

# 应用 Empower Stats 软件分析阈值效应

林林 陈常中 余晓丹

**【导读】** 生物医学研究中许多研究因素在一定范围内对结果变量无影响或有正效作用,超过某一阈值后,作用大小或/和方向可发生变化,称为阈值效应。在分析研究因素( $x$ )对结果变量( $y$ )的作用有无阈值效应时,可先通过平滑曲线拟合观察是否有分段线性关系,然后采用分段回归模型、LRT 检验和 Bootstrap 重抽样法进行阈值效应分析。美国 X & Y Solutions 软件公司开发的 Empower Stats 软件,设有阈值效应分析模块,可输入阈值后按所给阈值分段模拟数据,也可以不输入阈值,由软件自动确定最佳阈值模拟数据,并计算阈值置信区间。

**【关键词】** 阈值效应;分段式回归模型;LRT 检验;Bootstrap 重抽样

**The analysis of threshold effect using Empower Stats software** LIN Lin<sup>1</sup>, CHEN Chang-zhong<sup>2</sup>, YU Xiao-dan<sup>3</sup>. 1 Department of Epidemiology, School of Binzhou Medical College, Yantai 264003, China; 2 Dana-Farber Cancer Institute, Harvard Medical School, Boston, USA; 3 Shanghai Key Laboratory of Children's Environmental Health, Xinhua Hospital Affiliated to School of Medicine, Shanghai Jiaotong University

Corresponding author: YU Xiao-dan, Email: xdyu1108@163.com

**【Introduction】** In many studies about biomedical research factors influence on the outcome variable, it has no influence or has a positive effect within a certain range. Exceeding a certain threshold value, the size of the effect and/or orientation will change, which called threshold effect. Whether there are threshold effects in the analysis of factors ( $x$ ) on the outcome variable ( $y$ ), it can be observed through a smooth curve fitting to see whether there is a piecewise linear relationship. And then using segmented regression model, LRT test and Bootstrap resampling method to analyze the threshold effect. Empower Stats software developed by American X & Y Solutions Inc has a threshold effect analysis module. You can input the threshold value at a given threshold segmentation simulated data. You may not input the threshold, but determined the optimal threshold analog data by the software automatically, and calculated the threshold confidence intervals.

**【Key words】** Threshold effect; Segmented regression model; Likelihood ratio test; Bootstrap resampling method

生物医学研究中许多研究因素对结果变量的影响不是简单的直线关系,即在一定范围内无作用或有正效作用,超过某一阈值后作用大小或/和方向发生变化,称为阈值效应。以下介绍阈值效应分析原理及 Empower Stats 软件进行阈值效应分析的应用。

## 基本原理

1. 阈值效应检验:如果要检验研究因素  $x$  对结果变量  $y$  的作用,在  $k$  点存在转折。可以首先定义 2 个自变量  $x_1$ 、 $x_2$  (当  $x < k$ ,  $x_1 = x$ ,  $x_2 = 0$ ;  $x \geq k$ ,  $x_1 = 0$ ,  $x_2 = x$ ), 然后采用如下回归方程拟合数据<sup>[1]</sup>

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \varepsilon \quad (1)$$

当  $x < k$  时,式(1)简化为  $y = \alpha + \beta_1 x_1 + \varepsilon$ ; 当  $x \geq k$  时,式(1)简化为  $y = \alpha + \beta_2 x_2 + \varepsilon$ 。因此  $\beta_1$ 、 $\beta_2$  分别表示  $x < k$  与  $x \geq k$  两段的回归系数,即  $y$  随  $x$  变化的斜率。如果这两段斜率的差异无统计学意义,式(1)则简化为

$$y = \alpha + \beta x + \varepsilon \quad (2)$$

统计检验可以用 Wald 检验比较  $\beta_2 - \beta_1$  是否等于零,或用似然比检验(LRT)<sup>[2]</sup>比较式(1)与式(2);如 Wald 检验结果  $\beta_2 - \beta_1 \neq 0$ ,或 LRT 检验结果式(1)显著优于式(2),则表示  $x$  在阈值  $k$  前后对  $y$  的效应不同。

2. 确定阈值:一般通过平滑曲线拟合判断是否存在阈值<sup>[3]</sup>,观察危险因素  $x$  与结局变量  $y$  有无分段式的关系,然后采用最大似然法确定阈值。所谓最大似然法,即从数据中寻找一个  $k$  值,其所得出的式(1)给出最大的似然值。具体操作分两步。

第一步:从  $x$  的 5% 百分位数开始,按 5% 逐步递增,到 95% 百分位数,共 19 个点,分别赋  $k$  值为这 19 个点的  $x$  值,找出哪个百分位数给出最大的似然值,

记为  $p_1$ , 并分别找出  $p_1 - 4\%$  与  $p_1 + 4\%$  百分位数所对应的  $x$  值, 记为  $k_{\min}$ 、 $k_{\max}$ , 应将  $k$  值缩小到该范围内。

第二步: 用递归方法在  $k_{\min}$  至  $k_{\max}$  之间的所有观察到的  $x$  取值内, 找出哪个  $x$  作为  $k$  值给出最大似然值。具体方法是首先比较该范围内的 Q1 (25% 百分位点)、Q2 (50% 百分位点) 与 Q3 (75% 百分位点), 找出哪个位点作为  $k$  值所给出的模型似然值最大, 然后把范围缩小到该位点前后 25% 范围内, 这样每次递归剔除 50% 的  $x$  取值。最终得出能给出最大的模型似然值的  $k$  值。

3. 确定阈值置信区间 (CI): 采用 Bootstrap 重抽样方法确定阈值 CI<sup>[4]</sup>。即从现有数据中随机抽样, 重新抽取一个同样样本量的数据  $i$ 。抽出的个体再放回, 这样在抽取的数据中, 原数据中有些个体可能未被抽中, 也有些个体可能被抽出多次。对抽取的数据进行阈值分析, 记录所得阈值  $k_i$ 。这样重复随机抽样 1000 次, 计算出 1000 个阈值  $k_i$ 。再计算这 1000 个  $k_i$  的 2.5% 与 97.5% 百分位数, 即为所观察阈值的 95% CI。

4. Empower Stats 软件在阈值分析上的应用: 上述确定阈值及其 CI 的过程复杂, 工作量大, 且需很高的编程技巧。美国 X & Y Solutions 软件公司开发的 Empower Stats 软件, 设计有阈值效应分析模块, 用户不需自己编程就可方便快捷完成上述分析。在分析研究因素 ( $x$ ) 对结果变量 ( $y$ ) 的作用有无阈值效应时, 可先通过该软件的平滑曲线拟合模块, 观察  $x$  与  $y$  的曲线拟合是否有分段线性关系, 还可从曲线上观察出阈值。如有分段线性关系, 再使用阈值效应模块。可以输入阈值, 该软件按所给阈值分段拟合数据; 也可不输入阈值, 让软件按前述方法自动确定阈值, 并计算阈值的 CI。

### 实例分析

【例 1】 锰摄入过量可对机体产生不良影响。出生前低浓度的锰暴露对新生儿生长的影响尚不清楚。2008—2009 年 Yu 等<sup>[5]</sup> 对上海市 1377 对母婴进行多中心研究, 测定胎儿脐血中锰浓度、新生儿出生身长和 BMI 等指标。鉴于血锰浓度呈强偏态分布, 对锰浓度做常用对数变换后, 再进行分析。研究者使用 Empower Stats 软件, 首先绘制对数锰浓度 ( $\log Mn$ ) 与新生儿身长、BMI 的曲线拟合图, 发现新生儿 BMI 与  $\log Mn$  曲线分两段, 第一段随着  $\log Mn$  增加呈下降趋势, 后一段则呈上升趋势 (图 1), 根据曲线观察并结合其他数据分析结果,  $\log Mn$  定义转折点为 0.7  $\mu g/L$ 。第二步, 调用阈值效应分析模块分析胎儿脐带血中的

锰浓度与新生儿 BMI 关系 (表 1)。结果显示, 当脐血  $\log Mn < 0.7$  时 (即锰浓度  $< 5.0 \mu g/L$ ), 回归系数为  $-0.233$ ; 当脐血  $\log Mn \geq 0.7$  时 (即锰浓度  $\geq 5.0 \mu g/L$ ), 回归系数为  $0.115$ , 两回归系数差别的 Wald 检验,  $P < 0.001$ 。如果不用分段线性模型, 用一条线拟合数据, 得出回归系数为  $0.049$ ,  $P = 0.070$ , 但该直线回归系数显然不能正确反映胎儿脐血锰浓度与新生儿 BMI 关系。两模型比较的 LRT 检验  $P < 0.001$ 。

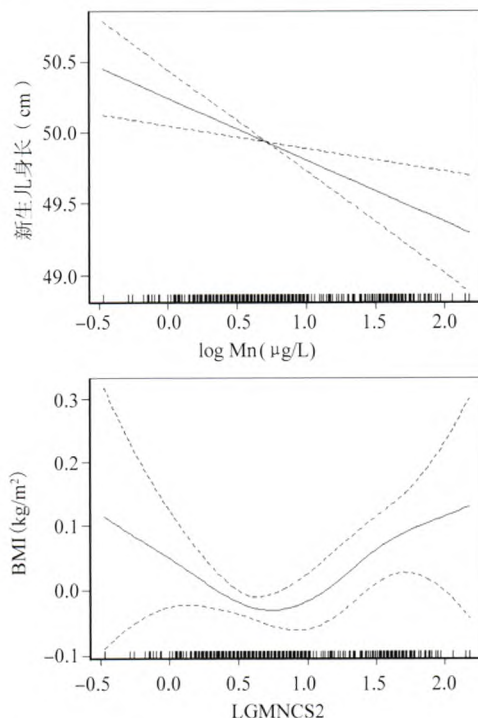


图 1 孕妇产前锰暴露与新生儿生长的剂量-反应关系平滑曲线拟合

表 1 新生儿脐血  $\log Mn$  与 BMI 关系的阈值效应分析

项 目	BMI		
	$\beta$	95%CI	P 值
模型 I			
一条直线回归系数	0.049	(-0.004, 0.101)	0.070
模型 II			
折点 ( $k$ )	0.700		
$< k$ 时回归系数 1	-0.233	(-0.407, -0.059)	0.009*
$> k$ 时回归系数 2	0.115	(0.049, 0.180)	$< 0.001$ *
回归系数 2 与 1 之差	0.348	(0.143, 0.552)	$< 0.001$ *
对数似然比检验		$P < 0.001$ *	
折点处 $y$ 预测值		2.670 (2.65, 2.70)	

注: \*  $P < 0.05$

【例 2】 黄爱群等<sup>[6]</sup> 1996 年对安徽省 795 名农村人群采用问卷调查分析家庭成员间年龄与 SBP/DBP 的关系。Empower Stats 阈值效应分析模块输出结果见图 2、表 2。当年龄  $< 39.6$  岁时, 年龄对 SBP 升高无显著作用 (回归系数 = 0.1,  $P = 0.274$ ); 当年龄  $\geq 39.6$  岁时, 年龄每增加 1 岁, SBP 升高 1.3 mm Hg (回归系数 = 1.3,  $P < 0.001$ )。比较



分段线性模型与一条线模型的似然比检验结果  $P < 0.001$ , 表明年龄对 SBP 升高作用的阈值为 39.6 (95%CI: 18.5 ~ 64.2) 岁。同理, 当年龄  $< 44.6$  岁时, 年龄每增加 1 岁, DBP 升高的风险增加 0.1 mm Hg (回归系数 = 0.1,  $P = 0.007$ ), 而当年龄  $\geq 44.6$  岁时, 年龄每增加 1 岁, DBP 升高 0.4 mm Hg (回归系数 = 0.4,  $P < 0.001$ )。比较分段线性模型与一条线模型的似然比检验结果  $P = 0.003$ , 表明年龄对 DBP 明显升高作用阈值为 44.6 (95%CI: 18.3 ~ 64.3) 岁 (表 2)。

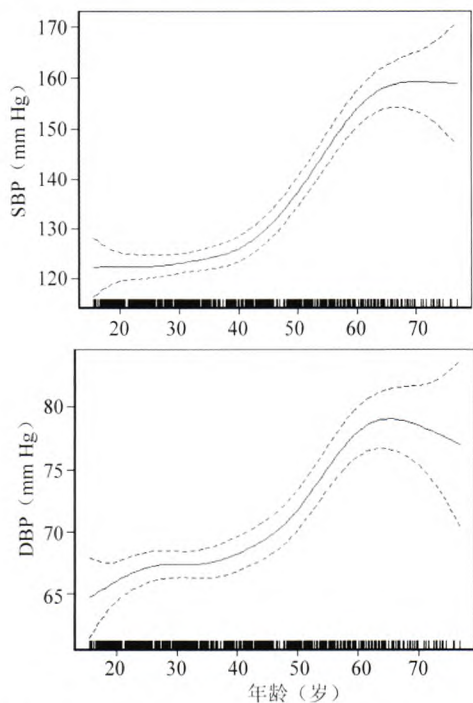


图 2 年龄与 SBP/DBP 关系的平滑曲线拟合

### 讨 论

现代医学研究中常应用阈值效应, 但经常被忽视。如本文例 1, 如不采用分段模型, 可得出新生儿 BMI 随脐血锰浓度增加的线性关系 ( $\beta = 0.049, P = 0.07$ ); 而采用分段线性模型则更确切反映两者的分段关系, 即当脐血锰浓度  $< 5.0 \mu\text{g/L}$ , 新生儿 BMI 随

脐血锰浓度增加而减小; 但脐血锰浓度  $\geq 5.0 \mu\text{g/L}$  时, 新生儿 BMI 随脐血锰浓度增加而增加。阈值效应分析也可用于确定高危人群。如本文例 2 中结果显示年龄对 SBP 升高作用的阈值为 39.6 (95%CI: 18.5 ~ 64.2) 岁; 年龄对 DBP 明显升高作用的阈值为 44.6 (95%CI: 18.3 ~ 64.3) 岁。根据该结果, 可选定  $> 39$  岁为高血压重点监测人群。

确定阈值最确切的方法是分别计算每个  $x$  (除去最小和最大的值以便分段模拟数据) 作为  $k$  值给出的模型似然值进行比较, 但计算量大且费时, 当样本量大又需要做 Bootstrap 重抽样计算 CI 时, 其可行性差。递归法首先比较 Q1 (25% 百分位点)、Q2 (50% 百分位点)、Q3 (75% 百分位点) 3 个点给出的模型似然值, 把范围缩小到最大点前后 25% 内, 每次递归缩小一半, 最终找出  $k$  值。直接用递归法效率最高, 当曲线拟合波动性小,  $k$  值比较居中时, 效果较好; 但当  $k$  值较偏向两端, 曲线拟合波动较大时, 不容易找出最佳  $k$  值。本文提出的分两步方法, 兼顾到  $k$  值偏向两端与曲线拟合波动较大的情况, 同时又充分利用到递归法的高效率。本文主要介绍两段线性模型分析  $x$  与  $y$  的关系, 所用原理与方法可扩展到分析  $x$  与  $y$  之间呈三段或多段线性关系的数据。

使用传统的如 SAS、R 等统计分析软件确定阈值及其 CI, 均需要通过复杂的编程才能实现。使用 Empower Stats 软件, 用户只要输入结果变量 ( $y$ )、研究因素 ( $x$ )、要调整的协变量 (如需要调整其他变量的作用)、分层变量 (如需要分层分析), 软件自动编写 R 程序, 实现上述分析过程, 输出内容中不仅包括所编写的 R 程序及其运算过程的详细记录, 并把最终分析结果整理成表格 (如例 1 所示), 方便用户阅读。该软件大大提高了用户数据分析能力与分析效率。

### 参 考 文 献

- [1] Harrell FE Jr. Regression modeling strategies: with applications to linear models, logistic regression, and survival analysis. New York: Springer-Verlag, 2001.
- [2] Chen XR. Advanced Mathematical Statistics. Hefei: China Science and Technology University Press, 2009. (in Chinese)  
陈希孺. 高等数理统计学. 合肥: 中国科学技术大学出版社, 2009.
- [3] Cleveland WS. Robust locally weighted fitting and smoothing scatterplots. J Am Stat Assoc, 1979, 74: 829-836.
- [4] Bradley E, Robert T. An introduction to the bootstrap. New York: Chapman & Hall Ltd, 1993.
- [5] Yu XD, Cao LL, Yu XG, et al. Elevated cord serum manganese level is associated with a neonatal high ponderal index. Environ Res, 2012, 121: 79-83.
- [6] Huang AQ, Liu X, Zhang WB, et al. The analysis of the risk factors about hypertension among rural population. Chin J Dis Control Prev, 1999, 3 (2): 82-83. (in Chinese)  
黄爱群, 刘学, 张文兵, 等. 农村地区居民高血压危险因素分析. 疾病控制杂志, 1999, 3 (2): 82-83.  
(收稿日期: 2013-07-10)  
(本文编辑: 张林东)

表 2 年龄与 SBP/DBP 的阈值效应分析

项 目	SBP			DBP		
	$\beta$	95%CI	P 值	$\beta$	95%CI	P 值
模型 I						
一条直线回归系数	0.8	(0.7, 0.9)	$< 0.001^*$	0.3	(0.2, 0.3)	$< 0.001^*$
模型 II						
折点 ( $k$ )	39.6	(18.5, 64.2)		44.6	(18.3, 64.3)	
$< k$ 时回归系数 1	0.1	(-0.1, 0.3)	0.274	0.1	(0.0, 0.2)	0.007*
$> k$ 时回归系数 2	1.3	(1.1, 1.5)	$< 0.001^*$	0.4	(0.3, 0.6)	$< 0.001^*$
回归系数 2 与 1 之差	1.1	(0.8, 1.5)	$< 0.001^*$	0.3	(0.1, 0.5)	0.003*
对数似然比检验	$P < 0.001^*$			$P = 0.003^*$		
折点处 $y$ 预测值	124.63 (122.03, 127.23)			69.33 (67.84, 70.83)		

注: \* 同表 1