

论水文计算中的假相关

丁 星 杨树林

(成都科技大学)

一、引 言

相关分析广泛地用于水文计算，但是，只有正确地使用这种分析技术才能得出有意义的结果。所谓“正确”，当然涉及的面相当广泛，本文仅就避免出现假相关这一点加以论述。

假相关一般指两个变量之间表现出的一种虚假的相关关系，而实际并不存在统计关系。早在1897年，皮尔逊^[1]就对假相关进行过研究，但是直到1965年，才由本森^[2]明确指出水文和水力学中的假相关现象。随后，叶菲耶维奇^[3]在其《水文概率和统计》一书中作了比较系统的论述。1977年赫安^[4]在《水文学的统计方法》一书中进一步论述了这一课题。我国金光炎工程师^[5]在1981年也发表文章指出水文中的假相关现象。最近，肯奈^[6]进一步分析了这一现象，并警告要当心自身的假相关。尽管对假相关的研究已有很长时间，但迄今对这一问题的认识还有些模糊不清，特别是在我国水文计算中尚未引起足够的重视，误用假相关的现象仍然存在。因此，系统地论述这一课题，澄清一些模糊的概念，指出水文中可能出现的各种假相关现象，找出避免的方法，无疑是非常必要的。

为叙述和分析方便，本文尝试将通常遇到的假相关分为隐假相关和显假相关。前者系指间接引起的假相关；后者系指明显的假相关。

二、隐 假 相 关

这里，隐假相关的含义和一般所说的自身假相关相同，是由于包含相同的变量造成的假相关。

设有两个变量 x 和 y ，相互间完全独立。现构成两个新变量 A 和 B 。

第一种情况： $A = x + y$, $B = x$;

第二种情况： $A = x \cdot y$, $B = x$;

第三种情况： $A = y / x$, $B = x$ 。

以后将会看到：尽管 x 、 y 之间相互完全独立，即 $\gamma_{xy} = 0$ ，但是， A 和 B 却显示出一定的相关性，即 $\gamma_{AB} \neq 0$ 。这种 A 和 B 间出现的相关关系常叫做假相关。实际上，这种提法并不严格，就 A 和 B 两变量而言， γ_{AB} 确是度量两者之间线性相互关系的客观指标，并非假的。之所谓“假”，仅仅是相对于 x 、 y 两变量的关系而言。下面结合水文中的实例分别讨论这三种情况。

1. 第一种情况 ($A = x + y$, $B = x$ 型)

对于变量 A 和 B ，其表征两者线性相关关系的定量指标——相关系数为^[6]：

$$\gamma_{ab} = \frac{n\sum A \cdot B - \sum A \cdot \sum B}{(\sum A^2 - (\sum A)^2)^{1/2} (\sum B^2 - (\sum B)^2)^{1/2}} \quad (1)$$

将 $A = x + y$, $B = x$ 代入(1)得：

$$\gamma_{AB} = \frac{1 + \gamma_{xy} \left(\frac{S_y}{S_x} \right)}{\left\{ 1 + \left(\frac{S_y}{S_x} \right)^2 + 2\gamma_{xy} \left(\frac{S_y}{S_x} \right) \right\}^{1/2}} \quad (2)$$

式中 γ_{xy} —— x 和 y 间的相关系数；

S_x —— x 的标准差；

S_y —— y 的标准差。

若 $\gamma_{xy} = 0$

$$\text{则 } \gamma_{AB} = \frac{1}{\left[1 + \left(\frac{S_y}{S_x} \right)^2 \right]^{1/2}} \quad (3)$$

式(3)可写成另一种形式

$$\gamma_{AB} = \frac{1}{\left[1 + \left(\frac{C_{y\bar{y}}}{C_{x\bar{x}}} \right)^2 \right]^{1/2}} \quad (4)$$

式中 $C_{y\bar{y}}$ —— y 的变差系数和平均值；

$C_{x\bar{x}}$ —— x 的变差系数和平均值。

如果变量 x 和 y 的均值和变差系数相等，那么，即使 $\gamma_{xy} = 0$ ，新变量 A 和 B 的相

关系数由式(4)计算得 $\gamma_{AB} = 0.707$ 。这种现象便看作是假相关，其假的程度与原始变量 x 和 y 的统计特性紧密有关。图1示 γ_{AB} 随 \bar{y}/\bar{x} 和 C_x/C_y 变化的情况。显然，在 \bar{y}/\bar{x} 一定的情况下， C_x/C_y 愈小， γ_{AB} 愈大。其原因在于 A 和 B 的变化主要由 x 的变化所制约。

由式(4)知， γ_{AB} 取决于原始变量 x 和 y 的一阶矩和二阶矩的统计特性， x 和 y 的三阶矩特性(偏态系数 C_s)是否影响 γ_{AB} 呢？即当原始变量呈偏态时，是否仍可用式(4)计算？为了回答这一问题，做了一系列统计试验，其结果列于表1的第二列。

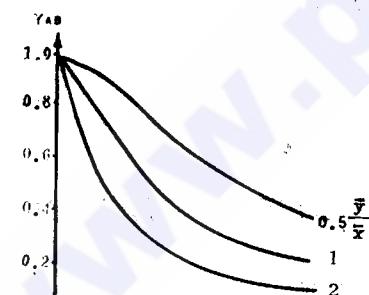


图1 $\gamma_{AB} \sim \bar{y}/\bar{x} \sim C_y/C_x$ 图

对比由式(4)和模拟系列计算的数值，考虑到模拟系列的抽样误差，可以认为两者无显著的差异，这就是说式(4)适用于偏态系列。

A 和 B 之所以相关，其根本原因是由于两者均含有相同的变量 x 。 A 和 B 之间的关系是客观的，但是利用 A 和 B 的这种相关信息，由 B (即 x) 估计出 A ，然后由 A 减去 x 而估出 y ，这就意味着 x 中包含了 y 的信息，即认为 x 和 y 之间有关系，实际上，两者之间完全无关。因此所谓“假”，就是由于 x 估计出与之完全无关的 y 。

水文计算中出现这种情况的例子屡见不鲜，现以北碚站年径流量为例说明。如图2所

表1 统计试验成果表

1							2		3		4	
总体参数值							A=x+y	B=x型	A=xy	B=x型	A= $\frac{y}{x}$	B=x型
x	y	C _x	C _y	C _{xy} /C _x	C _{xy} /C _y	γ _{xy}	γ _{AB}		γ _{AB}		γ _{AB}	
							模拟系列	式(4)	模拟系列	式(6)	模拟系列	式(8)
10	10	0.48	0.48	4	4	0	0.694	0.707	0.674	0.707	-0.535	-0.707
10	10	0.35	0.80	4	4	0	0.438	0.401	0.413	0.401	-0.336	-0.401
10	10	0.80	0.35	4	4	0	0.896	0.916	0.872	0.916	-0.626	-0.916
10	1000	0.48	0.48	4	4	0	0.002	0.01	0.674	0.707	-0.535	-0.707
1000	10	0.48	0.48	4	4	0	1.00	1.00	0.674	0.707	-0.535	-0.707
10	10	0.48	0.48	2	4	0	0.727	0.707	0.694	0.707	-0.610	-0.707
10	10	0.48	0.48	4	2	0	0.707	0.707	0.667	0.707	-0.521	-0.707
10	10	0.48	0.48	2	2	0	0.696	0.707	0.694	0.707	-0.604	-0.707
1000	1000	0.48	0.48	2	2	0	0.685	0.707	0.694	0.707	-0.604	-0.707
1000	1000	0.35	0.35	2	2	0	0.698	0.707	0.706	0.707	-0.651	-0.707
10	10000	0.01	0.90	4	4	0	0.018	0.001				
1000	10	0.90	0.01	4	4	0	1.00	1.00				

注：（1）模拟系列的相关系数是由模拟系列直接按相关系数的定义公式计算而得。

（2）模拟系列的随机模拟方法采用一般的舍选法。

示， x 为上游三站年径流量之和， y 为区间径流。 A 为北碚站的年径流量， $A = x + y$, $B = x$ 。由实测的28年资料得 $\gamma_{AB} = 0.99$ ，表明 A 和 B 之间确有关系。若直接计算 x 和 y 的相关系数，则 γ_{xy} 仅为 0.32，关系非常微弱，由 x 难以估计 y 。如通过 γ_{AB} 来估计 y ，显然其结果只是假象。

2. 第二种情况 ($A = y \cdot x$, $B = x$ 型)

对于变量 A 和 B 的相关系数 γ_{AB} ，在原始变量 x 和 y 的变差系数很小时，可以下式计算。

$$\gamma_{AB} = \frac{\gamma_{xy} C_x + C_y}{\sqrt{(C_y^2 + C_x^2 + 2\gamma_{xy} C_x C_y)^2}} \quad (5)$$

式中符号意义同前。若 $\gamma_{xy} = 0$ 。

则

$$\gamma_{AB} = \frac{1}{\sqrt{1 + (\frac{C_y}{C_x})^2}} \quad (6)$$

由式(6)知，在 $C_x = C_y$ 时， $\gamma_{AB} = 0.707$ 。这表明，即使 x 和 y 毫无关系， A 和 B 之间的线性相关系数也可达到 0.707，和第一种情况一样，这也是由于包含共同变量而引起的假相关。

式(6)是近似公式，在 C_x 和 C_y 较小时适用性较好，那么当 C_x 和 C_y 较大时，近似程度如何？另外在原始变量 x 和 y 具有偏态特性时，该公式适用性又如何？对这些问题，通过统计试验作了研究，其结果如表1第三列所示。总的来说，式(6)误差



图2 测站位置示意图

在 5% 以下，从适用的角度来看可以接受。

第二种情况的假相关在水文计算中可以找到许多例子。例如，为了展延含沙量 (y)，先建立流量 (x) 和输沙率 (A) 的关系，然后由流量展延输沙率，最后由输沙率转换成含沙量。这是典型的 $A = y/x$, $B = x$ 型的假相关。就流量和输沙率而言，两者的关系并非假的。黄河流域某站有 24 年观测资料，年径流量和输沙率之间的相关系数为 0.80，而年径流量和含沙量之间的相关系数则为 -0.10。这清楚地表明，含沙量绝对不可能由年径流量来展延。但是，用上面所说的方法，通过展延得到的输沙率转换成含沙量，其展延成果无疑是虚假的。

3. 第三种情况 ($A = y/x$, $B = x$ 型)

这种情况下的 γ_{AB} 计算式，在原始变量 x 和 y 变差很小时，可表示为^[6]

$$\gamma_{AB} = \frac{\gamma_{xy} C_y - C_x}{(C_x^2 + C_y^2 - 2\gamma_{xy} C_x C_y)^{\frac{1}{2}}} \quad (7)$$

若 $\gamma_{xy} = 0$

则 $\gamma_{AB} = \frac{-1}{(1 + (C_y/C_x)^2)^{\frac{1}{2}}} \quad (8)$

由式 (8) 知，当 $C_y = C_x$ 时， $\gamma_{AB} = -0.707$ 。

式 (8) 和式 (6) 一样是近似的，统计试验的结果（见表 1 第 4 列）表明，在一般情况下，误差较大，随变差系数的减小而减小，而随偏态系数的加大而加大，实际应用时，应予以注意。

属于 $A = y/x$, $B = x$ 型的假相关在水文计算中也时常碰到。例如，设降雨量为 y ，降雨历时为 x ，则降雨强度 $A = y/x$ 和降雨历时 $B = x$ 的关系便为这种类型的假相关。在 y 和 x 确是无关的情况下，若企图通过 A 和 B 的关系由历时 x 推估降雨量 y ，则结果很不可靠。又如设径流量为 y ，流域面积为 x ，则径流模数 $A = y/x$ 和流域面积 $B = x$ 的关系也隶属于这一类型的假相关。只要 x 和 y 之间无关系，通过 A 和 B ，间接由 x 估计 y 便是假相关的结果。

上面就隐假相关分别讨论了三种情况。还有很多其它情况，例如 $A = x - y$, $B = x$; $A = y \pm z$, $B = x \pm z$; $A = yz$, $B = xz$; $A = z/y$, $B = z/x$; 等类型。这些类型在水文计算中不很常见，若工作中遇此情况，可按类似上述的方法进行分析。为节省篇幅，在此不再赘述。

本小节讨论的隐假相关，主要涉及这样一类假相关：直接关系（表面上）是真实的，而间接关系（隐藏着）是虚假的。下小节讨论比较明显的假相关——显假相关。

三、显假相关

属于这种类型的假相关，可以举出下述三种情况：

1. 点簇假相关

图 3 为点簇假相关最典型的情况。尽管变量 x 和 y 没有关系，但计算出的相关系数

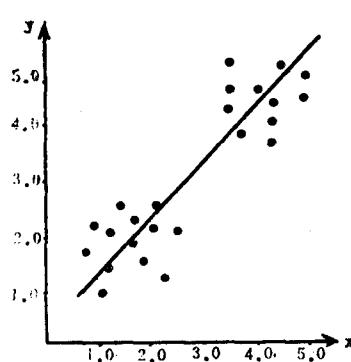


图3 点簇假相关示意图 对原始变量 x 和 y 作对数变换，是否造成假相关，不得一概而论，必须针对具体情况加以分析。现分以下几种情况论述

(1) 若原始变量 x 、 y 严格服从对数正态分布，则其相关系数 ρ_{xy} 和 $\rho_{\ln x, \ln y}$ 之间的理论关系为^[7]

$$\rho_{\ln x, \ln y} = \frac{1}{\sigma_{\ln x} \sigma_{\ln y}} \ln \left[\sqrt{(e^{\sigma_{\ln x}^2} - 1)(e^{\sigma_{\ln y}^2} - 1)} \rho_{xy} + 1 \right] \quad (9)$$

表 2 点簇假相关统计试验成果表

1										2		
总体参数										统计试验值		
\bar{x}_1	\bar{y}_1	\bar{x}_2	\bar{y}_2	C_x	C_y	C_{sx}/C_x	C_{sy}/C_y	γ_1	γ_2	$\hat{\gamma}_1$	$\hat{\gamma}_2$	γ_0
10	10	1000	1000	0.48	0.48	2	4	0.0	0.0	-0.006	-0.026	0.68
10	10	1000	1000	0.48	0.48	4	2	0.0	0.0	0.047	0.011	0.69
10	10	1000	1000	0.35	0.80	4	4	0.0	0.0	0.013	0.005	0.63
10	10	1000	1000	0.80	0.35	4	4	0.0	0.0	0.011	0.016	0.63
1	1	1000	1000	0.48	0.48	4	4	0.0	0.0	0.048	-0.054	0.69
100	100	1000	1000	0.48	0.48	4	4	0.0	0.0	0.048	-0.054	0.64

注：(1)对于变量 (x_1, \bar{y}_1) 和 (x_2, \bar{y}_2) 均分别模拟容量为40的样本50组，就每组计算相关系数，然后以50组平均得 γ_1 、 γ_2 ；

(2) 合并每两组系列计算相关系数，然后取平均得 γ_0 。

式中 $\sigma_{\ln x}$ 、 $\sigma_{\ln y}$ 分别为 $\ln x$ 、 $\ln y$ 的标准差。计算表明， $\rho_{\ln x, \ln y}$ 略大于 ρ_{xy} 。这就是说，在原始变量取对数后，其相关系数仅少许提高。如原始变量为微弱相关或不相关，通过对数变换一般不会显著提高相关系数，即不会造成明显的假相关。

(2) 若原始变量 x 、 y 互不相关且服从皮尔逊Ⅲ型分布，从这样的总体随机抽取样本，分别计算 $\hat{\rho}_{xy}$ 和 $\hat{\rho}_{\ln x, \ln y}$ 。统计试验表明， $\hat{\rho}_{\ln x, \ln y}$ 和 $\hat{\rho}_{xy}$ 相差不大，且 $\hat{\rho}_{\ln x, \ln y}$ 不一定都大于 $\hat{\rho}_{xy}$ 。表3是从大量统计试验中摘出的部分结果，该表清楚地说明，在所讨论的条件下，不会造成明显的假相关。

当然，由于对原始变量取对数，比例尺缩小（数量级减少），图上点据分散的程度看起来有所减小，给人一种假象，好象相关关系提高了。这单纯是肉眼上的错觉，若定量计算，变换前后的相关系数相差无几。

(3) 若原始变量 x 、 y 之间呈现出某种曲线相关关系，计算而得 $\hat{\rho}_{xy}$ 可能较小。此

却相当大，显然这是假的。设一对无关变量为 x_{11} 和 y_1 ，另一对无关变量为 x_2 和 y_2 ，如该两组变量的均值数量级相差很大，则往往引起点簇相关。

为了阐明这一点，进行了统计试验研究，结果列于表2。从该表可以明显看出，影响点簇假相关的因素是两组变量均值的差异。正是由于均值的显著差异性，才会出现分开的两组散乱点群，从而导致计算出的相关系数较大，造成假相关现象。

2. 对数变换假相关

表 3 对数转换部份统计试验成果

总体参数							样本容量	$\hat{\rho}_{xy}$	$\hat{\rho}_{\ln x \ln y}$
\bar{x}	\bar{y}	C_x	C_y	C_{sx}/C_x	C_{sy}/C_y	ρ_{xy}	n		
10	10	0.48	0.48	4	4	0	40	0.18	0.16
								0.29	0.15
								0.30	0.31
10	10	0.38	0.80	4	4	0	40	0.31	0.37
								0.39	0.30
								0.20	0.12
10	10	0.80	0.35	4	4	0	40	0.15	0.23
								0.18	0.15
								0.23	0.32
1000	10	0.48	0.48	4	4	0	40	0.29	0.23
								0.25	0.18
								0.19	0.23
10	10	0.48	0.48	2	2	0	40	0.32	0.31
								0.18	0.20
								0.33	0.35
10	10	0.48	0.48	2	2	0	40	0.11	0.16
								0.25	0.17
								0.19	0.23
1000	1000	0.35	0.35	2	2	0	40	0.28	0.23
								0.20	0.23
								0.55	0.45
10	10	0.90	0.30	4	4	0	40	0.35	0.41
								0.20	0.18
								0.30	0.35

时对原始变量作对数转换，计算而得的 $\hat{\rho}_{\ln x, \ln y}$ 可能相当大。对于这种情况，在我们看来不宜叫做假相关，因为原始变量之间确有关系，只不过这种关系并非线性而已。

综上所述，在一般情况下，对数变换不会导致明显的假相关。若变换后，相关系数大大提高，则原始变量可能存在着曲线相关关系。

水文计算工作中常常对变量取对数，这是无可非议的。要注意的是：若原始变量 x, y 的关系相当差（点据分散），通过对数变换缩小比尺使点据看起来比较密切，并以肉眼定出相关线，便会造成假相关。实际工作中应杜绝这种作法。

3. 辗转假相关

两变量 x, y 之间没有什么关系，但分别通过 x 和 y 同第三变量 z 相关，可能获得 x 和 y 之间的辗转关系。这似乎表明 x 中包含着 y 的信息，显然，这是假的。为了定量说明这种假相关的特性，我们作了辗转三次的统计试验研究，结果如表 4 所示。由该表清楚地看出：尽管 x 和 z , z 和 w , w 和 y 之间的相关系数均较高，而实际上 x 和 y 之间的相关系数却相当低。因此，就 x 和 y 两变量而言，这种辗转相关是假的。

表 4 辗转假相关统计试验成果表

总体参数						γ_{xz}	γ_{zw}	γ_{wy}	γ_{xy}
E_x	E_w	E_y	C_x	C_w	C_y	辗转相关系数值			相关系数
100	100	100	0.20	0.20	0.20	0.8	0.8	0.8	0.35
100	100	100	0.20	0.20	0.20	0.9	0.9	0.9	0.43
100	100	100	0.20	0.20	0.20	0.9	0.9	0.9	0.43
100	100	100	0.25	0.25	0.25	0.9	0.9	0.9	0.43
100	100	100	0.25	0.25	0.25	0.7	0.7	0.7	0.27

注 .. γ_{xy} 为容量是40的50个样本的平均值

辗转假相关的例子在水文计算中还是能见到的。例如，为了利用洪峰流量展延长历时的洪量，就可能误用辗转相关。松花江丰满站的资料表明，洪峰和历时十一日的洪量之间的相关系数仅为0.59。若由资料建立洪峰流量和三日洪量关系，三日洪量和七日洪量关系，进而七日洪量和十一日洪量关系，这些关系均较密切，通过它们由洪峰流量推求出的十一日洪量，就是典型辗转假相关的结果。

四、小 结

假相关是相关分析中出现的一种现象。它并非相关技术本身的问题，而是误用这种技术所引起的。若对两原始变量作变换或构成另一对新变量，则新变量之间的相关系数只能表明新变量本身之间的线性关系。在一般情况下，不能将新变量之间客观存在着的关系信息不加分析地用于原始变量。否则会导致假相关。

假相关可分为两大类型，一是隐假相关，二是显假相关。这两类假相关在水文计算中都会遇到，在假相关情况下获得的计算成果，显然是不可靠的，所以无论何时都应避免。

为了避免误用假相关，最关键的问题在于直接研究原始变量，就是说，不作任何数学转换和运算。

参 考 文 献

- [1]K. Pearson, On a form of spurious correlation which may arise when indices are used in the measurement of organs, Proc. R. Soc. London, 60, 1897
- [2]M. A. Benson, Spurious correlation in hydraulics and hydrology, J. Hydraul. Div. Am. Soc. Eng., 91, 1965
- [3]V. Yevjevich, Probability and Statistics in Hydrology Water Resources Publication, Fort Collins Co. U. S. A. 1972
- [4]C. T. Hearn, Statistical Method in Hydrology, The Iowa State University Press, 1977
- [5]金光炎, 用统计法计算设计洪水的若干问题, 《水文》, 1981, 第6期
- [6]B. C. Kenney, Beware of spurious self-correlation! Water Resources Research, Vol. 18, No. 4, 1982
- [7]J. R. Stedinger, Estimating correlation in Multivariate stream model, Water Resources Research, Vol. 17, No. 1 1981