

自动语言分析在精神分裂症和精神病临床 高危综合征中的应用

张丹 魏燕燕 王继军

200030 上海交通大学医学院附属精神卫生中心

通信作者: 王继军, Email: jijunwang27@163.com

DOI: 10.3969/j.issn.1009-6574.2022.01.001

【摘要】 语言反映个体的思维活动,因此语言障碍反映个体的思维障碍。自动语言分析是一种基于自然语言处理和机器学习的计算方法,主要用来处理和理解个体语言。常见分析指标有语义一致性和句法复杂性等,目前主要应用于识别精神分裂症和精神病临床高危综合征(CHR-P)以及预测CHR-P人群转化。研究表明,自动语言分析的优势在于敏感、准确和客观,且优于传统的量表评分。现就常见的自动语言分析指标及其应用进行综述。

【关键词】 精神分裂症; 精神病临床高危综合征; 自动语言分析; 自然语言处理; 综述

基金项目: 上海市自然科学基金(19ZR1445100); 国家自然科学基金(82001406); 上海市精神卫生中心基金(2020-FX-02)

Application of automated analysis of language in schizophrenia and clinical high-risk for psychosis

Zhang Dan, Wei Yanyan, Wang Jijun

Shanghai Mental Health Center, Shanghai Jiao Tong University School of Medicine, Shanghai 200030, China

Corresponding author: Wang Jijun, Email: jijunwang27@163.com

【Abstract】 Language reflects individual thinking, and language disturbance reflects thought disorder. Automated analysis of language (AAL) is a computational method based on natural language processing and machine learning, which is mainly used to process and understand individual language content. Common indicators include semantic coherence and syntactic complexity. AAL is mainly used to identify schizophrenia and clinical high-risk for psychosis (CHR-P), and predict the conversion of CHR-P. Research indicates that AAL is sensitive, accurate and objective, and better than clinical rating. This article reviews the common AAL indicators and related applications.

【Key words】 Schizophrenia; Clinical high-risk for psychosis; Automated analysis of language; Natural language processing; Review

Fund programs: Natural Science Foundation of Shanghai (19ZR1445100); National Natural Science Foundation of China (82001406); The Foundation of Shanghai Mental Health Center (2020-FX-02)

精神分裂症(schizophrenia)发病率高,具有高自杀率^[1]和高致残率^[2-3]。近年来,对于精神分裂症的临床研究热点前移至精神病临床高危综合征(clinical high-risk for psychosis, CHR-P)。CHR-P人群的临床表现为轻微的精神病性症状和社会功能损伤^[4],但是尚不符合精神分裂症标准。基于临床评估和随访,CHR-P人群在2年内转化为精神病的概率在30%左右^[5],这很难满足临床早期干预需求,亟待开发可以预测CHR-P临床转化的新技术。

语言反映精神活动结构和内容,语言分析常被

用来研究患者的思维障碍。较早的语言分析常采用手动语言分析(manual linguistic analyses),主要指临床访谈结合量表评分。研究者常用思维、语言和交流评估量表(scale for the assessment of thought, language, and communication, TLC)^[6]和儿童形式思维障碍量表(kiddie-formal thought disorder scale, K-FTDS)^[7]评估患者的思维障碍。综合前人研究发现,手动语言分析往往比较费时费力,主观性较强且准确度低^[8]。

随着人工智能技术的发展,自动语言分析(automated

analysis of language, AAL)在精神疾病领域的应用也愈加广泛,研究证明其效果可能优于手动语言分析^[9]。AAL是一种基于人工智能和自然语言处理的计算方法,主要用于处理和理解个体的语言内容^[8],其优势在于可敏感、客观和快速地提取语言特征,尤其是细微的语言障碍^[8-14]。这不仅有助于研究精神分裂症患者和CHR-P人群的思维障碍,也有助于预测CHR-P人群转化。现就语言采集方法、常见的指标及其应用进行介绍。

一、语言采集方法

目前语言采集方法分为三种,首先是自由访谈,要求被试在自然和放松的状态下谈论自己当前想到的事情,其目的是让被试尽可能多加谈论,进而获取丰富的语言信息^[5];其次是结构式访谈,要求被试完成语言任务,常见任务有故事游戏(story game)访谈^[7],要求被试复述所听到的故事、回答相应问题以及讲述一个新故事;最后是书面语言采集,要求被试完成一段情景式叙述描写^[12]。综合三种方法,自由访谈和结构式访谈应用较多,且针对CHR-P人群,往往需要更长的访谈时间,这是为了获取更多的语言信息进而探索细微的语言障碍。

二、AAL常用指标、分析技术及应用

AAL常见分析指标涉及语义、句法,最新研究也有涉及隐喻性和情感,具体见表1。接下来逐一介绍指标的意义、分析技术和在精神病领域的应用。

1. 语义一致性(semantic coherence): 语义一致性指个体语言信息有序衔接的程度^[18]。潜在语义分析(latent semantic analysis, LSA)是目前最常用的测量方法,优点在于灵活、客观和有效。在词汇习得理论的启发下,LSA的应用原理为词义是每个词与其他词之间关系的函数^[19]。计算机通过扫描相当大的语料库学习该词的含义,当两个词语同时出现的频率越高,则相似性越高(例如猫/狗和猫/铅笔)。LSA将每个单词映射到一个降维的向量空间,给每个单词分配一个相关联的单词向量,相邻单词向量之间的余弦值可以用来评估单词之间的相似性。将单词向量相加得到短语向量,通过测量相邻短语向量之间的余弦值可以测量语篇层面的语义一致性^[9]。低语义一致性提示精神分裂症的阳性思维障碍^[8]。Elvevåg等^[15]首次利用LSA分析语义一致性,并采用TLC量表评估患者的思维障碍,结果发现,患者的语义一致性低于健康对照,且语义一致性与思维障碍得分存在相关关系。另外,判别分

析结果显示,语义一致性区分精神分裂症患者和健康对照的准确率为82.4%。基于精神分裂症的遗传性,Elvevåg等^[20]在随后研究中要求患者分别和一级亲属以及陌生的健康被试自由交谈,发现语义一致性区分患者和健康对照(一级亲属和陌生的健康被试)的准确率达到86.8%。另外,不同年龄段的精神分裂症患者的低语义一致性现象可能是稳健的,因为Holshausen等^[21]发现老年精神分裂症患者的语义一致性也相对较低,且得分与适应性功能评分相关。最后,采用语义一致性预测CHR-P人群转化也是当前的研究热点^[9,14],研究者也在不同风险队列之间进行了交叉验证,表明语义一致性可能是预测CHR-P转化最有效的指标之一^[8]。

2. 语言连接性(language connectedness): 语言连接性指在不考虑内容和句法的情况下,单词之间的语序接近性^[8]。语音图解分析(speech graph analysis)^[11,13,17]常用来测量语言连接性,该方法将每个单词视为节点,单词之间的时间序列表示为有向边,即用语言创建图表量化不同个体的语言连接性差异^[13]。低语言连接性为精神分裂症的阴性思维障碍^[8]。首先,语言连接性可被应用于区分精神分裂症患者和健康对照,一项研究显示区分准确率达91.67%^[13]。其也可用于区分精神分裂症和躁狂症(敏感度为93.8%,特异度为93.7%),且区分效果优于传统的临床量表(敏感度为62.5%,特异度为62.5%)^[17]。其次,语言连接性与临床症状相关,包括阴性症状^[13],认知表现、思维障碍和使用功能磁共振成像在静止状态下测量的大脑连接障碍^[22]。另外,研究发现CHR-P人群的语言连接性介于健康对照和精神分裂症患者之间,语言连接性得分与CHR-P人群的临床转化相关^[23]。最后,研究发现针对健康被试,教育水平能较年龄更好地解释语言连接性和个体发展之间的关系,而对于精神分裂症患者,未发现教育水平有类似的作用^[11],这提示了精神分裂症患者语言发展轨迹可能出现了早期偏差。

3. 内容贫乏(poverty of content): 内容贫乏为典型的精神分裂症阴性症状,常用测量指标有句法复杂性(syntactic complexity)和语义密度(semantic density)。句法复杂性指语言的具体性^[9],常采用词性标注分析(part-of-speech tagging analyses, POS-Tag)技术进行测定。该技术首先根据语法功能对每个单词进行语法标记^[24],最后根据语法标记统计每个语法功能的使用频率。常见的统计指标有“that”

和“which”等补语词的频率以及具有语法规则的短句长度等,研究显示这些指标都可预测精神分裂症发作^[14]。语义密度是指一个句子中能够表达语义的成分(意义成分)数量,常使用 POS-Tag 技术和向量解包(vector unpacking)技术进行测定。向量解包技术原理在于将一个句子中的意义成分数量除以句中实词的数量,进而得到语义密度(范围0~1),且研究发现该指标有助于预测 CHR-P 人群转化^[10]。

4. 指称衔接(referential cohesion): 指称衔接是指不清楚或模棱两可的指称词的数量,是一种连接短语或者句子的语言特征^[12]。指称词包括代词、指示词和比较词,指称词的准确应用能表示上下文存在联系。研究者常采用 Coh-Metrix 工具(一种基于网络的计算语言分析工具)评估指称衔接。精神分裂症患者^[25]和 CHR-P 人群^[12]均被报道存在较低的指称衔接,且指称衔接也有助于预测精神分裂症预后效果和角色功能损害程度^[26]。

5. 其他指标: 近年来,有3个新指标被证实有助于预测 CHR-P 人群转化。首先是潜在内容(latent content),常用潜在内容分析(latent content analysis)技术进行测定,该技术将被试的语言样本与大型语料库进行比较识别语言表达中的潜在内容。其次是隐喻意义(metaphoricity),常采用隐喻识别算法(metaphor identification algorithm)进行测定。根据前人报道,精神分裂症患者可能会频繁使用一些具有隐喻意义的词语,如“手表”被称为“时间容器”^[6]。隐喻识别算法可以在大型隐喻语料库中学习词语的隐喻意义,进而将语言样本中的单词标记为字面意义或者隐喻意义^[16]。最后是情感分数(sentiment scores),采用自动情感分析(automated sentiment analysis)^[16]进行评估,该方法在词汇和短语层面进行情感评分,评分范围为1(非常消极)~5(非常积极)。

综上所述,随着 AAL 技术的发展,量化语言特征的指标种类愈加丰富,进而为研究精神疾病患者的思维障碍夯实了重要基础。

三、综合多种语言指标预测 CHR-P 人群转化的研究进展

精神分裂症发作前已经存在一段可以临床识别的时期,即前驱期(psychosis prodromal)^[27],此时为干预的最佳时机。处于前驱期的 CHR-P 人群在思想、知觉和交流方面表现出亚临床症状。预测 CHR 人群的临床结局对于早期干预意义重大,并且有助于早期预防 CHR-P 人群功能恶化并减少疾病迁延的

风险,因此相应预测技术和模型也成为该领域的研究热点。当前研究挑战在于如何在症状模糊和微妙的情况下来探测精神疾病的迹象。研究表明,综合多种语言指标预测精神分裂症发作的效果可能最佳^[8]。

结合语义特征和语法特征确定机器学习分类器是当前研究的趋势之一。Bedi 等^[14]采用开放式的基线访谈,结合 LSA、POS-Tag 和机器学习方法,确定了一个凸包(convex hull)分类器。分类器包括3个指标,分别为“相邻短语的语义一致性最小值”“短语长度”和“限定词使用的频率(如 which 和 that)”,结果表明,这3个指标与临床症状得分相关,而且分类器预测 CHR-P 人群转化的准确率达到 100%(34 名 CHR-P 个体,2.5 年后 5 人转化)。在此基础上,Corcoran 等^[9]修改了语言采集方法,即要求被试完成“故事游戏”访谈,创建逻辑回归分类器,确定了3个指标,包括“语义一致性最小值”“语义一致性方差”和“减少使用所有格代词的频率”。该分类器预测 CHR-P 人群转化的准确率为 83%(59 名 CHR-P 个体,2 年后 19 人转化),受试者工作特征(ROC)曲线下面积达到 0.87,交叉验证(跨队列;34 名 CHR-P 个体,2.5 年后 5 人转化)的准确率为 79%。

其他分类器也表现出较强的预测能力,“低语义密度”和“声音相关词汇(如语气词)的使用频率(患者隐秘地谈论声音相关的词语,可用来表征幻听的早期迹象^[10])”两个指标的分类器预测转化的准确率为 93%(30 名 CHR-P 个体,2 年后 7 人转化)。此外,Gutierrez 等^[16]通过分析 Bedi 等^[14]的样本数据得出了隐喻意义和情感分数特征,采用这两个指标以及性别和年龄创建了一个凸包分类器,结果发现区分首发精神分裂症和健康对照的准确率为 84%,预测 CHR-P 转化的准确率达到 97.1%(准确预测了 34 名 CHR-P 个体中 33 人的临床结局)。

综上所述,在预测转化方面的研究趋势有几点,综合不同指标,目的是涉及不同语言层面(例如同时涉及语义和语法);从基础指标中发展新的指标,如最小值、最大值和方差等;结合机器学习有助于敏感探测出更细微的语言特征。

四、总结与展望

AAL 技术如同“显微镜”般可快速、准确地识别语言的细微特征,并探测出语言和精神分裂症的隐秘联系。这不仅是了解精神分裂症风险及发作的有效途径,也有助于促进精神分裂症防治关口前移。

表1 AAL指标及相关应用

语言指标	AAL方法	应用
语义一致性	潜在语义分析	区分精神分裂症患者和健康被试 ^[15-16] ; 预测CHR-P人群转化 ^[9,16]
语言连接性	语音图解分析	区分精神分裂症患者和健康被试; 预测首发精神分裂症发作; 解释精神分裂症患者发病后6个月内阴性症状的变化 ^[13,17]
句法复杂性	词性标注分析	预测CHR-P人群转化 ^[9,14]
语义密度	向量解包技术	预测CHR-P人群转化 ^[10]
潜在内容	潜在内容分析	预测CHR-P人群转化 ^[10]
指称衔接	Coh-Metrix 工具	比较和区分CHR-P与健康对照差异 ^[12]
隐喻性	隐喻识别算法	区分首发精神分裂症患者与健康对照, 预测CHR-P人群转化 ^[16]
情感分数	自动情感分析	区分首发精神分裂症患者与健康对照, 预测CHR-P人群转化 ^[16]

注: AAL 自动语言分析; CHR-P 精神病临床高危综合征

未来研究可从几个方向考虑。(1)我国在精神病领域采用AAL分析语言障碍的研究还处于初级阶段,研究重点在于分析技术和语言特征的“汉化”,未来研究可聚焦汉语样本的特征分析。(2)将语言指标与已知的风险生物标志物结合分析,如认知指标、遗传学指标和神经影像指标,将有助于预测CHR-P人群多种临床结局,如功能不良、其他精神病发作、缓解和恢复,进一步探索临床转化的生物学机制。(3)综合考虑不同年龄段人群的语言特征,尤其是在儿童和青少年被试中进一步验证AAL的可重复性,确定语言特征的变异性来源和语言特征的变化轨迹。

利益冲突 文章所有作者共同认可文章无相关利益冲突

作者贡献声明 文章构思、文献收集、撰写及修订为张丹,文献收集、论文修订为魏燕燕,论文指导及审校为王继军

参 考 文 献

- [1] Sher L, Kahn RS. Suicide in schizophrenia: an educational overview [J]. *Medicina (Kaunas)*, 2019, 55(7): 361. DOI: 10.3390/medicina55070361.
- [2] Fakorede OO, Ogunwale A, Akinhanmi AO. Disability among patients with schizophrenia: a hospital-based study [J]. *Int J Soc Psychiatry*, 2020, 66(2): 179-187. DOI: 10.1177/0020764019894608.
- [3] Huang Y, Wang Y, Wang H, et al. Prevalence of mental disorders in China: a cross-sectional epidemiological study [J]. *Lancet Psychiatry*, 2019, 6(3): 211-224. DOI: 10.1016/s2215-0366(18)30511-x.
- [4] Cannon TD. How schizophrenia develops: cognitive and brain mechanisms underlying onset of psychosis [J]. *Trends Cogn Sci*, 2015, 19(12): 744-756. DOI: 10.1016/j.tics.2015.09.009.
- [5] Zhang TH, Li HJ, Woodberry KA, et al. Two-year follow-up of a Chinese sample at clinical high risk for psychosis: timeline of symptoms, help-seeking and conversion [J]. *Epidemiol Psychiatr Sci*, 2017, 26(3): 287-298. DOI: 10.1017/s2045796016000184.
- [6] Andreasen NC. Scale for the assessment of thought, language, and communication (TLC) [J]. *Schizophr Bull*, 1986, 12(3): 473-482. DOI: 10.1093/schbul/12.3.473.
- [7] Caplan R, Guthrie D, Fish B, et al. The Kiddie Formal Thought Disorder Rating Scale: clinical assessment, reliability, and validity [J]. *J Am Acad Child Adolesc Psychiatry*, 1989, 28(3): 408-416. DOI: 10.1097/00004583-198905000-00018.
- [8] Corcoran CM, Mittal VA, Bearden CE, et al. Language as a biomarker for psychosis: a natural language processing approach [J]. *Schizophr Res*, 2020, 226: 158-166. DOI: 10.1016/j.schres.2020.04.032.
- [9] Corcoran CM, Carrillo F, Fernández-Slezak D, et al. Prediction of psychosis across protocols and risk cohorts using automated language analysis [J]. *World Psychiatry*, 2018, 17(1): 67-75. DOI: 10.1002/wps.20491.
- [10] Rezaii N, Walker E, Wolff P. A machine learning approach to predicting psychosis using semantic density and latent content analysis [J]. *NPJ Schizophr*, 2019, 5(1): 1-12. DOI: 10.1038/s41537-019-0077-9.
- [11] Mota NB, Sigman M, Cecchi G, et al. The maturation of speech structure in psychosis is resistant to formal education [J]. *NPJ Schizophr*, 2018, 4(1): 1-10. DOI: 10.1038/s41537-018-0067-3.
- [12] Gupta T, Hespos SJ, Horton WS, et al. Automated analysis of written narratives reveals abnormalities in referential cohesion in youth at ultra high risk for psychosis [J]. *Schizophr Res*, 2018, 192: 82-88. DOI: 10.1016/j.schres.2017.04.025.
- [13] Mota NB, Copelli M, Ribeiro S. Thought disorder measured as random speech structure classifies negative symptoms and schizophrenia diagnosis 6 months in advance [J]. *NPJ Schizophr*, 2017, 3(1): 1-10. DOI: 10.1038/s41537-017-0019-3.
- [14] Bedi G, Carrillo F, Cecchi GA, et al. Automated analysis of free speech predicts psychosis onset in high-risk youths [J]. *NPJ Schizophr*, 2015, 1(1): 1-7. DOI: 10.1038/npjschz.2015.30.
- [15] Elvevåg B, Foltz PW, Weinberger DR, et al. Quantifying incoherence in speech: an automated methodology and novel application to schizophrenia [J]. *Schizophr Res*, 2007, 93(1/3): 304-316. DOI: 10.1016/j.schres.2007.03.001.
- [16] Gutierrez ED, Cecchi GA, Corcoran C, et al. Using automated metaphor identification to aid in detection and prediction of first-episode schizophrenia [C]//Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, 2017: 2923-2930.
- [17] Mota NB, Vasconcelos NA, Lemos N, et al. Speech graphs provide a quantitative measure of thought disorder in psychosis [J]. *PLoS One*, 2012, 7(4): e34928. DOI: 10.1371/journal.pone.0034928.

- [18] Hansson K, Bååth R, Löhndorf S, et al. Quantifying semantic linguistic maturity in children[J]. J Psycholinguist Res, 2016, 45(5): 1183-1199. DOI: 10.1007/s10936-015-9398-7.
- [19] Landauer TK, Dumais ST. A solution to Plato's problem; the latent semantic analysis theory of acquisition, induction, and representation of knowledge[J]. Psychol Rev, 1997, 104(2): 211-240. DOI: 10.1037/0033-295x.104.2.211.
- [20] Elvevåg B, Foltz PW, Rosenstein M, et al. An automated method to analyze language use in patients with schizophrenia and their first-degree relatives[J]. J Neurolinguistics, 2010, 23(3): 270-284. DOI: 10.1016/j.jneuroling.2009.05.002.
- [21] Holshausen K, Harvey PD, Elvevåg B, et al. Latent semantic variables are associated with formal thought disorder and adaptive behavior in older inpatients with schizophrenia[J]. Cortex, 2014, 55: 88-96. DOI: 10.1016/j.cortex.2013.02.006.
- [22] Palaniyappan L, Mota NB, Oowise S, et al. Speech structure links the neural and socio-behavioural correlates of psychotic disorders[J]. Prog Neuropsychopharmacol Biol Psychiatry, 2019, 88: 112-120. DOI: 10.1016/j.pnpbp.2018.07.007.
- [23] Spencer TJ, Thompson B, Oliver D, et al. Lower speech connectedness linked to incidence of psychosis in people at clinical high risk[J]. Schizophr Res, 2021, 228: 493-501. DOI: 10.1016/j.schres.2020.09.002.
- [24] Santorini B. Part-of-speech tagging guidelines for the Penn Treebank Project. (3rd revision, 2nd printing) [J]. Beatrice Santorini; Beatrice Santorini, University of Pennsylvania, 1990.
- [25] Ditman T, Goff D, Kuperberg GR. Slow and steady: sustained effects of lexico-semantic associations can mediate referential impairments in schizophrenia[J]. Cogn Affect Behav Neurosci, 2011, 11(2): 245-258. DOI: 10.3758/s13415-011-0020-7.
- [26] Bearden CE, Wu KN, Caplan R, et al. Thought disorder and communication deviance as predictors of outcome in youth at clinical high risk for psychosis[J]. J Am Acad Child Adolesc Psychiatry, 2011, 50(7): 669-680. DOI: 10.1016/j.jaac.2011.03.021.
- [27] Nelson B, Yuen HP, Wood SJ, et al. Long-term follow-up of a group at ultra high risk ("prodromal") for psychosis: the PACE 400 study[J]. JAMA Psychiatry, 2013, 70(8): 793-802. DOI: 10.1001/jamapsychiatry.2013.1270.

(收稿日期: 2021-08-06)

(本文编辑: 赵金鑫)

· 消息 ·

欢迎订阅2022年《神经疾病与精神卫生》杂志

《神经疾病与精神卫生》杂志是神经、精神科学及精神卫生领域的学术性期刊,国内外公开发行人,2006年被中国科学技术信息研究所收录为中国科技论文统计源期刊(中国科技核心期刊)。本刊坚持党的出版方针和卫生工作方针,遵循学科发展规律,以提高杂志质量、扩大社会效益为使命,及时反映科学研究的重大进展,更好地促进国内外学术交流。主要读者对象为广大神经科学、精神科学及精神卫生领域中从事基础、临床医学、教学、科研的工作者及学生。报道内容包括相关各学科领先的教学、科研成果及临床诊疗经验。主要栏目有专家论坛(述评)、论著、学术交流、短篇报道、综述、病例报告、会议纪要、国内外学术动态等。

《神经疾病与精神卫生》杂志国内邮发代号为82-353,由北京市邮政局发行;国外发行代号BM1690,由中国国际图书贸易总公司发行。每期定价15.00元,全年180.00元。欢迎直接通过本社订阅。

银行汇款: 开户行: 中国建设银行建华支行 户名: 《神经疾病与精神卫生》杂志社

账号: 23001626251050500949

联系电话: (010)83191160 传真: (010)83191161