

主成分分析与因子分析的 异同和 SPSS 软件

——兼与刘玉玫、卢纹岱等同志商榷

林海明 张文霖

ABSTRACT

Between the Principal Component Analysis and the Factor Analysis. This paper puts forward the difference and the Identity, which advances positive proposal to some users of this two methods essentially.

关键词：主成分分析；因子分析；混淆；出错；避免

设 $X = (X_1, \dots, X_p)'$ 为标准化的随机向量 ($p \geq 2$), R 为相关系数矩阵, $F_m = (F_1, \dots, F_m)'$ 为主成分向量, $Z_m = (Z_1, \dots, Z_m)'$ 为因子向量, $m \leq p$, 为方便, 因子、因子估计、因子得分用同一记号。

一、问题的提出

主成分分析与 R-型因子分析是多元统计分析中的两个重要方法, 同是降维技术, 应用范围十分广泛, 但通过流行甚广的 SPSS 软件调用这两种方法的过程命令, 有些使用者容易出现混淆性错误, 如《统计研究》2003 年第 12 期发表的论文《经济全球化程度的量化研究》(以下称《刘文》)、电子工业出版社 2002 年 9 月出版的《SPSS for Windows 统计分析(第二版)》(以下称《卢书》)就是这种情况。是什么原因造成这些错误呢? 主成分分析与 R-型因子分析到底有何异同呢?

经过对一些论文和一些 SPSS 软件教科书仔细分析、比较我们发现出错的主要原因在于有些使用者和 SPSS 软件教科书作者对怎样用 SPSS 软件得出主成分分析与 R-型因子分析的结果掌握不全面, 对主成分分析与 R-型因子分析异同的认识不透彻。

经过仔细查证出现的错误有:

使用主成分分析时: ①叙述主成分分析概念出错; ②主成分 F_i 求解出错, 如 $F_m = A'_m X$ 中 $A'_m A_m \neq I_m$ (I_m 为单位矩阵, A_m 的意义见表 1); ③找不到主成分 F_i 的命名依据, 对主成分 F_i 命名出错; ④某变量 X_k 被丢失; ⑤对 A_m 错误地进行旋转; ⑥错误地进行回归求 F_i ; ⑦错误地把因子分析法(含初始因子分析法)当作主成分分析法。

使用因子分析时: ①将因子分析的思想叙述为主成分分析的思想; ②因子 Z_i 的命名出错, 如用因子得分函数对因子 Z_i 进行命名; ③某变量 X_k 被丢失; ④将主成分或因子错误地表示为 $B'_m X$ (B_m 的意义见表 1); ⑤不知相关系数矩阵特征值 λ_i 与因子贡献 v_i 的区别, 如综合因子得分函数 $Z_{\text{综}} = \sum_{i=1}^m (v_i/p) Z_i$ 中的 v_i 错误地取为特征值 λ_i 。

二、主成分分析与 R-型因子分析数学模型的异同比较

相同之处: 主成分分析与 R-型因子分析都是对协差阵的逼近, 都是打算降维解释数据集。具体为指标的正向化, 指标的标准化(SPSS 软件自动执行), 通过相关系数矩阵判断变量间的相关性, 求相关系数矩阵的特征值和特征向量, 主成分间、因子间线性无关, 用累计贡献率 ($\geq 85\%$)、变量不出现丢失确定主成分、因子个数 m , 前 m 个主成分与前 m 个因子对 X 的综合贡献相同、是最大化的, 命名依据都是主成分、因子与变量的相关系数。

不同之处: 方差, 最大化方向, 所处的坐标系(标准正交性), 应用上侧重等不同见表 1。

主成分分析与因子分析定量上不同的显著性标志是方差。事实上, $\text{Var}F_i > (<) \text{Var}Z_i = 1$, 即 F_i 的取值范围比 Z_i 的取值范围大(小); 通常 $\text{Var}F_{\text{综}} > \text{Var}Z_{\text{综}}$, 即 $F_{\text{综}}$ 的取值范围比 $Z_{\text{综}}$ 的取值范围大, 这些都肯定了主成分分析与因子分析的计量值、评价体系不同。

结论: 主成分分析与因子分析两种方法方差、最大化

表 1 主成分分析与 R—型因子分析的不同

区别项目	主成分分析数学模型:	R—型因子分析数学模型:
表达式与系数矩阵	$F_m = A' mX, A_m = (a_{ij})_{p \times m} = (\alpha_1, \alpha_2, \dots, \alpha_m), R\alpha_i = \lambda_i \alpha_i, \lambda_i, \alpha_i$ 是相应的特征值和单位特征向量, $\lambda_1 \geq \dots \geq \lambda_m \geq 0$ 。	$X = B_m Z_m + \epsilon$ (ϵ 为特殊因子), 因子载荷矩阵 $B_m = (b_{ij})_{p \times m} = B_m C, B_m = (\sqrt{\lambda_1} \alpha_1, \sqrt{\lambda_2} \alpha_2, \dots, \sqrt{\lambda_m} \alpha_m)$ 为初始因子载荷矩阵 (λ_i, α_i 同左)。
因变量方差最大化	F_i 依次达到信息贡献(方差)最大化, $\text{Var} F_i = \lambda_i$ 。	Z_i 没有达到方差最大化, $\text{Var} Z_i = 1$ 。
矩阵方差最大化旋转	无, 旋转后就不是主成分了, 因为 $\text{Var} F_i \neq \lambda_i$ 。	有, $C = (c_{ij})_{m \times m}$ 为 B_m 方差最大正交旋转矩阵, B_m 达到方差最大化。
因变量对 X 的贡献	特征值 λ_i 。	$v_i = \sum_{k=1}^p b_{ki}^2, v_i \neq \lambda_i$, 通常 $\lambda_1 > v_i$ 。
相关系数	$r_{X_i F_j} = \sqrt{\lambda_j} a_{ji}$ 。	$r_{X_i Z_j} = b_{ji}$ 。
命名依据	用 $\sqrt{\lambda_j}(a_{1j}, \dots, a_{pj})$ 式中系数绝对值大的对应变量对 F_j 命名, 有时命名清晰性低。	将 B_m 的第 j 列绝对值大的对应变量归为 Z_j 一类并由此对 Z_j 命名, 命名清晰性高(精细)。
回归过程	无。	有, 因子得分函数 $Z_m = B_m' R^{-1} X$
标准正交性	是, $A_m' A_m = I_m$ (判据之一)。	非, $B_m' R^{-1} (B_m' R^{-1})' \neq I_m, B_m' B_m \neq I_m$
综合评价函数及方差	$F_{\text{综}} = \sum_{i=1}^m (\lambda_i / \kappa) F_i, \text{Var} F_{\text{综}} = (\sum_{i=1}^m \lambda_i^3) / \kappa^2, \kappa = p$ 或 $\lambda_1 + \dots + \lambda_m$, 通常 $\text{Var} F_{\text{综}} > \text{Var} Z_{\text{综}}$, 即 $F_{\text{综}}$ 的取值范围通常比 $Z_{\text{综}}$ 大。	$Z_{\text{综}} = \sum_{i=1}^m (v_i / \kappa) Z_i, v_i \neq \lambda_i$ (判据之一) $\text{Var} Z_{\text{综}} = (\sum_{i=1}^m v_i^2) / \kappa^2$ (旋转后因子贡献从 λ_i 变为 v_i , 因此权数应取为 $v_i / \kappa, \kappa = p$ 或 $v_1 + v_2 + \dots + v_m$ 。
应用上侧重	信息贡献影响力综合评价。	成因清晰性的综合评价。

* 取初始因子的方法为主成分法。

方向不同, 直接导致主成分值、因子得分值、综合评价值和应用侧重上不同, 综合评价应该分开进行, 混淆在一起是不同计量值交替错误。

三、避免出错的方法步骤

1 主成分分析法和 SPSS 软件应用时一对一的正确步骤:

- (1) 指标的正向化;
- (2) 指标数据标准化(SPSS 软件自动执行);
- (3) 指标之间的相关性判定: 用 SPSS 软件中表“Correlation Matrix(相关系数矩阵)”判定;
- (4) 确定主成分个数 m : 用 SPSS 软件中表“Total Variance Explained(总方差解释)”的主成分方差累计贡献率 $\geq 85\%$ 、结合表“Component Matrix(初始因子载荷阵)”中变量不出现丢失确定主成分个数 m 。
- (5) 主成分 F_i 表达式(这是 SPSS 软件及其教科书中没完善的地方): 将 SPSS 软件中表“Component Matrix”中的第 i 列向量除以第 i 个特征根的开根后就得到第 i 个主成分 F_i 的变量系数向量(在“transform \rightarrow compute”中进行计算), 由此写出主成分 F_i 表达式。用 $F_m = A_m' X$ 的 $A_m' A_m = I_m$ 检验之。
- (6) 主成分 F_i 命名: 用 SPSS 软件中表“Component

Matrix”中的第 i 列中系数绝对值大的对应变量对 F_i 命名(有时命名清晰性低)。

(7) 主成分与综合主成分(评价)值(这是 SPSS 软件及其教科书中没完善的地方): 综合主成分(评价)公式

$$F_{\text{综}} = \sum_{i=1}^m (\lambda_i / p) F_i$$
 (在“transform \rightarrow compute”中进行计算), λ_i / p 在 SPSS 软件中表“Total Variance Explained”下“Initial Eigenvalues(主成分方差)”栏的“% of Variance(方差率)”中。 $\text{Var} F_{\text{综}} = (\sum_{i=1}^m \lambda_i^3) / p^2$ 。

(8) 检验: 综合主成分(评价)值用实际结果、经验与原始数据做聚类分析进行检验(对有争议的结果, 可用原始数据做判别分析解决争议)。

(9) 综合实证分析。

2 因子分析法和 SPSS 软件应用时一对一的正确步骤:

- (1) ~ (3) 步骤同主成分分析步骤。
- (4) 确定因子个数 m : 用 SPSS 软件中表“Total Variance Explained”特征值累计贡献率 $\geq 85\%$ 、结合表“Rotated Component Matrix(旋转后因子载荷阵)”中变量不出现丢失确定因子个数 m 。
- (5) 求因子载荷矩阵 B_m : SPSS 软件中表“Rotated Component Matrix”。

量的比重综合指标;第三主成分 F_3 仅与 X_3 十分显著正相关,所以我们可以称之为外国分支机构比重指标;而第四主成分与变量没有明显的相关性,因此不对其进行命名。从这里也可以看出前三个主成分包含了全部的指标所具有的大部分信息。

四个主成分的表达式还不能从输出窗口中直接得到,因为“Component Matrix”是指初始因子载荷矩阵,为了得到四个主成分的表达式,以便求主成分值,还需进一步操作:将前四个因子载荷矩阵输入到数据编辑窗口(为变量 B_1 、 B_2 、 B_3 、 B_4),然后利用“Transform→compute”,在对话框中输入“ $A_1=B_1/\text{SQR}(6.049)$ ”,即可得到主成分系数向量 A_1 。同理,可得到 A_2 、 A_3 、 A_4 。于是,四个主成分表达式如下(这里的 ZX_i 是 X_i 的标准化数据):

$$F_1=0.1653ZX_1+0.2424ZX_2-0.0596ZX_3+0.364ZX_4+0.2495ZX_5+0.3357ZX_6+0.1113ZX_7+0.2584ZX_8+0.2516ZX_9+0.2244ZX_{10}+0.2659ZX_{11}+0.2707ZX_{12}+0.3507ZX_{13}+0.2961ZX_{14}+0.2355ZX_{15}$$

$$F_2=0.3341ZX_1-0.3016ZX_2+0.0064ZX_3-0.138ZX_4+0.3163ZX_5-0.0512ZX_6-0.2602ZX_7+0.29192ZX_8+0.2914ZX_9+0.3176ZX_{10}-0.2865ZX_{11}-0.284ZX_{12}$$

$$F_4=-0.297ZX_1-0.251ZX_2+0.4973ZX_3+0.1406ZX_4-0.0593ZX_5+0.4498ZX_6+0.3046ZX_7+0.2032ZX_8+0.1569ZX_9-0.1766ZX_{10}-0.2528ZX_{11}-0.2421ZX_{12}+0.1805ZX_{13}-0.1084ZX_{14}-0.1376ZX_{15}$$

应用这一线性组合计算出各主成分值,最后利用综合主成分函数($\kappa=p=15$):

$$F_{\text{综}}=0.40327F_1+0.38754F_2+0.07621F_3+0.0584F_4=0.1979zx_1-0.0189zx_2+0.066zx_3+0.0885zx_4+0.2218zx_5+0.1218zx_6-0.0251zx_7+0.2321zx_8+0.2241zx_9+0.2172zx_{10}-0.0064zx_{11}-0.0032zx_{12}+0.1002zx_{13}+0.022zx_{14}+0.2095zx_{15}$$

可以求得各个国家世界经济全球化程度的综合主成分值(见表4)。

五、主成分分析与(初始)因子分析的实证比较

《刘文》表2的结果为初始因子分析结果(经仔细验算确认),现将其与主成分分析结果表4进行比较。

主成分分析与初始因子分析的命名依据都是初始因子载荷矩阵表3的相应列,《刘文》对初始因子分析的命名准确性不够,致使相应经济分析有些偏离实际,如《刘文》中“中国参与经济全球化程度总体水平很低,但对生产与贸易全球化依存度及投资全球化依存度很高”并不显现。实际结果表4中为:在华外国分支机构占世界全部外国分支机构的比重很高,表明中国参与经济全球化进程正受到

世界各国的高度关注。

表4 主成分、综合主成分值

国家	F ₁	排名	F ₂	排名	F ₃	排名	F ₄	排名	F _综	排名
美国	3.29	3	6.07	1	1.46	2	-0.80	14	3.74	1
英国	4.45	2	0.98	4	-1.76	16	2.17	1	2.17	2
德国	1.40	4	1.34	3	-0.25	5	-0.23	8	1.05	3
日本	0.44	6	1.85	2	-0.25	6	-1.23	16	0.81	4
法国	0.87	5	0.46	5	-0.52	14	0.45	4	0.52	5
新加坡	5.27	1	-6.26	16	1.18	3	-0.95	15	-0.27	6
意大利	-0.61	8	0.11	6	-0.54	15	-0.65	13	-0.29	7
加拿大	-0.43	7	-0.47	12	-0.31	11	0.00	7	-0.38	8
中国	-2.18	14	0.05	7	3.00	1	1.83	2	-0.52	9
巴西	-1.91	13	-0.05	8	-0.43	12	0.14	6	-0.81	10
澳大利亚	-1.36	10	-0.92	14	-0.30	10	0.22	5	-0.91	11
韩国	-1.69	12	-0.45	11	-0.27	7	-0.61	12	-0.92	12
墨西哥	-1.67	11	-0.68	13	0.02	4	-0.30	9	-0.95	13
新西兰	-0.98	9	-1.73	15	-0.28	8	0.73	3	-1.05	14
俄罗斯	-2.34	15	-0.19	10	-0.30	9	-0.36	10	-1.06	15
印度	-2.56	16	-0.10	9	-0.46	13	-0.39	11	-1.13	16

表4中主成分 F_i 、 $F_{\text{综}}$ 的值与《刘文》表2中因子 f_i ($=Z_i$)、 $F(=Z_{\text{综}})$ 的值全部不等,这是二者函数方差不同造成的,这里, $\text{Var}F_{\text{综}}=1.87$, $\text{Var}Z_{\text{综}}=0.32$, $\text{Var}F_{\text{综}}>\text{Var}Z_{\text{综}}$ 。

由表4与《刘文》的表2对比可知:部分国家参与经济全球化程度综合主成分值排名中,中国的排名相差较大,在本文表4中,中国排第9,而在《刘文》中国排第6;新加坡、意大利、加拿大、韩国在本文中表4分别排第6、7、8、12,而在《刘文》中分别排第7、8、9、13;墨西哥在本文中排第13,而在《刘文》中排第12。

通过表4可将综合主成分结果在等距 $d=(3.74+1.13)/4=1.2175$ 下可分为四类国家。

- 第一类国家:综合主成分值取值范围为[2.523 3.74]。
- 第二类国家:综合主成分值取值范围为(1.305 2.523]。
- 第三类国家:综合主成分值取值范围为(0.086 1.305]。
- 第四类国家:综合主成分值取值范围为(-1.13 0.086]。

通过样品的综合主成分值取值可以确定样品的类别,如美国的综合主成分值为3.74是第一类国家,英国的综合主成分值为2.17是第二类国家;但在《刘文》表2中,美国的综合值为1.57,在此只能划分为第二类国家,英国的综合值为0.9,在此只能划分为第三类国家。如果将表4中美国、英国、中国、巴西、澳大利亚、韩国、墨西哥、新西兰、俄罗斯、印度的综合主成分值在《刘文》表2中来确定样品的类别,结果是这些国家不在《刘文》表2的取值范围[-0.49 1.57]内。即不同定量值会带来混乱。

以上可看出:主成分分析与因子分析的实证结果是有差异的,定量值全部不同,不能混用。

日本向知识经济的战略转型^{*}

刘彦

ABSTRACT

Intellectual Property strategy is an important strategy selection as Japan entered the 21th. Century. The key problem for economic depression in Japan is having not established the modern industry structure in time after the traditional manufacturing moving out. The important reason for Japan's economic miracle is the hitchhike R&D input structure and the special R&D structure mainly depends on big company. During the period of transferring to knowledge economy, the traditional dominant reason for miracle becomes the restricting reason, for competition in knowledge economy is focusing on the intellectual property right system and the original creation in R&D rather than hitchhike.

关键词: 知识财产战略; 知识产权制度; 战略选择

2003年7月,日本总理内阁直属的知识财产战略本部正式推出了《创造、保护及应用知识财产推进计划》(以下简称《推进计划》)。就提出和实施国家知识产权战略来说,日本并不是第一个。令人关注的是,日本战后依靠成功的技术引进,实现了长达20年的经济增长“奇迹”,并相继提出“技术立国”和“科学技术创造立国”。日本已是世界技术强国之一,为何又提出知识财产立国?从1997年4月“21世纪知识财产委员会”向政府提出政策研究报告算起,日本为这一战略的制定做了长达6年的准备。考虑到日本经历了20世纪90年代长期的经济衰退和历届政府的多次经济与行政改革,知识财产战略显然是日本进入21世纪对多方比较权衡后的重大战略选择。日本能否成功将受到世人瞩目。

一、日本经济衰退之谜^①

《推进计划》认为,日本战后通过技术引进,成功地建立了称雄世界的制造体系。产业竞争力居世界领先水平。日本20世纪90年代泡沫经济崩溃以来,至今未能摆

脱经济衰退,原因之一是在剧烈变化的国际环境中,日本仍然沉浸在以往的成功经验之中,黠守陈旧的制度,未能对以往“日本模式”进行大胆的变革。

日本战后大量引进欧美先进技术,每年引进技术迅速增加,20世纪60年代平均每年引进1000多项,20世纪70年代前半期平均每年引进2000多项,1949~1975年,日本用不到60亿美元引进了25000多项技术,几乎引进了20世纪70年代前所有重要的产业技术。技术引进使日本产业技术水平与欧美国家的差距迅速缩小。20世纪50年代初,日本科学技术落后于美国20世纪20~30年,到20世纪70年代,在大部分产业领域,日本已基本消除了

*基金项目:科学技术部科技兴贸行动计划资助项目(2004EE660003)。

① 波特在《日本还有竞争力吗》一书中从另一个角度提出过日本经济之谜:日本既有高度竞争力的产业,也有不具竞争力的产业,好象同时存在两个日本。

参考文献

- [1] 刘玉玫、张芄. 经济全球化程度的量化研究. 统计研究, 2003. 12.
- [2] 卢纹岱. spss for windows 统计分析. 电子工业出版社, 2002. 9.

作者简介

林海明,男,广东商学院统计学系副教授,相继获得

武汉大学、湖南大学理学(数学)学士、硕士学位,从事多元统计学应用等研究。邮编:510320 地址:广东省广州市广东商学院统计学系。电话:020-84096763。邮箱:linhml@yahoo.com.cn.

张文霖,男,22岁,广东商学院统计学系学生。地址:广东商学院601信箱。邮编:510320。邮箱:zwenlin@126.com. 电话:020-84094459.