构建多维度AI审计框架思考

◎ 文/ 林慧涓 陈宋生

一、人工智能对财会领域的影响: 机会和威胁

近年来深度学习技术的崛起、神经网络架构的改进、计算能力的增强以及互联网数据的可用性,共同推动人工智能领域的技术革新。这些人工智能技术在财务和会计方面已经展示出巨大的潜力和价值。

首先,在数据分析方面,AI和 机器学习算法能够处理大量复杂 的财务数据,提供更及时的财务报 告。这不仅提高财务报告的质量, 也使得财务报告使用者的决策更 快速

其次,在风险评估和管理方面, AI技术能够分析历史数据和市场趋势,预测潜在的财务风险,从而助力企业更有效地进行资本配置和风险规避。例如 AI 算法能够识别信贷风险、市场风险和操作风险等多种类型的财务风险,为企业提供更全面的风险评估。

最后,在自动化和效率提升方面,AI技术也发挥重要作用。通过自动化财务流程,如发票处理、支付和收款等,AI不仅降低人力成本,还提高财务操作的准确性和效率。此外,AI还能够自动执行复杂的财务模型和算法,从而在预算编制、资金分配和投资决策等方面提供更高水

平的自动化支持。在合规和审计方面,AI技术能够自动检测财务报告和交易记录中的不规范和异常行为,从而提高财务合规和审计的效率。不仅在财经方面有着广泛的的应用前景,而且在医疗诊断、代码生成、语言翻译等多个领域体现其优越的性能。特别是像 ChatGPT 这样的大型语言模型,在某些基准测试上已经接近人类的性能水平。

然而,AI技术的广泛应用也带来 伦理和社会的挑战。Weidinger et al. 指出,大型语言模型如 ChatGPT 存 在加剧社会偏见的风险。这些模型 在输出内容中可能无意中强化对特 定社会群体的偏见性描述,从而加 剧现有社会的刻板印象。此外,这 些模型还可能被用于恶意目的,包 括生成高度复杂和个性化的欺诈行 为或进行大规模的欺诈活动。同 时,AI的架构存在数据泄露和敏感 信息推断的潜在风险,这严重威胁 个人隐私。因此,从伦理角度对AI 严格分析不仅需要扩展现有的治理 框架,还需要在风险评估方法、性能 基准和伦理框架方面取得革命性的 进展。正是这些挑战需要一种跨学 科的方法,将数据科学、管理学、伦 理学等多个学科联合起来,以全面 审视和解决这些挑战。特别是在人 工智能审计方面,需要构建一套全 面而有效的审计机制,以确保AI技 术的健康、可持续发展。

二、人工智能审计的必要性

随着人工智能(AI)技术在各个领域的广泛应用,其潜在的风险和挑战日益凸显。这些风险不只局限于技术层面,还涉及社会、伦理、法律和经济等多方面。因此,对AI进行全面和深入的审计成为迫切需求,有必要通过分析未经审计的AI模型可能带来主要问题,探讨开展人工智能审计的必要性。AI模型的主要问题:

(一)数据偏见

大型语言模型通常基于大量的 文本数据进行训练。如果这些数据 包含某种形式的偏见(如性别偏见、 种族偏见等),模型可能会学习并在 输出中反映这些偏见。这不仅可能 误导用户,还可能加剧社会不平等 和歧视,从而对社会和伦理方面产 生不良影响。

(二)不准确的信息传播

由于AI模型是基于其训练数据 进行预测和决策的,如果训练数据 包含错误或过时的信息,模型的输 出也可能是不准确的。这在医疗、 法律或金融等关键领域尤为危险, 因为不准确的信息可能导致严重的 后果,包括误诊、误判或财务损失。

(三)安全性问题

未经审计的 AI 模型可能更容易受到对抗性攻击。这些攻击旨在

通过输入特定的数据来误导模型, 从而产生不准确或出现误导性的输 出。这不仅威胁到模型的可靠性和 有效性,而且可能导致更广泛的安 全风险。

(四)法律和伦理风险

AI模型在生成内容时可能触犯 法律和伦理规定,如侵犯版权或涉嫌诽谤。使用这些未经审计的模型 的个人或组织可能因此面临法律责 任。此外,模型可能生成不道德或 不合法的内容,从而引发更广泛的 伦理和社会问题。

(五)可解释性和透明度

由于未经审计的 AI 模型通常 缺乏可解释性,用户难以理解模型 的决策逻辑和依据。这不仅影响用 户对模型的信任,还可能导致错误 或不合理的决策。

(六)社会和文化影响

AI模型可能会生成与特定文化 或社会观念不一致的内容,从而引 发文化冲突或误解。这不仅可能导 致模型在某些文化或社会环境中的 应用受限,还可能加剧社会分裂和 冲突。

因此,审计在AI模型的治理体系中被认为是一个重要且务实的工具。这一观点不仅得到学界的重视,而且受到高科技企业和政策制定者的重视。例如,欧洲委员会定在考虑将大型语言模型归类为"高风险AI系统",这意味着设计这些模型的技术供应商必须接受"独立第三方参与的符合性评估",即另一种名为审计的活动。

AI 审计的早期研究主要集中在 探索和缓解算法应用中可能产生的

歧视及其他不良影响,特别是在决 策支持系统以及对个体和组织的影 响方面。近年来,研究领域逐渐引 入"基于伦理的自动决策系统审计" (EBA)这一概念。EBA被构想成一 种结构化流程,用于评估某实体(无 论是现存还是历史)的行为是否与 相关原则或规范相符。与此相对, Brown et al.(2021)提出的伦理算法 审计更侧重于评估算法对利益相关 者产生的实际影响,并据此识别特 定算法的情境或属性。两者主要区 别在于伦理算法审计更侧重于实际 影响,而EBA更注重原则和规范的 一致性。此外,伦理算法审计(Brown et al., 2021)明确将算法作为审计的 主要对象,而没有保留被审计实体 的开放性。

以大型语言模型审计为例,AI 审计可细分为模型审计、治理审计 和应用审计三个关键阶段,以全面 而系统地评估与人工智能相关的各 方面风险。

首先,模型审计专注于预训练 但尚未发布的大型语言模型技术属 性,以揭示可能反映在历史不公正 的训练数据集中偏见或其他歧视性 做法。该审计应由所有设计和传播 大型语言模型的技术供应商强制执 行,并产生一份详细报告以概述模 型的属性和局限性。

其次,治理审计主要针对技术 供应商在设计和传播语言模型过程 中的管理和决策机制。该审计阶段 由独立的第三方或内部机构执行, 旨在核实技术供应商自我报告的准 确性,并通常只能访问有限的内部 流程.其主要目标是评估语言模型 的设计和传播流程,尽管它无法预 见所有随着AI系统与复杂环境随 时间互动而出现的风险。

最后,应用审计则是针对基于 语言模型的具体应用,包括两个。 阶段:功能性审计和影响审计。 的能性审计评估应用是否本身致; 功能性审计评估应用是否本身致; 的审计则更关注应用输出如何。 这是 种审计类型各自审计与语言模型 并有的用户群体和自然环境。 。 数型 等的审计类型各自审计与语言相相关 ,并在实践中相实,并在实践中相关 以中 可能带来的风险。

在AI审计的概念框架和评估 方法中,应综合考虑基于原则和影 响的两种方法,而不是过早地界定 该领域。以当前发展迅速的大型语 言模型为例,从治理角度看,大型语 言模型带来方法论和规范方面的挑 战。它们通常分两个阶段开发与应 用。首先是在大型、非结构化文本 语料库上进行预训练,其次在较小、 特定任务的数据集上进行微调。这 使得在没有特定情景下很难评估大 型语言模型。此外,它们在特定任 务上或在规模上的表现可能是不可 预测的。因此,治理和技术审计都 有局限性。治理审计不能预见所有 随着AI系统与复杂环境随时间互 动而出现的风险。并且,历史上专 注于在明确定义环境中的特定功 能,没有能力捕捉大型语言模型等 AI模型在许多下游应用中带来的社 会和伦理风险的全貌。因此,尽管 设计针对AI模型的审计程序面临 诸多实践性和概念性的困难,但这

27

些困难不应成为不开展AI审计的 借口。相反,这些问题应视为警示, 即不能期望单一的审计机制能全面 捕捉与AI模型相关的所有伦理风 险。作为一种治理机制,审计能助 力技术提供商识别并可能预防风 险,影响AI模型的持续(重新)设计, 并丰富有关技术政策的公共讨论。 尽管面临多重挑战,审计在确保AI 模型的准确性、可靠性和伦理性方 面仍具有不可或缺的核心作用。

三、人工智能审计的框架构思

在当今高度数字化的商业环境 中,如何避免使用人工智能(AI)带来 的风险是一个相对较新且日益重要 的研究领域。为此有必要从AI审 计的定义、实践应用,以及在全球范 围内的法律和政策方面构建人工智 能审计框架。

(一)人工智能审计的定义与范畴

审计是一种治理机制,它可以 用于监控行为和绩效,并且在诸如 财务会计和工人安全等领域促进程 序规则性和透明度的悠久历史。因 此,就像财务交易可以针对真实性、 准确性、完整性以及合法性进行审 计一样,AI审计中,AI系统的设计 和使用也可以针对技术安全性、法 律合规性或与预定义的伦理原则一 致性进行审计。AI审计是一种系统 性和独立的活动,它涉及获取和评 估与特定实体(无论是AI系统、组 织、流程,或是这些元素的任何组 合)的行为或属性相关的证据,并将 评估结果传达给利益相关方。它强 调审计活动的系统性和独立性,突 出在高度依赖数据和算法的现代商

业环境中,审计在治理机制中的重 要性。被审实体不局限于算法,而 是开放的,可以是一个AI系统、一 个组织、一个过程,或者是它们的任 何组合.根据其侧重点,AI审计可分 为狭义和广义两种。狭义的人工智 能审计是以影响为导向,重点在于 探测和评估AI系统对不同输入数 据的输出。广义的人工智能审计则 是以过程为导向,侧重于评估整个 软件开发流程和质量管理流程。

(二)实践应用层面: AI 审计的 应用场景

纽约市的AI审计法(NYC Local Law 144)要求使用 AI 进行就业决策 的公司必须接受独立审计。此外, 全球范围内的专业服务公司,如普 华永道(PwC)、德勤(Deloitte)、毕马威 (KPMG)和安永(EY),也提供AI审计 服务,以帮助企业评估和管理与AI 应用相关的风险。在法律与政策背 景方面,在全球范围内,AI审计逐渐 受到法律和政策的关注。尽管全球 范围内人工智能(AI)审计的实践与 学术研究呈逐渐上升的趋势,但该 领域仍然是一个复杂且至关重要的 研究领域。它不仅跨越多个学科领 域和研究方法论,而且受到各类法 律和政策环境的深刻影响。在我 国,该领域的研究处于关键阶段。

(三)AI审计框架

随着关于如何审计AI项目的 指导方针日益普及,多个国际组织 和政府已发布一系列有助于内部审 计功能的AI审计框架。

COBIT框架。最新版本的COBIT 框架,即COBIT2019,由信息系统审 计与控制协会(ISACA)在2018年发

布,以取代其前身COBIT5。作为一 个"综合性"框架,COBIT在企业信 息和技术的治理与管理方面得到了 国际认可。该框架包括几乎所有IT 领域的流程描述、期望结果、基础实 践和工作产品,因此非常适合作为 审计AI启用项目时内部审计功能 的初始起点。然而, COBIT2019的 综合性也是其弱点之一。由于覆盖 范围广泛,该框架可能缺乏针对特 定 AI 应用场景的深入指导。此外, 其复杂性可能导致实施难度增加。

COSO ERM 框架。由赞助组织 委员会(COSO)在2017年更新的COSO ERM 框架包括五个组成部分和20个 原则,为内部审计提供一种集成和全 面的风险管理方法。COSO ERM 的 风险管理方法可以为AI的治理提供 指导,并有效地管理其相关风险,以造 福组织。然而,COSO ERM 主要侧重 于风险管理,可能忽视AI治理中的 其他关键方面,如伦理和社会责任。

美国审计署(GAO)AI框架。GAO 开发并于2021年6月发布的人工智 能问责框架旨在"帮助管理者确保 AI在政府程序和流程中的负责任使 用"。尽管这一AI审计框架主要关 注政府使用AI的受托责任,但由于 它基于现有的控制和政府审计标 准,主要适用于政府组织,可能不适 用于私营企业或非营利组织。此 外,它可能缺乏对特定AI技术或应 用的深入分析。

国际内部审计帅协会(IIA)人工 智能审计框架。由IIA发布的人工 智能审计框架包括三个主要组成 部分和七个元素,有助于内部审计 在短期、中期或长期内评估、理解和 传达人工智能对组织创造价值能力 的影响。该框架主要侧重于审计, 可能不足以全面地解决 AI 治理的 复杂性和多维性。

四、人工智能审计的 建设方向和建议

在AI审计的复杂体系中,内部 审计与外部审计各自扮演着不可或 缺的角色,共同构建一个多维度的 审计框架。内部审计作为企业自我 评估和监控的重要机制,其核心价 值不仅局限于模型和算法的性能与 安全性,而且延伸至对企业内部流程 和决策机制的深度审查。从信息透 明度和准确性的角度,内部审计有助 于确保模型设计和性能的自我评估 是准确和可靠的,这一点在企业商业 决策和战略规划中具有至关重要的 作用。通过内部审计,企业能更精准 地了解到模型的优缺点,从而做出更 明智和合理的商业决策。此外,内部 审计员由于能够访问企业的内部流 程,因此能全面评估模型从设计到 应用的全过程,这不仅有助于识别 和管理潜在的技术和安全风险,还 确保模型的设计和实施与企业的整 体战略和目标是一致的。

然而,内部审计也面临一系列 挑战和局限性。内部审计需要高度 重视企业的商业利益和伦理风险, 以避免与被审计对象之间存在潜在 的利益冲突。例如,如果一个企业 的AI模型是由高级管理层直接推动 或批准的,内部审计员可能会面临来 自上级的压力,隐瞒或忽略模型存在 的问题。特别是在面临激烈市场竞 争和快速技术发展的环境中,任何对 模型性能的负面评价都可能影响企 业的市场地位和竞争力。此时,内部 审计很难独立发表审计意见,特别是 当审计结果与企业的商业目标不一 致时。因此,审计中缺乏第三方的问 责机制可能会导致决策者忽视或淡 化那些可能威胁到商业利益的审计 改进建议,这可能影响模型的性能和 可靠性,还可能引发一系列伦理和法 律问题。例如,如果内部审计报告指 出,AI模型存在数据偏见,但是要解 决这一问题,将降低模型性能,那么 管理层可能忽视或淡化这类问题,以 维护企业利益。如果存在性别或种 族偏见的信贷评估模型用于实践中, 可能触发社会不满,还可能导致企 业面临法律风险。

与内部审计相辅相成,外部审 计作为第三方独立评估机制,其核 心价值体现在全方位、多维度的审 查能力。在模型性能评估方面,外 部审计不仅对模型在特定任务和数 据集上的准确性进行深入评估,而 且对模型在不同文化和背景上的适

用性进行评价。这一点明确了审计 师需要具备跨学科的知识体系和视 野。从安全性维度出发,外部审计 致力于深度识别和评估模型可能存 在的安全漏洞,这包括对抗性攻击、 数据泄露和未经授权的访问等问 题。这体现了审计工作的严谨性和 全面性。在伦理和合规性方面,通 过先进的数据分析和算法检测技 术,外部审计可识别模型是否存在 性别、种族、文化或其他形式的偏见 和歧视,并据此提出针对性的改进 建议。这一环节不仅体现审计工作 的社会责任感,也是对模型公平性 和合规性的有力保证。除了技术和 伦理两个主要维度外,模型的可解 释性和透明度也是外部审计的重要 组成部分。这涵盖了模型决策逻辑 的可解释性、模型训练和应用过程 的透明度,以及与模型相关的所有 文档和元数据的完整性和准确性, 有助于提升模型的社会可接受度和 信任度。从更为宏观的社会和文化 影响角度来看,外部审计还需关注 模型可能对社会结构、公众舆论及 信息传播等方面产生的长期和短期 影响。这不仅是对模型影响力的全 面评估,也是对其可持续发展性和 社会责任的一种体现。

因此,内部审计与外部审计应 当相互补充,与其他治理机制协同 作用,以实现更全面和有效的AI模 型审计。这不仅要求审计员具备跨 学科的知识和视野,还要求其能全 面评估和管理模型在技术、商业、伦 理和社会等多个方面的风险和影 响。这一多维度的审计框架为AI 模型的全周期治理提供了一个全面 而深入的研究基础。

五、人工智能审计可能带来的影响

AI审计作为综合性的治理机 制,其核心目标不会仅局限于提升模 型在技术方面的准确性与可靠性,而 是更广泛地涵盖社会、伦理、法律及 经济等多个维度的标准与期望。从 技术角度出发,AI审计通过精密的 数据分析和模型验证,能有效地识别 并修正模型的潜在缺陷与误差,从而 增强模型在具体应用场景中的准确 性和可靠性,促使其整体性和稳健性 的提升。在社会与伦理层面,AI审 计具有减缓数据偏见和歧视现象的 功能。通过对模型的数据处理和决 策机制进行严格审查,审计活动能揭 示潜在的偏见或歧视,并据此提出针 对性的纠正措施,对缓解模型在实际 应用中可能引发的社会不平等和不 公具有至关重要的意义。从透明度 和可解释性方面考虑,AI审计通常 会产出翔实的分析报告,以明确模型 的运行逻辑和决策依据,这不仅增加 模型的透明度和可解释性,也有助于 提升用户及其他利益相关方对模型 的信任度。在法律和合规性方面,AI 审计发挥着风险规避和合规保障的 作用。通过审计,能确保模型和算法 的运行符合相关法律和规定,如数据 保护法和版权法,从而降低AI应用 可能导致的法律风险,并为构建更完 善的合规体系提供坚实的基础。从 伦理和社会责任视角看,AI审计确 保模型在设计和应用阶段严格遵循 公平、透明和责任等伦理原则,这不 仅提升模型的社会接受度,还能促进 AI的开发和应用实践。

在经济层面,经过审计的AI模 型通常表现出更高的运行效率和准 确度,从而有助于降低运营成本和 提高经济效益。此外,通过减少潜 在的法律风险和提升用户及其他利 益相关方的信任,模型的经济价值得 以进一步放大。从多元利益相关者 的视角来看,AI审计在企业、政府、 消费者和学术界等多个层面都具有 不可忽视的价值,对企业而言,AI审 计不仅是品牌信誉建设的有效手段, 还是一种战略性资产,能够提供更高 的投资同报率,吸引更多资本注入, 进而推动企业的技术创新和市场拓 展。对政府和监管机构来说,AI审 计生成的数据和分析报告为政策制 定提供科学依据,有助于更准确地把 握AI技术的社会影响,从而制定更 合理和有效的法规与政策。此外,AI 审计在识别和纠正模型偏见及歧视 方面具有重要作用,这不仅促进社会 公平和稳定,也是构建和谐社会的必 要条件。从消费者和公众角度,AI 审计增强模型的透明度和可解释性, 有助于保障消费者的知情权和选择 权,同时也增强公众对AI应用的整 体信任度。在学术界和研究人员层 面,AI审计作为一种研究工具,可以 深入探究AI模型在社会、伦理和法 律等方面的复杂影响,这不仅推动 了AI伦理和社会影响等领域的学 术研究,还可能促进不同学科和领 域之间的交叉合作与知识整合。

六、结论与展望

AI审计的实施,有助于确保人工智能模型、算法及其应用在技术、伦理和社会应用层面达到社会的期

望。通过分析数据和检验模型,AI审 计能有效提升模型在实际应用中的 准确性和可靠性,进而优化运行效率 和降低运营成本。然而,面对人工智 能,尤其是强人工智能(AGI)的快速发 展,审计领域将遭遇更为复杂和多维 的挑战。从审计方法论角度看,传统 审计手段可能不完全适用于AI审计, 因为这些方法通常针对更为简单和专 门化的系统。因此,学术界和实践界 需要多学科合作,共同研发适用于AI 的新型审计方法和工具。总体而言, AI审计将面临多层次、多维度和跨学 科的技术与伦理等多方面的挑战,这 需要全球范围内的企业、高校与研究 机构的通力合作和技术创新,方能保 证AI模型的安全性、可靠性和伦理合 规性。综合考虑技术、法律、伦理和 社会因素,AI审计不仅是技术进步的 必要条件,还是实现AI有效应用的 关键途径,有助于构建一个更为安 全、公平和可持续的AI生态系统。

基金项目:北京理工大学珠海学院实验教学示范中心"智能财经人才实验教学示范中心"(2023006ZLGC);教育教学改革一重点项目"基于知识图谱的智能会计核心课程群改革计划"(2023008ZLGC);北京理工大学珠海学院校级课题"预算绩效管理对粤港澳大湾区企业高质量发展的影响研究"(XK-2023-013);线下一流课程"管理会计模拟"(2023012YLKC)。

作者简介:林慧涓,美国注册管理会 计师,北京理工大学珠海学院财务管理 系主任、副教授;陈宋生,北京理工大学 管理与经济学院教授、博士生导师,全国 会计领军人才,北京市优秀人才。

原载《会计之友》(太原),2023.23.32~37