

# 健康咨询对话式人工智能应用实践指南研究

马骋宇<sup>1\*</sup> 吴海锋<sup>2</sup> 张毓辉<sup>3,4</sup> 顾雪非<sup>5</sup> 韩优莉<sup>1</sup>

1. 首都卫生管理与政策研究基地/首都医科大学公共卫生学院 北京 100069
2. 抖音集团旗下小荷健康医学中心 海南海口 571925
3. 海南省卫生健康委员会 海南海口 570203
4. 中国卫生经济学会 北京 100191
5. 国家卫生健康委卫生发展研究中心 北京 100032

**【摘要】**人工智能技术发展为公众获取健康信息提供了新途径。目前国内外已出现健康咨询类对话式人工智能的研究及实践探索,但对其发展中的关键问题,尚缺乏统一意见。在此背景下,本研究在广泛检索国内外文献的基础上,经过多轮专家咨询,针对健康咨询对话式人工智能的定义、特点、应用场景、评价要素、评价方法、展望与挑战等六大问题开展研究,形成健康咨询对话式人工智能应用实践指南,为推动人工智能技术在医疗健康领域的创新应用和规范发展提供参考。

**【关键词】**对话式人工智能;健康咨询;实践指南

中图分类号:R197 文献标识码:A doi:10.3969/j.issn.1674-2982.2025.04.009

## Research on practice guideline of health conversational artificial intelligence

MA Cheng-yu<sup>1</sup>, WU Hai-feng<sup>2</sup>, ZHANG Yu-hui<sup>3,4</sup>, GU Xue-fei<sup>5</sup>, HAN You-li<sup>1</sup>

1. Capital Health Management and Policy Research Base, School of Public Health, Capital Medical University, Beijing 100069, China
2. Xiaohe Health Medical Center, Haikou Hainan 571925, China
3. Health Commission of Hainan Province, Haikou Hainan 570203, China
4. Chinese Society of Health Economics, Beijing 100191, China
5. China National Health Development Research Center, Beijing 100032, China

**【Abstract】**The development of artificial intelligence (AI) technology has provided new avenues for the public to access health information. While both domestic and international research and practical applications of health conversational AI have emerged, there remains a lack of consensus on the key issues surrounding their development. Based on a comprehensive review of domestic and international literature and multiple rounds of expert consultations, the study focuses on six issues: the definition, characteristics, application scenarios, evaluation elements, evaluation methods, as well as prospects and challenges of AI-driven health consultation systems, and develops the practice guideline for health conversational artificial intelligence. The findings aim to provide scientific reference for promoting the innovative application of AI technology in the healthcare field.

**【Key words】**Conversational artificial intelligence; Health consultation; Practice guideline

构建“以人为本”的医疗卫生服务体系,其核心在于全面、深入地理解和满足患者的多元化服务需求,良好的沟通被视为其中重要的一环。有效的沟

通机制不仅是医患之间信息传递的重要桥梁,更是建立信任、理解与支持的关键。随着人口老龄化趋势的不断加剧,公众对更好的健康和更优质的医疗

\* 作者简介:马骋宇(1981年—),女,博士,教授,主要研究方向为智慧医疗、互联网医疗。E-mail:machengyu@ccmu.edu.cn  
通讯作者:韩优莉。E-mail:hanyouli@ccmu.edu.cn

卫生服务需求日益高涨,这对于医疗卫生服务体系的高效、精准和人性化水平提出了更高要求。然而,我国医疗卫生服务资源配置的不均衡导致患者和健康信息寻求者难以获得与医生进行充分沟通交流的机会。随着人工智能技术特别是大语言模型(Large Language Model, LLM)的快速发展,健康咨询对话式人工智能具备的学习、理解和整合大规模相关知识并通过自然对话进行交互的能力,为满足患者和公众的健康信息与服务需求提供了新策略和新途径。<sup>[1]</sup>相关应用在促进健康信息传播、提升医疗卫生服务质量和效率,为患者提供优质、便捷、高效且可负担的医疗卫生服务方面,具有巨大的应用潜力和社会价值。

然而,人工智能在医疗卫生服务领域的应用也伴随着有效性、安全性、伦理道德以及法律合规性等一系列挑战。<sup>[2-3]</sup>2023 年 7 月,国家互联网信息办公室、国家发展和改革委员会等七部门联合发布《生成式人工智能服务管理暂行办法》<sup>[4]</sup>,为人工智能的技术应用和安全性提供了明确的指导。在此背景下,为了更好地规范健康咨询对话式人工智能相关技术和产品的应用场景,提供科学评估框架,在借鉴国内外相关研究及实践应用的基础上,开展关于健康咨询对话式人工智能实践应用指南的研究,旨在为推动人工智能技术在医疗卫生服务领域的创新应用,保障用户的合法权益,增强应用场景的适应性,制定科学、系统、全面的医疗健康大模型评估标准提供参考。

## 1 指南的发起及目标人群

本实践指南研究由首都卫生管理与政策研究基地、国家卫生健康委卫生发展研究中心、中国卫生经济学会青年卫生经济委员会和抖音集团旗下小荷健康医学中心联合发起,启动时间为 2024 年 8 月,实践指南文本已在国际实践指南注册与透明化平台注册(注册号:PREPARE-2024CN1204)。

本实践指南使用者为健康咨询对话式人工智能平台及应用的开发者、运营者、医务人员、卫生系统管理者与决策者、各级行政主管部门以及相关领域的科研工作者。目标人群为使用相关应用及平台的健康信息寻求者,包括患者、亚健康人群及健康人群。

## 2 指南研制方法

### 2.1 专家遴选

指南研制专家的人选标准为对健康咨询对话式人工智能熟悉且从事相关工作,具有丰富理论与实践经验的领域专家,最终遴选专家 22 人,包括医学人工智能专家 3 人、医院管理专家 3 人、临床专家 2 人、卫生政策研究者 5 人、法学及伦理学研究者 3 人、指南方法学专家 2 人和医疗大模型平台开发者 4 人。专家负责指南的审阅、评价和修改。

### 2.2 文献检索

在中国知网、万方数据库、PubMed 和 Web of Science 数据库中对相关研究进行检索。中文检索关键词为“健康咨询”“聊天机器人”“会话代理”“对话式人工智能”“大语言模型”“生成式人工智能”“医疗卫生”“医疗健康”“医学”;英文检索关键词为“health consultation”“chatbot”“chatterbot”“Conversational AI”“ChatGPT”“AI chat agent”“Large Language Model”“generative AI”“NLP”“Machine Learning”“Deep Learning”“transformer model”“health”“healthcare”“medical”“medicine”。利用以上关键词,采用布尔组合检索策略(AND, OR 和 NOT),检索了 2004—2024 年发表的中英文文献,通过阅读标题和摘要排除不相关文献,并阅读其余全文。本指南根据文献检索证据撰写指南草案,并根据专家意见对指南进行修改与完善。

### 2.3 专家咨询

专家咨询采用改良的 Delphi 法<sup>[5]</sup>,进行两轮咨询。共包括 3 个环节:(1)针对指南初稿开展第 1 轮在线问卷调查,征求专家意见,评估是否达成重要性共识。(2)根据专家反馈意见对指南初稿进行修改,并发放第 2 轮在线问卷调查,对反馈意见进行持续修改直到达成共识。(3)通过组织会议,向参与专家反馈研制结果。

问卷调查中,邀请专家对健康咨询对话式人工智能发展中面临问题的重要性进行评分,并说明原因及修改意见。专家意见采用李克特五级量表收集,包括“完全不重要”“比较不重要”“不确定”“比较重要”和“非常重要”。指南条目草案满足重要性评分“非常重要” $\geq 80\%$ ,记为“通过”,不满足记为“待定”。指南起草组在第 1 轮问卷调查结束后征询专家意见,并修改相关条目陈述,将相关内容制订为

第2轮问卷。第2轮问卷调查的专家意见采用李克特三级量表收集,包括“同意”“不同意”“不确定”。指南条目如满足同意 $\geq 90\%$ ,记为“通过”,不满足记为“未通过”。指南起草组在第2轮调查后根据专家意见对指南草案进行进一步修改。

### 3 结果

#### 3.1 专家咨询

本研究共进行两轮专家咨询,第一轮邀请专家16人,第二轮邀请专家22人。22位专家中,高级职称专家占比90.90%,10年以上医疗健康相关工作经验的专家占比81.82%。专家工作单位均衡分布于医疗机构、大学院校和研究机构(表1)。

表1 咨询专家基本信息

| 基本信息     | 第一轮(n=16) |       | 第二轮(n=22) |       |
|----------|-----------|-------|-----------|-------|
|          | 人数        | 占比(%) | 人数        | 占比(%) |
| 性别       |           |       |           |       |
| 男        | 6         | 37.50 | 7         | 31.82 |
| 女        | 10        | 62.50 | 15        | 68.18 |
| 年龄(岁)    |           |       |           |       |
| ≤40      | 6         | 37.50 | 9         | 40.91 |
| 41~50    | 9         | 56.25 | 13        | 59.09 |
| >50      | 1         | 6.25  | 0         | 0.00  |
| 教育水平     |           |       |           |       |
| 本科       | 1         | 6.25  | 1         | 4.55  |
| 硕士       | 5         | 31.25 | 6         | 27.27 |
| 博士       | 10        | 62.50 | 15        | 68.18 |
| 职称       |           |       |           |       |
| 中级       | 1         | 6.25  | 2         | 9.10  |
| 副高       | 7         | 43.75 | 10        | 45.45 |
| 正高       | 8         | 50.00 | 10        | 45.45 |
| 工作年限(年)  |           |       |           |       |
| ≤10      | 3         | 18.75 | 4         | 18.18 |
| 11~15    | 4         | 25.00 | 6         | 27.27 |
| 16~20    | 5         | 31.25 | 8         | 36.37 |
| >20      | 4         | 25.00 | 4         | 18.18 |
| 工作单位     |           |       |           |       |
| 医疗机构     | 4         | 25.00 | 5         | 22.73 |
| 大学院校     | 6         | 37.50 | 7         | 31.82 |
| 研究机构     | 5         | 31.25 | 6         | 27.27 |
| 其他       | 1         | 6.25  | 4         | 18.18 |
| 研究方向     |           |       |           |       |
| 卫生管理与政策学 | 7         | 43.75 | 10        | 45.45 |
| 计算机科学    | 1         | 6.25  | 2         | 9.09  |
| 医学伦理学    | 1         | 6.25  | 1         | 4.55  |
| 卫生法学     | 2         | 12.50 | 2         | 9.09  |
| 医院管理     | 2         | 12.50 | 3         | 13.64 |
| 其他       | 3         | 18.75 | 4         | 18.18 |

两轮专家积极系数均为1.00,表明专家积极性较高。第1轮咨询专家16人,问卷包括9个条目,其中通过5项,分别为“健康咨询对话式人工智能的特点”“用户体验评价”“伦理与安全评价”“智能评价方法”“手工评价方法”;待定4项,分别为“健康咨询对话式人工智能的概念”(非常同意占比68.75%,11/16)、“健康咨询对话式人工智能的应用场景”(62.50%,10/16)、“健康咨询对话式人工智能模型技术水平评价”(68.75%,11/16)、“健康咨询应用有效性评价”(75.00%,12/16)。指南起草组根据专家意见修改了指南草案的条目表述,第2轮咨询专家22人,通过问卷对修改后的4个条目进行调查,结果全部通过,所邀请专家对9个条目达成一致意见。

#### 3.2 指南研制主要结果

##### 3.2.1 健康咨询对话式人工智能的概念

健康咨询对话式人工智能是指利用大语言模型、智能算法和医学知识库等人工智能技术,模拟人类医师与健康信息寻求者进行智能对话,提供健康教育与知识科普、症状自诊自查、健康生活方式指导与疾病预防、用药咨询与指导、健康管理、健康教育、就医导航服务等个性化健康信息服务,是对话式人工智能在医疗卫生服务领域应用的新阶段。

##### 3.2.2 健康咨询对话式人工智能的特点

健康咨询对话式人工智能融合了生成式人工智能的对话能力和健康咨询的医学专业性,与传统健康信息获取渠道(如搜索引擎等)相比,具备以下特点:

(1)知识专业性。集成大量专业医学和健康数据,如最新的研究成果、医学指南和常见健康问题的解答等,训练数据来源广泛且可不断扩展,涵盖多个医学领域和全生命周期的健康问题。

(2)自主学习性。健康咨询对话式人工智能可自动捕捉用户输入的重点内容和关键信息,通过数据驱动学习,不断识别并分析用户的真实需求。

(3)实时交互性。健康咨询对话式人工智能可以迅速响应用户提问,在短时间内处理大量咨询请求,并支持连续多轮问答,通过结构化地展示决策逻辑与分析路径,使健康咨询过程清晰可见、有理有据,有效提升健康咨询服务效率。

(4)响应个性化。根据用户的个人健康数据、历史记录,健康咨询对话式人工智能能够提供个性化的健康咨询服务。

(5)情感智能化。健康咨询对话式人工智能具

备情感识别能力,能感知用户的情绪状态,并据此调整回复的语气与内容,为用户提供情感慰藉。

(6)多模态集成性。健康咨询对话式人工智能支持文本、图像、音视频等多模态、跨模态数据分析,精准识别用户的就医需求和健康问题,满足用户多元化的医疗健康信息需求。

(7)便捷可及性。用户可以通过移动通讯、电脑等设备,随时随地获取健康咨询服务。尤其对于偏远地区或行动不便的用户,健康咨询对话式人工智能可以成为获取医疗信息与服务的重要渠道。

### 3.2.3 健康咨询对话式人工智能的应用场景

健康咨询对话式人工智能面向公众提供各类健康咨询服务,满足用户个性化的健康信息与服务需求,具体应用场景包含但不限于以下内容:

(1)健康教育与知识科普。通过互动问答、图文、视频等多种形式,为用户提供健康知识、疾病预防、健康生活方式等方面的科普信息,提升用户的健康意识和自我管理能力和<sup>[6-7]</sup>,促进其健康素养的提升。

(2)常见症状自诊自查。通过理解用户的症状描述,为用户提供常见病、多发病的疾病自诊与报告解读,帮助用户准确识别和理解自身症状,提供潜在健康问题的分析,帮助用户合理分诊,及时寻求专业医疗帮助,确保其在必要时病情可以得到有效处理。

(3)健康生活方式指导与疾病预防。针对肥胖、吸烟、饮酒等不良生活方式,提供科学指导,鼓励用户制定切实可行的健康目标,并通过持续跟踪和反馈,帮助用户逐步形成健康的生活方式,预防疾病发生,全面提升健康水平。<sup>[8]</sup>

(4)个性化健康管理及康复支持。为亚健康人群、慢性病或康复期患者提供全面的健康管理信息支持,包括体检报告解读、健康风险预测、日常护理建议、康复锻炼计划等,以满足用户的多元化需求,实现个性化健康管理和康复指导。<sup>[9]</sup>

(5)用药安全与咨询。解答用户对药物的使用、剂量、可能的副作用及药物相互作用等方面的问题,提供用药安全指导,帮助用户正确理解医嘱。<sup>[10-11]</sup>

(6)就医导航服务。根据用户的病情和需求,协助其匹配到适合的医疗卫生服务资源,如医生、医院或特定健康服务,提供便捷的预约挂号和互联网医疗服务,提高用户就医连续性、便捷性和准确性。

(7)紧急事件的应对处理。在突发疾病、意外伤害等紧急情况下,为用户提供院前急救处理建议和

指导,使用户可以迅速判断情况的紧急程度。与专业急救平台建立联动机制,必要时启动危机干预流程,确保用户获得及时有效的援助。

### 3.2.4 健康咨询对话式人工智能的评价要素

健康咨询对话式人工智能在实践应用中,需进行持续的测试和评价,控制可能带来的风险,以确保系统安全有效使用。结合国内外文献研究结果,归纳健康咨询对话式人工智能的评价要素,主要包括以下四类。

#### (1)模型技术水平评价

模型的数据来源和算法性能是衡量健康咨询大模型内在质量与能力的基石,直接关系到模型能否准确、高效地分析并完成健康咨询任务。

数据质量评估。健康咨询对话式人工智能的应用需要遵循《生成式人工智能服务管理暂行办法》和《互联网信息服务算法推荐管理规定》<sup>[12]</sup>等文件要求,完成算法备案和大模型备案。在数据采集阶段,应加强数据资源管理,保证训练数据的多样性和准确性,所采集数据应能够客观反映真实世界的医疗健康信息需求。模型训练应使用具有合法来源的数据,确保模型在维护用户隐私和数据安全方面的合规性。

模型能力评估。在健康咨询对话式人工智能的模型训练阶段,应开展持续的模型优化和算法验证,通过模拟测试和实际应用测试,评估模型在健康科普、常见病自查自诊、慢性病患者健康管理、康复期患者康复指导、用药咨询等场景下的适配性和稳定性。适配性是指模型能否准确理解并回应各种医疗健康信息服务需求,包括常见疾病咨询、药物使用建议等。稳定性则是指模型在不同时间、不同环境下能否保持一致的性能表现。

在应用阶段,需关注健康咨询对话式人工智能的功能性、性能、有益性、公平性和鲁棒性等特征。<sup>[13]</sup>功能性是指系统能否满足用户的信息与服务需求,提供准确、有用的信息与服务。性能则是指系统的响应速度、处理效率等,需确保用户在使用过程中获得流畅的服务体验。有益性是指系统提供的健康咨询建议是否对用户有益。公平性是指系统对不同用户群体的服务是否一致,避免出现歧视或偏见。鲁棒性则是指系统在面对异常输入或攻击时能否保持正常运行,确保用户数据的安全和隐私。

系统安全性评价。在模型训练阶段,应加强训练数据的安全性,重点关注训练数据的来源安全、内

容安全、标注安全以及模型安全,以确保健康咨询对话式人工智能在应用中能够展现出高度的准确性和可靠性。<sup>[14]</sup>在部署阶段,需评估健康咨询对话式人工智能的系统安全与抵御外部攻击的能力。这包括评估系统的安全防护措施是否完善,能否有效抵御各种网络攻击和数据泄露风险。同时,应建立监测系统,实时监控系统的运行状态和用户数据的安全情况,及时识别并防范潜在风险。此外,还需接受相关管理机构(如卫生健康部门、数据保护机构等)的评估和监管,确保系统的合规性和安全性。

### (2) 应用有效性评价

健康咨询对话式人工智能在应用过程中,应结合医学专业知识,提供专业的临床解释,由具有医学背景的专业人员参与模型的设计和解释,以确保模型的决策与临床实践相符合。

健康问题识别能力。健康咨询对话式人工智能能够准确理解并提取用户描述的健康问题,在与用户交互过程中能够提供多轮流畅的智能交互。具备多模态信息识别能力,能够准确处理多种类型的输入信息(如文本、语音、图像)。采集信息完整,能够全面收集疾病初步判断所需要的用户信息,包含环境、心理、社会、生理、健康行为等多个方面。健康咨询对话式人工智能的问询内容应与用户症状和疾病密切相关,能够针对需要问询的内容进行提问,避免过度问询或提出不相关、冗余的问题。

健康问题推理能力。健康咨询对话式人工智能能够根据用户的主诉信息给出最可能的健康问题判断结果,能够覆盖常见病病种。推理结果的可解释性较好,能够清晰、明确地解释其决策过程和依据。健康咨询对话式人工智能能够使用医学专业术语进行回复,对生成答案具备解释和推理能力,能给出支撑答案的可信证据。具备处理不完整或噪声数据的能力,如当用户输入信息包含少量错误或不规范数据时,仍然可以生成合理一致的结果。对于少见或隐性的用户问题,如疑难杂症、罕见病等问题,具备一定的敏感性,能够为用户提示潜在的疾病风险。

健康咨询建议生成能力。健康咨询对话式人工智能所提供的健康咨询建议应基于证据,符合我国医学标准和最佳实践,内容全面、无遗漏。能够根据用户的个体特征(如病史、年龄、性别、过敏史等)以及用户群体(如儿童、孕妇、老年人等)特点,提供定制化的健康咨询信息与服务。建议所生成的内容具有可操作性,可以根据用户所提供的经济能力、地理

位置等信息,结合医疗资源可得性等条件给出合理建议。也可为用户提供后续的医疗服务资源和接口,如线下就医建议、用药提醒、康复期患者健康管理计划,以及未受商业影响的医疗资源接口信息(如医疗机构在线预约挂号平台链接)等。

### (3) 用户体验评价

健康咨询对话式人工智能的应用还需关注用户的交互体验,如信息可靠性、交互友好性和服务响应性等。可以通过用户反馈和满意度调查等方式评估用户体验,据此对系统进行改进和优化。

健康咨询信息的可靠性。健康咨询对话式人工智能所提供的健康咨询信息是可靠的、准确的。可以根据用户的年龄、性别、文化程度以及健康素养等特点,提供个性化的建议和解决方案。建立有效的用户反馈渠道,及时收集并反馈用户的意见和建议,不断优化系统的功能和性能。

以人为本的响应性。在用户体验设计中,应秉持以人为本的理念,加强对用户的人文关怀,模拟人性化的关怀与互动,提升用户的体验感和满意度。健康咨询对话式人工智能在与用户交互的过程中,应保持礼貌和尊重,避免使用冒犯性或冷漠的语言。健康咨询对话式人工智能能够耐心倾听用户的问题和需求,提供及时的服务响应和充分的沟通交流,确保用户感受到被重视和关注。

人机交互的友好性。为确保用户轻松获取所需的健康信息与服务,健康咨询对话式人工智能的对话应具备高度的可读性,避免使用复杂的语句或专业术语。在研发和应用中,可以采用更友好的人机交互方式,如语音、图片、视频等,以更好地服务于老年人、儿童、残疾人或文化程度较低的人群。此外,建议健康咨询对话式人工智能关注多语言支持和本土化内容,以避免因方言等问题而引起的歧义,确保生成咨询建议准确无误。

### (4) 伦理与安全评价

在指导开发者、用户和监管机构改进和监督此类技术的设计和使用过程中,需要关注伦理性和合规性要求<sup>[15-17]</sup>,并建立相应的风险防范机制。

恪守伦理准则。在医疗卫生服务领域,应用人工智能需符合人类的普适性价值观和职业伦理道德准则。2023年国家人工智能标准化总体组、全国信标委人工智能分委会发布的《人工智能伦理治理标准化指南》指出,医疗领域的人工智能应以人类的健康需求为核心,确保人类用户可以掌控医疗决策过

程和对健康咨询对话式人工智能系统的控制,自行选择采纳人工智能建议的方式和程度。2024 年世界卫生组织(WHO)发布的《卫生领域人工智能的伦理与治理:多模态大模型指南》,提出六项伦理原则:①保护人类的自主性;②增进人类福祉、安全和公共利益;③确保透明、可以解释和可以理解;④培养责任感和实行问责制;⑤确保包容性和公平;⑥推广反应迅速且可持续的人工智能。为此,系统的开发方应通过建立相应的机构伦理委员会,确保其在算法设计、训练数据选择、模型生成和优化、服务提供、健康咨询建议生成等全过程均不存在歧视或偏见,保证系统的公平正义性和伦理规范性,体现“以人为本、智能向善、造福人类”的伦理价值准则。

遵守合规性要求。以健康咨询对话式人工智能为代表的人工智能技术的开发和使用,必须严格遵守所在国家及地区的数据安全和医疗安全等政策法规要求。在数据合规使用方面,平台开发者需确保个人隐私及数据安全,实施数据脱密脱敏处理,并在隐私政策中明确告知用户数据的收集、使用、存储、处理和收益政策,约定数据合理使用范畴,保障用户权益。在医疗合规性方面,系统应当明确标注并提醒用户,其提供的内容仅供初步参考,不可直接视为或替代专业的诊断结论。防范出现绝对化诊断、错误诊断、危急重症延误治疗等情况。警惕功效夸大、危害夸大、违背科学事实、伪造健康信息等问题。不允许出现涉及医疗广告、可能引发舆情风险的内容,以及其他任何形式的医疗违法行为。

建立风险防范机制。健康咨询对话式人工智能在对话过程中遇到患者自残自伤甚至自杀倾向等高风险行为时,平台应及时采取预警提醒、患者回访、危机干预等措施。平台应对相关功能及应用设定年龄限制,例如 OpenAI 要求其用户必须年满 18 周岁,或者年满 13 周岁并经过监护人同意才能使用相关服务,旨在保护未成年人的合法权益。

### 3.2.5 健康咨询对话式人工智能的评价方法

为确保健康咨询对话式人工智能生成内容的质量和效果,需要构建科学、全面的质量评价体系及标准,涵盖上述的质量评价要素。目前国内外已经出现一些评估医疗大模型性能的标准化测试,即基准测试(Benchmark),通过使用预定义的数据集、任务和评估指标,对模型在特定任务上的表现进行量化评估,以便比较不同模型之间的性能差异。不同的基准测试需设计一套独立的评估架构与评估体系,

通常包括不同的评估指标,对指标体系的评估多采用智能评估和人工评估两类。

#### (1) 智能评估

已有研究通过应用程序编程接口(API)实现模型生成数据的接入或上传,并基于统计方法进行自动量化评估。常见的定量评估指标,例如精准率(Precision)、召回率(Recall)、F1 分数(F1 Score)可用于测量模型的准确性;BLEU(Bilingual Evaluation Understudy)、ROUGE(Recall-Oriented Understudy for Gisting Evaluation)可用于评估生成文本的质量;可读性(Readability)关注生成文本对人类用户的友好程度。随着人工智能代理(AI Agent)技术的不断发展,人们开始探索更为复杂和动态的智能评估方法。例如,谷歌利用多代理系统构建自我博弈模拟环境,通过患者代理、医生代理和评价代理,共同生成医患模拟对话场景,覆盖多种疾病状况和专科领域。通过内置的自动反馈机制不断优化对话质量和诊断能力,实现病史采集和诊断对话的性能提升。<sup>[18]</sup>

#### (2) 人工评估

人工评估的评估者包括医学专家、临床医生、研究人员以及用户等,通常利用李克特量表对模型的回答进行量化。评估者就临床诊断的准确性和全面性、同理心等难以客观计算的指标,根据模型表现给予相应的评分。使用李克特量表或类李克特量表将复杂的评估标准转化为可量化的分数。<sup>[18]</sup>然而,人工评估的过程中可能缺乏统一的量化标准,且各个研究所用的李克特量表差异较大,这可能导致评估结果的解释复杂化,从而会对临床决策产生影响。因此,在进行人工评估时,确保专家评分的标准化和系统性尤为重要。

## 4 建议

目前,我国健康咨询对话式人工智能的应用与评估尚处于起步阶段,随着技术的进步,人工智能的应用场景将进一步拓展,模型性能将进一步优化,评价体系和方法也将不断丰富与深化。本实践指南结合已有研究提出以下建议,并坚持与时俱进的原则,定期进行证据更新和完善,为相关领域的理论研究和实践探索提供有价值的参考框架。

### 4.1 加快推动健康咨询对话式人工智能评价指标体系的实践应用

国内学者已针对医疗健康大模型应用和评估展

开初步探索,并取得一些成果,但未来在评价指标体系的细化,评价方法的更新迭代等方面还需进一步完善,例如,2024年互联网医疗健康产业联盟发布了《医疗健康行业大模型成熟度评估模型 第1部分:健康咨询》标准,未来仍需要进一步推动健康咨询对话式人工智能相关应用及产品评价体系的深化和标准的落地执行。在此过程中,应根据政策引领方向,充分考虑应用场景、潜在风险、影响效应等因素,开展分级分类评估,为健康咨询对话式人工智能健康可持续发展奠定基础。

#### 4.2 拓宽健康咨询对话式人工智能的应用领域及场景模式

当前,健康咨询对话式人工智能系统多被集成于在线咨询平台或移动应用程序中,未来可融入医疗机构诊疗服务流程,以医疗机构为主体提供健康咨询服务,提升信息咨询与服务的连续性。在诊前阶段,成为用户了解自身健康状况、疾病自诊自查、智能分诊导诊的有效工具;在诊后阶段,通过提供健康咨询、用药指导、个人健康档案管理和智能随访等服务,为用户提供全方位的医疗健康信息与服务支持,持续提升用户作为健康第一责任人的理念,改善诊疗效果。此外,健康咨询对话式人工智能还可以作为智慧家居生态系统的一部分,通过提供健康管理和咨询服务,在居家养老、智能护理等方面发挥其应用价值。

#### 4.3 强化健康咨询对话式人工智能的伦理和安全治理规范

健康咨询对话式人工智能在健康咨询和健康管理中的应用可以有效提升用户体验、提高服务效能,但随之而来的医疗安全及伦理风险问题对现行治理体系提出巨大挑战。虽然,我国针对人工智能领域的立法监管进行了初步尝试和实践探索,但治理过程仍处于初级阶段。例如,由于无法对数据来源进行合法合规性审查,可能导致个人信息与商业秘密的泄露。大模型幻觉、生成虚假信息等问题<sup>[19]</sup>,也引起学术界与业界的广泛担忧与思考。未来,实现健康咨询对话式人工智能与人类的价值对齐,提升数据可靠性和模型可解释性,仍需要进一步探索。

#### 4.4 探索开展应用效果的卫生经济学评价

健康咨询对话式人工智能是患者获取医疗健康信息与服务的新渠道。随着相关应用的持续深化,健康咨询对话式人工智能将对患者的疾病认知与健

康行为产生影响,并有望在重塑现有医患沟通模式、推动医疗卫生服务资源优化配置等方面发挥潜力。为此,未来应探索性开展针对应用效果及影响效应的评价研究,为推动健康咨询对话式人工智能在医疗卫生服务领域的广泛应用与健康发展提供证据支持。

### 致谢

感谢专家组的无私支持,他们的专业贡献为本实践指南的形成和完善提供了重要指导(按姓氏笔画排序)。

王丹(抖音集团旗下小荷健康医学中心)、孔愨(北京市丰台区方庄社区卫生服务中心)、孙鑫(国家癌症中心/国家肿瘤临床医学研究中心/中国医学科学院北京协和医学院肿瘤医院)、李筱永(首都医科大学医学人文学院)、李萌(抖音集团旗下小荷健康医学中心)、王慧超(抖音集团旗下小荷健康医学中心)、杨学来(中日友好医院)、吴浩(首都医科大学全科医学与继续教育学院)、闵栋(中国信息通信研究院)、邱月(清华大学医疗管理学院)、邱英鹏(国家卫生健康委卫生发展研究中心)、武雅文(中国信息通信研究院)、郭珉江(中国医学科学院医学信息研究所)、郭蕊(首都医科大学公共卫生学院)、龚楠(北京市百瑞律师事务所)、戚森杰(首都医科大学附属北京中医医院)、彭迎春(首都医科大学医学人文学院)、鲍薇(中国电子技术标准化研究院)

同时也感谢秘书组的辛勤工作,为本实践指南的顺利完成提供保障。

刘昊鹏(首都医科大学公共卫生学院)、廖委真(首都医科大学公共卫生学院)、王依(首都医科大学公共卫生学院)

**作者贡献:**马骋宇、韩优莉负责文章的构思、设计与起草;吴海锋、顾雪非负责论文的修改;张毓辉负责论文的质量控制。所有作者同等贡献。

**作者声明本文无实际或潜在的利益冲突。**

### 参 考 文 献

- [1] YANG R, TAN T F, LU W, et al. Large language models in health care: Development, applications, and challenges [J]. Health care science, 2023, 2(4): 255-263.
- [2] PARK Y J, PILLAI A, DENG J, et al. Assessing the research landscape and clinical utility of large language models: a scoping review [J]. BMC medical informatics

- and decision making, 2024, 24(1): 72.
- [3] NAZI Z A, PENG W. Large Language Models in Healthcare and Medical Domain: A Review [J]. *Informatics*, 2024, 11(3): 57.
- [4] 国家网信办, 国家发展改革委, 教育部, 等. 生成式人工智能服务管理暂行办法 [EB/OL]. (2023-07-10) [2024-11-28]. [https://www.gov.cn/zhengce/zhengceku/202307/content\\_6891752.htm](https://www.gov.cn/zhengce/zhengceku/202307/content_6891752.htm)
- [5] HSU C C, SANDFORD B A. The Delphi technique: making sense of consensus[J]. *Practical Assessment Research & Evaluation*, 2007, 12(10): 1-8.
- [6] NOH E, WON J, JO S, et al. Conversational agents for body weight management: systematic review [J]. *J Med Internet Res*, 2023, 25: e42238.
- [7] ANISHA S A, SEN A, BAIN C. Evaluating the potential and pitfalls of ai-powered conversational agents as humanlike virtual health carers in the remote management of noncommunicable diseases: scoping review [J]. *J Med Internet Res*, 2024, 26: e56114.
- [8] SINGH B, OLDS T, BRINSLEY J, et al. Systematic review and meta-analysis of the effectiveness of chatbots on lifestyle behaviours[J]. *Npj Digit Med*, 2023, 6(1): 118.
- [9] LI Y, LIANG S, ZHU B, et al. Feasibility and effectiveness of artificial intelligence-driven conversational agents in healthcare interventions: a systematic review of randomized controlled trials[J]. *Int J Nurs Stud*, 2023, 143: 104494.
- [10] REIS Z S N, PAGANO A S, RAMOS DE OLIVEIRA I J, et al. Evaluating Large Language Model: Supported Instructions for Medication Use: First Steps Toward a Comprehensive Model [J]. *Mayo Clinic Proceedings: Digital Health*, 2024, 2(4): 632-644.
- [11] NAYAK A, VAKILI S, NAYAK K, et al. Use of voice-based conversational artificial intelligence for basal insulin prescription management among patients with type 2 diabetes: a randomized clinical trial [J]. *JAMA Netw Open*, 2023, 6(12): e2340232.
- [12] 国务院. 互联网信息服务算法推荐管理规定[EB/OL]. (2022-03-01) [2024-11-28]. [https://www.gov.cn/zhengce/2022-11/26/content\\_5728941.htm](https://www.gov.cn/zhengce/2022-11/26/content_5728941.htm)
- [13] ABBASIAN M, KHATIBI E, AZIMI I, et al. Foundation metrics for evaluating effectiveness of healthcare conversations powered by generative ai[J]. *Npj Digit Med*, 2024, 7(1): 82.
- [14] 许志伟, 李海龙, 李博, 等. AIGC 大模型测评综述: 使能技术、安全隐患和应对[J]. *计算机科学与探索*, 2024, 18(9): 2293-2325.
- [15] 夏光辉, 曹艳林, 陈炳澍, 等. 大模型人工智能技术在医疗服务领域应用的专家共识[J]. *中国卫生法制*, 2023, 31(5): 124-126.
- [16] FOURNIER-TOMBS E, MCHARDY J. A medical ethics framework for conversational artificial intelligence [J]. *J Med Internet Res*, 2023, 25: e43068.
- [17] SEDLAKOVA J, TRACHSEL M. Conversational artificial intelligence in psychotherapy: a new therapeutic tool or agent? [J]. *Am J Bioeth*, 2023, 23(5): 4-13.
- [18] TU T, PALEPU A, SCHAEKERMANN M, et al. Towards conversational diagnostic AI [EB/OL]. (2024-01-11) [2024-11-28]. <https://arxiv.org/abs/2401.0565>
- [19] CLUSMANN J, KOLBINGER F R, MUTI H S, et al. The future landscape of large language models in medicine [J]. *Communications Medicine*, 2023, 3(1): 141.

[收稿日期:2025-02-20 修回日期:2025-03-28]

(编辑 赵晓娟)