

累积和控制图法在传染病暴发探测中的应用

张洪龙 赖圣杰 李中杰 兰亚佳 杨维中

【导读】 近年来,为提高传染病暴发早期发现能力,许多学者开展了基于各种传染病监测数据源的暴发早期探测方法研究,期望通过采用适当的数学算法,对疾病监测数据进行分析,从而早期发现疾病与症状发生的异常增加或聚集性。本文介绍了一种国外广泛应用于传染病监测数据异常探测的统计过程控制方法——累积和(CUSUM)控制图法,对其基本原理与特征、方法设计要点,以及在传染病暴发探测中应用的典型案例进行了深入分析,从而为在我国传染病预警领域推广应用该方法提供参考与借鉴。

【关键词】 累积和; 传染病; 探测

Application of cumulative sum control chart algorithm in the detection of infectious disease outbreaks ZHANG Hong-long¹, LAI Sheng-jie¹, LI Zhong-jie¹, LAN Ya-jia², YANG Wei-zhong¹. 1 Chinese Center for Disease Control and Prevention, Beijing 102206, China; 2 West China School of Public Health, Sichuan University

Corresponding author: YANG Wei-zhong, Email: yangwz@chinacdc.cn

This work was supported by grants from the National Science and Technology Key Project on Building the Platform of Infectious Disease Surveillance (No. 2009ZX10004-201), the National Science and Technology Support Projects for the "Eleventh Five-Year Plan" of China (No. 2008BAI56B00) and the Cooperation Project Between Chinese Ministry of Health and World Health Organization from 2010 to 2011 on the Research of Outbreak Detection Temporal and Spatial Model Development (No. WPCNH1002405).

【Introduction】 In recent years, for improving the ability of early detection on infectious disease outbreak, many researchers study the disease outbreak detection algorithms, based on many disease surveillance data, expecting to detect the abnormal increasing and cluster of disease and symptom at an early stage by adopting appropriate algorithm. This paper introduces a cumulative sum control chart method, one of statistical process control algorithms widely used in foreign countries and describes its basic principle and characteristic, key points of design, typical examples in application of disease outbreak detection of cumulative sum method, with expect to provide reference for its application in studies of disease outbreak early warning in China.

【Key words】 Cumulative sum; Infectious disease; Detection

早期发现传染病暴发的迹象,并及时采取应对措施,可有效防止其进一步扩散,减少疾病对社会和经济造成的损失^[1]。近年来全球广泛开展了基于各种传染病监测数据源的暴发早期探测方法研究,期望通过采用适当的数学算法,分析疾病监测数据,从而早期发现疾病与症状发生的异常趋势或聚集性^[2-4]。目前用于疾病发生异常探测的方法较多,根据分析数据项的类型分为时间模型、空间模型和时空模型三大类^[5]。统计过程控制(Statistical Process

Control)作为一种成熟且易于掌握的时间模型方法,在公共卫生领域应用较多。常见的统计过程控制图包括休哈特控制图、移动平均(moving average, MA)控制图、指数加权移动平均(exponentially weighted moving average, EWMA)控制图、累积和(cumulative sum, CUSUM)控制图四类,其中CUSUM在探测疾病监测数据微小变化时具有较好的及时性和灵敏度,越来越多地被采用^[3,6-8]。本研究在查阅大量国外相关文献的基础上,对CUSUM控制图方法在传染病暴发探测中应用的典型案例进行分析,介绍其基本原理、应用类型、方法设计等内容,从而为在我国传染病预警领域推广应用该方法提供参考与借鉴。

基本原理

CUSUM是一种在传统休哈特控制图方法基础上改良的方法。休哈特控制图法是由美国贝尔电话

DOI: 10.3760/cma.j.issn.0254-6450.2010.12.019

基金项目: 国家科技重大专项(2009ZX10004-201); "十一五"国家科技支撑计划(2008BAI56B00); 中国-世界卫生组织2010-2011合作项目(WPCNH1002405)

作者单位: 102206 北京, 中国疾病预防控制中心(张洪龙、赖圣杰、李中杰、杨维中); 四川大学华西公共卫生学院(兰亚佳)

通信作者: 杨维中, Email: yangwz@chinacdc.cn

实验室 Shewhart 博士在 1924 年首先提出,最早用于产品质量控制,其原理是将引起质量波动的原因分为偶然和异常两种因素,并利用控制图中上下控制限($\bar{x} \pm 3s$)和中心线(\bar{x})进行质量控制:正常情况下观察数据在偶然因素影响下围绕中心线波动,当观察数据超出了控制限或在控制限内排列非随机时,即表明质量出现了异常波动^[9]。休哈特控制图法的优点是方法简便,对观察数据突然出现的较大变化较为敏感,而对一些缓慢变化的微小偏移则难以探测到^[9,10]。为进一步提高控制图法的敏感性,剑桥大学 Page^[11]于 1954 年提出了 CUSUM 方法。其原理是通过不断累加计算观察值与期望值的差值,逐渐放大数据出现的波动,从而更加快速、灵敏地探测到休哈特控制图法无法识别的微小异常情况^[12]。该方法适用于分析正态分布、平稳且不自相关的数据^[11]。

CUSUM 方法的计算公式:

$$S_0 = 0$$

$$S_t = \max[0, S_{t-1} + (X_t - r - K)] \quad (1)$$

式中 t 为时间, X_t 为观察值, r 为期望值, K 表示观察值大于期望值的最小偏移量, S_0 是最初累积和值, S_t 是当前期累积和值, S_{t-1} 是上一期累积和值。该方法的两个重要参数是阈值 H 和参考值 K , 其值的选取直接影响方法的探测能力^[13]。在该控制图中,从监测开始到发现异常的过程中,所需平均时间为平均运行长度 (ARL), 是评价控制图效果的一个重要指标。从无异常到发出错误警报所需时间为 ARL_0 , 其值越大越好; 从出现异常到发出正确警报所需时间为 ARL_1 , 其值越小越好。通常先固定 ARL_0 , 然后计算出满足要求的不同 H 和 K 参数组合, 再分别计算这些参数组合所对应的 ARL_1 , ARL_1 值最小时的 H 和 K 为最优参数组合, 此时控制图效果最佳。当 S_t 值超过阈值参数 H 时即表明出现异常情况。该公式仅表示正向变化的探测, 即关注观察值的异常增加。当需要探测负向变化(观察值异常降低)时, 公式中的“max”替换成“min”, 即当 S_t 值低于阈值参数 H 时即表明出现异常情况。在传染病暴发探测应用中, 一般仅探测正向变化的异常情况。

实例分析

CUSUM 方法最初应用于工业领域, 经过几十年的发展, 仍是产品质量控制的主要方法之一。在医学领域中由于探测疾病数据的异常变化与制造业中发现产品质量变化的原理相似, CUSUM 也适用于疾病暴发的早期探测^[6,8,14]。早在 1981 年 CUSUM 就已

在英格兰和威尔士用于流感监测数据的分析^[15]。

在传染病暴发探测领域应用中, CUSUM 方法的主要设计环节有阈值参数 H 和参考值参数 K 的设定^[8]、期望值的计算方式和基线数据范围的确定^[6]等。根据可利用的历史数据时间长短情况, CUSUM 可以选择不同长度的基线数据计算期望值。

1. 基于长期基线数据的 CUSUM 方法: 考虑到疾病的季节效应, 基于长期基线数据的方法利用当前数据和历史同期数据进行比较^[1]。自 1995 年美国疾病预防控制中心 (CDC) 的沙门菌暴发探测算法 (SODA) 就是利用 CUSUM, 以每周沙门菌感染病例合计数为观察值, 并假定观察值服从正态分布, 分析国家沙门菌监测系统的数据^[16], 其公式:

$$S_0 = 0$$

$$S_t = \max[0, S_{t-1} + (X_t - \mu_0) / \sigma_x - K] \quad (2)$$

式中 X_t 为当前观察值, 期望值 μ_0 为过去 5 年同期每周病例数的均数, σ_x 为过去 5 年同期每周病例数的标准差。根据公式 (2) 计算 S_t 值, 当 S_t 小于阈值 H 时, 直接累加计算下一个 S_{t+1} 值; 当 S_t 大于阈值 H 时, 表明监测数据出现异常增加^[14], 需要引起注意, 同时重新设定 S_t 为 0^[6,12]。该研究发现, 当 $K=1, H=0.5$ 时, ARL_1 的值最小, 方法效果最佳。

SODA 最初备选了 3 种期望值的算法: 过去 5 年同期 5 周病例数 ($X_2, X_3, X_4, X_{11}, X_{14}$) 的均数、过去 5 年同期 5 周病例数 ($X_2, X_3, X_4, X_{11}, X_{14}$) 的中位数和过去 5 年同期及其前后摆动 1 周的 15 周 ($X_1, X_2, X_3, \dots, X_{15}$) 病例数的均数, 当前数据和基线数据的比较见表 1。该研究最终发现采用过去 5 年同期 5 周病例数的均数计算期望值, 其错误预警信号较少、特异度最高。

表 1 SODA 备选的当前数据和基线数据比较

年份	上周	当前周	下周
2009		X_0	
2008	X_1	X_2	X_3
2007	X_4	X_5	X_6
2006	X_7	X_8	X_9
2005	X_{10}	X_{11}	X_{12}
2004	X_{13}	X_{14}	X_{15}

注: X 为沙门菌周病例合计数, 假设当前数据为 2009 年某周观察值 X_0 。

此外, 为避免在病例数较少期间, 该方法经过不断累积运算后可能生成过多假的异常信号, SODA 中还依据经验设置了方法开始运算的基准值, 即仅当观察值 > 5 时才进行累积和运算, 从而可进一步减少假信号的数量。

对于基于长期基线数据的 CUSUM 方法, 基线

数据的范围可根据需要进行适当调整。如为了增加基线数据的代表性,也有研究采用过去5年同期及其前后摆动6周(共65周)的病例数的均数作为期望值^[6]。此外,为了排除历史中疾病暴发对基线数据的影响,有研究仅从没有暴发事件发生的时期中选择基线数据^[15],也有文献建议选择中位数代替均数计算期望值^[1]。在传染病监测中,一般长期历史数据的范围最长采用过去5年,如果选择更长时间的基线数据,一方面当地人口学特点可能发生较大变化^[3],另一方面疾病的流行或变化趋势也会影响基线的稳定,5年以上的那部分数据其作用就被抵消了^[17]。

近年来基于长期基线数据的CUSUM方法也被应用于美国各州的公共卫生实验室信息系统,进行多种疾病的病原体监测^[18]。美国CDC利用CUSUM方法曾探测到几起沙门菌肠炎暴发^[14],1995年5月还成功发现了一起斯坦利沙门菌引起的大暴发^[19]。

2. 基于短期基线数据的CUSUM方法:在实际工作中,依赖于长期历史数据的统计方法在某些情况下并不适用,而只能采用基于短期基线数据的方法,例如在很多发展中国家,监测系统建立时间还不长,有些疾病仅有很少(甚至没有)历史数据^[11,14,20];又如在大型公共活动中,大量人群短期内聚集在特定场所^[21],疾病监测需要信息的快速反馈,而常常只有几天的历史数据可以利用^[3]。美国在早期异常报告系统(Early Aberration Reporting System, EARS)中建立了仅需要短期基线数据的一套异常探测方法,该方法命名为C1、C2和C3^[3]。

C1、C2和C3方法均以上述CUSUM公式(2)为基础,以日为单位计算期望值,其中C3灵敏度最高,C2次之,C1最低。考虑到数据的周末效应,C1、C2和C3的基线数据时间长度均为7d^[21-23]。各方法采用了不同时间范围的基线数据计算均数和标准差:C1采用第-1至第-7天作为基线数据范围;C2和C3均采用第-3至第-9天作为基线数据范围,即基线数据和当前数据之间设置了2d的缓冲期,从而减少近期数据可能对基线数据造成的影响^[24,25](图1)。对于C1和C2,设定上一个累积和值 S_{t-1} 为0,当前值大于基线数据的均数加上3倍标准差时即表明有异常数据产生。如果在特殊时期要求灵敏度更高,可以降低阈值选择当前值大于基线数据的均数加2倍标准差^[21]时视为出现异常数据(故也有文献指出,C1和C2方法并非真正意义的累积和方法,而是一种变化的休哈特方法^[22])。C3的累计和值

是过去3d的C2累积和值相加(3d的累积时间是由美国CDC与州及地方卫生部门实践经验而得^[21]),当C3累积和值大于阈值2时,表明数据异常。

C1、C2和C3方法建立后,在美国成功探测出几起具有重要公共卫生意义的传染病暴发,包括西尼罗河病暴发和流感流行季节的开始^[23]。北京市流感监测系统利用该方法,以周为单位分析流感监测数据,成功探测到北京市2007—2008年流感流行季节的开始时间^[26]。

此外,为提高当前数据与基线数据的可比性,美国CDC公共卫生信息网的事件早期探测系统(BioSense)采用了改进的C2方法^[7],将监测数据分为工作日和周末进行分别计算,即如果当日是工作日,则仅以过去的工作日(7d)作为基线数据范围计算期望值和标准差;如果当日是周末,则仅以过去的周末(7d)作为基线数据范围计算期望值和标准差。

基于短期基线数据的探测方法有效地满足了缺乏足够历史数据的分析需要,国外有研究通过选用模拟数据进行测试,结果表明基于短期基线数据的方法具有与基于长期基线数据的方法同样好的灵敏度和特异度^[23]。除了较为常用的将近期7d的数据作为基线数据的CUSUM方法外,为了增加基线数据的稳定性,也有研究将过去14d或28d作为基线数据范围^[25]。

讨 论

相比于其他控制图方法,CUSUM法原理简单、易于实现,且能灵活应用于不同时间范围的基线数据,个别方法还提供了成熟的免费软件,如EARS开发的SAS版和EXCEL版的C1、C2和C3软件^[27]。因此,该方法在传染病暴发的早期探测研究领域具有很好的应用前景。同时,许多研究者也指出,CUSUM法最初应用于工业领域的产品数据分析,而疾病监测数据的类型与特征更为复杂,其往往并非独立的随机变量,也不一定服从正态分布,因此CUSUM法在传染病暴发探测应用中也面临着一些挑战^[28,29]。对于监测数据不能直接应用CUSUM的

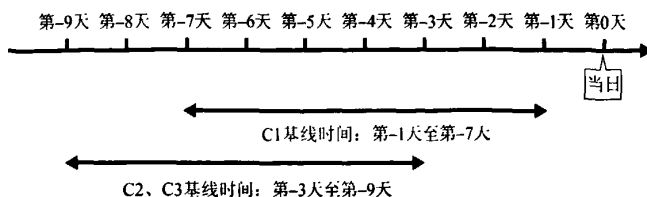


图1 美国EARS的C1、C2和C3方法基线范围示意图

情况,有学者提出可先采用时间序列模型计算一个预测值,然后再对实际观察值与预测值的差值进行累积和运算^[10,22,30]。此外,CUSUM擅于发现数据中小的异常波动,而在现实中疾病监测数据的变化形式未知,故其探测能力也会受到影响^[9]。CUSUM在方法学上面临的挑战,要求使用者应在充分分析监测数据的特征基础上,对数据进行适当的处理和转化,科学设定基线数据,并采用适当的评价标准对方法的各种参数进行优选,从而确保方法达到最优效果。目前,CUSUM法在我国用于传染病暴发探测的实例并不多见,如何针对不同来源的监测数据,合理设计及应用还有待进一步研究和实践。

参 考 文 献

- [1] Stroup DF, Williamson GD, Herndon JL, et al. Detection of aberrations in the occurrence of notifiable diseases surveillance data. *Stat Med*, 1989, 8(3):323-329.
- [2] Heffernan R, Mostashari F, Das D, et al. Syndromic surveillance in public health practice, New York City. *Emerg Infect Dis*, 2004, 10(5):858-864.
- [3] Hutwagner LC, Thompson W, Seeman GM, et al. The bioterrorism preparedness and response Early Aberration Reporting System (EARS). *J Urban Health*, 2003, 80(2 Suppl 1):i89-96.
- [4] Lombardo J, Burkorn H, Elbert E, et al. A systems overview of the Electronic Surveillance System for the Early Notification of Community-Based Epidemics (ESSENCE II). *J Urban Health*, 2003, 80(2 Suppl 1):i32-42.
- [5] Kulldorff M, Heffernan R, Hartman J, et al. A space-time permutation scan statistic for disease outbreak detection. *PLoS Med*, 2005, 2(3):e59.
- [6] O'Brien SJ, Christie P. Do CuSums have a role in routine communicable disease surveillance? *Public Health*, 1997, 111(4):255-258.
- [7] Centers for Disease Control and Prevention. BioSense real-time hospital data user guide application version 2.11, 2007 [2010-07-22]. http://www.cdc.gov/biosense/files/CDC_BioSense_BioSense_Hospital_Data_User_Guide_V2.11.pdf.
- [8] Carpenter TE. Evaluation and extension of the cusum technique with an application to Salmonella surveillance. *J Vet Diagn Invest*, 2002, 14(3):211-218.
- [9] Shmueli G, Fienberg SE. Current and Potential Statistical Methods for Monitoring Multiple Data Streams for Bio-Surveillance// Wilson AG, Wilson GD, Olwell DH. *Statistical Methods in Counterterrorism*. Springer, 2005:109-140.
- [10] Williamson GD, Weatherby HG. A monitoring system for detecting aberrations in public health surveillance reports. *Stat Med*, 1999, 18(23):3283-3298.
- [11] Page SE. Continuous Inspection Schemes. *Biometrika*. 1954(41):100-115.
- [12] Wong W, Moore AW. Classical Time-Series Methods for Biosurveillance//Wagner MM, Moore AW, Aryel RM. *Handbook of Biosurveillance*. London: Elsevier Academic Press, 2006: 217-234.
- [13] Wang X, Zeng D, Seale H, et al. Comparing early outbreak detection algorithms based on their optimized parameter values. *J Biomed Inform*, 2010, 43(1):97-103.
- [14] Hutwagner LC, Maloney EK, Bean NH, et al. Using laboratory-based surveillance data for prevention: an algorithm for detecting Salmonella outbreaks. *Emerg Infect Dis*, 1997, 3(3):395-400.
- [15] Tillet HE, Spencer IL. Influenza surveillance in England and Wales using routine statistics. *J Hyg(Lond)*, 1982, 88(1):83-94.
- [16] Janes GR, Hutwagner LC, Jr WC, et al. *Descriptive Epidemiology: Analyzing and Interpreting Surveillance Data// Teutsch SM, Churchill RE. Principles and practice of public health surveillance*. 2nd ed. Oxford: Oxford University Press, 2000: 112-167.
- [17] Stroup DF, Wharton M, Kafadar K, et al. Evaluation of a method for detecting aberrations in public health surveillance data. *Am J Epidemiol*, 1993, 137(3):373-380.
- [18] Martin SM, Bean NH. Data management issues for emerging diseases and new tools for managing surveillance and laboratory data. *Emerg Infect Dis*, 1995, 1(4):124-128.
- [19] Mahon BE, Ponka A, Hall WN, et al. An international outbreak of Salmonella infections caused by alfalfa sprouts grown from contaminated seeds. *J Infect Dis*, 1997, 175(4):876-882.
- [20] Simonsen L, Clarke MJ, Stroup DF, et al. A method for timely assessment of influenza-associated mortality in the United States. *Epidemiology*, 1997, 8(4):390-395.
- [21] Hutwagner LC, Thompson WW, Seeman GM, et al. A simulation model for assessing aberration detection methods used in public health surveillance for systems with limited baselines. *Stat Med*, 2005, 24(4):543-550.
- [22] Jr Fricker RD, Hegler BL, Dunfee DA. Comparing syndromic surveillance detection methods: EARS' versus a CUSUM-based methodology. *Stat Med*, 2008, 27(17):3407-3429.
- [23] Hutwagner LC, Browne T, Seeman GM, et al. Comparing aberration detection methods with simulated data. *Emerg Infect Dis*, 2005, 11(2):314-316.
- [24] Burkorn HS, Elbert Y, Feldman A, et al. Role of data aggregation in biosurveillance detection strategies with applications from ESSENCE. *MMWR*, 2004, 53 Suppl: S67-73.
- [25] Tokars JI, Burkorn H, Xing J, et al. Enhancing time-series detection algorithms for automated biosurveillance. *Emerg Infect Dis*, 2009, 15(4):533-539.
- [26] Yang P, Duan W, Lv M, et al. Review of an influenza surveillance system, Beijing, People's Republic of China. *Emerg Infect Dis*, 2009, 15(10):1603-1608.
- [27] Centers for Disease Control and Prevention. EARS: Download the latest versions of EARS. [2010-07-22]. <http://www.bt.cdc.gov/surveillance/ears/downloads.asp>.
- [28] Montgomery DC. *Introduction to Statistical Quality Control*. New York: John Wiley & Sons, 1991.
- [29] De M. *Instruction to Statistical Quality Control*. 4th ed. New York: Wiley, 2001.
- [30] Reis BY, Mandl KD. Time series modeling for syndromic surveillance. *BMC Med Inform Decis Mak*, 2003, 3:2.

(收稿日期:2010-07-31)

(本文编辑:张林东)