

# 随机死亡率模型的改进与预测\*

张志强 杨帆

**【摘要】**文章将多变点检测方法应用于人口死亡率预测,并对年龄别死亡率的偏差进行主成分提取,利用变点检测法分别估计了主要主成分得分随时间变化的最优变点个数及位置,据此对主成分得分进行分段线性回归拟合,从最后一段回归模型外推主成分得分的预测值,得到死亡率预测值;同时利用发达国家 1951~2010 年连续 60 年死亡率数据,对改进的 PC 模型与经典 Lee-Carter 进行比较研究,结果表明,改进的 PC 模型在死亡率预测的精度和稳定性方面均优于经典 Lee-Carter 模型,多变点检测方法提高了死亡率模型的预测精度。研究结果显示,基于奇异值分解的经典 Lee-Carter 模型中的时间因子和基于特征值分解的经典 PC 模型中的第一主成分得分反映出了几乎一致的死亡率变化趋势;经典 PC 模型中的第二主成分主要综合了队列效应对死亡率的影响。

**【关键词】**随机死亡率 主成分分析 变点检测 队列效应 分段线性回归

**【作者】**张志强 厦门大学经济学院统计系,教授;杨帆 厦门大学经济学院,硕士研究生。

## 一、引言

20 世纪以来,人口死亡率持续下降和预期寿命不断延长成为世界人口发展最显著的特征,死亡率的这种非预期降低导致的长寿风险已成为各国政府及个人所面临的一类日益严重的社会风险。Koissi 等(2006)研究显示,多数国家依据经验数据对人口死亡率下降的预测普遍存在低估,导致各国养老保障事业的发展滞后于人口老龄化进程,同时也给养老金成本核算带来巨大风险,严重影响各类养老计划的可持续发展,因此,提高死亡率的预测精度应引起人口和统计学界的高度重视。

死亡率预测是人口预测的基础,也是寿险产品、养老产品定价和养老金财务计划的基础, Lee 等(1992)提出的 Lee-Carter 模型(以下简称经典 LC 模型)堪称随机死亡率预

\* 本文为国家社会科学基金重大项目“大数据与统计学理论的发展研究”(编号:13&ZD148)的阶段性成果。

测模型的基准,随后经典 LC 模型得到了广泛的应用和改进,主要研究成果大致可分为两类,一类是从死亡率的影响因素出发考虑对经典 LC 模型的改进和扩展及各死亡率模型的短期和中期预测精度的比较,如 Booth 等(2005)将经典 LC 模型与 LM 模型(Lee 等,2001)和 BMS 模型(Booth 等,2002a、2002b)进行了比较研究;Brouhns 等(2002a、2002b)、Delwarde 等(2007)在 Lee 等(1992)提出的方法中嵌入了 Poisson 回归模型,更好地解释了死亡率的变化;Renshaw 等(2006)、Debonneuil(2010)在经典 LC 模型中增加了队列效应,反映了队列效应对年龄别死亡率的影响;随着统计方法的发展,非参数统计法及处理函数性数据的统计方法被应用,为提高死亡率预测精度提供了新途径(Currie 等,2004;Hyndman 等,2007);死亡率预测的深入研究及综述也不断涌现(Booth 等,2006、2008;Shang 等,2010;Cairns 等,2006、2011)。国内学者主要利用 LC 模型对中国人口死亡率进行预测,陈秉正、祝伟(2009)针对部分年份中国城市分性别人口死亡率数据缺失的状况,运用死亡人数服从 Poisson 分布的 LC 模型进行死亡率预测;李志生、刘恒甲(2010)选择 LC 模型对中国人口的死亡率数据进行拟合和预测,探讨了 LC 模型在中国的适用性和表现形式;王晓军、任文东(2012)讨论了有限数据下 LC 模型在人口死亡率预测中的应用;吴晓坤、王晓军(2014)利用 Poisson 最大似然估计法建立中国人口年龄别死亡率的 LC 模型,在最大似然估计的基础上附加再抽样方法对模型参数、死亡率及其他相关变量进行估计和预测;王志刚等(2016)从理论上推导了 LC 模型预测区间,并给出了中国人口死亡率的预测区间。另一类是基于主成分分析思路构建死亡率预测模型。主成分分析的思路是基于特征值分解,按照数据方差最大方向调整数据的主成分分析降维方法,这与基于奇异值分解的经典 LC 模型不同,虽然这两种方法都是以提取出年龄别死亡率矩阵这一重要的特征为目的。

本文将解释变量为主成分得分的对数中心死亡率回归模型(Logarithmic Central Mortality Regression Model)称为经典 PC 模型,该模型最大的优点是结构简单,模型参数估计便利,最早的工作可追溯到 Bell 等(1991)。基于奇异值分解的经典 LC 模型可视为经典 PC 模型含第一主成分的情形(这里指经典 LC 模型中的时间因子和经典 PC 模型中的第一主成分得分反映了几乎一致的死亡率变化趋势),经典 LC 模型对死亡率的影响因素考虑不足,必然会导致经典 LC 模型的随机误差项存在异方差,但是,经典 PC 模型通过增加主成分的个数,异方差问题会迎刃而解,进而增强死亡率预测结果的稳定性。经典 PC 模型也存在两个缺点,一直阻碍经典 PC 模型的发展及广泛应用。一是各主成分得分的预测问题,二是各主成分的人口学意义不够明确。随着函数性数据主成分分析方法的提出,人们再度关注经典 PC 模型(Hyndman 等,2007;Shang 等,2011),但仍没有解决经典 PC 模型遇到的上述两个问题,为此本文将针对经典 PC 模型的这两个问题进行研究。

事实上,经典 PC 模型各主成分得分时间序列在各国的实际数据中其平稳性往往不能满足,对于短期预测来说,采用分段回归模型有其独特的优势,但其预测的精度取决于最优分段数及分段点的位置。随着统计变点检测理论的推进,本文首次采用多变点检测法客观地估计了各主成分得分随时间变化的最优分段数及分段点的位置,为利用分段回归模型进行各主成分得分短期预测奠定基础。从各国主成分得分变化来看,第一主成分得分下降趋势明显,但下降的速度不同,利用变点检测法可以检测到速度改变之处,进而提高第一主成分得分预测的精度;对于其他主成分得分预测,变点检测法更能体现其优势,因为其他主成分得分波动不定,利用变点检测法可获得最优分段进而保证最后一段回归分析时使用了最多的经验信息,从而提高短期预测的稳定性和精度,再将主成分得分的预测值代入经典 PC 模型中即可获得死亡率的预测值,实现提高死亡率预测精度的目的。本文利用发达国家 1951~2010 年的中心死亡率数据,对改进的 PC 模型与经典 LC 模型进行比较,并对经典 PC 模型中主要主成分进行人口学解释。

## 二、数据来源与经典 PC 模型的改进

### (一) 数据来源

本文选取美国、加拿大、英国、法国、丹麦、挪威、日本 7 个发达国家 1951~2010 年 0~99 岁 5 岁组分性别中心死亡率数据,这样就产生了 21 个年龄组(其中 0~4 岁组分为 0 岁组和 1~4 岁组),每组有 60 个观察值,同时,选取中国 1995~2012 年 0~84 岁 5 岁组分性别中心死亡率数据,即 17 个组,每组有 18 个观察值,以此作为本研究的基础数据。其中发达国家的数据直接源于世界人口死亡率数据库(HMD),中国的数据根据相应年份《中国人口统计年鉴》、《中国人口和就业统计年鉴》及《中国统计年鉴》计算得到。

选择这些国家一是这些发达国家的死亡率数据记载完整,记录了长达 60 年的人口死亡率变化;二是主成分的方差解释比例较高(见表 1)。从表 1 可以看出,并不是所有的发达国家的第一主成分的方差解释比例都超过 85%,有的发达国家的第一主成分的方差解释比例会低至 70%左右,选取丹麦、挪威两个国家作为代表,主要是这两个国家男性的第一主成分的方差解释比例分别为 72.76%、71.61%,与中国男性的第一主成分的方差解释比例相当;不过多数发达国家的第一主成分的方差解释比例都会超过 85%,有的甚至高达 95%以上。例如,美国、加拿大、英国、法国男性的第一主成分的方差解释比例分别为 85.81%、91.50%、89.32%、92.12%;日本男性这一比例达 95.98%。

### (二) 经典 PC 模型的改进

#### 1. 经典 PC 模型

首先引入 3 个符号  $\mu_x$ 、 $\sigma_{xt}$  及  $m_{xt}, \mu_x$  代表年龄组  $x$  的对数中心死亡率的平均值,表示

该年龄组的对数中心死亡率随时间变化的一般水平;用  $\ln m_{xt}$  代表  $t$  年年龄组  $x$  的对数中心死亡率,即  $\mu_x = \frac{1}{T} \sum_{t=t_0}^{T+t_0-1} \ln m_{xt}$ ,其中  $t_0$  为初始观察年, $T$  为观察年数; $\sigma_{xt}$  代表随着时间的推移各种影响中心死亡率的因素引发  $\ln m_{xt}$  的波动,即  $\sigma_{xt} = \ln m_{xt} - \mu_x$ ,称之为中心死亡率的偏差。通过对  $\mu_x$  和  $\sigma_{xt}$  含义的解释可见,有效提取  $\sigma_{xt}$  的影响因素是建立随机死亡率模型的关键,因此本文采用 Bell 等(1991)提出的模型(即经典 PC 模型)建立方程,即:

$$\ln m_{xt} = \mu_x + \sum_{j=1}^J \beta_{xj} PC_{jt} + \varepsilon_{xt} \tag{1}$$

式(1)中, $J$  表示主成分的个数, $PC_{jt}$  表示第  $j$  个主成分随时间  $t$  变化的主成分得分, $\beta_{xj}$  表示  $PC_{jt}$  对  $\ln m_{xt}$  的偏效应; $\varepsilon_{xt}$  表示  $t$  年时  $x$  岁对数中心死亡率的残差,并且  $\varepsilon_{xt}$  是一个均值为 0,方差为  $\sigma_\varepsilon^2$  的白噪声。

将各年龄组作为变量, $T$  作为样本容量,在 0.05 的显著性水平下 Bartlett 球形检验表明各国的  $\sigma_{xt}$  均适合进行主成分分析。表 1 给出了 7 个发达国家和中国  $\sigma_{xt}$  第一和第二主成分的累计方差解释比例。

表 1 各国  $\sigma_{xt}$  第一和第二主成分的累计方差解释比例 %

国 家	男 性		女 性		国 家	男 性		女 性	
	第一主成分	第二主成分	第一主成分	第二主成分		第一主成分	第二主成分	第一主成分	第二主成分
美 国	85.81	93.66	93.42	97.50	日 本	95.98	97.87	96.64	98.93
加 拿 大	91.50	95.72	95.80	98.01	丹 麦	72.76	84.51	82.77	89.74
英 国	89.32	96.12	94.64	97.09	挪 威	71.61	85.69	82.94	88.13
法 国	92.12	95.66	96.45	98.18	中 国	70.23	83.31	73.89	83.56

从表 1 可以看出,女性的  $\sigma_{xt}$  前两个主成分的累计方差解释比例在各个国家均高于男性;且并不是所有的国家的第一主成分的方差解释比例都会超过 85%,如丹麦、挪威和中国的第一主成分对  $\sigma_{xt}$  的方差解释比例均不足 75%,这表明仅用第一主成分不足以表达  $\sigma_{xt}$  的信息,加入第二个主成分后,表 1 显示,丹麦、挪威和中国的  $\sigma_{xt}$  前两个主成分的累计方差解释比例在 85%左右,而美国、加拿大、英国、法国、日本的  $\sigma_{xt}$  前两个主成分的累计方差解释比例均超过 90%,有的甚至超过 95%,说明选取前两个主成分可以基本保留原来  $\sigma_{xt}$  的信息,因此,本文讨论  $J=2$  的情形。

采用经典 PC 模型预测中心死亡率关键要确保主成分得分预测的精度,Bell 等(1991)在假设各主成分得分的时间序列是平稳的基础上,利用时间序列法根据各主成分得分来完成预测。事实上,各主成分得分时间序列其平稳性在各国的实际数据中往往不能满足,因此本文在经典 PC 模型中引入带有变点检测的分段线性回归法使死亡率预



测精度得到提高,以下称之为改进的 PC 模型。

## 2. 基于变点检测的分段线性回归的主成分得分预测

Zeileis 等(2003)介绍了线性回归模型中的变点检测方法,设一般的线性回归模型为:

$$y_i = x_i^T \beta_i + u_i \quad (i=1, 2, \dots, n) \quad (2)$$

式(2)中,  $y_i$  为因变量在时刻  $i$  的观察值,  $x_i$  为  $k$  维自变量向量, 方程中回归元的个数为  $k$  个, 常数项为 1。  $\beta_i$  是  $k$  维系数向量,  $\beta_i$  可能随着时间的变化而有所不同。线性回归中变点检测的原假设为  $\beta_i$  不随时间发生变化, 此时线性回归中不存在变点; 备择假设意为至少有 1 个系数会随时间发生变化。若假定回归中存在  $m$  个变点, 那么线性回归模型的形式为:

$$y_i = x_i^T \beta_j + u_i \quad (i=i_{j-1}+1, \dots, i_j; j=1, \dots, m+1) \quad (3)$$

式(3)中,  $j$  是分段指数,  $\xi_{m,n} = \{i_1, \dots, i_m\}$  表示变点的集合,  $i_0=0, i_{m+1}=n$ , 实际上变点很少是外生给定的, 而是未知的且需要从数据中进行估计。

给定  $m$  个分段点  $i_1, \dots, i_m$ , 可以得到  $\beta_j$  的估计值, 则最小残差平方和为:

$$RSS(i_1, \dots, i_m) = \sum_{j=1}^{m+1} rss(i_{j-1}+1, i_j) \quad (4)$$

式(4)中,  $rss(i_{j-1}+1, i_j)$  表示第  $j$  个分段回归中的最小残差平方和。最优变点的估计问题就是找到变点  $\hat{i}_1, \dots, \hat{i}_m$ , 使对于所有满足  $i_j - i_{j-1} \geq n_h \geq k$  的  $(i_1, \dots, i_m)$  能够最小化目标函数,  $\hat{i}_1, \dots, \hat{i}_m = \arg \min_{(i_1, \dots, i_m)} RSS(i_1, \dots, i_m)$ , 其中  $n_h = [nh]$ , 一般设定  $h=0.1$  或  $h=0.5$ 。为了得到

$\hat{i}_1, \dots, \hat{i}_m = \arg \min_{(i_1, \dots, i_m)} RSS(i_1, \dots, i_m)$ , 通常采用阶数为  $O(n^m)$  的网格搜索法, 但这种方法在

$m>2$  的情形下计算非常繁琐。为此, 联结子样本的分段算法被提出用以改善这种繁琐的计算方法。但这些分段算法并不能保证找到全局最小的  $(i_1, \dots, i_m)$ 。全局最小的  $(i_1, \dots, i_m)$  是通过一种动态规划算法来获取。对于任意的变点个数  $m$ , 这种算法的阶数均为  $O(n^2)$ 。Bai 等(2003)基于 Bellman 准则的思路提出了普通最小二乘回归情形下应用动态规划算法估计分段回归模型的方法, 这种方法确定的最优分段点的集合满足递归式:

$$RSS(\zeta_{m,n}) = \min_{mn_h \leq i \leq n-n_h} [RSS(\zeta_{m-1,i}) + rss(i+1, n)] \quad (5)$$

式(5)中, 如果  $i$  是分段点集合  $\zeta_{m,n} = \{i_1, \dots, i_m\}$  中的最后一个变点, 对于每一个  $i$ , 递归式可以知道  $i$  之前的最优子集。递归式的结果可以从满足  $j-i \geq n_h$  的  $rss(i, j)$  的一个三角矩阵中得到, 递归式(5)可以通过递归关系  $rss(i, j) = rss(i, j-1) + r(i, j)^2$  进行计算, 其中  $r(i, j)$  表示样本从  $i$  处开始, 时刻  $j$  时的递归残差。

本文根据上述普通最小二乘回归情形下应用动态规划算法估计分段回归模型, 来

表 2 英国和丹麦  $PC_{1t}$  的分段回归的参数估计结果

分段时间 (年)	男 性		分段时间 (年)	女 性	
	截距项	斜率系数		截距项	斜率系数
英国					
1951~1959	673.2069	-0.3415	1951~1959	944.8210	-0.4801
1960~1978	323.7066	-0.1628	1960~1978	366.2907	-0.1845
1979~1998	459.0404	-0.2317	1979~2010	531.2613	-0.2682
1999~2010	840.9149	-0.4228			
丹麦					
1951~1960	263.1336	-0.1331	1951~1961	717.5923	-0.3640
1961~1977	177.9239	-0.0890	1962~1977	546.6314	-0.2764
1978~1994	416.8840	-0.2094	1978~1995	272.4891	-0.1375
1995~2010	1272.6480	-0.6384	1996~2010	910.1975	-0.4573

确定前两个主成分得分的分段线性回归中最优变点的个数及其位置,再借助 R 软件中的 breakpoints 函数对前两个主成分得分分段线性回归的参数进行估计。

由于中国数据不能满足变点检测方法的要求,本文以英国、丹麦为例,这两个国家

前两个主成分得分的分段线性回归结果如表 2、表 3 和图 1、图 2 所示,其中图 1 显示了英国和丹麦分性别的分段回归模型与  $PC_{1t}$  的真实值的拟合效果。通过分段拟合可以看出  $PC_{1t}$  随时间下降速度的变化,从人口学意义来看,分段回归模型可以直接反映死亡率水平的变化,Amanian 等(2013)对此也给出了分析。图 2 显示了英国和丹麦分性别的分段回归模型与  $PC_{2t}$  真实值的拟合效果。通过分段拟合可以显示  $PC_{1t}$  和  $PC_{2t}$  的变化速度,依据最后分段回归模型外推得到  $PC_{1t}$  和  $PC_{2t}$ ,预测精度明显提高。

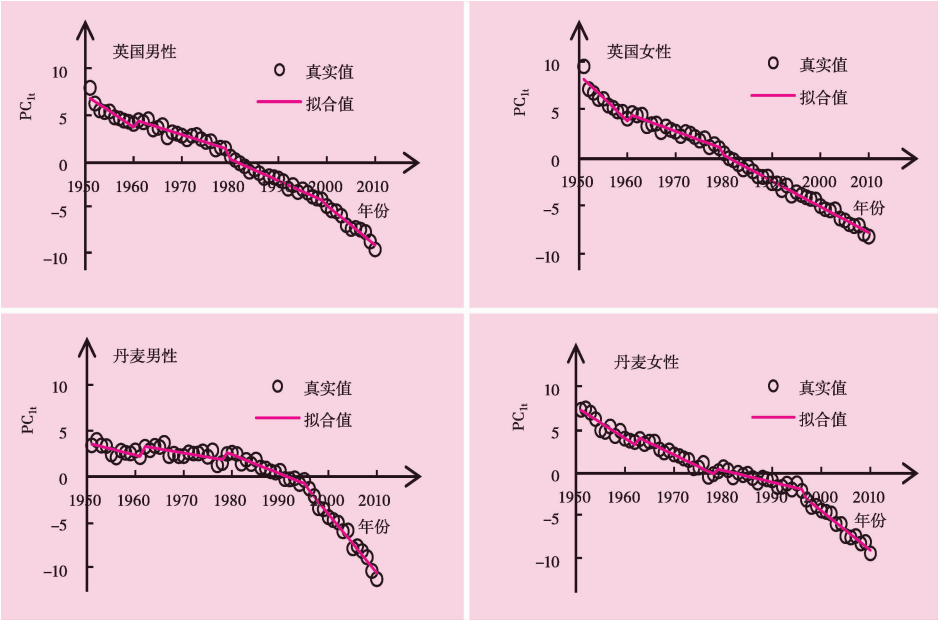


图 1 英国和丹麦分性别  $PC_{1t}$  的分段线性回归

注:纵坐标  $PC_{1t}$  表示  $t$  年年龄组  $x$  的对数中心死亡率的偏差的第一主成分得分。

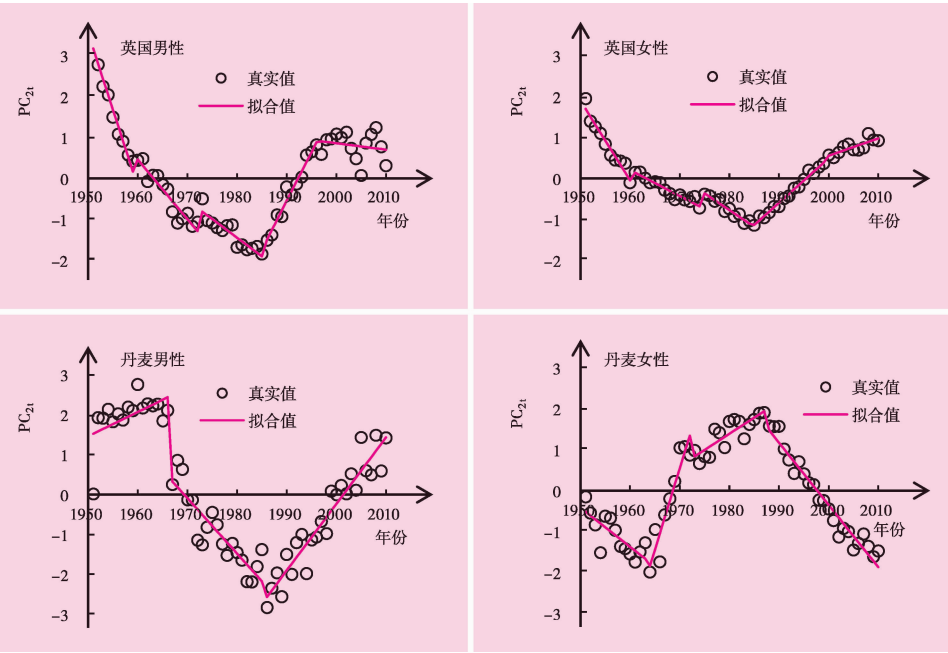


图 2 英国和丹麦分性别  $PC_{2t}$  的分段线性回归

注:纵坐标  $PC_{2t}$  表示  $t$  年年龄组  $x$  的对数中心死亡率的偏差的第二主成分得分。

从变点检测确定的分段线性回归的结果看,英国男性 $PC_{1t}$ 的最优分段点为 3 个,即得分数据被分为 4 段,可以用最后一段数据即 1999~2010 年的数据外推得到  $PC_{1t}$ ;  $PC_{2t}$  的最优分段点为 4 个,即得分数据被分为 5 段,可以用最后一段

表 3 英国和丹麦  $PC_{2t}$  的分段回归的参数估计结果

分段时间 (年)	男 性		分段时间 (年)	女 性	
	截距项	斜率系数		截距项	斜率系数
英国					
1951~1958	735.2333	-0.3752	1951~1959	380.2200	-0.1940
1959~1971	286.3193	-0.1458	1960~1973	124.4035	-0.0634
1972~1984	178.8456	-0.0911	1974~1983	160.9400	-0.0817
1985~1995	-473.6893	0.2378	1984~1998	-216.9286	0.1087
1996~2010	32.3674	-0.0158	1999~2010	-83.9721	0.0423
丹麦					
1951~1965	-117.5324	0.0610	1951~1962	186.2010	-0.0957
1966~1984	272.7513	-0.1385	1963~1971	-782.9358	0.3977
1985~2010	-333.1085	0.1664	1972~1986	-155.4078	0.0792
			1987~2010	303.7832	-0.1521

数据即 1996~2010 年的数据外推得到  $PC_{2t}$ 。英国女性  $PC_{1t}$  的最优分段点为 2 个,即得分数据被分为 3 段,可以用最后一段数据即 1979~2010 年的数据外推得到  $PC_{1t}$ ;  $PC_{2t}$  的最优分段点为 4 个,即得分数据被分为 5 段,可以用最后一段数据即 1999~2010 年的数据外推得到  $PC_{2t}$ 。同理可以获得丹麦的  $PC_{1t}$ 、 $PC_{2t}$  的最优分段数和各分段回归的参数估计(见表 2、表 3),这为采用经典 PC 模型预测英国和丹麦两性死亡率奠定了基础。

3. 改进的 PC 模型与经典 LC 模型的年龄别中心死亡率预测结果比较

本文将获得的英国和丹麦两国分性别第一和第二主成分得分 2011 年的预测值分别对应代入经典 PC 模型,就可以预测 2011 年各年龄组的中心死亡率,即为利用改进的 PC 模型获得的 2011 年各年龄组的中心死亡率的预测值,并将预测结果与各年龄组中心死亡率的实际数据进行比较。同时本文给出了经典 LC 模型的预测结果,这里经典 LC 模型参数估计采用矩阵奇异值分解法,对所估计的  $\kappa_t$  使用 ARIMA 模型拟合预测,进而获得经典 LC 模型预测的各年龄组的中心死亡率,其预测效果对比如图 3 所示,从中可见改进的 PC 模型的预测精度优于经典 LC 模型。

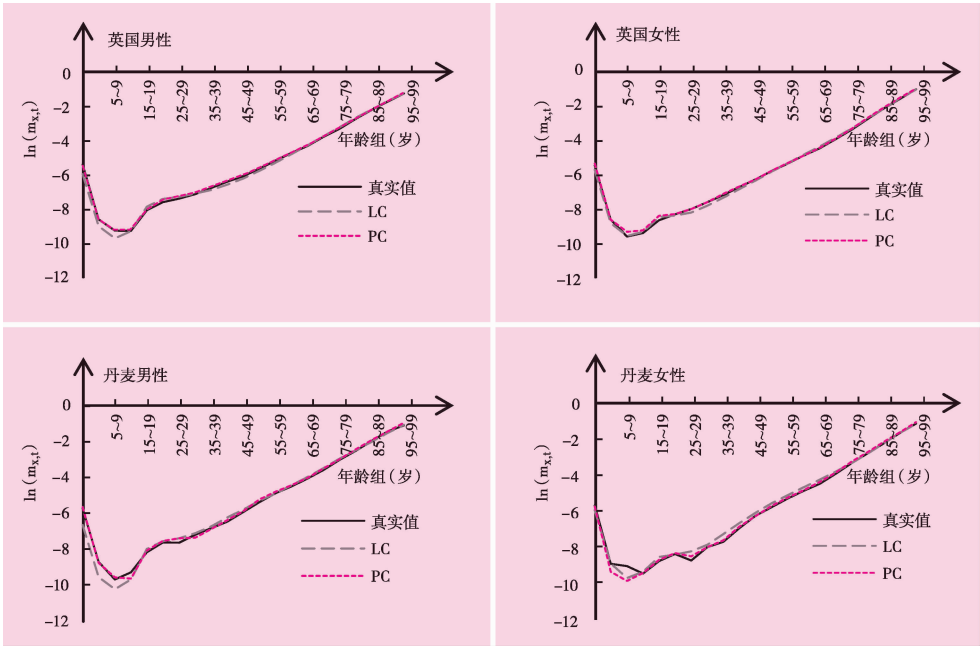


图 3 2011 年英国和丹麦分性别实际死亡率与预测死亡率比较

三、主成分的人口学解释

(一)第一主成分的人口学意义

1. 经典 LC 模型

经典 LC 模型表达式为:  $\ln m_{xt} = \mu_x + \beta_x \kappa_t + e_{xt}$ , 其中  $m_{xt}$  表示  $t$  年  $x$  岁的中心死亡率,  $\mu_x$  表示  $x$  岁的中心死亡率取对数后的平均值,  $\kappa_t$  表示时间因子, 反映了中心死亡率随时间变动的趋势,  $\beta_x$  表示变化  $\kappa_t$  时, 对数中心死亡率变化的敏感度,  $e_{xt}$  表示  $t$  年  $x$  岁对数中心死亡率的残差, 并且  $e_{xt}$  是一个均值为 0, 方差为  $\sigma_e^2$  的白噪声。对比经典 PC 模型和经典 LC 模型的结构可见, 比较两个模型关键是比较  $\kappa_t$  与  $PC_{1t}$  的关系。



2. 经典 LC 模型中  $\kappa_t$  与经典 PC 模型中的  $PC_{1t}$  的关系

本文在  $\sigma_{xt}$  的第一主成分的方差解释比例超过 85% 的国家中选择英国和日本作为代表,在不足 85% 的国家中选择丹麦为代表,为了显示  $\kappa_t$  与  $PC_{1t}$  关系的稳定性,对历史数据缺失严重的中国也进行了分析。虽然中国的历史数据缺失,但在图 4 中仍显示出  $\kappa_t$  与  $PC_{1t}$  具有一致性的规律。根据英国、日本、丹麦 1951~2010 年连续 60 年 0~99 岁 5 岁组分性别中心死亡率数据和中国 1995~2012 年连续 18 年 0~84 岁 5 岁组分性别中心死亡率数据,采用主成分分析方法即可获得  $\sigma_{xt}$  的第一主成分在各年的得分,即  $PC_{1t}$ ;经典 LC 模型中的参数  $\kappa_t$  估计采用奇异值分解法即可获得,将  $\kappa_t$  与  $PC_{1t}$  随时间推移而变化

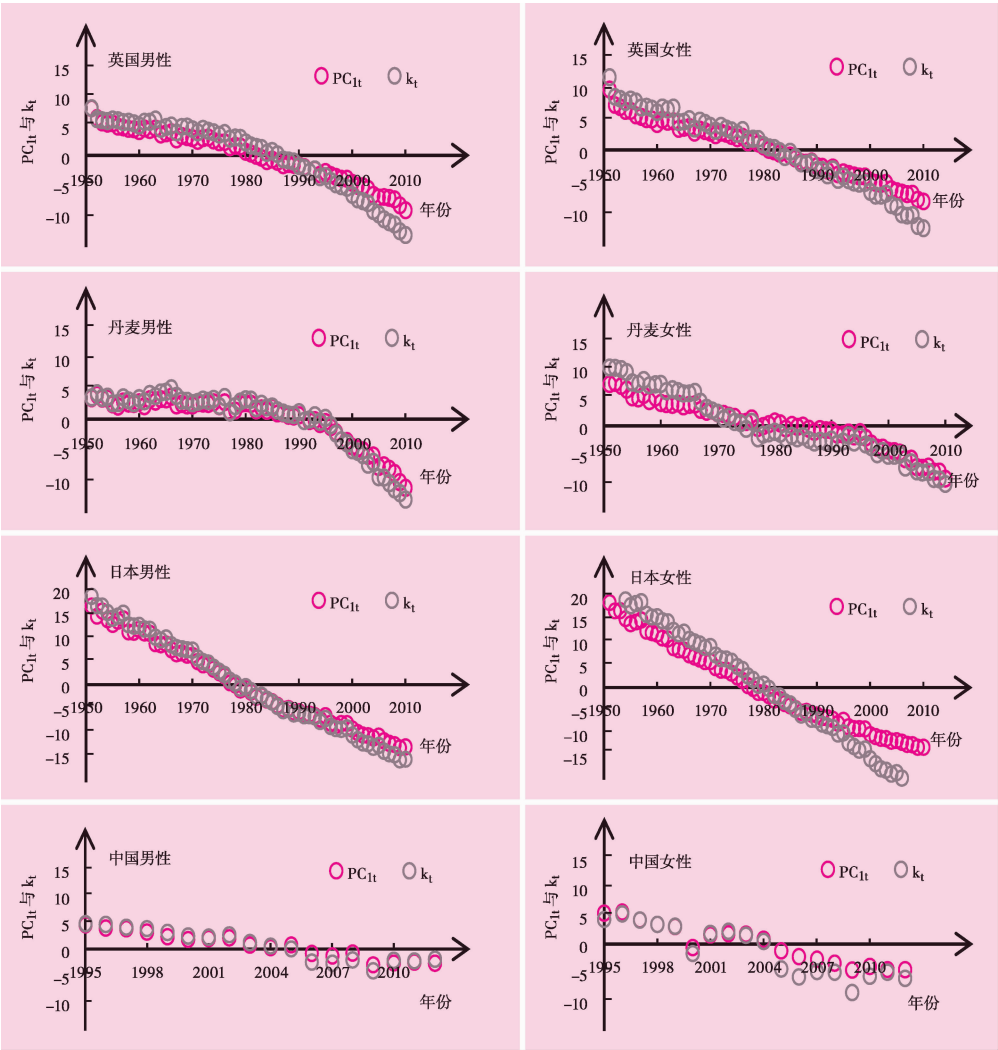


图 4 英国、丹麦、日本、中国  $PC_{1t}$  与  $\kappa_t$  变化趋势

注: $PC_{1t}$  表示  $t$  年  $x$  岁的对数中心死亡率的偏差的第一主成分得分; $\kappa_t$  表示时间因子,反映中心死亡率随时间变动趋势。

的结果呈现在同一坐标系里可见，随时间变动  $\kappa_t$  与  $PC_{1t}$  具有几乎相同的变化趋势(见图4)。事实上,其他国家也有相同的规律,说明  $\sigma_{xt}$  的第一主成分的人口学意义主要综合了各年龄组的中心死亡率随时间变化而下降的信息,这是年龄别死亡率变化的第一大特征,这一解释从第一主成分的线性组合的系数均较大也可以得到支持(见表4)。

(二)第二主成分的人口学意义

为了探索第二主成分的人口学意义,本文首先计算了死亡率的改善率(MIR),其计算公式为:

$$MIR=1-\frac{m_{xt}}{m_{x(t-1)}} \tag{6}$$

式(6)中, $m_{xt}$ 表示  $t$  年  $x$  岁的中心死亡率, $m_{x(t-1)}$ 表示  $t-1$  年  $x$  岁的中心死亡率, $MIR$  的数值越大说明死亡率的改善程度越大,若  $MIR$  为负值,则说明死亡率恶化了。为了能够清楚地显示死亡率改善率的变化,需要使用单岁组的中心死亡率数据,但是丹麦的原

表 4 英国和日本第一和第二主成分的线性组合系数

年龄组 (岁)	英国男性		英国女性		日本男性		日本女性	
	第一 主成分	第二 主成分	第一 主成分	第二 主成分	第一 主成分	第二 主成分	第一 主成分	第二 主成分
0	0.9819	0.0224	0.9873	0.0032	0.9958	-0.0695	0.9958	-0.0675
1~4	0.9896	0.0173	0.9903	-0.0218	0.9819	-0.2075*	0.9808	-0.2359*
5~9	0.9864	-0.0485	0.9883	-0.0605	0.9950	-0.2685*	0.9897	-0.2502*
10~14	0.9877	-0.0620	0.9792	-0.0420	0.9950	-0.2141*	0.9888	-0.2186*
15~19	0.9217	-0.3011	0.9484	0.0416	0.9778	0.0227	0.9678	-0.1634
20~24	0.9332	0.2073	0.9507	0.1645	0.9753	-0.0678	0.9662	-0.1154
25~29	0.7591	0.5793*	0.9398	0.3062*	0.9672	-0.0864	0.9740	-0.1156
30~34	0.7142	0.6838*	0.9500	0.2846*	0.9712	-0.1507	0.9773	-0.2072
35~39	0.8823	0.3836*	0.9784	0.2481*	0.9856	-0.1185	0.9842	-0.1702
40~44	0.9644	0.0440	0.9854	0.0445	0.9898	-0.0553	0.9926	-0.1163
45~49	0.9772	-0.0960	0.9837	-0.0809	0.9917	-0.0481	0.9953	-0.0853
50~54	0.9820	-0.1396	0.9752	-0.1661	0.9905	-0.0668	0.9961	-0.0628
55~59	0.9810	-0.1724	0.9628	-0.2459	0.9934	-0.0447	0.9978	-0.0197
60~64	0.9765	-0.2008	0.9584	-0.2488	0.9935	0.0012	0.9972	0.0494
65~69	0.9692	-0.2207	0.9689	-0.1845	0.9948	0.0414	0.9939	0.0965
70~74	0.9614	-0.2359	0.9823	-0.1158	0.9926	0.0964	0.9865	0.1560
75~79	0.9686	-0.2037	0.9934	-0.0513	0.9832	0.1681	0.9716	0.2329
80~84	0.9734	-0.1617	0.9955	-0.0420	0.9735	0.2151	0.9538	0.2983
85~89	0.9850	-0.0423	0.9922	-0.0090	0.9687	0.2385	0.9402	0.3378
90~94	0.9720	0.0581	0.9802	0.0239	0.9747	0.2106	0.9322	0.3561
95~99	0.9223	0.2402	0.9360	0.1802	0.9063	0.3526	0.9267	0.3588

注:\* 表示第二主成分的线性组合系数较大的年龄组。

始数据中部分单岁组的死亡人数为零,导致其中心死亡率也为零,这样死亡率改善率(MIR)会出现部分无意义的情形。另外,中国单岁组的中心死亡率数据缺失严重,因此,本文通过热度图显示英国、日本男女死亡率的改善率(见图 5、图 6)。

从图 5 中最明显的那条深色斜线来看,英国 1920 年左右的出生人口存在一个明显的队列效应,对应的观察期为 1951~2010 年,这代人 1951 年时年龄为 30 岁左右,其死亡率受到队列效应的影响,这一现象通过主成分的线性组合的系数也被表现出来。主成分是原来变

量的线性组合,在这个线性组合中各变量的系数绝对值大者表明该主成分主要综合了这些变量的特点。从表 4 可以看出,第二主成分的线性组合的系数在 30 岁左右年龄组最大,表明第二主成分的人口学意义主要综合了队列效应对死亡率的影响。同样从图 6 中最明显的那条深色斜线来看,日本 1945 年左右的出生人口存在一个明显的队列效应,对应的观察期为 1951~2010 年,这代人 1951 年时年龄为 5 岁左右,其死亡率受到队列效应的影响。从表 4 可见第二主成分的线性组合的系数在 5 岁左右的年龄组很大,说明第二主成分的人口学意义主要综合了队列效应对死亡率的影响。

四、结 语

本文在经典 PC 模型中引入带有变点检测的分段回归方法,有效地解决了各主成分得分的短期预测,进而提高了死亡率的短期预测精度;通过  $\kappa_t$  与  $PC_{1t}$  的比较发现二者反映出几乎一致的死亡率变化趋势;通过研究各国队列效应对死亡率的影响,发现了第

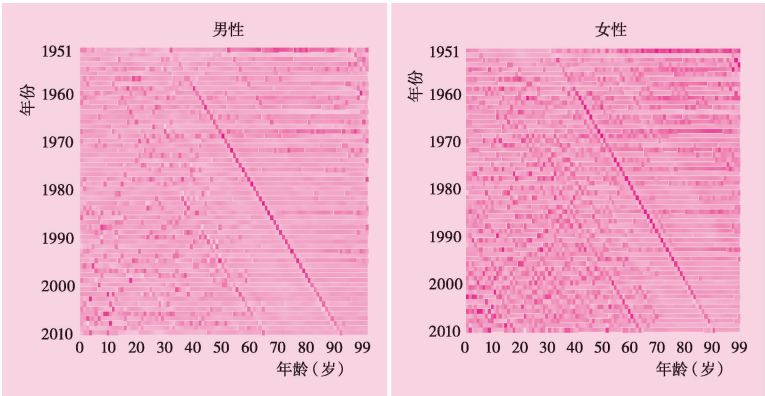


图 5 英国人口死亡率的改善率热度图  
注:图中颜色越深表示死亡率改善率越小。

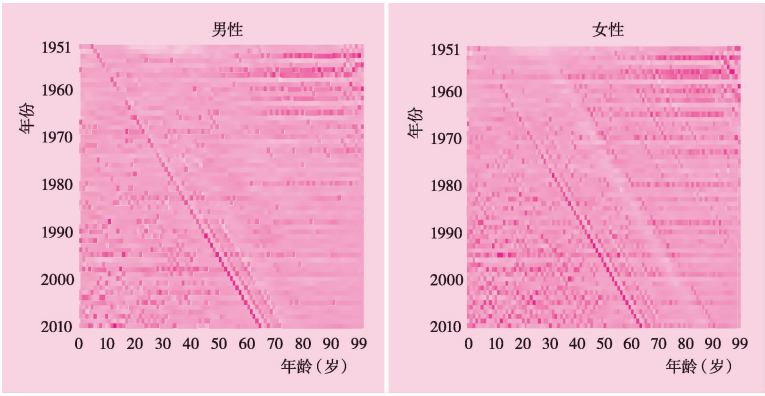


图 6 日本人口死亡率的改善率热度图  
注:同图 5。

二主成分的人口学意义,主要主成分含义的获得大大增强了经典 PC 模型的人口学解释力。研究还发现,由于经典 LC 模型中的随机干扰项存在异方差,使该模型对第一主成分的方差解释比例不足 85% 的国家不适用,因为主成分的方差解释比例越高,原来对数中心死亡率与其均值的偏差中保留的信息越多,那么中心死亡率的预测精度和稳定性就会越好。但多数发达且死亡率数据记载完整的国家的第一主成分的方差解释比例均超过 85%,有的甚至会高达 95% 以上,这与经典 LC 模型得到广泛应用的情形相一致。

采用带有变点检测的分段线性回归方法进行短期预测不仅有统计学意义而且也有人口学意义,对主成分得分采用带有变点检测的分段线性回归模型拟合相当于考虑了死亡率水平的时间可变性,并且第二主成分得分随时间推移而发生的变化使死亡率的极限值不再为零,实现了在不损失模型参数估计便利性的情况下,提高了模型的预测精度。

需要注意的是对主成分得分采用带有变点检测的分段线性回归模型拟合,要求中心死亡率原始数据的时间跨度足够长,否则只能主观确定变点,这也是本文没有分析中国情形的原因。本文论证了改进的 PC 模型在死亡率预测的精度和稳定性方面均优于经典 LC 模型,但改进的 PC 模型在个别国家的个别年龄组的死亡率预测结果也可能出现不及经典 LC 模型的预测结果的现象,例如,丹麦女性 5~9 岁年龄组的中心死亡率的预测结果(见图 3),这也表明改进的 PC 模型仍有不足,这是今后需要进一步研究完善的问题。

#### 参考文献:

1. 陈秉正、祝伟(2009):《中国城市人口死亡率的预测》,《数理统计与管理》,第 4 期。
2. 李志生、刘恒甲(2010):《Lee-Carter 死亡率模型的估计与应用——基于中国人口数据的分析》,《中国人口科学》,第 3 期。
3. 王晓军、任文东(2012):《有限数据下 Lee-Carter 模型在人口死亡率预测中的应用》,《统计研究》,第 6 期。
4. 王志刚等(2016):《Lee-Carter 模型的理论分布和区间预测》,《数理统计与管理》,第 3 期。
5. 吴晓坤、王晓军(2014):《中国人口死亡率 Lee-Carter 模型的再抽样估计、预测与应用》,《中国人口科学》,第 4 期。
6. Amania F., Kazemnejadb A., Habibic R.(2013), Change Point Detection in Trend of Mortality. *Canadian Journal on Computing in Mathematics, Natural Sciences, Engineering and Medicine*. 4(1):75-80.
7. Bai J., Perron P.(2003), Computation and Analysis of Multiple Structural Change Models. *Journal of Applied Econometrics*. 18(1):1-22.
8. Bell W., Monsell B.(1991), Using Principal Components in Time Series Modelling and Forecasting of Age-specific Mortality Rates. In: *1991 Proceedings of the American Statistical Association*. Alexandria, Virginia, American Statistical Association, 154-159.
9. Booth H., Maindonald J., Smith L.(2002a), Applying Lee-Carter under Conditions of Variable Mortality Decline. *Population Studies*. 56(3):325-336.
10. Booth H., Maindonald J., Smith L.(2002b), Age-Time Interactions in Mortality Projection: Applying Lee-Carter to Australia. Working Papers in Demography No.85, The Australian National University.

11. Booth H., Tickle L., Smith L. (2005), Evaluation of the Variants of the Lee-Carter Method of Forecasting Mortality: A Multi-country Comparison. *New Zealand Population Review*. 31(1):13-34.
12. Booth H., Jong P.D., Hyndman R.J., Tickle L. (2006), Lee-Carter Mortality Forecasting: A Multi-country Comparison of Variants and Extensions. *Demographic Research*. 15(9):289-310.
13. Booth H., Tickle L. (2008), Mortality Modelling and Forecasting: A Review of Methods. *Annals of Actuarial Science*. 3(1-2):3-43.
14. Brouhns N., Denuit M., Vermunt J.K. (2002a), A Poisson Log-Bilinear Regression Approach to the Construction of Projected Lifetables. *Insurance: Mathematics and Economics*. 31(3):373-393.
15. Brouhns N., Denuit M., Vermunt J.K. (2002b), Measuring the Longevity Risk in Mortality Projections. *Bulletin of the Swiss Association of Actuaries*. 2:105-130.
16. Cairns A.J., Blake D., and Dowd K. (2006), A Two-Factor Model for Stochastic Mortality with Parameter Uncertainty: Theory and Calibration. *Journal of Risk and Insurance*. 73(4):687-718.
17. Cairns A.J., Blake D., Dowd K., Coughlan G.D., Epstein D., and Khalaf-Allah M. (2011), Mortality Density Forecasts: An Analysis of Six Stochastic Mortality Models. *Insurance: Mathematics and Economics*. 48(3):355-367.
18. Currie I.D., Durban M., Eilers P. (2004), Smoothing and Forecasting Mortality Rates. *Statistical Modelling*. 4(4):279-298.
19. Debonneuil E. (2010), A Simple Model of Mortality Trends Aiming at Universality: Lee Carter+Cohort. *Quantitative Biology*.
20. Delwarde A., Denuit M., Eilers P. (2007), Smoothing the Lee-Carter and Poisson Log-bilinear Models for Mortality Forecasting: A Penalized Log-likelihood Approach. *Statistical Modeling*. 7(1):29-48.
21. Hyndman R.J., Ullah M.S. (2007), Robust Forecasting of Mortality and Fertility Rates: A Functional Data Approach. *Computational Statistics & Data Analysis*. 51(10):4942-4956.
22. Koissi M.C., Shapiro A.F., Högnäs G. (2006), Evaluating and Extending the Lee-Carter Model for Mortality Forecasting: Bootstrap Confidence Interval. *Insurance: Mathematics and Economics*. 38(1):1-20.
23. Lee R.D., Carter L.R. (1992), Modeling and Forecasting U.S. Mortality. *Journal of the American Statistical Association*. 87(419):659-671.
24. Lee R., and Miller T. (2001), Evaluating the Performance of the Lee-Carter Method for Forecasting Mortality. *Demography*. 38(4):537-549.
25. Renshaw A.E., Haberman S. (2006), A Cohort-based Extension to the Lee-Carter Model for Mortality Reduction Factors. *Insurance: Mathematics and Economics*. 38(3):556-570.
26. Shang H.L., Hyndman R.J., and Booth H. (2010), A Comparison of Ten Principal Component Methods for Forecasting Mortality Rates. Monash Econometrics and Business Statistics Working Papers, 8/10, Monash University.
27. Shang H.L., Booth H., and Hyndman R. (2011), Point and Interval Forecasts of Mortality Rates and Life Expectancy: A Comparison of Ten Principal Component Methods. *Demographic Research*. 25(5):173-214.
28. Zeileis A., Kleiber C., Krämer W., Hornik K. (2003), Testing and Dating of Structural Changes in Practice. *Computational Statistics & Data Analysis*. 44(1-2):109-123.

(责任编辑:朱 萍)