

纵向数据分析方法

刘红云 孟庆茂

(北京师范大学心理学院, 北京 100875)

摘要 纵向研究方法是心理学研究领域的一种重要方法。近年来,国外在纵向研究数据分析方法上取得了一系列理论和应用上的进展。文章对此方法进行了简要的回顾,并重点阐述了最近发展起来的纵向研究的方法:多层线性模型和潜变量增长曲线模型,并在此基础上对几种常用的方法进行了比较。

关键词 纵向研究, 多层线性模型, 潜变量增长曲线模型。

分类号 B841

1 引言

纵向研究在心理学中处于特殊的地位,这一研究主要用来分析一段时间或某几个时间点总体的平均增长趋势和个体之间的差异。也就是说,对于纵向研究设计,主要关心两个问题,一个是描述总体的平均增长趋势,另一个是用来描述不同个体之间增长趋势的差异。纵向研究与横向研究相比,最大的优点是纵向研究设计可以合理地推论变量之间存在的因果关系。从方法论的角度讲,要想得出变量之间的因果关系,原因变量和结果变量之间至少要满足下列 3 个条件^[1]:(1) 假设存在因果关系的原因变量和结果变量之间是相关的;(2) 从时间上来讲,原因变量在前,结果变量在后;(3) 在所考虑的模型中,其他原因变量对结果变量的影响能够被控制或排除。可见从方法论的角度来讲,横向研究永远不可能满足上述的第二个条件,所以要从横向研究数据本身探索变量之间的因果关系,几乎是不可能的。正是由于纵向研究有这一显著优点,所以在心理研究中,多用纵向研究探讨数据之间的因果关系和分析事物的增长规律。

近年来,随着社会科学研究方法的快速发展,提供了一系列分析变量增长趋势的统计方法。其中概括起来,主要有以下几种:(1) 重复测量的方差分析(repeated measures analysis of variance);(2) 时间序列分析(time series analysis);(3) 潜变量增长曲线模型(latent growth curve model);(4) 多层线性模型(hierarchical linear model)。

上面常用的几种方法各有优缺点,前面两种方法主要是解决总体平均发展趋势的问题,而后两种方法除了对总体平均增长趋势进行分析外,同时注重个体发展趋势之间的差异。因此,从心理学纵向研究方法的进展而言,纵向研究的问题,逐渐由以往的注重总体平均趋势的发展过渡到综合考虑总体平均趋势和个体发展差异的系统分析的问题^[2]。

2 传统纵向数据分析方法综述

2.1 重复测量方差分析

重复测量的方差分析在实际中有非常广泛的应用,其中的一个作用就是用来分析重复测量实验设计(又称被试内设计、混合设计等)得来的数据。该方法通过把总的变异分解为被试内和被试间两部分,对被试

的平均增长趋势进行分析,可以通过多项式比较分析线性增长趋势和非线性增长趋势。如果研究中我们只关心不同时间点的平均数间是否存在差异,可以用单变量方差分析解决这一问题。但是值得注意的是应用重复测量的方差分析时,必须满足协方差矩阵球形(sphericity)的假设条件,也就是说,MANOVA要求所有重复测量的总体的方差相等并且所有重复测量总体之间的协方差也相等。如这一条件不满足,那么得到的F检验统计量的值正偏,拒绝虚无假设的概率增大,也就是说如果观测变量协方差矩阵球形假设条件不满足,传统重复测量的方差分析的统计检验力降低,F检验犯第一类错误的概率增大。另外,MANOVA不能用来处理依时间变化的协变量对因变量的影响。关于重复测量方差分析的详细介绍在大多数的统计资料中都有较详细的介绍,这里不再重复。

用于重复测量的方差分析的软件有很多,最基本的有SAS、SPSS和Statistics等,另外,这一方法还可看成是后面介绍的LGM和HLM的特例,也可用SEM和HLM软件进行分析。

2.2 时间序列分析

时间序列分析是对纵向研究数据进行分析的另外一类非常重要的统计分析技术,它在许多领域都有十分重要的应用,尤其在预测和控制应用方面,有着其它方法不可比拟的优点。时间序列分析以回归分析为基础,目的在于测定时间序列中存在的长期趋势、季节性变动、循环波动及不规则变动,并进行统计预测。为了对时间序列中不同的变化趋势进行分析,主要有两大类模型:经典模型(Kinetic Model)和动态模型(Dynamical Model),经典模型是将时间序列 $\{x_t, t \in T\}$ 看作是时间的函数: $x_t = f(t)$;而动态模型是将t时刻的观测看成是t时刻前观测值(可以与t时刻的观测类型相同,也可以不同)的函数: $x_t = f(x_{t-1}, x_{t-2}, \dots)$,通常所说的AR, ARMA, ARMA模型都属于这一类,这里为了便于和其他几种方法比较,我们只简单介绍第一类类型模型。对于第一类类型的模型常用的模型有:加法模型,即假定各构成部分对时间序列的影响是相互独立的,这时可以将时间序列表示为: $x_t = T + C + S + I$,其中,T、S、C、I分别代表时间序列中存在的长期趋势、季节性变动、循环波动及不规则变动。另一类是乘法模型,即假设各组成部分对时间序列的影响均按比例变化,从而可以把时间序列表示为: $x_t = T \times C \times S \times I$ 。除上面的加法模型和乘法模型外,还有其它混合模型,不再一一列举。进行时间序列分析,如果我们要测定长期趋势(可以是直线的也可以是非直线的),可以通过移动平均法、时距扩大法或数学模型法,剔除时间序列中循环波动C、季节性变动S及不规则变动I,使得时间序列的长期增长趋势显现出来。

对于时间序列中的第二类模型,在实际中有许多应用,模型分类也比较复杂,需要对时间序列的平稳性进行分析,并且要求研究者有较高的数学素养。另外,由于时间序列分析往往要求较多的连续观测时间点,所以在心理学和教育学中用的不是很多。

目前常用的统计软件SAS、SPSS和BMDP都含有时间序列分析过程,可以对常见的几种时间序列模型进行统计分析。

3 纵向数据分析方法新进展

3.1 潜变量增长曲线模型

潜变量增长曲线模型是用于固定情形(fixed occasion)纵向研究数据的一种统计分析方法,也就是说,该方法适用于在某几个固定时间点观测得来的纵向研究资料。在潜变量增长曲线模型中,用潜变量来描述总体的平均增长趋势和依时间变化的情况^[3,4]。基本模型可以用下图表示(图1):

图1描述的是含有五个测试时间点的潜变量增长模型, $y_{1i}, y_{2i}, y_{3i}, y_{4i}, y_{5i}$ 分别表示第i个被试的5次测量结果,上述模型可以表示为:

$$y_{ii} = \pi_{0i} + \pi_{1i}T_{ii} + e_{ii} \quad i = 1,2,3,4,5; \quad i = 1,2\Lambda, n \quad (1)$$

$$\pi_{0i} = \beta_{00} + \beta_{01}Z_i + u_{0i} \quad (2)$$

$$\pi_{1i} = \beta_{10} + \beta_{11}Z_i + u_{1i} \quad (3)$$

其中 π_{0i} , π_{1i} 分别表示截距和斜率, 在上面的模型中, 这一截距和斜率为随机参数, (2) 和 (3) 进一步解释上述截距 π_{0i} 和斜率 π_{1i} 的变化。

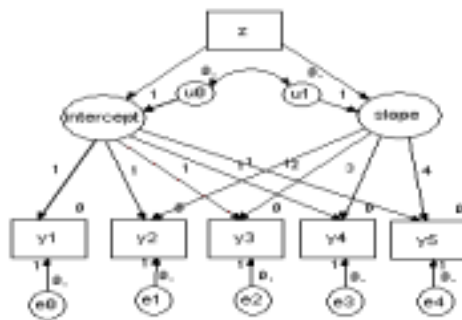


图1 潜变量增长模型结构图

从上面模型的描述可以看出潜变量增长曲线模型同时考虑因素的平均值和方差, 也就是说, 潜变量增长曲线模型不仅分析了总体的发展趋势, 而且可以分析总体之间存在的差异。

事实上, 在上述的潜变量模型中, 只是简单地定义了线性增长模型, 在实际中, 我们可以不固定斜率测量的因素载荷 (如在图 1 中让固定为 2, 3, 4 的斜率载荷自由估计) 得到增长曲线模型, 我们还可以定义测量误差之间的不同关系 (如限定测量误差相等, 误差间存在一阶自相关, 二阶自相关等等)。有关潜变量增长曲线的更详细的和深入的介绍, 可以参看 Duncan 的著作^[1]。

潜变量增长曲线模型可以用协方差结构模型 (SEM) 软件进行分析, 常用的软件有 Lisrel^[3], AMOS^[5], EQS^[6] 和 MPLUS^[7]等。

3.2 多层线性模型

多层线性模型是用于分析具有嵌套结构特点数据的一种统计分析技术, 近年来在教育、管理等领域有相当广泛的应用。当对相同的观测对象进行重复测量时, 可以将这些重复测量的数据本身看成是具有嵌套结构特点的。如对生长发育期儿童身高和体重变化情况的追踪调查等, 可将这些重复测量数据构造出一个两水平的层次结构, 其重复测量或测量点为水平 1 的单位, 观测个体为水平 2 的单位^[8], 这时就可用多层分析的方法对纵向数据进行分析。

对于重复测量的数据, 用层次分析法描述数据之间的关系, 对应的两水平重复测量模型, 可以用下式表示 (下面只给出最简单的一种多层模型形式, 实际上, 可以进一步考虑更多的不同水平预测变量和更复杂的随机残差之间的关系):

$$\text{水平 1 (重复测量): } Y_{ii} = \pi_{0i} + \pi_{1i}T_{ii} + e_{ii}$$

$$\text{水平 2 (个体): } \pi_{0i} = \beta_{00} + \beta_{01}Z_i + u_{0i}$$

$$\pi_{1i} = \beta_{10} + \beta_{11}Z_i + u_{1i}$$

$$\text{合并模型可以表示为: } Y_{ii} = \beta_{00} + \beta_{10}T_{ii} + \beta_{01}Z_i + \beta_{11}Z_iT_{ii} + u_{0i} + u_{1i}T_{ii} + e_{ii}$$

从上面的模型中可以看出,与潜变量增长曲线模型类似,多层分析不仅可以分析总体上个体随时间的变化(截距 β_{00} 和斜率 β_{10}),而且可以将个体之间增长的差异进行分析(截距的差异 u_{0i} ,斜率的差异 u_{1i}),并将这一差异的原因进行解释(β_{01} 解释截距的差异和 β_{11} 解释斜率的差异)。

可以在上述模型中包含更多的水平 1 的随机误差。这主要是由于在重复测量的模型,测量与测量之间往往是相关的而不是独立的(如在个体水平上的多次测量,由于具有相同的个体特征和测量间的相互影响,存在的测量误差(第一水平的随机误差)之间的“自相关”。

对于多层线性模型的数据分析,可以采用专门的软件进行分析,常用的用于多层分析的统计分析软件有:HLM^[8],MLn^[9],VARCL^[10],SAS 和 Mplus^[7]等。

3.3 应用举例

为了使读者对这两种方法有一个初步的认识,下面通过一个简单的例子来说明多层线性模型和潜变量增长曲线模型在分析纵向研究数据时的应用。

数据:随机抽取 155 名婴儿,测量其出生时的体重,然后在婴儿 3 个月,9 个月,15 个月和 21 个月每隔半年测量一次体重,目的在于分析婴儿 2 岁以前体重的增长情况,以及出生时的体重对婴儿体重的影响。

对于上面的问题,我们可以用传统的协方差分析(出生时的体重为协变量),对平均体重增长的情况进行分析,这里我们不再重复这一传统分析方法的结果。主要就多层线性模型和潜变量增长曲线模型在分析这类数据时的应用进行简单介绍。我们首先假设简单的线性增长模型,对于多层线性模型和潜变量增长模型分别用 HLM 和 Lisrel 对上面的数据进行分析,得到结果如下(见表 1)。

表 1 多层线性模型和潜变量增长模型参数估计结果

| 固定部分: | 多层线性模型 | | | 潜变量增长模型 | | |
|--------------|--------|-------|-----------|---------|-------|---------|
| | 系数 | 标准误 | t | 系数 | 标准误 | t |
| β_{00} | 8.779 | 1.029 | 8.532** | 8.780 | 1.030 | 8.524** |
| β_{01} | 0.387 | 0.079 | 4.898** | 0.389 | 0.080 | 4.862** |
| β_{10} | 1.862 | 0.297 | 6.269** | 1.863 | 0.300 | 6.210** |
| β_{11} | 0.134 | 0.093 | 1.438 | 0.132 | 0.089 | 1.483 |
| 随机部分: | 方差 | 自由度 | 卡方值 | 方差 | 标准误 | t |
| u_0 | 2.321 | 153 | 198.851** | 2.332 | 0.563 | 4.144** |
| u_1 | 2.709 | 153 | 239.246** | 2.710 | 0.512 | 5.303** |
| e | 6.142 | | | 6.321 | | |

注: * $p < 0.05$, ** $p < 0.01$

从固定部分参数估计结果可看出,婴儿 2 岁以前的体重有明显的线性增长趋势(HLM: $\beta_{10}=1.862$, $t=6.269$; LGM: $\beta_{10}=1.863$, $t=6.210$); 婴儿出生时的体重对婴儿 3 个月后的平均体重有显著影响,出生时体重较重的婴儿,3 个月时的平均体重也较重(HLM: $\beta_{01}=0.387$, $t=4.898$; LGM: $\beta_{01}=0.389$, $t=4.862$); 但是婴儿出生时的体重对体重的增长速度没有显著影响(HLM: $\beta_{11}=0.134$, $t=1.438$; LGM: $\beta_{11}=0.132$, $t=1.483$)。随机部分参数估计结果表明,婴儿出生后 3 个月时的体重存在显著的个体间差异(HLM: $\text{Var}(u_0)=2.321$, $\chi^2=198.85$; LGM: $\text{Var}(u_0)=2.332$, $t=4.144$); 婴儿 2 岁前体重的增长速度也存在显著的个体间差异(HLM: $\text{Var}(u_1)=2.709$, $\chi^2=239.246$; LGM: $\text{Var}(u_1)=2.710$, $t=5.303$)。两种方法得到的结果在误差

允许的范围类相同。

显然,这两种方法不仅可以就个体平均发展趋势进行分析,而且可以分析个体间的差异;不仅可以回答平均水平是否存在差异的问题,而且可以回答发展趋势是否存在差异的问题;另外,在实际应用中,可以根据随机部分参数估计的结果,进一步考虑个体层次的预测变量对可能导致这一差异的原因进行分析。

4 纵向数据分析方法述评

上面介绍的用于纵向研究的常用方法,各有优缺点,简述如下:

重复测量的方差分析主要用来比较均值间的差异,一般不对增长的变异情况进行分析,也就是说,重复测量的方差分析主要用来描述总体的平均增长趋势,而不关注个体增长曲线存在的差异,有计算简单,易于理解等优点。最主要的缺点是不能就个体之间存在差异的原因进行分析和解释,数据中的缺失值不能得到精确的估计,在数据缺失量较大时,分析所用数据信息损失较大。另外,重复测量方差分析不能处理分段间距不等或测量次数不等的的数据。

时间序列分析是一类很有用的分析数据随时间变化趋势的统计技术,在自然科学和社会科学各个领域都有非常重要的应用价值,但是由于其理论比较复杂、要求测试的时间点相对具有连续性和要求较多的测试时间点等特点,所以在心理学和教育学的研究中用的不是特别普遍。

采用多层分析的方法处理重复测量数据与时间变量之间的关系,在多层结构中,可以对非平衡测量数据得到参数的有效估计,因此用多层分析法处理重复测量的数据,不要求所有的观测个体有相同的观测次数,在纵向调查研究中,由于各种各样的原因,被试个体观测值部分缺失的情况时有发生,因此多层分析法处理缺失数据而不影响参数估计精度的这一特征,使得多层分析法处理在处理纵向观测数据时,比传统多元重复测量方法有很大的优势。与传统的用于处理多元重复测量数据的方差分析和回归分析方法相比,多层分析法至少具有以下优点^[8]:多层分析法通过考虑测量水平和个体水平不同的差异,明确表示出个体在水平 1(不同测量点)的变化情况,因而对于数据的解释(个体随时间的增长趋势)是在个体与重复测量交互作用基础上的解释,即不仅包含了不同测量点的差异,而且包含了个体之间存在的差异;多层分析法对数据资料较传统多元重复测量方法有较低的要求,对于重复测量的次数和重复测量之间的时间跨度都没有严格的限制,不同个体可以有不同的测量次数,测量与测量之间的时间跨度也可以不同;多层分析模型可以定义重复观测变量之间复杂的协方差结构,并且对所定义的不同的协方差结构进行显著性检验,在多层分析模型中,通过定义第一水平和第二水平的随机变异来解释个体随时间的复杂变化情况;当数据满足传统多变量重复测量模型对数据的要求和假设时,层次分析法得到与传统固定效应多元重复测量模型相同的参数估计和假设检验结果;用多层分析模型可以考虑更高一层的变量(如不同地区儿童)对个体增长的影响。但是多层分析模型也有缺点,首先用于多层分析模型的参数估计方法较传统估计参数的方法要复杂得多^[4,10],而且与后面介绍的 LGM 方法相比也不能处理变量之间间接的影响关系和处理复杂的观测变量和潜变量之间的关系。

潜变量增长曲线模型(LGM)^[11]可以直接处理变量之间复杂的因果关系,即不仅可以对变量之间直接的影响关系进行分析,而且可以将变量之间间接的因果关系进行分析;另外,由于潜变量结构模型是基于协方差结构模型的理论,所以不仅可以分析观测变量之间的关系,而且可以在考虑测量误差的基础上对潜变量之间的因果关系进行考察;上面介绍的多层分析模型只能分析变量之间的直接因果路径,对于潜变量之间关系的分析要比 LGM 复杂得多,并且在测量模型上也有更多的限定条件。LGM 模型可以简便地处理变量测量误差(残差)之间的关系,而不必限定残差之间相互独立,如可以直接定义类似于 AR(1)和 ARMA

模型中所要求的残差之间的关系类型；用 HLM 虽然没有残差之间相互独立的要求，但是用现有的多层分析软件定义起来要比 LGM 复杂得多。LGM 的另外一个优点是，因为 LGM 分析可以采用标准的用于 SEM 的分析软件，所以可以得到模型整个拟合的情况，并且可以根据提供的修正指数对模型进行修改。LGM 不仅就个体的发展轨迹进行描述，而且可以分析个体之间存在的差异以及存在差异的原因；LGM 不仅可以对给定的增长趋势进行检验，而且在观测时间点多于两点的情况下可以对个体随时间变化的趋势类型（如直线或曲线）进行探索。LGM 可以分析依时间变化的预测变量对因变量的影响，并且可以用类似于 SEM 中多样本比较的方法对多个样本之间的差异进行检验，可以有效处理缺失值。但是 LGM 也有如下缺点，因为 LGM 用 SEM 的基本原理对变量之间的关系进行分析，所以为了得到可靠的分析和检验结果，往往要求比较大的样本容量；对于所有个体的评估要求测试时间间隔相同，如果个体的变化随时间变化趋势不是很明显，LGM 方法与传统方法相比没有明显的优势。

5 纵向数据分析方法应用前景

在心理学应用中，对于纵向研究的资料，我们往往不仅对个体增长的平均趋势感兴趣，而且希望分析个体之间增长存在的差异。作为综合分析方法，应当能够同时解决这两个问题。潜变量增长曲线模型和多层分析模型是在传统分析方法基础上发展起来的综合分析的统计技术，这两种方法可以同时解决上面提到的两个问题。从国外纵向研究的发展趋势来看，这两种方法近年来越来越受到重视，其原因不仅是因为这两种方法是一种新的统计分析技术，更重要的是他们可以帮助我们发现事物发展的更深一层的规律，可以对个体之间的发展变化进行进一步的分析和解释，为理论研究提供更加有意义的实证研究的成果。国内心理学研究中，多层分析方法处于起步阶段，潜变量增长曲线模型还没有见有介绍。

在心理学研究中，人们已经不满足于现象的描述（横断数据资料的分析）和简单的差异的检验，要对人类心理现象发展的内在心理机制进行研究，把握事物发展的内在规律，纵向研究必然越来越受到研究者的重视。随着纵向研究方法的应用，用于纵向数据分析的综合统计分析技术——潜变量增长曲线模型和多层分析法必然受到研究者们的青睐。

层次潜变量增长模型^[11-13]是将层次线性模型和潜变量增长曲线模型结合起来的一种统计分析技术，它试图将两种方法的优点结合起来更加合理地解决实际问题。但这一方法目前尚处于发展阶段，在理论和应用上都还不是特别成熟，也不能完全融入上面两种方法的优点，因此，对于此方法的发展也必然成为今后理论研究者 and 应用研究者共同关注的一个问题。

参考文献

- [1] Duncan T E, Duncan S C, Strycker L A. An Introduction to Latent Variable Growth Curve Modeling: Concepts, Issues, and Applications. New Jersey, London: Lawrence Erlbaum Associates, 1999. 12~65
- [2] Raudenbush S W, Chan W. Growth curve analysis in accelerated longitudinal designs. Journal of Research in Crime and Delinquency, 1992, 29: 387~411
- [3] Jöreskog K G, Sörbom D. Lisrel 8: Structural Equation Modeling with the SIMPLIS command language. Chicago: Scientific Software International, 1993. 12~145
- [4] Liang K Y, Zeger S L. Longitudinal data analysis using generalized linear models. Biometrika, 1986, 73: 13~22
- [5] Arbuckle J L. AMOS for windows, analysis of moment structures Version 3.5. Chicago IL: Smallwaters, 1995. 25~75
- [6] Bentler P M, Wu E. EQS structural equations program manual. Encino, CA: Multivariate software, 1995. 1~10
- [7] Múthen B, Múthen J. Mplus user's guide: <http://www.statmodel.com/Mplus> Mplus Version 2.1.

- [8] Bryk A S, Raudenbush S W. Hierarchical Linear Models: Applications and Data Analysis Methods. Newbury Park,CA: Sage Publication, 1992. 12~52
- [9] Rasbash J, Browne W, Goldstein H, Yang M, Plewis I, Healy M, Woosdouse G, Draper D. A user's guide to MLwin. London: Institute of Education, 1999. 1~225
- [10] Longford N T. A fast scoring algorithm for maximum likelihood estimation in unbalanced mixed models with nested effects. *Biometrika*, 1987,74: 817~827
- [11] Duncan S C, Duncan T E. A multilevel latent growth curve analysis of adolescent substance use. *Structural Equation Modeling*, 1996, 3: 323~347
- [12] Múthen B. Multilevel factor analysis of class and student achievement components. *Journal of Educational Measurement*, 1991, 28: 338~354
- [13] Múthen B. Multilevel covariance structure analysis. *Sociological methods and Research*, 1994, 22: 376~398

A Review on Longitudinal Data Analysis Method and It's Development

Liu Hongyun Meng Qingmao

(Department of Psychology, Beijing Normal University, Beijing 100875)

Abstract : Longitudinal method is one of the central topics in psychology. A series of theoretical and application advances have been made recently. In this article, these advances are reviewed briefly and overlapping Hierarchical Linear Model(HLM) and Latent Growth Curve Model(LGM) are mainly discussed. In addition, the difference of several longitudinal methods are discussed briefly.

Key words: Longitudinal research, Hierarchical Linear Model, Latent Growth Curve Model.