

# Discrimination on Application of Principal Component Analysis in Environmental Assessment\*

Xiaoli Lu, Dingsheng Zhong

Environmental College, Jiangsu University, Zhenjiang  
Email: xiao0820li@126.com, zhongds@gmail.com

Received: May 4th, 2012; revised: May 21st, 2012; accepted: May 29th, 2012

**Abstract:** Principal component analysis (PCA) has been widely applied to environmental assessment. Taking the water environmental quality assessment as an example, whether the method can be applied to environmental quality assessment has been discussed. And the results showed that it was obviously wrong in many cases. We analyze the reasons why it was wrong in these cases from two aspects including the physical meanings and operation process of PCA. The conclusions were as follows: 1) PCA is a Mathematics method about numerical analysis of covariance essentially, but comprehensive evaluation is the judgment of the value of the evaluation system. When the numerical difference has not constituted value judgments on the comprehensive evaluation, it is wrong to apply PCA to environmental valuation forcibly; 2) The variance contribution rate only represents the variance information of the index system, not all of the information. It may cause the removal of the important index, if researchers select important indexes only according to variance contribution rate. We must consider other factors to select important index, such as, Natural and Geographical conditions, Social and Economic status and so on. So it can not be taken as a mainly yardstick to judge whether the value is high or low, and also, it is not appropriate to regard the variance contribution and the weight coefficient as the same in comprehensive evaluation; 3) PCA is a good mathematical method which has physical meaning itself, but we can't quote its physical meaning blindly. Whether a method could be applied to environmental quality assessment or not, if and only if the numerical difference of the evaluation index system constitutes on the value differ of things being evaluated. Specifically speaking, it is difficult to get a simple conclusion of when it is appropriate or not. We must make a detail analysis according to concrete problems.

**Keywords:** Environmental Assessment; Principal Component Analysis; Numerical Difference; Value Judgment; Variance Contribution

## 主成分分析法在环境评价中的应用辨析\*

卢小丽, 钟定胜

江苏大学环境学院, 镇江  
Email: xiao0820li@126.com, zhongds@gmail.com

收稿日期: 2012年5月4日; 修回日期: 2012年5月21日; 录用日期: 2012年5月29日

**摘要:** 主成分分析法在环境评价中已被广泛应用, 以水环境质量评价为例, 探讨该方法应用于环境评价的可行性。结果表明, 在不少情形下主成分分析法在环境评价中的应用明显有误。从物理意义和操作流程两个方面分析错误原因, 得出以下结论: 1) 主成分分析实质主要做协方差的数值分析, 综合评价则是对评价体系的价值判断, 当围绕指标体系的协方差所得数值差异不构成对该体系价值判断时, 强行应用主成分分析法进行环境价值评价是一种错误的做法; 2) 方差贡献率只代表指标体系的方差信息, 而非全部信息, 仅依据方差贡献率的大小判断信息量载荷大小的观点不妥, 仅根据方差贡献率的大小甄选指标可能会造成重要指标被错误剔除; 不能简单拿方差贡献率作为价值高低的评价尺

\*资助信息: 江苏大学高级人才专项资助(07JDG060)。

度,在综合评价中,将方差贡献率作为各指标的权重系数不恰当;3)主成分分析法本身是个好的数学方法,但不能盲目引申该数学方法的物理意义,当且仅当评价指标体系中各个指标的协方差所得数值差异构成对其所评价事物的价值判断时,才可能可以采用主成分分析法进行价值评价,具体何时恰当,何时不恰当,必须具体问题具体分析,难以简单下结论。

**关键词:** 环境评价; 主成分分析; 数值差异; 价值判断; 方差贡献率

## 1. 引言

在评价类的科学研究中,由于研究对象的复杂性和多面性,往往需要对其进行全面的、综合的定性或定量评价。环境评价即是对环境状况按照一定的标准和方法给予定性和定量的说明与描述<sup>[1]</sup>。其核心在于通过各种方法来量化地衡量环境质量的“价值”高低,得出其综合效用或综合水平,从而揭示环境的质量高低及其发展变化的特征。常用的环境评价方法有简单指数法、分级加权平均法、综合污染指数法、模糊数学法、普通概率统计法、主成分分析法等。主成分分析法自1901年美国统计学家 Pearson<sup>[2]</sup>首次将其应用于生物学理论研究中,后经 Hotelling<sup>[3]</sup>等学者进一步扩充完善和推广,成为一种全新的统计方法。由于该方法具有消除各指标不同量纲的影响,以及不受主观因素影响等优势,已被广大国内外学者应用于经济、社会、生态、环境等领域,现已成为环境评价中的常用方法之一。然而实践证明,主成分分析法在环境评价中的应用存在诸多问题,本文将对此做详细的分析。

## 2. 主成分分析法的原理和方法

### 2.1. 基本思想

本主成分分析法又称为主分量分析法,是利用降维的思想,把多个反映研究对象各方面信息、具有相关性的指标,利用数学变换的方法转变成几个不相关的新变量<sup>[4]</sup>,且新的指标体系仅用较少的新变量即可反映原指标的大部分信息。

### 2.2. 基本原理

借助正交变换,将分量相关的原始随机向量转变成分量不相关的新的随机向量。在代数上,表现为将原始随机向量的协方差阵变换成对角阵;在几何上,表现为将原坐标系变换成新的正交坐标系,使之

指向样本点散布最开的  $P$  个正交方向,然后对多维变量系统进行降维处理,使之由一个较高的维数转换成低维变量系统,再通过构造适当的函数,进一步把低维系统转化成一维系统<sup>[4]</sup>。

### 2.3. 基本步骤

主成分分析法的具体操作步骤<sup>[4]</sup>如下:1)根据评价对象,构建尽可能全面且合理的评价指标体系(设选取  $P$  个指标)。将指标正向化处理,即将逆向指标和适度指标转化为正向指标;2)将  $P$  项指标的原始数据标准化。设  $P$  项指标原始数据矩阵为  $\mathbf{X} = (x_{ij})_{m \times p}$ , 标准化矩阵为  $\mathbf{Y} = (y_{ij})_{m \times p}$ , 为了消除不同的量纲和数量级对评价结果的影响,普遍采用 Z-score 标准化公式:  $y_{ij} = \frac{1}{s_j}(x_{ij} - \bar{x}_j)$ 。其中,  $\bar{x}_j$  为  $j$  项指标的均值,  $s_j$  为  $j$  项指标的方差;3)计算  $P$  项指标的相关矩阵  $\mathbf{R}$ , 求  $\mathbf{R}$  的  $p$  个特征值  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$  及对应的正特征向量  $u_i = (u_{i1}, u_{i2}, \dots, u_{ip})$ ,  $i = 1, 2, \dots, p$ ; 4)提取主成分。设方差贡献率为  $\sigma_i = \lambda_i / \sum_{k=1}^p \lambda_k$ , 根据前  $m$  个主成分的累积贡献率  $\sum_{i=1}^m \sigma_i$ , 当其达到一定数值时(一般取  $\geq 85\%$ ), 取  $m$  个主成分  $F_i = u_{i1} \sum x_1 + u_{i2} \sum x_2 + \dots + u_{ip} \sum x_p$ ,  $i = 1, 2, \dots, m$ 。进而得到综合评价函数  $F = \sum_{i=1}^m \sigma_i F_i$ ; 5)根据综合评价函数算出每个评价单位的结果并对其进行排名和分析。

## 3. 主成分分析法的应用

近十几年,主成分分析法在环境综合评价中的应用非常普遍,国内外均有相关的应用研究,但以国内的应用最多,在应用中开发的方法变种也最多。

Susan<sup>[5]</sup>在可持续农业生态系统管理的土壤质量评估工具的设计中,应用主成分分析法进行指标选

取,以方差贡献率  $\geq 85\%$  作为准则,根据各个主成分的得分大小及总的相关系数值提取重要的指标。Ying Ouyang<sup>[6]</sup>利用主成分分析法对美国佛罗里达州的圣约翰斯河的水质进行了综合评价和站点筛选。C. Sarbu<sup>[7]</sup>等在经典主成分分析法的基础上,提出了模糊主成分分析法,并应用于多瑙河水质的评价。Uttam Kumar Mandal<sup>[8]</sup>等利用主成分分析法评价灌溉水对石灰性粘土的影响。鲍卫锋<sup>[9]</sup>等运用主成分分析法进行流域水环境等各方面的计算和评价。杜敏<sup>[10]</sup>、张继承<sup>[11]</sup>等认为将主成分分析法与其他各类评价方法相比较具有以下优势:1) 所确定的指标权数是基于数据分析而得到的指标间的内在结构关系,具有较好的客观性;2) 它能有效消除不相关指标的影响,从而可进行有针对性的定量化评价;3) 得出的综合指标主成分间相互独立,不仅简化了评价体系,而且减少了信息的交叉和冗余。

在主成分分析法被广泛应用的同时,也有很多学者意识到主成分分析法应用于综合评价存在很多缺陷。不少学者针对该方法进行改进和开发新的方法变种,比如阎慈琳<sup>[12]</sup>认为:只要第一主成分的方差贡献率  $\geq 75\%$ ,只取第一主成分评价效果更好;若第一主成分贡献率不够大,可作方差最大正交旋转等。姚焕玖<sup>[13]</sup>从对逆向指标进行线性变换、用对数化法对指标进行无量纲化及评价指标的选择这三个方面进行改进,其中论证了正向化中倒数法改变 Pearson 线性相关系数,从而导致特征值和特征向量的改变。洪素珍<sup>[14]</sup>对阎慈琳提出质疑并论证得出:1) 即便第一主成分的贡献率达到了某一定值,也未必能用于综合评价,要想第一主成分能用于做综合评价,则第一主成分做综合评价的得分值必须是有序的;2) 主成分能否用于做综合评价与其方差贡献率大小无关,采用正交旋转同样不可行,因为它改变了原始数据的线性关系。

也有学者对该方法的某些应用进行了否定。徐龙封<sup>[15]</sup>提出该方法评价过程中并未最大限度地发掘样本反映的质量信息,也未排除评价指标间的相关重叠性;冯利华<sup>[16]</sup>举例论述了该方法对环境质量的分辨能力不高等局部失效问题;俞立平<sup>[17]</sup>等就期刊评价中应用此方法的失效问题进行研究,他认为,虽然主成分分析法应用中提供了一定的统计检验方法,但是这种检验不可靠,要采用该方法,除了同类可比及它本身

提供的统计检验外,必须添加一个条件,即其结果与所有评价指标回归后的指标系数为正。

从上述文献可看出,不少学者已经意识到主成分分析法应用于综合评价存在不足,并进行了改进。尽管有些改进措施的确能降低主成分分析法失效的可能性,然而,笔者认为如果该方法的适用性如此狭窄,那么它在应用于环境综合评价中就可能是存在重大缺陷的,甚至可能是错误的。笔者通过进一步深入研究,发现许多情形下的应用确实是错误的。下文将从主成分分析法的物理意义和操作流程两方面入手分析其在环境评价中失效的原因。

## 4. 环境评价中主成分分析法的应用辨析

### 4.1. 围绕协方差所得数值差异不等于价值差异

主成分分析法作为一种统计方法,是一种简化数据集的技术。它通过线性变换,将数据变换到一个新的坐标中,使得样本数据投影的第一大方差在第一个坐标,第二大方差在第二坐标上,依次类推。主成分分析法被应用于综合评价时,第一坐标被称为第一主成分,第二坐标被称为第二主成分,然后根据主成分的特征向量及各主成分贡献率线性组合成综合评价函数。有不少文献将该方法应用于环境综合评价,以主成分分析的计算结果当作衡量环境综合评价中价值高低的依据,对此做法,本文提出质疑,质疑的核心理由是:本文认为评价体系中指标的协方差所得数值差异不等于价值差异,主成分分析法虽然是一个好的数学方法,但应用于多数环境评价中,主成分分析法的计算结果都不具备对被评价对象的价值高低(即优劣程度)进行评价的物理意义。为充分说明这一问题,本文通过以下几个例子进行论证。

比如有些研究者将主成分分析法应用于南淝河流域水环境质量<sup>[18]</sup>,该文章将5个断面的原始数据与《地表水环境质量标准(GB3838-2002)》中II类~V类水体的标准值结合,应用主成分分析法得到5个断面和四类水体的综合得分,评价出5个断面的水质情况。然而,就5个断面原始数据本身而言,仅以TN这个指标即可判断屯溪路桥、合钢二厂下游、板桥码头、施口小学四个断面均属于劣V类水质,而该文应用主成分分析法综合分析得到的结果却是屯溪路桥、合钢二厂下游、板桥码头为V类水质,施口小学为IV类水质。

显然,该方法未能准确判断各断面水质污染的程度,分析结果严重不合理。与该论文的问题类似的研究还有不少,在此不再一一剖析。为更好地论证本文的核心观点:在大多数环境评价中,主成分分析法的计算结果都不具备对被评价对象的价值高低进行评价的物理意义,即评价体系中指标的协方差所得数值差异不等于价值差异,本文构造了以下两个更为简洁明了的例子进一步进行论证。

案例 1: 本文参照地表水环境质量标准(GB3838-2002)中各级标准,构造了以下数据情景(采用构造数据的方式举例,理由是: 1) 主成份分析法用在环境评价中,并非只在河流湖泊的水质评价中有应用,它也同样可以被拿来对水样的水质评价,任何水样,只要可以提供(或配制)出来,就可以拿来对水质评价; 2) 本文所构造的数据从物理和化学的角度来说,没有理由会导致这种水样配制不出来。下同此理,见表 1。

应用 SPSS16.0 软件包,根据主成分分析步骤,对数据进行 Z-score 标准化,计算各项指标的相关系数矩阵  $R$ , 求取  $R$  的特征值及对应的特征向量,提取前两个主成分的方差贡献率  $\sum_{i=1}^2 \sigma_i \geq 98.64\%$ , 综合评价结果对照单因子评价结果见表 2。

表 2 中由单因子评价法可得: 1 区明显重金属超标, 属劣 V 类水质, 2 区、3 区、4 区、5 区分别为 II

类、III类、IV类、V类水质。主成分分析结果却显示: 5 区比 3 区 4 区反而更好, 与单因子评价法的结果相反, 这显然是一个难以被接受的评价结果, 出现这种现象, 要么是主成分分析法的结果有问题, 要么是单因子评价法的结果有问题, 二者必居其一。那么, 究竟该如何来评判二者的对错呢? 本文认为, 尽管单因子评价法并非完美, 但作为一个物理意义简单、明确和相对粗略的评价方法, 也早已被有关领域广泛认同和采用, 是一个常用的经典方法, 其结果的合理性和科学性毋庸置疑, 即使其评价结果有瑕疵(尤其是在做差别较小的样本的区分时, 结果容易迟钝甚至丧失细节上的分别能力), 这也仅仅是细节上的技术性问题, 而非根基上的对错问题, 单因子评价法的评价结果若有大问题, 除非是其依据的环境质量标准设置出了大问题, 不同指标的环境质量标准在同一级别上对环境的影响程度迥异。相比起来, 主成分分析法在环境评价中的意义却非常堪疑, 最为重大的疑点/缺陷在于, 该方法至始至终均未引入各个环境指标的环境影响程度的判别数据, 仅仅是就各指标的协方差所得数值差异情况进行的数学分析, 尽管这种分析本身是有一定的物理意义的, 但与实质要做的工作——环境质量的优劣评价——有明显的差别。而为何主成分分析法无法或难以直接引入各个环境指标的环境影响程度的判别数据, 其根源与所研究对象的线性与非线性这个

Table 1. Performance about 9 environmental indicators of 5 regions (unit: mg/L)  
表 1. 5 个区的 9 个环境指标状况(单位: mg/L)

地区	COD	BOD <sub>5</sub>	砷	锌	六价铬	铅	铜	总氮	总磷
1	15.00	3.00	1.00	3.00	1.00	2.00	3.00	0.20	0.01
2	15.00	3.00	0.05	0.05	0.01	0.01	0.01	0.50	0.10
3	20.00	4.00	0.05	1.00	0.05	0.05	1.00	1.00	0.20
4	30.00	6.00	0.10	1.00	0.05	0.05	1.00	1.50	0.30
5	40.00	10.00	0.10	2.00	0.05	0.10	1.00	2.00	0.40

\*数据解释: 关于是否需要数据进行数据的正向化, 有些学者提出不同的观点<sup>[13]</sup>, 不过, 笔者认为正向化并不是主成分分析法失效的根源所在, 关于这一点, 限于篇幅, 本文暂不进行深入探讨。为简化计算, 本文所选指标均为越小越好型指标, 故不需要再做正向化。

Table 2. Results comparison of PCA and evaluation of single index about 5 regions  
表 2. 5 个区主成分分析与单因子分析结果对比

地区	第一主成分得分值	第二主成分得分值	主成分综合得分	主成分分析结果排名	单因子评价结果
1	4.19568	0.88307	3.10238	5	劣 V 类
2	-0.30012	-2.33324	-0.93197	1	II 类
3	-0.43436	-0.85873	-0.56118	4	III 类
4	-1.29984	0.24919	-0.79761	3	IV 类
5	-2.16136	2.059715	-0.81162	2	V 类

问题也有关,关于这一点,笔者将在后文进行进一步论述。

基于上述分析,笔者认为主成分分析的结果存在问题,出现了明显的优劣判别的错误。

案例 1 中所反映的问题并非笔者第一个发现,已有不少学者在研究中发现了类似的问题,如徐龙封<sup>[15]</sup>、冯利华<sup>[16]</sup>和南英子<sup>[19]</sup>等等,但在为何出现这种问题以及该如何解决这个问题上,笔者并不认同他们的观点。比如,徐龙封<sup>[15]</sup>针对主成分分析应用于综合评价中所得综合评价函数的多种情况做了详细的分析,认为如果要采用第一主成分做综合评价,样本点必须是占绝对优势的流向,这样的评价才可能是有效的。本文并不认同这一观点,理由是用主成分分析法来进行评价,其结果有效的前提是主成分分析的分析结果具有对被分析事物的优劣进行评判的物理意义,至于样本是否占有绝对优势,这一点并不影响主成分分析本身的数学意义,更不会改变主成分分析结果的物理意义。与此同理,采用变换法的“改进”形式仅仅只是通过数学变换,使得样本的分布和流向看起来更漂亮,但同样不能改变该方法的数学意义和物理意义,恰恰相反,仅仅通过变换法就使得样本集的分布和流向发生巨大变化这一事实说明样本的分布和流向是可以通过各种物理意义不明确的数学变换方法来进行剧烈改变的,这就更说明样本的分布信息和流向信息与我们要进行的优劣评价是两回事。冯利华<sup>[16]</sup>就主成分分析应用于环境评价存在的诸多失效问题进行了详细的举例,比如主成分分析对环境质量的分辨能力不高,再比如污染物超标倍数增大,所得评价结果环境质量等级反而更好,但遗憾的是,该文仅仅

发现了主成分分析在环境质量评价中的失效现象,对于为何会出现这些现象的本质原因并未详细论述,而这个问题正是本文论述的核心内容。

此外,有学者提出通过 KMO 检验(即样本充分度检验)和 Bartlett 检验(即球形检验)来解决主成分分析法的失效问题。比如南英子<sup>[19]</sup>在其论文中论述了主成分分析法进行 KMO 检验和 Bartlett 检验的必要性。对此,我们认为,通不通过检验都无关紧要,是否如此,可以通过案例 2 来进行论证。

案例 2: 本文参照地表水环境质量标准(GB3838—2002)中各级标准,构造了另一个数据情景,见表 3。

数据先进行 KMO 检验和 Bartlett 检验,检验结果 KMO 值为  $0.714 > 0.5$ , Bartlett 值为 109.734,  $P < 0.000$ ,通过了统计检验,再根据主成分分析步骤,提取前两个主成分的方差贡献率  $\sum_{i=1}^2 \sigma_i \geq 88.00\%$ , 结果见表 4。

与例 1 类似,由单因子评价法可得: 2 区属 I 类水, 3 区属 II 类水, 6 区属 III 类水, 8 区属 IV 类水, 1 区铅含量超标, 4 区水质中锌指标含量超标, 7 区砷指标含量超标, 9 区重金属指标和氮、磷指标均严重超标, 10 区氮磷指标含量超标, 均属劣 V 类水; 然而, 按照主成分分析法却会得出这样的结论: 属劣 V 类水的 1 区排在了第一, 属 IV 类水的 8 区却排在倒数第四, 这个结果显然错误。由此可见, 对评价而言, 应用该方法时通不通过 KMO 检验和 Bartlett 检验都无关紧要, 并没有解决问题的根本, 通过了检验并不意味着它的物理意义就符合环境评价的物理意义, 就能够拿来作价值判断。

Table 3. Performance about 9 environmental indicators of 10 regions (unit: mg/L)  
表 3. 10 个区的 9 个环境指标状况(单位: mg/L)

地区	COD	BOD5	砷	锌	六价铬	铅	铜	总氮	总磷
1	7.00	1.00	0.01	0.05	0.01	2.00	0.01	0.01	0.01
2	14.00	2.00	0.01	0.05	0.01	0.01	0.01	0.10	0.01
3	15.00	3.00	0.05	1.00	0.01	0.01	1.00	0.05	0.02
4	15.00	3.00	0.05	15.00	0.01	0.01	1.00	0.05	0.02
5	18.00	3.50	2.34	2.56	2.31	2.84	3.15	0.80	0.18
6	20.00	4.00	0.05	1.00	0.05	0.05	1.00	1.00	0.20
7	25.00	5.50	0.85	2.00	0.03	0.02	1.00	1.26	0.25
8	30.00	6.00	0.10	2.00	0.05	0.05	1.00	1.50	0.30
9	40.00	10.00	5.26	20.15	2.56	6.15	7.81	6.54	6.98
10	28.00	10.00	0.20	0.30	0.45	0.50	0.42	5.46	9.86

Table 4. Results comparison of PCA and evaluation of single index about 10 regions  
表 4. 10 个区主成分分析与单因子分析结果对比

地区	第一主成分得分值	第二主成分得分值	主成分综合得分	主成分分析结果排名	单因子评价结果
1	-1.93932	-0.88779	-1.50915	1	劣V类
2	-1.91404	-0.21346	-1.36511	2	I类
3	-1.57986	-0.20585	-1.13233	3	II类
4	-1.02971	-0.70383	-0.84492	4	劣V类
5	0.83606	-1.75104	0.25019	8	劣V类
6	-1.12487	0.21809	-0.73780	5	III类
7	-0.55114	0.46094	-0.29505	6	劣V类
8	-0.43704	0.82634	-0.14748	7	IV类
9	6.38801	-0.77742	4.27640	10	劣V类
10	1.35190	3.03402	1.50526	9	劣V类

此外,我们还可以从单个地区比较看出主成分分析法的评价结果不恰当。比如,4区、6区的评价结果与单因子评价结果对比一下我们不难发现,两种方法评价结果完全对立,原因何在?由案例1分析中我们认为,单因子评价法一般不会出现根基上的对错问题,除非是其依据的环境质量标准设置出了大问题,否则,即使出问题,也应只有小的不够精细的问题,相比之下,最为可能是主成分分析评价结果出了问题。那么,案例2是否也是如此呢?我们可以肯定一个简单事实:劣V类水质不可能优于III类水质,一个重金属污染严重的地区其综合环境质量不应该高于III类水质地区环境质量,所以,肯定是主成分分析法出现问题。那为何主成分分析会出现这样的结果呢?结合主成分分析步骤和基本原理可知:主成分分析是在保证总方差不变的情况下进行线性变换,使第一大方差在第一坐标,被称为第一主成分,第二大方差在第二坐标,被称为第二主成分,依次类推。这样做主观上已经认定方差信息就是指标的所有信息,而实际上,方差只是表示一批数据波动的大小,反映不了指标的所有信息,这就直接导致主成分提取的信息与实际评价存在偏差。这也说明主成分分析仅是一种数学工具,是根据协方差所得数值差异来进行的机械化的数学分析,无法充分考虑各个指标本身相对重要程度及其真实的物理意义。环境质量评价则不然,它是对被评价事物的价值进行优劣判断,从而揭示出其中的发展规律。显然,在环境质量评价中,主成分分析法这种纯粹的数学分析的物理意义与被分析事物的优劣进行价值评判的物理意义不相符合。

## 4.2. 根据累计方差贡献率来压缩信息值得商榷

主成分分析应用中根据方差贡献率的大小提取主成分,一般提取累计方差贡献率大于85%的主成分,得到原始数据的绝大部分信息,从而达到压缩数据维数,简化模型的目的。由于其简便性,该方法的应用形式多样化。然而,“根据累计方差贡献率来压缩信息”这一做法是否恰当,笔者将在下文逐一分析。

### 4.2.1. 方差贡献率不能作为提取信息量多少的唯一准则

主成分分析中所谓的信息就是指标的变异性,通常用标准差和方差表示<sup>[20]</sup>。环境评价中的信息则涵盖较广,包括方差信息、各个指标的毒理性程度和机理、指标与指标的作用关系等。然而主成分分析法不少应用中,仅仅以方差贡献率大小代表各指标所含信息量的多少,提取的主成分以所占方差贡献率大小为依据,贡献率大就取,反之,舍弃。实际上,评价中所得的方差贡献率仅是提取指标体系的方差信息。方差,通俗点讲,就是用来衡量一批数据波动大小(即这批数据偏离平均数的大小),而这并不能全面反映指标间信息(事实上,在经过线性变换后,各个指标的方差的大小可以被显著改变)。如果仅以方差贡献率大小选取指标信息,极有可能会造成重要信息的遗漏。因此,仅以方差贡献率提取信息量的做法不恰当。

### 4.2.2. 根据载荷量剔除体系中不重要的指标不妥

有文献利用主成分载荷矩阵进行指标筛选,认为该矩阵能够反映各个指标在各个主成分上的重要程度,可作为指标重要性的评判依据。对此做法,笔者

同样认为值得商榷，主成分载荷矩阵在一定程度上确实具有选择的功能，可以对指标的重要与否起到一定的借鉴作用。但应该具体情况具体分析。我们知道：载荷量是由特征值的平方根乘以特征向量得到。倘若主成分分析应用于指标体系的物理意义尚不符合，由此得到的特征值物理意义就更加不明确了，所得的载荷量的物理意义也就更加模糊，如此执意而为，只会造成指标筛选有误。为充分说明，举以下例子进行论证。

比如运用该方法对评价指标体系进行筛选优化<sup>[21]</sup>，该指标优化运用特征向量的大小作为筛选指标的标准，选取特征向量大、去除特征向量小的指标，从而达到优化指标的目的。为了便于说明问题，本文直接用表3数据，采用该方法<sup>[21]</sup>的变种——载荷量——来进行指标体系指标的筛选，分析得到主成分初始因子载荷矩阵图(见图1)。

如果初始因子载荷矩阵代表了每个指标分别在各个主成分上的载荷量，反映了每个指标在主成分上的重要程度，那么第一主成分上，COD、BOD<sub>5</sub>、砷、六价铬、铅、铜、总氮都有较高载荷量(80%以上)，因此第一主成分就可以用这几个指标来表征。同理，可判断第二主成分上表征的指标。

按照提取第一主成分中载荷量 80%和第二主成分中载荷量 60%的方法选取指标，锌指标将被删除。事实上，锌指标属于重金属污染指标，它的重要性不亚于铜等指标，是衡量4区污染状况的关键指标。删除该指标，4区重金属污染状况会被忽略，这显然是一个不可接受的重大错误。筛选指标的目的在于将不重

要的指标筛选出局，但是绝不允许出现重要的指标被错误删除的情况出现。所以，如何有效的应用主成分分析法来进行指标甄选，清楚的了解所研究对象的指标的物理意义与主成分分析法的物理意义之间的联系和区别是极为重要的前提。

在应用主成分分析法进行指标筛选的研究中，还有一类这样的变种——应用主成分分析法进行样本筛选，对于这一类的应用，本文同样持高度怀疑的态度。本文以主成分分析法对美国佛罗里达州的圣约翰斯河的水质进行综合评价<sup>[6]</sup>为例进行分析。该文采用主成分分析法进行监测站点的优化筛选。对于该文的这一做法，本文持质疑的观点。本文认为，决定一个监测站点是否应该保留或删除，首先应该考虑该站点所处位置的自然地理和社会经济状况。比如，如果该站点位于饮用水取水口、位于重要鱼类繁殖区等等，则不论该站点的水质变化是否剧烈、对全河流的水质变化的方差贡献率是否大，都不应被剔除而应被优先设置和保留。类似的可能还有其他需要优先考虑的因素。只有当这些需要优先考虑的因素都已经被考虑过了之后，仍然需要剔除一些站点，才可以考虑用主成份分析的方法来进行进一步的优选。

#### 4.2.3. 综合评价函数的线性系数不能作为指标的权重

不少学者将主成分分析得到的综合评价函数的线性系数作为加权处理的权重系数，这些权重系数反映了各个变量在综合评价中的相对重要性。从该方法的基本操作可知，综合评价函数的线性系数是由方差贡献率与特征向量的乘积所得，是由数学变换的过程

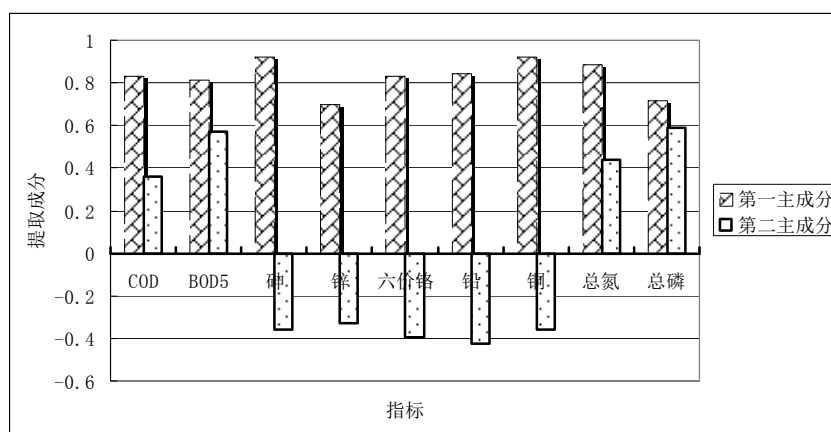


Figure 1. Component matrix picture of PC about 10 regions  
图 1. 10 个区域主成分因子载荷图

中生成的，不能进行人为调整。这表面上看起来比较客观，但是实际上主成分分析对原始变量与分量的关系都是由线性关系处理，这样往往导致与现实关系的偏差，两者在实际评价中不可混为一谈。事实上，实际应用中往往可能会碰到两者重要性相当但线性系数相差巨大，或者两者的重要性大小与其系数大小正好相反等异常情况。这一点通过例 1 就可以得到验证。

例 1 中 5 个区的 9 个指标，得到相关系数矩阵后，求得如下特征值：

$$\lambda_1 = 6.062, \lambda_2 = 2.815, \lambda_3 = 0.093, \lambda_4 = 0.03.$$

而前两个主成分的方差贡献率已达到 98.64%，则选取两个主成分为：

$$F_1 = -0.290x_1 - 0.272x_2 + 0.379x_3 + 0.239x_4 + 0.384x_5 \\ + 0.383x_6 + 0.328x_7 - 0.334x_8 - 0.354x_9$$

$$F_2 = 0.417x_1 + 0.431x_2 + 0.204x_3 + 0.477x_4 + 0.193x_5 \\ + 0.195x_6 + 0.331x_7 + 0.322x_8 + 0.288x_9$$

综合评价函数为：

$$F = -0.063x_1 - 0.049x_2 + 0.319x_3 + 0.310x_4 + 0.318x_5 \\ + 0.319x_6 + 0.324x_7 - 0.128x_8 - 0.148x_9$$

可以看出： $x_1$ COD， $x_2$ BOD<sub>5</sub> 的系数值远远小于  $x_8$  总氮， $x_9$  总磷的系数值。然而，从环境保护的长期实践经验来看，COD、BOD<sub>5</sub> 指标的重要性至少与总氮、总磷相当，甚至在环境保护开展的早期阶段中，COD、BOD<sub>5</sub> 指标重要性远大于总氮、总磷，也是污染处理的首要对象之一，只是近年来，伴随富营养化的出现和人们对环境质量要求的不断提高，总氮、总磷指标的重要性才开始突显，但还是远达不到其重要性远超 COD、BOD<sub>5</sub> 的状况。进一步对线性系数与相关矩阵分析可知，综合评价函数中的线性系数只由数值间的相关程度决定，但数值间相关程度并不代表指标间的重要程度。在实际的研究中，极有可能出现某个指标与其他指标的相关性很低，但是这个指标却是最重要的关键性指标的现象。因此，由以上分析可以断定，函数综合评价中的系数作为指标权重具有不合理性。

#### 4.3. 小结

上文的分析表明，主成分分析法的评价结果经常会出现明显的错误，其错误的根源是错误地引申了主成分分析法的物理意义。水环境质量评价的本质在于

价值判断，这种“价值”判断，以笔者目前的经验来看，主要有两大类方法：1) 通过特殊的指标(如经济价值、产量、等标污染负荷等)进行直接度量；2) 通过专家经验、问卷调查、头脑风暴等方式进行度量，从而得到具有物理意义的“价值”量用以评价和对比。而主成分分析法并未引入各个环境指标的环境影响程度的判别数据，仅仅就各个指标的协方差所得数值差异进行数值分析，这种数值分析与被分析事物的优劣分析具有物理意义上的不同。主成分分析法本身是个好的具有一定物理意义的数学方法，但应用过程中不能盲目引申该数学方法的物理意义，只有当评价指标体系中各个指标的协方差所得数值差异构成对其所评价事物的价值判断时，才可能可以采用主成分分析法进行价值评价(优劣评价)，至于具体何时恰当，何时不恰当，必须具体问题具体分析，难以简单下结论。

其次，从方法论的角度分析，也可以大致判断出直接将主成分分析法用于环境质量优劣评价不恰当，理由如下：主成分分析法是一种线性分析方法，提取的相关系数矩阵也只是反映指标间的线性相关程度。它不具备对非线性指标体系的分析能力，除非事先已经对该非线性系统进行了非常有针对性且恰当的线性简化。而水环境系统是一个由多指标组成的复杂非线性系统，比如重金属等指标，它们逐级分类，每个等级对环境的危害越来越大，其具体数值与其环境危害呈典型的非线性关系。从这个角度来说，用一个线性分析方法去分析一个非线性系统，而且在分析的过程中并未考虑该非线性系统线性简化的做法，不恰当。

## 5. 结论和探讨

主成分分析法被认为是一种简单易行的评价方法，它以少数的综合变量取代原有的多维变量，使数据结构大为简化，避免主观随意性而得到了广泛的推崇。但大量的事实和模拟分析均表明，主成分分析法在环境评价中的应用存在着诸多问题，学者们众说纷纭。本文通过举例论证了主成分分析法在相当多的情形下应用于环境评价是失效的，并分析了失效的原因。进而指出了方差贡献率作为提取信息的唯一准则是不可行的；运用载荷量剔除不重要指标存在不合理性；方差贡献率不可作为价值评价权重，由此确定的



综合评价函数的线性系数作为指标权重系数同样存在类似的问题。总之，当体系中指标的协方差所得数值差异不构成对该体系的价值判断时应用主成分分析法是失效的，当且仅当体系中各个指标的数据差异构成对所研究对象的价值判断时，采用主成分分析法进行事物的优劣评价才有可能有效。

如何引申一个数学方法的物理意义是一个需要慎重处理的科学问题，盲目引申极有可能导致大面积的错误应用。本文并非要否认主成分分析法在环境评价中的应用，只是认为主成分分析法在许多情形下的应用有误，对于主成分分析法的适用条件还有待于进一步的研究。

## 参考文献 (References)

- [1] 李祚泳, 丁晶, 彭荔红. 环境质量评价原理与方法[M]. 北京: 化学工业出版社, 2004: 2-35.
- [2] K. L. Pearson III. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine Series 6*, 1901, 2(11): 559-572.
- [3] H. Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, 1933, 24(6): 417-441.
- [4] 何晓群. 现代统计分析方法与应用[M]. 北京: 中国人民大学出版社, 1998: 281-315.
- [5] S. S. Andrews, C. R. Carroll. Designing a soil quality assessment tool for sustainable agroecosystem management. *Ecological Applications*, 2001, 11(6): 1573-1585.
- [6] Y. Ouyang. Evaluation of river water quality monitoring stations by principal component analysis. *Water Research*, 2005, 39(12): 2621-2635.
- [7] C. Sarbuu, H. F. Pop. Principal component analysis versus fuzzy principal component analysis: A case study: The quality of danube water (1985-1996). *Talanta*, 2005, 65(5): 1215-1220.
- [8] U. K. Mandala, D. N. Warringtonb, A. K. Bhardwaj, et al. Evaluating impact of irrigation water quality on a calcareous clay soil using principal component analysis. *Geoderma*, 2008, 144(1-2): 189-197.
- [9] 鲍卫锋, 黄介生, 孔祥元. 基于主成分分析法的流域水循环效应[J]. *武汉大学学报(工学版)*, 2007, 40(2): 29-33.
- [10] 杜敏. 基于主成分分析法的环境质量综合指数研究[D]. 四川大学, 2006.
- [11] 张继承. 基于RS/GIS的青藏高原生态环境综合评价研究[D]. 吉林大学, 2008.
- [12] 阎慈琳. 关于用主成分分析做综合评价的若干问题[J]. *数理统计与管理*, 1998, 17(2): 22-25.
- [13] 姚焕玫, 黄仁涛, 甘复兴等. 用改进的主成分分析法对东湖的水质污染进行评价[J]. *武汉大学学报(信息科学版)*, 2005, 30(8): 732-735.
- [14] 洪素珍. 如何有效利用主成分分析中的主成分[D]. 华中师范大学, 2008.
- [15] 徐龙封. 对主成分分析法应用的思考[J]. *中国统计*, 1994, 5: 39-40.
- [16] 冯利华. 主成分分析在环境质量评价中的失效问题[J]. *数学实践与认识*, 2005, 35(6): 12-16.
- [17] 俞立平, 潘云涛, 武夷山. 学术期刊评价中主成分分析法应用悖论研究[J]. *情报理论与实践*, 2009, 32(9): 84-87.
- [18] 盛周君, 孙世群, 王京城等. 基于主成分分析的河流水环境质量评价研究[J]. *环境科学与管理*, 2007, 32(12): 172-175.
- [19] 南英子. 实证分析中运用主成分分析法应注意的几个问题[J]. *统计与决策*, 2009, 21: 155-156.
- [20] 刘大海, 李宁, 晁阳. SPSS15.0 统计分析从入门到精通[M]. 北京: 清华大学出版社, 2008: 157-268.
- [21] 苏喜军. 基于主成分分析法的区域创业环境评价指标体系研究[J]. *河南社会科学*, 2010, 18(3): 85-87.