	应用领域	具体应用	特点 (优点/缺点)
决策树	倾向性得分匹配	抗生素治疗和儿科患者死亡率及肾	优点:可直接划分试验组和对照组;可以用来处理分类、
		毒性的因果关系[40]	有序、连续和缺失数据
			缺点:对于异常值和变量的单调变换不敏感;在建模平滑
			函数和主效应方面存在困难; 可能存在过拟合
广义增强模型	倾向性得分匹配	在缓刑时相关治疗对于青少年药物	优点: 计算速度快; 可以分析非线性效应和交互作用项;
		滥用的疗效[41]	可用于拟合平滑模型
			缺点: 仅能分析有限个因素间的交互作用
深度学习	倾向性得分匹配	低出生体重早产儿中婴儿认知测试	优点: 能够处理高维、非线性/非平行治疗分配; 可以与
		得分的影响因素[31]	传统方法结合使用;可通过调整损失函数来估计因果效
	直接匹配	专家就诊对儿童认知发展的影响[19]	应,避免了模型过拟合
	工具变量	机票价格与客户是否购买的关系[37]	缺点: 只适用于不存在潜在混杂因素的情况
神经网络	倾向性得分匹配	探究抗糖尿病药物对高水平糖化血	优点:能够处理协变量较多的情况;高度自动化;拟合速
		红蛋白的影响 (HbAlc) <sup>[33]</sup>	度较快
			缺点: 在少数情况下,难以保证所有协变量平衡
套索回归	工具变量	各种社会经济因素对病毒传播的影	优点:避免了模型过拟合的情况
		哨[42]	缺点:难以分离出每一项干预的影响
搜索算法	直接匹配	探索加入某职业培训项目的人参与	优点:较传统方法更加高效;可解决协变量平衡问题
		工作类型的影响因素[20]	缺点: 当治疗分配的估计概率接近0或1时模型效果不佳
生成式对抗网络	倾向性得分匹配	探索与抗生素耐药性相关基因[30]	优点:可解决样本量不足和权重不稳定的问题;适用于高维数据
			缺点:无法克服由于数据不平衡而产生的偏差

## 5 结语

机器学习方法应用于因果推断领域, 在控制混 杂因素进而得到合理的因果效应方面发挥着重要作 用。本文介绍因果推断、混杂因素定义、混杂因素 对因果推断的影响及常见的混杂因素识别方法。重 点介绍3种控制混杂因素的方法,即基于机器学习 的样本匹配方法、基于机器学习的倾向性评分方 法、基于机器学习的工具变量法,以及这些方法如 何提升估计因果效应的能力。在基于机器学习的倾 向性评分方法中, 重点介绍基于树模型的倾向性评 分方法、基于深度学习的倾向性评分方法和基于神 经网络的倾向性评分方法。在因果推断混杂因素控 制方面, 机器学习较传统的简单线性模型更具有灵 活性,能够用来分析暴露和结局间的非线性关系, 从而更加全面地评估暴露和结局间的因果关系,提 升预测准确性。在分析高维数据时, 机器学习方法 能够更好地控制和识别混杂因素,提升暴露组和对 照组间的匹配效果,实现"随机化"。此外,许多 机器学习算法利用模型集成的方法来降低单一方法 可能存在的预测偏差,使得预测结果更加精准。但 是机器学习方法在因果推断中控制混杂因素尚存在 不足,包括机器学习的过拟合、训练样本集构建、 结果可解释性等问题。机器学习模型通过纳入大量 参数和复杂的非线性关系提升预测能力,这可能导 致模型出现过拟合的问题, 使得机器学习模型进行 外部验证时效果不佳。若加入惩罚项来消除模型中 的过拟合问题容易造成因果效应估计量的一致性被 破坏, 当前机器学习领域有学者提出双重机器学 习、样本分割等方法来尝试解决该问题。此外,机 器学习模型是一个"黑箱"模型,研究者通常难以 了解模型构建过程,导致预测结果缺乏可解释性。 因此在利用机器学习方法定义因果关系时应更加谨 慎,需要考虑多种证据后才能得出确切结论。机器 学习算法作为一种统计推断的方法,需要在正确的 因果推断框架内使用才能达到最佳效果。在因果推 断领域,良好的试验设计方案是正确估计因果效应 的关键所在。

## 参考文献

- 黄丽红,赵杨,王陵,等.获得现实世界证据的因果推 断统计学思考[J].中国临床医学,2021,28 (5): 738-743.
- 2 詹思延.流行病学(第八版)[M].北京:人民卫生 出版社,2017.
- 3 陶秋山,李立明.基于虚拟事实理论的病因效应模型「J].中华流行病学杂志,2014,23 (1):60-62.
- 4 Pearl J. Causal Diagrams for Empirical Research [J]. Biometrika, 1995, 82 (4): 669-688.
- 5 Rex B K. Principles and Practice of Structural Equation Modeling (Fourth Edition) [M]. New York: Guilford Press, 2016.
- 6 Pearl J. Theoretical Impediments to Machine Learning With Seven Sparks from the Causal Revolution [C]. Marina Del Rey: 11th ACM International Conference on Web Search and Data Mining, 2018.
- 7 Rubin D B. Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies [J]. Journal of Educational Psychology, 1974, 66 (5): 688 701.
- 8 Howards P. An Overview of Confounding. Part 1: The Concept and How to Address It [J]. Acta Obstetricia Et Gynecologica Scandinavica, 2018, 97 (4): 394 399.
- 9 胡永华, 耿直. 关于混杂概念的讨论 [J]. 中华流行病 学杂志, 2001, 22 (6): 459-461.
- 10 Setoguchi S, Schneeweiss S, Brookhart M A, et al. Evaluating Uses of Data Mining Techniques in Propensity Score Estimation: a Simulation Study [J]. Pharmacoepidemiology and Drug Safety, 2008, 17 (6): 546-555.
- Mansournia M A, Jewell N P, Greenland S. Case Control Matching: Effects, Misconceptions, and Recommendations [J]. European Journal of Epidemiology, 2018, 33 (1): 5-14.
- 12 Su X, Tsai C L, Wang H, et al. Subgroup Analysis via Recursive Partitioning [J]. Journal of Machine Learning Research, 2009 (10): 141-158.
- 13 Chipman H A, George E I, Mcculloch R E. BART: Bayesian Additive Regression Trees [J]. The Annals of Applied Statistics, 2010, 4 (1): 266-298.
- 14 Green D P, Kern H L. Modeling Heterogeneous Treatment Effects in Survey Experiments with Bayesian Additive Regression Trees [J]. Public Opinion Quarterly, 2012, 76 (3): 491-511.
- 15 Leonti M, Cabras S, Weckerle C S, et al. The Causal De-

- pendence of Present Plant Knowledge on Herbals contemporary Medicinal Plant use in Campania (Italy) Compared to Matthioli (1568) [J]. Journal of Ethnopharmacology, 2010, 130 (2): 379 391.
- 16 Athey S, Imbens G. Recursive Partitioning for Heterogeneous Causal Effects [J]. Proceedings of the National Academy of Sciences of the United States of America, 2016, 113 (27): 7353-7360.
- Wager S, Athey S. Estimation and Inference of Heterogeneous Treatment Effects using Random Forests [J]. Journal of the American Statistical Association, 2017, 113 (523): 1228-1242.
- 18 Athey S, Tibshirani J, Wager S. Generalized Random Forests [J]. The Annals of Statistics, 2019, 47 (2); 1148-1178.
- 19 Schwab P, Linhardt L, Karlen W. Perfect Match: a Simple Method for Learning Representations for Counterfactual Inference with Neural Networks [EB/OL]. [2021 12 20]. https://arxiv.org/abs/1810.00636.
- 20 Diamond A, Sekhon J S. Genetic Matching for Estimating Causal Effects: a General Multivariate Matching Method for Achieving Balance in Observational Studies [J]. Review of Economics and Statistics, 2013, 95 (3): 932 - 945.
- 21 Louizos C, Shalit U, Mooij J, et al. Causal Effect Inference with Deep Latent - Variable Models [J]. 31st Annual Conference on Neural Information Processing Systems (NIPS), 2017 (30): 6449 - 6459.
- 22 Rosenbaum P R, Rubin D B. The Central Role of the Propensity Score in Observational Studies for Causal Effects [J]. Biometrika, 1983, 70 (1): 41-55.
- 23 Shalit U, Johansson F D, Sontag D. Estimating Individual Treatment Effect: Generalization Bounds and Algorithms [J]. Sydney 34th International Conference on Machine Learning, 2017 (70): 3076 - 3085.
- 24 Cook E F, Goldman L. Asymmetric Stratification. An Outline for an Efficient Method for Controlling Confounding in Cohort Studies [J]. American Journal of Epidemiology, 1988, 127 (3): 626-639.
- 25 Mccaffrey D F, Ridgeway G, Morral A R. Propensity Score Estimation with Boosted Regression for Evaluating Causal Effects in Observational Studies [J]. Psychological Methods, 2004, 9 (4): 403-425.
- 26 Mccafdrey D F, Griffin B A, Almirall D, et al. A Tutorial on Propensity Score Estimation for Multiple Treatments Using Generalized Boosted Models [J]. Statistics in Medicine,

- 2013, 32 (19): 3388 3414.
- 27 Zhao P, Su X, Ge T, et al. Propensity Score and Proximity Matching Using Random Forest [J]. Contemporary Clinical Trials, 2016 (47): 85 - 92.
- Yang L, Wu H, Jin X, et al. Study of Cardiovascular Disease Prediction Model Based on Random Forest in Eastern China [J]. Scientific Reports, 2020 (10): 5245.
- 29 Di Castelnuovo A, Bonaccio M, Costanzo S, et al. Common Cardiovascular Risk Factors and In – Hospital Mortality in 3 894 Patients with COVID – 19: Survival Analysis and Machine Learning – based Findings from the Multicentre Italian CORIST Study [J]. Nutrition, Metabolism, and Cardiovascular Diseases, 2020, 30 (11): 1899 – 1913.
- 30 Ghosh S, Boucher C, Bian J, et al. Propensity Score Synthetic Augmentation Matching Using Generative Adversarial Networks (PSSAM GAN) [J]. Computer Methods and Programs in Biomedicine Update, 2021 (1): 10.
- 31 Alaa A M, Weisz M, Van Der Schaar M. Deep Counterfactual Networks with Propensity Dropout [EB/OL]. [2022 01 10]. https://arxiv.org/pdf/1706.05966.pdf.
- 32 Shi C, Blei D M, Veitch V. Adapting Neural Networks for the Estimation of Treatment Effects [EB/OL]. [2022 01 10]. https://arxiv.org/pdf/1906.02120.pdf.
- 33 Belthangady C, Stedden W, Norgeot B. Minimizing Bias in Massive Multi – arm Observational Studies with BCAUS: Balancing Covariates Automatically Using Supervision [J]. BMC Medical Research Methodology, 2021 (21): 190.
- 34 Kuang K, Cui P, Zou H, et al. Data Driven Variable Decomposition for Treatment Effect Estimation [ J ]. IEEE Transactions on Knowledge and Data Engineering, 2020: 1.
- 35 Brookhart M A, Wang P, Solomon D H, et al. Evaluating Short – term Drug Effects Using a Physician – specific Prescribing Preference as an Instrumental Variable [J]. Epidemiology (Cambridge, Mass.), 2006, 17 (3): 268 – 275.
- 36 Bargagli Stoffi F J, De Witte K, Gnecco G. Heterogeneous Causal Effects with Imperfect Compliance: a Bayesian Machine Learning Approach [EB/OL]. [2022 01 10]. https://arxiv.org/pdf/1905.12707v2.pdf.
- 37 Hartford J, Lewis G, Leyton Brown K, et al. Deep IV: a Flexible Approach for Counterfactual Prediction [EB/OL]. [2022 – 01 – 10]. https://www.cs.ubc.ca/~jasonhar/slides/deepiv.pdf.

(下转第33页)

## 参考文献

- 国家自然科学基金委员会. 资助格局 [EB/OL]. [2021 08 14]. http://www.nsfc.gov.cn/publish/portal0/jg-sz/08/default.htm#02.
- 2 刘仲林. 交叉科学时代的交叉研究 [J]. 科学学研究, 1993 (2): 11-18.
- 3 Rafols I, Meyer M. How Cross disciplinary is Bionanotechnology? Explorations in the Specialty of Molecular Motors [J]. Scientometrics, 2007, 70 (3): 633 650.
- 4 Porter A L, Roessner J D, Cohen A S, et al. Interdisciplinary Research: Meaning, Metrics and Nurture [J]. Research Evaluation, 2006, 15 (3): 187-195.
- 5 Porter A L, Rafols I. Is Science Becoming more Interdisciplinary? Measuring and Mapping Six Research Fields over Time [J]. Scientometrics, 2009, 81 (3): 719-745.
- 6 商宪丽. 基于 LDA 的交叉学科潜在主题识别研究——以数字图书馆为例 [J]. 情报科学, 2018, 36 (6): 57-62.
- 7 阮光册, 夏磊. 学科间交叉研究主题识别——以图书情报学与教育学为例 [J]. 情报科学, 2020, 38 (12): 152-157.
- 8 徐庶睿, 卢超, 章成志. 术语引用视角下的学科交叉测度——以 PLOS ONE 上六个学科为例 [J]. 情报学报, 2017, 36 (8): 809-820.
- 9 Xu H, Guo T, Yue Z, et al. Interdisciplinary Topics of Information Science: a Study Based on the Terms Interdisciplinarity Index Series [J]. Scientometrics, 2016, 106 (2): 583-601.
- 10 Dong K, Xu H, Luo R, et al. An Integrated Method for In-

- terdisciplinary Topic Identification and Prediction: a Case Study on Information Science and Library Science [J]. Scientometrics, 2018, 115 (2): 849 868.
- 11 韩正琪, 刘小平, 寇晶晶. 基于 Rao Stirling 指数和 LDA 模型的领域学科交叉主题识别——以纳米科技为 例「J]. 情报科学, 2020, 38 (2): 116-124.
- 12 叶春蕾. 基于 Web of Science 学科分类的主题研究领域 跨学科态势分析方法研究 [J]. 图书情报工作, 2018, 62 (2): 127-134.
- 13 李长玲, 高峰, 牌艳欣. 试论跨学科潜在知识生长点及其识别方法 [J]. 科学研究, 2021, 39 (6): 1007-1014.
- Otte E, Rousseau R. Social Network Analysis: a Powerful Strategy, also for the Information Sciences [J]. Journal of Information Science, 2002, 28 (6): 441-453.
- 15 Web of Science. Research Areas [EB/OL]. [2021 11 04]. http://images.webofknowledge.com//WOKRS534DR3/help/zh\_ CN/WOS/hp\_ research\_ areas\_ easca. html.
- 16 王旻霞, 赵丙军. 跨学科知识交流网络结构特征研究 [J]. 情报科学, 2016, 34 (5): 46-50, 104.
- 17 Leydesdorff L, Rafols I. Indicators of the Interdisciplinarity of Journals: Siversity, Centrality, and Citations [J]. Journal of Informetrics, 2011, 5 (1): 87-100.
- 18 黄水清,张俊,阎素兰.黄金分割法在学科及机构评价中的应用「J〕.图书情报工作,2012,56(22):33-36.
- 19 Review M T. MIT Technology Review Presents 10 Breakthrough Technologies of 2021 [EB/OL]. [2021 11 15]. https://www.technologyreview.com/press releases/mit technology review presents 10 breakthrough technologies of 2021/.

## (上接第26页)

- 38 Singh A, Hosanagar K, Gandhi A. Machine Learning Instrument Variables for Causal Inference [EB/OL]. [2021 12 20]. https://www.researchgate.net/publication/342899660 \_ Machine\_ Learning\_ Instrument\_ Variables\_ for\_ Causal \_ Inference.
- 39 Belloni A, Chen D, Chernozhukov V, et al. Sparse Models and Methods for Optimal Instruments with an Application to Eminent Domain [ J ]. Econometrica, 2010, 80 (6): 2369 - 2429.
- 40 Tamma P D, Turnbull A E, Harris A D, et al. Less Is More; Combination Antibiotic Therapy for the Treatment of

- Gram negative Bacteremia in Pediatric Patients [J]. JA-MA Pediatrics, 2013, 167 (10): 903 910.
- 41 Mccaffrey D F, Ridgeway G, Morral A R. Propensity Score Estimation with Boosted Regression for Evaluating Causal Effects in Observational Studies [J]. Psychological Methods, 2004, 9 (4): 403-425.
- 42 Qiu Y, Chen X, Shi W. Impacts of Social and Economic Factors on the Transmission of Coronavirus Disease 2019 (COVID - 19) in China [J]. Journal of Population Economics, 2020, 33 (4): 1127 - 1172.