

三种时间序列模型在尘肺发病预测中的适用性研究

赵俊琴, 李建国, 赵春香

(河北省疾病预防控制中心, 河北 石家庄 050021)

摘要: **目的** 对基于时间序列的三种预测模型即自回归滑动平均混合模型 (ARIMA)、灰色模型 (GM)、广义回归神经网络模型 (GRNN) 进行尘肺发病预测的适用性比较。**方法** 选用河北省 1954—2015 年 62 年的尘肺发病数据, 前 54 年数据用来拟合预测, 后 8 年数据来比较三种模型的预测效果; 采用预测误差 (prediction error, PE)、平均绝对误差 (mean absolute error, MAE) 和平均相对误差 (mean relative error, MRE) 评价拟合效果。**结果** GM (1, 1) 的预测结果较差, ARIMA 的 MAE 和 MRE 是三种模型中最小的, 其短期预测的 PE 也最低; 三种方法长期预测的 PE 都比较大, 比较而言 GRNN 的长期预测结果最好。**结论** ARIMA 适用于尘肺发病的短期预测, GRNN 适用于长期预测。

关键词: 尘肺发病预测; 时间序列; 自回归滑动平均混合模型 (ARIMA); 灰色模型 (GM); 广义回归神经网络模型 (GRNN); 模型比较

中图分类号: R135.2; R195.1 文献标识码: A 文章编号: 1002-221X(2017)03-0168-04 DOI: 10.13631/j.cnki.zggyyx.2017.03.002

Applicability study on three time series models in incidence prediction of pneumoconiosis

Zhao Junqin, Li Jianguo, Zhao Chunxiang

(Hebei Provincial Center for Disease Prevention and Control, Shijiazhuang 050021, China)

Abstract: Objective The applicability of three prediction models based on time series namely autoregressive integrated moving average model (ARIMA), gray model (GM) and generalized regression neural network model (GRNN) in incidence prediction of pneumoconiosis was compared. **Methods** The pneumoconiosis incidence data of Hebei Province from 1954 to 2015 were collected, the first 54 years of data was used for fitting predictions, the last 8 years of data was used for comparing the prediction effect of the three models, the prediction error (PE), mean absolute error (MAE) and mean relative error (MRE) were used as well in the study to evaluate the fitting effect. **Results** The results showed that prediction effect of GM (1, 1) was poor; the MAE and MRE of ARIMA were the smallest among three methods, its PE was also the lowest in short-term prediction; meanwhile, the PE of three methods were all larger in long-term prediction, but in comparison, GRNN's long-term prediction effect was the best in these models. **Conclusion** The results suggested that ARIMA is more suitable for short-term incidence prediction of pneumoconiosis, and GRNN seems more suitable for long-term incidence prediction of pneumoconiosis.

Key words: incidence prediction of pneumoconiosis; time series; autoregressive integrated moving average model (ARIMA); gray model (GM); generalized regression neural network model (GRNN); model comparison

时间序列是指按照时间顺序排列的一组随机变量, 其在特定时间段上的观测样本可以视为随机过程的一次实现, 称为样本序列。由于时间的不可重复性, 时间序列通常仅有一次实现, 即只能在唯一可观测到的样本序列的基础上来推测时间序列的总体特征^[1]。科学预测尘肺发病趋势是准确判定尘肺防治效果的重要环节。随着预测理论及预测技术的发展与完善, 将有越来越多的统计理论、预测方法及适宜模型被应用于尘肺发病的预测和预警。本研究利用河北省历年尘肺数据, 旨在通过比较不同时间序列的预测模型, 探索适用于尘肺发病预测的适宜方法, 分析尘

肺发病趋势、预测未来尘肺发病水平及其规律, 为制定防治措施提供理论依据。

近年来, 有学者将时间序列方法应用于尘肺病的发病预测, 使用最多的是灰色模型 (gray model, GM), 并且预测效果较好^[2~4]。自回归滑动平均混合模型 (autoregressive integrated moving average model, ARIMA) 在传染病领域应用较多, 近期在尘肺发病预测中也有应用^[5]。神经网络模型中的时间序列方法——广义回归神经网络模型 (generalized regression neural networks model, GRNN), 在职业病领域的应用未见报道。本文将介绍 GRNN, 并与以上两种模型比较, 分析其在尘肺病发病预测应用上的实用性, 同时使用实际数据做模型拟合与结果验证, 结论更为可信、可靠。

1 原理与方法

1.1 原理

收稿日期: 2017-01-17; 修回日期: 2017-04-10

基金项目: 河北省卫生计生委医学科学研究重点课题, 河北尘肺病流行规律与防治对策研究 (20130089)

作者简介: 赵俊琴 (1987—), 女, 硕士, 医师, 研究方向: 劳动卫生与职业病。

通信作者: 赵春香, 主任医师, E-mail: hbxyy0911@126.com。

1.1.1 ARIMA ARIMA 又称 Box-Jenkins 模型, 由自回归 (AR) 模型与滑动平均 (MA) 模型组合而成。应用 ARIMA 模型时要求时间序列是零均值的平稳随机序列, 但大部分医学现象随着时间的推移, 表现出偶尔的上升或下降趋势, 为非零均值的非平稳时间序列, 因此在应用该模型前需对时间序列进行零均值化 (对均值不为零的时间序列中的每一项数据值都减去该序列的平均数, 构成一个均值为零的新的时间序列) 和差分平稳化 (对均值为零的非平稳的时间序列进行差分, 使之成为平稳的时间序列, 一般经过一阶差分或二阶差分后都可以平稳化) 处理。

ARIMA 的计算过程分为三个阶段。第一阶段为模型的识别: 根据自相关系数和偏自相关系数确定 ARIMA (p, d, q) 3 个参数, 先确定差分次数 d , 使时间序列实现平稳性; 然后确定自回归阶数 p , 即确定每个值对前第 p 个值的依赖程度; 最后确定滑动平均阶数 q , 即确定 q 个前干扰被用于平均 (在滑动平均过程中, 每一个值都是由当前干扰以及前一个或多个干扰的均值决定的)。第二阶段为模型中参数估计: 根据第一步给定的参数 (p, d, q) 初始值, 进行迭代计算获得拟合值和预测值。第三阶段为诊断与预测: 估计结果的残差应该是随机的, 即认为残差序列是白噪声序列, 又称之为白噪声检测, 常用 Box-Ljung Q 统计量进行检验, 检验零假设为在滞后 n 阶没有明显证据说明残差是非零自相关的, n 应选择约 $1/4$ 样本量 (但不应多于 $50^{[6]}$) 的时点。如果模型诊断合格则进行预测, 反之从第一阶段开始重复。

1.1.2 GM 灰色系统理论认为, 一切随机量都是在一定范围内、一定时段上变化的灰色量和灰过程。GM 是对灰色量进行处理, 将杂乱无章的原始数列, 变成比较有规律的时间序列数据, 即以数找数的规律, 再建立动态模型。对原始数据的处理有两个目的, 一是为建立模型提供中间信息, 二是将原始数据的波动性弱化。在灰色系统理论建模中, 一般用累加生成运算 (accumulated generating operation, AGO) 生成新的有规律的数列。GM 较常用的预测模型为用一次累加生成序列建立的一阶微分方程和一个变量的灰色模型, 记为 GM (1, 1) 模型。其预测模型为累加一次数列。

$$\hat{x}(k+1) = [x^0(1) - \frac{u}{a}] e^{-ak} + \frac{u}{a}$$

$$k=0, 1, 2, \dots, n$$

其中 a 表示发展参数, u 表示灰色作用量, k 为累加个数, $x^0(1)$ 表示原始数列。

模型建立后, 用后验差比值检验法对模型的拟合结果进行检验, 即求估计的残差方差与真实数列方差的比值 C 。 C 越小越好, 一般要求 $C \leq 0.45$, 最大不超过 $0.65^{[7]}$ 。其判定标准为 $C < 0.35$, 拟合精度等级为好; $C < 0.50$, 为合格; $C < 0.65$, 为勉强合格; $C \geq 0.65$, 为不合格。

1.1.3 GRNN GRNN 是径向基 (radical basis function, RBF) 神经网络模型的一个分支, 具有高度的容错性 (即在不影响预测的前提下, 允许误差存在) 和鲁棒性 (即在一定参数下, 预测能力保持稳定), 因此适用于非线性和不稳定数据的处理, 在逼近能力和学习速度上较 RBF 神经网络好, 网络最后收敛于样本量积聚较多的优化回归面, 并且在样本数据较少时, 预测效果也较好^[8]。

GRNN 由一个 RBF 网络层和一个线性网络层组成。其预测过程分为两个阶段: 一是训练阶段, 将样本分为训练样本和测试样本, 使用训练样本对网络进行训练, 使误差达到预定值, 获得预测模型。这过程中有一个扩展参数, 即 RBF 的平滑参数 spread, spread 越大, 函数拟合越平滑, 但是逼近误差会变大, 需要的隐藏神经元也越多, 计算也越大。反之, 函数的逼近会越精确, 但是逼近过程会不平滑, 网络的性能差, 会出现过适应现象。因此可通过调整 spread 取值来优化模型, 其最佳值可通过交叉验证法获取。二是预测阶段, 根据预测模型利用测试样本进行预测。整个过程计算快, 不需要循环的训练过程, 需要计算的参数仅有一个 spread 值。

1.2 统计方法

1.2.1 资料来源 以河北省 1954—2015 年累计 62 年间诊断的尘肺病病例为研究对象。

1.2.2 统计方法 先利用 1954—2007 年连续 54 年的年尘肺病发病数建立时间序列模型, 再以 2008—2015 年连续 8 年数据验证模型效果。模型拟合及预测效果采用预测误差 (prediction error, PE)、平均绝对误差 (mean absolute error, MAE) 和平均相对误差 (mean relative error, MRE) 进行评价。各自计算公式如下, 其中 $i=1, 2, \dots, n$ 表示预测个数。

$$PE = \text{实际发病数} - \text{预测发病数} \quad (1)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |PE| \quad (2)$$

$$MRE = \frac{1}{n} \sum_{i=1}^n \frac{|PE|}{\text{实际发病数}} \times 100\% \quad (3)$$

ARIMA 模型和 GM (1, 1) 模型采用 R 软件实现, GRNN 模型采用 Matlab 软件实现。

2 结果

2.1 拟合模型与预测

2.1.1 ARIMA 模型经过差分定阶和参数估计、根据 AIC 值选择合适的模型、拟合结果及白噪声诊断, 最终选择 ARIMA(1,1,1), 即 1 阶自回归 1 阶差分 1 阶滑动平均模型, 模型 AIC 值为 757.71, 残差 Box-Ljung 检验, 滞后 13 阶 χ^2 值为 12.619, $P = 0.4777$, P 值大于 0.05 不拒绝 H_0 , 可认为在滞后 13 阶没有明显证据说明残差是非零自相关的, 参数通过检验。

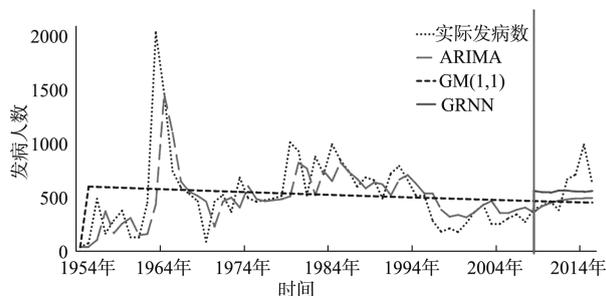
2.1.2 GM(1,1) 发展参数 $a = 0.0048$, 灰色作用量参数 $u = 602.3676$, 对模型进行后验差比值检验, 得 C 值为 0.707, 预测精度为不合格。

2.1.3 GRNN 使用 1954—2007 年共 54 年的数据作为 GRNN 模型的训练样本, 其中 1954—1980 年共 27 年的数据作为网络的输入, 1981—2007 年共 27 年的数据作为网络的输出, 然后根据训练的网络, 用 2008—2015 年共 8 年的数据作为测试样本来预测 2008—2015 年的发病数。经过反复交叉验证, spread = 0.4 时预测误差最小。

2.2 结果比较

图 1 显示三种模型中 ARIMA 的拟合效果和预测效果最好, 大致能反映河北省尘肺病发病的变化趋

势, 且其在预测开始阶段与实际发病曲线较为吻合。GM(1,1) 的拟合曲线和预测曲线为平滑曲线, 与实际发病曲线不一致。GRNN 的预测曲线也较为平滑, 与实际发病曲线不一致。



注: 竖线左侧为拟合部分, 右侧为预测部分 (GRNN 只有预测曲线)

图 1 三种方法的拟合及预测曲线

从表 1 可以看出, ARIMA 对 2008—2010 年发病预测的 PE 小, 短期预测结果较为准确, 整体来看 MAE 和 MPE 也是三种模型中最小的。GM(1,1) 对 2008—2011 年的发病预测 PE 比 GRNN 的小, 总的 MAE 较 GRNN 的大, 可见其对尘肺病的长期预测结果较差。三种方法的长期预测结果都不理想, 从 2012 年开始 GRNN 的 PE 是三种模型中最小的, 比较而言 GRNN 的长期预测结果较好。

表 1 三种方法的预测结果

年份	实际 发病数	ARIMA		GM(1,1)		GRNN	
		预测发病数	PE	预测发病数	PE	预测发病数	PE
2008	391	362	29	466	-75	563	-172
2009	421	417	4	464	-43	549	-128
2010	464	450	14	461	3	543	-79
2011	382	470	-88	459	-77	565	-183
2012	669	483	186	457	212	564	105
2013	714	490	224	455	259	556	158
2014	993	495	498	453	540	552	441
2015	625	498	127	451	174	561	64
MAE		146.25		172.88		166.25	
MRE		20.51%		25.05%		28.97%	

2.3 发病预测

2016 年尘肺病发病数已知为 990 例, 以 1954—2016 年连续 63 年的数据预测 2017—2020 年的尘肺病发病数。2017—2019 年、2020 年的尘肺病发病数分别使用 ARIMA 模型和 GRNN 模型预测, 结果显示 2017—2020 年的尘肺病预测发病数分别为 812、704、638、414 例。

3 讨论

尘肺病是我国乃至我省目前危害劳动者健康的主要职业病, 制约了我国生产力发展, 影响了劳动者身心健康及家庭、社会的稳定。早期预防尘肺病是职业病防治工作的重要任务, 探索准确预测尘肺病发病情况的手段与方法, 是提前预警的一种有效而重要的措施。尘肺病发病受多种因素影响, 如产业结构调整、

企业管理模式的变化等。构建尘肺发病预测的模型有很多方法,但各有所长。时间序列分析因突出时间序列暂不考虑外界因素影响,数据易获取,但当外界发生较大变化时,往往会有较大偏差。本次研究显示,ARIMA的拟合及预测效果基本能反映河北省尘肺病发病数波动的特点,适用于尘肺发病的短期预测,且其短期预测精度较高,但随着预测时间的延长,预测误差会相应的增大。三种模型比较而言GRNN对尘肺发病的长期预测精度相对较高,虽然其预测曲线不能完全反映尘肺病发病的波动特点,但基本能预测出发病趋势的平均水平。GM(1,1)的拟合曲线为平滑曲线,与河北省尘肺病发病的波动现象不符。有研究认为GM是短、中期预测法,不适用于长期趋势预测^[9];GM(1,1)适用于时间序列呈指数型趋势的方法,并且要求发病呈单一上升或下降趋势,不适宜对发病波动较大的疾病进行预测,此外如果数据离散程度越大,波动性越强,则预测精确度较差^[10]。本次GM(1,1)对尘肺发病做预测时,C值为0.707,预测精度为不合格,GM不适用于河北省尘肺病发病预测,这与有些文献报道的灰色模型适用于尘肺病预测的结论不同^[2-4]。

本研究通过河北省1954—2015年实际尘肺病数据验证,认为基于时间序列的方法中ARIMA适用于尘肺发病的短期预测,GRNN适用于长期预测。虽然ARIMA对数据要求相对较高,需要较多的连续序列数据(有文献建议不少于30个^[11]),但相对基于影响因素的预测方法来说,数据易获取;GRNN则对数据无要求,学习能力与容错能力都较好,且这两种方法使用软件计算都较简便。在今后尘肺发病预测中可采用ARIMA预测短期(如近二三年)尘肺发病例数,用于尘肺病防治效果评价;采用GRNN预测长期尘肺发病趋势,用于指导尘肺病防治策略。不同的

预测方法提供不同的有用信息,且其预测精度也不同,下一步将试探索这两种模型对尘肺发病的组合预测效果。

本研究预测结果显示,“十三五”期间河北省尘肺病发病数仍处于较高水平,特别是2017年、2018年。提示尘肺病在“十三五”期间依然是我省职业病防治重点。

参考文献:

- [1] 黄红梅. 应用时间序列分析 [M]. 北京: 清华大学出版社, 2016: 1.
- [2] 杨永利, 刘宏志, 孙世奎, 等. 灰色系统GM(1,1)模型对大隆矿煤工尘肺发病趋势的预测 [J]. 职业与健康, 2007, 23(18): 1581-1583.
- [3] 马兰, 陈彦龙. 应用数列模型对煤矿尘肺发病趋势的预测 [J]. 南通医学院学报, 2003, 23(3): 282-283.
- [4] 史善富, 魏春龙. 应用灰色数列模型预测南京市尘肺病发病危险度 [J]. 现代医药卫生, 2014, (24): 3838-3840.
- [5] 钟清, 隋怡, 庞燕, 等. 南京市新发尘肺自回归求积移动平均模型拟合预测 [J]. 中华劳动卫生职业病杂志, 2014, 32(3): 211-213.
- [6] 孙振球, 徐勇勇. 医学统计学 [M]. 3版. 北京: 人民卫生出版社, 2010: 394.
- [7] 宁宣熙, 刘思峰. 管理预测与决策方法 [M]. 北京: 科学出版社, 2009: 126.
- [8] 胡晓媛, 吴娟, 孙庆文, 等. ARIMA模型与GRNN模型对肺结核发病率预测的对比研究 [J]. 第二军医大学学报, 2016, 37(1): 115-119.
- [9] 党耀国. 灰色预测与决策模型研究 [M]. 北京: 科学出版社, 2009: 45.
- [10] 庞艳蕾, 张惠兰, 李向云, 等. 灰色模型GM(1,1)和ARIMA在拟合全国婴儿、5岁以下儿童死亡率中的应用 [J]. 中国卫生统计, 2015, 32(3): 461-463.
- [11] 陈远方, 张熳, 王小莉, 等. ARIMA模型和BP神经网络模型在我国乙型肝炎发病预测中的应用 [J]. 江苏预防医学, 2015, 26(3): 23-26.

· 声 明 ·

关于网络上出现假冒“中国工业医学杂志网站”及在线投稿的声明

中国工业医学杂志网站名称 <http://www.zggyyx.com>, 可在线投稿; 也可发至邮箱 zggyyx80@126.com 投稿。目前, 网络上出现的假冒“中国工业医学杂志网站”及在线投稿系统与本刊无关, 望广大作者和读者认真鉴别, 谨防受骗。

本刊编辑部