

广义极值分布参数估计方法的对比分析*

陈子燊¹, 刘曾美^{1,2}, 路剑飞¹

(1. 中山大学水资源与环境系, 广东 广州 510275;
2. 华南理工大学水利水电工程系, 广东 广州 510640)

摘要: 简介了广义极值分布的 3 种参数估计方法: 极大似然 (ML)、线性矩 (LM) 和间隔最大积 (MPS) 方法的特点和计算方法, 采用历年马口月最大径流量和广州日最大降水量作为广义极值分布不同参数估计方法的实证分析例子。分析结果表明, 两实例各自 3 种参数估计方法得到的 3 个参数值较为接近, 各种拟合优度检验结果表明两实例均服从广义极值分布, 但 MPS 参数估计推算的设计值与观测值拟合更好。

关键词: 广义极值分布; 参数估计方法; 拟合优度检验; 实证分析

中图分类号: TV122 **文献标志码:** A **文章编号:** 0529-6579 (2010) 06-0105-05

Comparative Analysis of Parameter Estimation Methods of Generalized Extreme Value Distribution

CHEN Zishen¹, LIU Zengmei^{1,2}, LU Jianfei¹

(1. Department of Water Resource and Environment, Sun Yat-sen University, Guangzhou 510275, China;
2. Department of Water Conservancy and Hydropower Engineering,
South China University of Technology, Guangzhou 510640, China)

Abstract: Three parameter estimation methods of generalized extreme value distribution function were briefly introduced in the paper, which included the maximum-likelihood estimation, the linear moment estimation and the maximum product of spacing estimation. Two demonstration examples including monthly maximum runoff at Makou Station and daily maximum precipitation in Guangzhou over the past years were analyzed by the three parameter estimation methods of GEV distribution. The results indicated that three parameters obtained by three different estimation methods were very close. Several goodness fit tests showed that two examples were obeyed the generalized extreme value distribution. And the designed values predicted by maximum product of spacing estimation were better fitted with the measured values.

Key words: generalized extreme value distribution; parameter estimation methods; goodness fit test; demonstration analysis

政府间气候变化专门委员会 (IPCC) 第一工作组 2007 年编写的第四份评估报告指出, 全球暖化愈趋明显, 水文气象极端事件发生的频率可能会增大^[1]。准确地度量极值事件发生的概率及其分位数, 预测极值事件可能造成的危害, 已经成为工程风险管理、控制和决策研究的重要问题。通常对

水文气象极端事件频率分析采用两种方法: 一是选用某些统计分布函数拟合水文气象要素的累积概率及其相应的重现水平, 如显示了较大实用性并得到广泛应用的皮尔逊 III 型分布^[2-4], 但其是根据一定代表性的经验得出的分布曲线, 缺乏严格的概率理论依据; 另一种方法是应用极值分布理论, 分别使

* 收稿日期: 2009-12-09

基金项目: 国家自然科学基金重点基金资助项目 (50839005); 2009 年广东水利科技创新与推广项目

作者简介: 陈子燊 (1952 年生), 男, 教授, 博士生导师; E-mail: eesczs@mail.sysu.edu.cn

用 Weibull、Gumbel、Frechet 三种极值分布函数之一对水文气象极值序列加以拟合计算^[5-9]。近年来,许多研究人员采用了适用性更强的广义极值分布 (generalized extreme-value distribution, 简称 GEV) 理论,解决了只能用一种极值分布函数类型的局限性,广义极值 (GEV) 分布已经在水文气象极端事件研究领域得到了较多的应用^[10-15]。

频率分析统计推断的精度除了取决于样本和选用的概率分布模型外,参数估计方法对计算结果具有重要影响。因此,概率分布模型的参数估计是极值统计理论研究的核心内容之一。参数估计的常用方法有极大似然 (ML)、矩法 (MOM)、概率权值矩 (PWM)、线性矩 (LM)、最小二乘 (LS) 和较新提出的间隔最大积 (MPS) 等估计方法。本文拟采用广义极值分布的极大似然法、线性矩法和间隔最大积 3 种参数估计法,对降水和径流极值序列做实证分析,并通过拟合优度检验加以对比。

1 广义极值分布

20 世纪 30 年代, Fisher 和 Tippett^[16]对独立同分布的极大值渐近分布进行理论研究时提出了 3 种极值分布:

极值 I 型 (Gumbel) 分布:

$$F_X(x) = P(X < x) = \exp\left[-\exp\left(\frac{x-\mu}{\sigma}\right)\right],$$

$$-\infty < x < +\infty;$$

极值 II 型 (Fréchet) 分布:

$$F_X(x) = P(X < x) = \begin{cases} \exp\left[-\left(\frac{x-\mu}{\sigma}\right)^\xi\right], & x > \mu \\ 0, & x \leq \mu \end{cases};$$

极值 III 型 (Weibull) 分布:

$$F_X(x) = P(X < x) = \begin{cases} \exp\left[-\left(\frac{x-\mu}{\sigma}\right)^\xi\right], & x < \mu \\ 1, & x \geq \mu \end{cases}$$

式中, ξ, μ, σ 分别为形状参数、位置参数和尺度参数。

Jenkinson^[17]、Coles^[18]根据极值分布理论,证明当极值的渐近分布存在且为非退化时可以将 3 种类型的经典极值分布发展为一种统一的具有 3 参数的极值分布函数——广义极值分布 (简记为 GEV)。设 X_1, \dots, X_m 是服从 GEV 分布的独立随机变量,则分布函数 F_X 为:

$$F_X(x) = P(X < x) = \begin{cases} \exp\left\{-\left[1 - \xi\left(\frac{x-\mu}{\sigma}\right)\right]^{1/\xi}\right\}, & \xi \neq 0 \\ \exp\left[-\exp\left(\frac{x-\mu}{\sigma}\right)\right], & \xi = 0 \end{cases}$$

当 $\xi \rightarrow 0$ 为极值 I 型,即 Gumbel 分布; $\xi < 0$ 为极值 II 型,即 Fréchet 分布, $\xi > 0$, 为极值 III 型,即 Weibull 分布。在计算得到分布函数参数后,即可对给定的频率 F 求解其对应的分位数 x_F :

$$x_F = \begin{cases} \mu + \frac{\sigma}{\xi} \{1 - [-\ln(F)]^\xi\}, & \xi \neq 0 \\ \mu - \sigma \ln[-\ln(F)], & \xi = 0 \end{cases}$$

2 参数估计

2.1 极大似然估计 (MLE)

极大似然估计法是 Fisher 于 1922 年提出的一种点估计方法。由于极大似然估计法使得参数估计结果在总体上反映样本的统计信息,具有良好的统计性质,因此得到了广泛重视与普遍应用^[12-15]。

参数估计: 设 $\{x_i\}$ 为服从 GEV 分布的独立同分布集合,当 $\xi = 0$ 时,其 n 个观测 $\{x_1, x_2, \dots, x_n\}$ 的对数似然函数为: $\ln[L(\theta|x)] = -n \ln(\sigma) +$

$$\sum_{i=1}^n \left[\left(\frac{1}{\xi} - 1 \right) \ln(y_i) - (y_i)^{1/\xi} \right]$$

式中, $\theta = (\mu, \sigma, \xi)$, $y_i = [1 - (\xi/\sigma)(x - \mu)]$ 。

由对数似然函数对 θ 求一阶导数,令 $\frac{d \ln L(\theta)}{d \mu} =$

$$0, \frac{d \ln L(\theta)}{d \sigma} = 0, \frac{d \ln L(\theta)}{d \xi} = 0, \text{ 求得的似然方程组}$$

为非线性方程组,可采用 Newton-Raphson 迭代方法得到参数的极大似然估计值。

2.2 线性矩估计 (LME)

Hosking 等^[19]指出,小样本参数的极大似然估计量是很不稳定的,推荐使用线性矩估计法。线性矩估计法同传统矩估计法基本原理一样,都是令样本矩与总体矩相等,得到总体参数的估计。Hosking^[20]曾就线性矩法在 P-III 分布的参数估计算法与传统矩参数估计方法在统计性能方面的差异进行了对比分析计算,统计试验结果表明,线性矩法确实具有良好性能,明显优于传统矩法。

参数估计: 设随机变量 X 的分布函数为 $F = (x, \theta)$, $x_{n,n} \leq \dots \leq x_{1,n}$, n 是样本的次序统计量,称

$$\lambda_r = r^{-1} \sum_{k=0}^{r-1} (-1)^k \binom{r-1}{k} EX_{k+1,r} \quad r = 1, 2, \dots$$

为 r 阶 L 矩。前 3 阶样本 L 矩的无偏估计为:

$$\hat{\lambda}_1 = \sum_{i=1}^n x_{i,n}/n; \hat{\lambda}_2 = \sum_{i>j} (x_{i,n} - x_{j,n})/(n(n-1));$$

$$\hat{\lambda}_3 = \sum_{i>j>k} 2(x_{k,n} - 2x_{j,n} + x_{i,n})/(n(n-1)(n-2))$$

Hosking 等^[19]定义线性矩的偏态系数为: $\hat{\tau}_3 =$

$\hat{\lambda}_3/\hat{\lambda}_2$ ，并给出了 GEV 分布参数的线性矩估计的近似方法：当 $-0.5 < \hat{\tau}_3 \leq 0.5$ 时， $\xi \approx 7.8590z + 2.9554z^2$ ，其中 $z = 2/3(3 + \hat{\tau}_3) - \ln(2)/\ln(3)$ ；尺度参数 $\beta: \beta = \lambda_2 \xi [(1 - 2^{-\xi})\Gamma(1 + \xi)]$ ；位置参数： $\mu = \lambda_1 - \sigma [1 - \Gamma(1 + \xi)]/\xi$ 。

2.3 间隔最大积估计 (MPSE)

间隔最大积估计 (Maximum Product of Spacing Estimation) 方法是 Cheng 和 Amin^[21] 提出并由 Ranneby^[22] 推荐的一个新的参数估计方法。Cheng 和 Amin^[21] 证明了 MPS 和 ML 估计都同样具有渐近充分性、一致性和有效性，但极大似然参数估计方法在许多情况下由于似然函数的无界性或不存在局部极大值而失效，间隔最大积估计可适用于任何连续的单变量分布的参数估计，因而具更好的可辨识性。与包括极大似然估计在内的其它参数估计方法相比较，MPS 特别适合于三参数分布函数，如广义极值分布、广义 Pareto 分布、威布尔分布、对数正态分布、3 参数的 Γ 分布曲线的参数估计，并能提供更好的稳健性、一致性、有效性等统计估计量^[23]。

参数估计：设 $\{x_i\}$ 为服从 GEV 分布的独立同分布集合，其递增排列的 n 个观测 $\{x_1, x_2, \dots, x_n\}$ 的累积分布函数为： $F(x; \theta^0): \theta^0 \in \Theta, \Theta \subseteq R^k (k \geq 1)$

设 $\theta \in \Theta$ 是 θ^0 的估计量，则 $F(x; \theta)$ 是 $F(x; \theta^0)$ 的估计量。定义一阶间隔为：

$$D_i(\theta) = F(x_i; \theta) - F(x_{i-1}; \theta),$$

$$D_1(\theta) = F(x_1; \theta), D_{n+1}(\theta) = 1 - F(x_n; \theta)$$

式中，规定附加的数据点 $x_0 = -\infty, x_{n+1} = \infty$ ，因此， $F(x_0; \theta) = 0; F(x_{n+1}; \theta) = 1$

设 $G_n(\theta)$ 为 $D_i(\theta)$ 的几何平均， $S_n(\theta)$ 为 $G_n(\theta)$ 的自然对数，即

$$G_n(\theta) = \left(\prod_{i=1}^{n+1} D_i(\theta) \right)^{1/(n+1)},$$

$$S_n(\theta) = \ln G_n(\theta) = \frac{1}{n+1} \sum_{i=1}^{n+1} \ln D_i(\theta)$$

类似于极大似然估计方法，MPS 法通过似然函数可得到参数 θ 的估计值。

3 实证分析

3.1 实例

实例 1：马口水文站为珠江三角洲西江进入河网区的国家重点监测站。马口站多年平均径流量 2 277 亿 m^3 ，占珠江径流总量的 77.1%；年内径流十分集中，汛期（4-9 月）的径流量占全年径流

总量的 77.14%。枯期的径流量占全年径流总量的 22.86%。样本序列为 1959-2007 年历年的月最大径流量。

实例 2：广州位于珠江三角洲中部，多年平均年降水量约 1 736 mm，近年来强降水过程引发的内涝趋于严重。样本序列为 1951-2008 年历年的日最大降水量。

3.2 参数估计结果

两极值序列的广义极值分布 3 种参数估计方法与参数估计的结果值见表 1，广义极值分布的不同重现水平的分位数与标准误差计算结果见表 2。

表 1 广义极值分布的参数估计结果

Table 1 The results of parameter estimation of generalized extreme value distribution

极值实例	参数估计方法	形态参数	尺度参数	位置参数
广州日最大降水序列	LM	-0.132	34.235	99.256
	ML	-0.109	34.628	100.005
	MPS	-0.116	37.018	99.593
马口月最大径流序列	LM	0.142	103.327	424.990
	ML	0.170	103.773	427.077
	MPS	0.151	111.465	423.999

表 1 显示，广州日最大降水序列和马口月最大径流序列的广义极值分布的 3 种参数估计方法推算参数值比较接近。广州日最大降水序列的广义极值分布的形态参数为负值表明，日最大降水序列属于极值 II 型分布，即 Frechet 分布，属于负偏右长尾，分布形态同皮尔逊 III 型密度曲线。

马口月最大径流序列的广义极值分布的形态参数为正值。序列属于极值 III 型分布，即 Weibull 分布，属于右短尾分布。理论分析表明，极值序列上限为： $x \leq \mu + \sigma/\xi$ 。由参数估计推算的马口月最大径流序列的上界极限值介于 1036.5 亿 ~ 1163.7 亿 m^3 。

三种参数估计推算的设计频率分位值说明，广州日最大降水和马口月最大径流序列的相同重现期由 MPS 法推算的分位值最大，线性矩法推算结果最小，而极大似然法介于二者之间。

3.3 拟合优度检验

采用了以下几个拟合优度检验方法：

(1) PPCC 检验法

根据实测样本序列排序后观测值 x_i 和经验频率分布 P_{ei} 计算的与广义极值分布相应的分位值 \hat{x}_i ，求二者的相关系数 r ：

表 2 广义极值分布的设计频率分位值
Table 2 The generalized extreme-value quantile estimation

重现期/a	广州日最大降水/mm		
	LM	ML	MPS
100	316.1	306.7	324.3
50	274.2	268.3	282.1
20	223.8	221.4	230.8
10	189.0	188.3	194.7

重现期/a	马口月最大径流/亿 m ³		
	LM	ML	MPS
100	773.1	756.4	791.5
50	733.0	721.0	750.4
20	673.3	667.0	688.5
10	621.8	619.1	634.5

$$r = \frac{\sum_{i=1}^n (x_i - x_m)(\hat{x}_i - \hat{x}_m)}{[\sum_{i=1}^n (x_i - x_m)^2 \sum_{i=1}^n (\hat{x}_i - \hat{x}_m)^2]^{\frac{1}{2}}}$$

(2) 均方根误差

均方根误差计算公式:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n \left(\frac{\hat{x}_i - x_i}{x_i} \right)^2}$$

(3) χ^2 拟合检验

检验假设, H_0 : 样本的总体与该理论分布无区别。取显著水平 $\alpha = 0.05$, 若 $\chi^2 < \chi^2_{1-\alpha}$, 接受原假设 H_0 ; H_1 : 若 $\chi^2 \geq \chi^2_{1-\alpha}$ 则拒绝 H_0 。计算统计量: $\chi^2 = \sum_{i=1}^k (n_i - np_i)^2 / np_i$; n_i : 实际频数, $i = 1, 2, \dots, k$; np_i : 理论频数。

(4) 柯尔莫哥洛夫拟合检验

原检验假设, $H_0: X \sim F_0(x, \mu, \sigma, \xi)$ 。计算柯尔莫哥洛夫统计量:

$$D_n = \sup_x |F_n(x) - F_0(x)| = \max_i |F_n(x_{(i)}) - F_0(x_{(i)})|,$$

$|F_n(x_{(i+1)}) - F_0(x_{(i)})|$
 $F_n(x)$: 样本经验分布; $F_0(x)$: 广义极值分布。
 取显著水平 $\alpha = 0.05$, 若 $D_n > D_{n,\alpha}$, 则拒绝 H_0 。
 广州日最大降水和马口月最大径流广义极值分布累积频率曲线 (图 1) 显示, 各种参数估计结果的累积分布与经验分布符合良好, 尤其是作为概率分布曲线的上尾部由 MPS 参数估计法推算的马口月最大径流更符合实测数据。各种拟合优度检验结果见表 3。从表 3 可看出, 柯尔莫哥洛夫检验、 χ^2

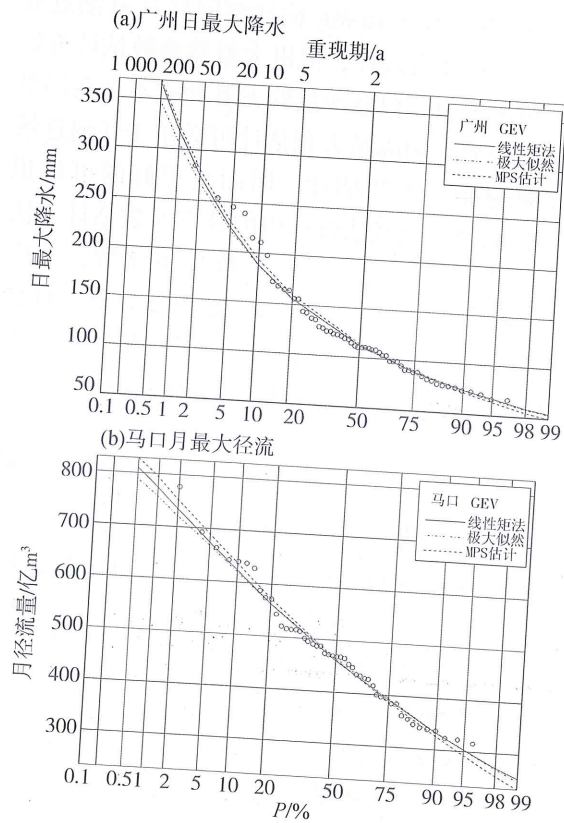


图 1 广州日最大降水与马口月最大径流广义极值分布三种参数估计
 Fig. 1 Three parameter estimation of generalized extreme value distribution for daily maximum precipitation in Guangzhou and monthly maximum runoff at Makou Station

表 3 拟合优度检验结果
 Table 3 The results of goodness-fit test

实例	参数估计法	柯尔莫哥洛夫检验		χ^2 检验		RMSE	PPCC
		$D_{n,\alpha}$	D_n	$\chi^2_{1-\alpha}$	χ^2		
广州日最大降水	线性矩	1.000	0.052	18.265	1.000	13.6	0.993
	极大似然	1.000	0.052	20.377	1.000	12.2	0.992
	MPS	0.999	0.069	13.997	1.000	1.6	0.992
马口月最大径流	线性矩	0.995	0.083	26.479	0.993	7.7	0.992
	极大似然	0.995	0.083	28.624	0.984	5.6	0.991
	MPS	0.946	0.104	25.121	0.996	5.6	0.992

检验和 PPCC 相关系数计算结果表明广州日最大降水和马口月最大径流序列均完全服从广义极值分布。不同参数估计结果对比表明, MPS 参数估计方法计算的 RMSE 明显小于其余二者估计结果, 说明拟合最优, 可以作为广州日最大降水和马口月最大径流序列广义极值分布推算的最佳设计值。

4 结 论

本文分析了广义极值分布理论的 3 种参数估计方法的特点和计算方法, 采用马口月最大径流量和广州日最大降水量作为广义极值分布不同参数估计方法的两个实例, 对分析结果有以下结论:

1) 三种参数估计方法的 3 个参数值较为接近, 广州日最大降水的形态参数为负值, 属于极值 II 型分布。马口月最大径流的形态参数为正值, 属于极值 III 型分布;

2) 拟合优度检验结果表明两实例均服从广义极值分布, 均可选择作为统计推断的分布函数类型;

3) MPS 参数估计推算的设计值与观测值拟合更好。

参考文献:

- [1] IPCC. IPCC Fourth Assessment Report (AR4) [M]. Cambridge: Cambridge University Press, 2007.
- [2] 杨远东, 王辉, 杨树佳, 等. 皮尔逊 III 型分布三参数估计新方法[J]. 水资源研究, 2007, 28(3): 18-28.
- [3] 周芬, 郭生练, 肖义, 等. P-III 型分布参数估计方法的比较研究[J]. 水电能源科学, 2003, 21(3): 10-13.
- [4] 黄振平, 王春霞, 马军建. P-III 型分布的适应性与水文设计值的误差分析[J]. 水文, 2002, 22(5): 21-24.
- [5] 张秀芝. Weibull 分布参数估计方法及其应用[J]. 气象学报, 1996, 54(1): 108-116.
- [6] 蔡敏, 丁裕国, 江志红. L-矩估计方法在极端降水研究中的应用[J]. 气象科学, 2007, 27(6): 597-603.
- [7] 段忠东, 周道成. 极值概率分布参数估计方法的比较研究[J]. 哈尔滨工业大学学报, 2004, 36(12): 1605-1609.
- [8] 曹兵, 王义刚, YOU Zaijin. 三种计算设计波高方法的比较[J]. 海洋工程, 2006, 24(4): 75-80.
- [9] 赵伟, 杨永增, 于卫东, 等. 长期极值统计理论及其在海洋环境参数统计分析中的应用[J]. 海洋科学进展, 2003, 21(4): 471-476.
- [10] 金光炎. 广义极值分布及其在水文中的应用[J]. 水文, 1998(2): 9-15.
- [11] 陈元芳, 李兴凯, 陈民, 等. 可考虑历史洪水信息的广义极值分布线性矩法的研究[J]. 水文, 2008, 28(3): 8-13.
- [12] PRESCOTT P, WALDEN A T. Maximum-likelihood estimation of the parameters of the three-parameter generalized extreme-value distribution from censored samples [J]. J Stat Comput Simul, 1983, (6): 241-250.
- [13] MARTINS E S. Generalized maximum-likelihood generalized extreme-value quantile estimators for hydrologic data [J]. Water Resources Research, 2000, 36(3): 737-744.
- [14] 刘聪, 秦伟良, 江志红. 基于广义极值分布的设计基本风速及其置信限计算[J]. 东南大学学报: 自然科学版, 2006, 36(2): 331-334.
- [15] 陈兴旺. 广义极值分布理论在重现期计算的应用[J]. 气象与减灾研究, 2008, 31(4): 52-54.
- [16] FISHER R A, TIPPETT L H. Limiting forms of the frequency distribution of the largest or smallest member of a sample [J]. Proc Cambridge Philos Soc, 1928, 24: 180-190.
- [17] JENKINSON A F. The frequency distribution of the annual maximum (or minimum) values of meteorological elements [J]. Q J R Meteorol Soc, 1955, 81: 158-171.
- [18] COLES S. An introduction to statistical modeling of extreme values [M]. New York: Springer Verlag, 2001: 36-78.
- [19] HOSKING J R M, WALLIS J R, WOOD E F. Estimation of the generalized extreme value distribution by the method of probability-weighted moments [J]. Technometrics, 1985, 27(3): 251-261.
- [20] HOSKING J R M. L-moments: Analysis and estimation of distributions using linear combinations of order statistics [J]. J Roy Statist Soc: Ser B, 1990, 52: 105-124.
- [21] CHENG R C H, AMIN N A K. Estimating parameters in continuous univariate distributions with a shifted origin [J]. J Roy Statist Soc: Ser B, 1983, 45(3): 394-403.
- [22] RANNEBY B. The maximum spacing method: an estimation method related to the maximum likelihood method [J]. Scand J Statist, 1984, 11: 93-112.
- [23] MAGNUS E. Alternatives to maximum likelihood estimation based on spacings and the Kullback-Leibler divergence [J]. Journal of Statistical Planning and Inference, 2008, 138: 1778-1791.