

ARIMAX多元时间序列模型在国家财政收入预测中的应用

陈 珊

上海工程技术大学管理学院, 上海

收稿日期: 2023年8月14日; 录用日期: 2023年9月18日; 发布日期: 2023年9月25日

摘 要

自1978年以来, 中国的改革开放政策推动了财政收入的快速增长, 并且出现了收入增长超过经济增长的局面。一个国家的财政收入是受多种因素共同影响的, 其中国民生产总值是最重要原因之一。本文结合国家财政收入的结构特点, 考虑国民生产总值对财政收入影响的前提下, 按时间序列方法对国家财政收入建立了ARIMA和ARIMAX模型, 并代入了近两年国家财政收入统计数据进行了检验, 并对检验结果进行了比较分析。

关键词

国家财政收入, 国民生产总值, ARIMA, ARIMAX模型

Application of ARIMAX Multivariate Time Series Model in State Revenue Forecasting

Shan Chen

School of Management, Shanghai University of Engineering and Technology, Shanghai

Received: Aug. 14th, 2023; accepted: Sep. 18th, 2023; published: Sep. 25th, 2023

Abstract

Since 1978, China's reform and opening-up policy has pushed the rapid growth of fiscal revenues and a situation where revenue growth exceeds economic growth has emerged. A country's fiscal revenue is affected by a combination of factors, among which the gross national product is one of

the most important reasons. In this paper, combining the structural characteristics of national fiscal revenue and considering the impact of GNP on fiscal revenue, we established ARIMA and ARIMAX models for national fiscal revenue according to the time-series method, and substituted the national fiscal revenue statistics of the last two years for testing, and made a comparative analysis of the test results.

Keywords

State Revenues, Gross National Product, ARIMA, ARIMAX Modeling

Copyright © 2023 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

国家财政收入是政府用于满足经济发展的必要条件[1]，它不仅反映了政府的经济实力，还体现了政府的社会责任感，它既可以用于维护社会稳定，又可以用于促进经济发展，从而确保社会的可持续发展。财政收入是政府践行其职能的资金保障，是保证国家有效运转的经济基础，因此随着我国经济实力的提升，国家财政收入也进入高速增长阶段。财政通过税收和各种收费等形式来筹集扩大资金，而筹集的资金又以各种形式的政府支出推动国民经济发展。在国家一系列财政活动中，财政收入和财政支出处于核心地位，而财政收入作为财政支出的先决条件更是重中之重。

因此，深入研究国家财政收入未来演变方向，为政府部门合理优化财政预算提供重要依据，这对促进国民经济稳定协调发展有着重要的现实意义。对于财政收入预测研究，国外起步早，美英等国早期采用的预测方法多样主要以时间序列预测模型为主，常用移动平均法以及自回归滑动平均混合模型等[2]。而国内的研究起步晚早期只是简单的模型，随着经济的高速发展，政府部门的大力推进，在学术界掀起热潮，其研究方法逐渐多样起来，集中在五大类：第一类多元回归模型，毛琴等[3]利用逐步回归得到影响财政收入的显著变量并通过建立多元线性回归模型进行预测与分析。姜昕[4]等通过多元线性回归方法发现税收收入对中国财政收入影响最大。第二类时间序列模型，郑鹏辉等人[5]使用自回归移动平均 ARIMA 模型对国家财政收入进行预测，总结得出此模型适合短期预测，长期预测误差会增大。第三类神经网络，随着人工智能兴起，神经网络运用于各个领域，包括金融领域，李伟[6]运用神经网络分别对财政收入和财政支出进行了预测，发现对财政支出的预测度较高。第四类灰色系统预测模型，连强[7]建立多因素灰色模型来预测河南省的财政收入，通过关键因素分析，建议加大公共服务投入促进财务收入。第五类多模型组合应用，赵海华[8]以及刘茂如等[9]都采用将灰色预测模型与神经网络相结合的方法对安徽省的财政收入进行预测分析，刘茂如等[9]还通过 Lasso 回归法筛选出财政收入的主要影响因素。

通过国内研究现状分析可以看出当前很多学者通过多种研究方法对财政收入进行预测分析，筛选各种影响财政收入的关键因素来建立模型，时间序列模型也被广泛应用于此，但大多都采用 ARIMA 模型，但当前研究很少从 GDP 这一关键影响入手分析预测，从两个经济指标的定义可以看出，一个国家的财政收入情况与国民生产总值有着极大的关联，只有整个社会创造了财富，政府才能获得更多的财政收入。因此通过加入国内生产总值(GDP)协变量建立更精确的时间序列模型 ARIMAX 模型来预测财政收入是十分有必要的。

2. 理论模型方法

2.1. ARIMA 时间序列模型

如果时间序列 $\{X_t\}$ 是一个非平稳序列, 通过对其进行 d 阶的差分运算, 可以使其成为一个平稳的时间序列, 那么就称 $\{X_t\}$ 是一个具有阶 p, d, q 的求和自回归移动平均模型, 简称 ARIMA(p, d, q) 模型:

$$\Phi(B)\nabla^d X_t = \Theta(B)\varepsilon_t$$

其中 $\nabla^d = (1-B)^d$, B 为延迟算子; ∇ 为差分运算; d 为差分阶数; $\Phi(B)$ 为残差序列的移动平均系数多项式, $\Theta(B) = 1 - \theta_1 B - \dots - \theta_q B^q$ 这是移动平均系数多项式。

2.2. 多元时间序列 ARIMAX 模型

比起 ARIMA 模型只能针对一维时间序列进行预测分析, 带有输入变量的 ARIMA 模型即 ARIMAX 模型是现在较为常用的多元时间序列分析模型, 其建模过程如下: 1) 平稳性检验。检验响应序列 y_t 和输入变量序列 $x_{1t}, x_{2t}, \dots, x_{kt}$ 是否都是平稳序列, 不平稳则进行差分直至两个序列都是平稳序列, 如果两序列都是平稳的, 则进行下一步。2) 构建回归模型, 考察两序列相关系数, 构建响应序列和输入变量序列之间回归模型, 如式(1)所示。3) 拟合残差序列 ε_t , 使用 ARIMA 模型继续提取残差序列 ε_t 中的相关信息, 最终得到动态回归模型 ARIMAX 模型, 如式(2) (3)所示。4) 模型预测。用构建的模型预测未来趋势。

$$y_t = \mu + \sum_{i=1}^k \frac{\Theta_i(B)}{\Phi_i(B)} B^{l_i} x_{it} + \varepsilon_t \quad (1)$$

$$\varepsilon_t = y_t - \left[\mu + \sum_{i=1}^k \frac{\Theta_i(B)}{\Phi_i(B)} B^{l_i} x_{it} \right] \quad (2)$$

$$\varepsilon_t = \frac{\Theta(B)}{\Phi(B)} a_t \quad (3)$$

其中 $\Phi_i(B)$ 为第 i 个输入序列的移动平均系数多项式; $\Theta_i(B)$ 为第 i 个输入序列的自回归系数多项式; ε_t 为回归残差序列; B 为延迟算子; l_i 为第 i 个输入变量 x_{it} 的延迟阶数, $\Phi(B)$ 残差序列自回归系数多项式; $\Theta(B)$ 为残差序列移动平均系数多项式; a_t 为零均值白噪声序列。

3. 基于 ARIMA 模型的全国财政收入的预测与分析

3.1. 选取数据

Table 1. National revenues (in 100 million yuan) from 1978 to 2021

表 1. 1978~2021 年的全国财政收入(单位: 亿元)

年份	国家财政收入	年份	国家财政收入	年份	国家财政收入
1978	1132.26	1993	4348.95	2008	61330.35
1979	1146.38	1994	5218.1	2009	68518.3
1980	1159.93	1995	6242.2	2010	83101.51
1981	1175.79	1996	7407.99	2011	103874.43
1982	1212.33	1997	8651.14	2012	117253.52
1983	1366.95	1998	9875.95	2013	129209.64
1984	1642.86	1999	11444.08	2014	140370.03

Continued

1985	2004.82	2000	13395.23	2015	152269.23
1986	2122.01	2001	16386.04	2016	159604.97
1987	2199.35	2002	18903.64	2017	172592.77
1988	2357.24	2003	21715.25	2018	183359.84
1989	2664.9	2004	26396.47	2019	190390.08
1990	2937.1	2005	31649.29	2020	182913.88
1991	3149.48	2006	38760.2	2021	202538.88
1992	3483.37	2007	51321.78		

如表 1 所示, 本文使用 1978 年至 2019 年的全国财政收入数据进行模型的构建, 数据来源于锐思数据库, 然后用所建立的 ARIMA 模型预测未来五年的全国财政收入数据, 2021 年和 2020 年数据用来检验预测误差。

3.2. 平稳性和纯随机检验

建立模型的第一步就是对研究的时间序列进行平稳性检验, 首先对整体数据画出时间序列, 对于有明显趋势的数据, 观察时序图就可以判断是不是平稳时间序列, 如果数据是平稳时, 直接用 ARMA 模型进行拟合; 当数据不平稳时, 先采用差分的方法将其进行处理为平稳时间之后再对非平稳时间序列进行数据进行处理, 使之变成平稳的时间序列。1978 年至 2021 年全部数据时序图如下:

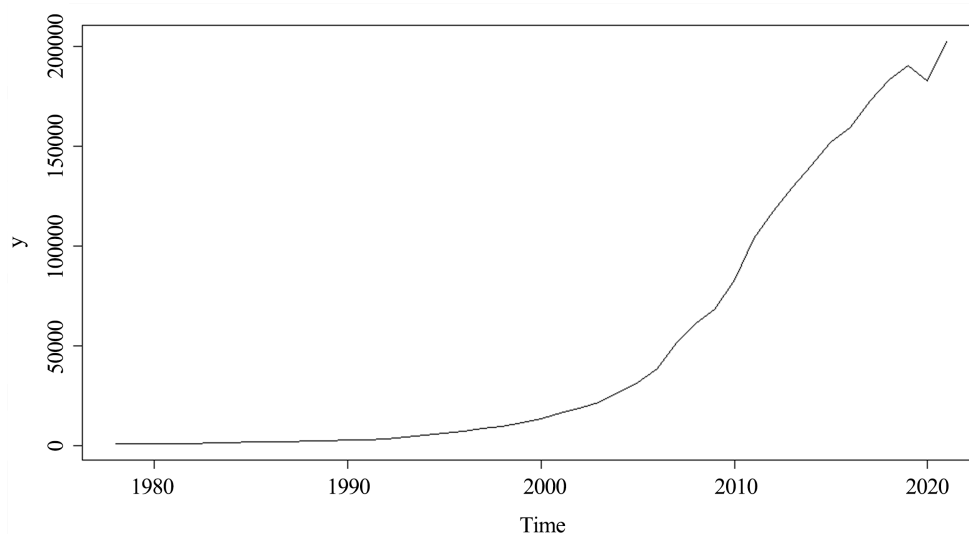


Figure 1. Time series of national revenues from 1978 to 2021

图 1. 1978~2021 全国财政收入时序图

从图 1 时序图上看, 显然这个数据是增长趋势型的数据, 处于非平稳数据, 需要进行差分运算, 才能用 ARMA 模型进行拟合。

首先我们对截取的 1978~2019 年的数据进行一阶差分运算, 一阶差分后的时序图如下图 2 所示。从图中可以看出, 一阶差分后的时序图仍然具有增长趋势, 也就是说一阶差分后的数据仍然是不平稳的, 并没有充分提取这个序列长期趋势的信息, 还需要再进一步差分。

二阶差分后的时序图如图 3，此时时序图中的曲线没有向上增长的趋势，说明二阶差分比较充分的提取了国家财政收入的长期趋势信息，从时序图来看，初步判定二阶差分后的数据是平稳，但由于图像观察太具有主观性，我们还是要进一步进行单位根检验——ADF 检验，ADF 检验通过后还需要进行纯随机检验，也就是白噪声检验。

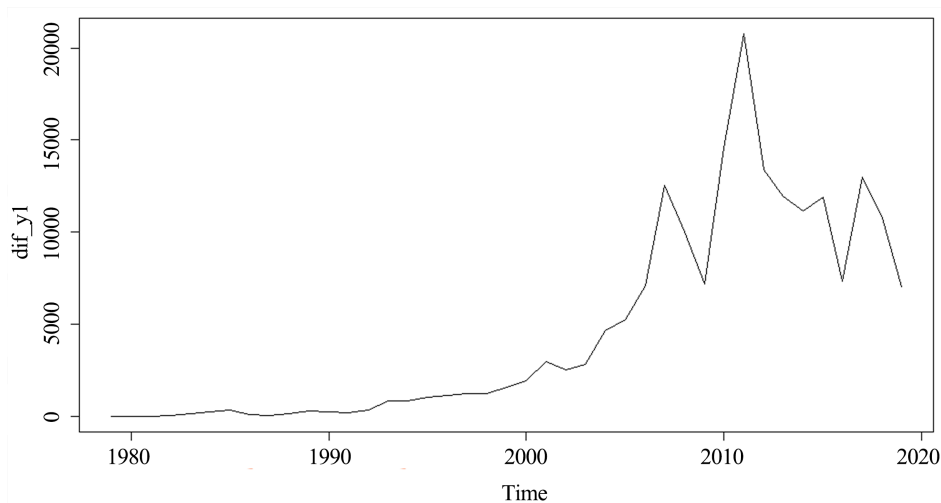


Figure 2. State revenue data post first order differential time series plot from 1978 to 2019
图 2. 1978~2019 国家财政收入数据一阶差分后时序图

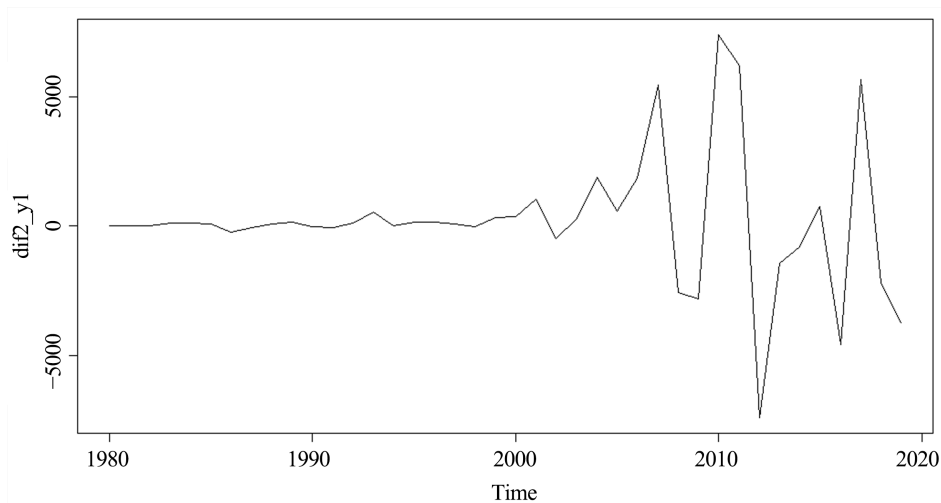


Figure 3. State revenue data post second order differential time series plot from 1978 to 2019
图 3. 1978~2019 国家财政收入数据二阶差分后时序图

由 ADF 检验可知，以上 3 种大类的前 3 个子类型的序列结构的 p 值均小于显著性水平 ($\alpha = 0.05$)，所以可以认为该序列是显著平稳的。白噪声的检验可知延迟 6 阶的 LB 统计量的 P 值大于显著性水平 α ，所以看该序列拒绝纯随机性原假设，是平稳非白噪声序列。没有检验其他延迟阶数是因为如果平稳序列短期延迟阶数都不存在显著相关关系，通常长期延迟就更不会存在显著的相关关系。

3.3. ARIMA 模型建立

首先做出国家财政收入二阶差分序列的自相关图和偏自相关图，如图 4 和图 5 所示。

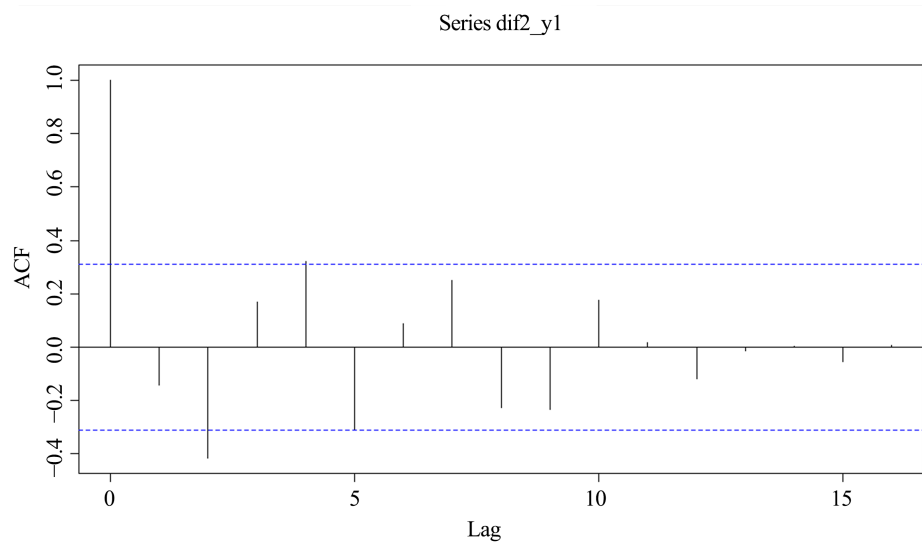


Figure 4. Serial autocorrelation coefficients after second-order differencing of state revenues from 1978 to 2019

图 4. 1978~2019 年国家财政收入二阶差分后序列自相关系数

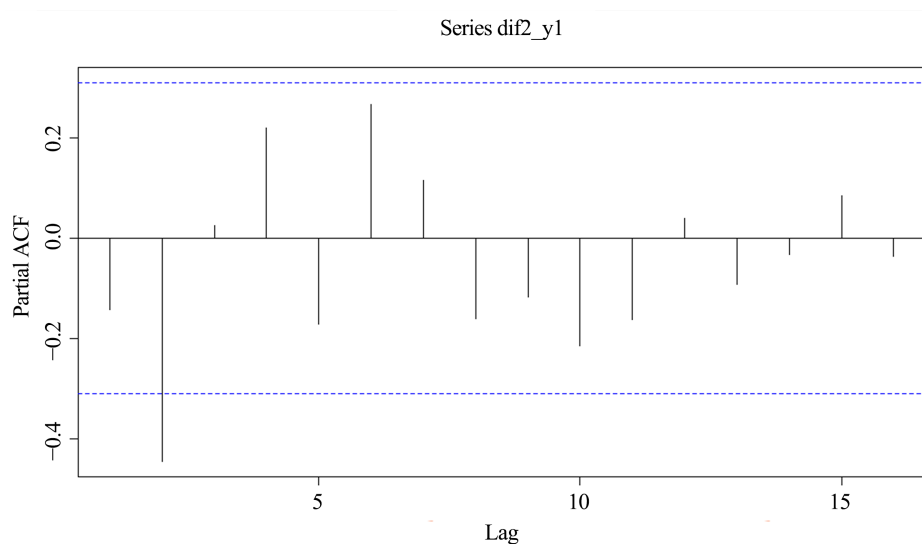


Figure 5. Partial autocorrelation coefficients of the second-order post-differential series of state revenues from 1978 to 2019

图 5. 1978~2019 年国家财政收入二阶差分后序列偏自相关系数

由样本自相关图 4 和样本偏自相关图 5 可知，二阶差分后的全国财政收入数据的偏自相关函数显示在延迟 4 阶以后，自相关系数都落在了 2 倍的标准差范围内，并在此范围内波动、自相关函数呈现较强的拖尾性显示在延迟 2 阶以后，自相关系数全部衰减在 2 倍的标准差范围内，并在此范围内波动，依据经验可以对全国财政收入序列构建 ARIMA(2,2,4) 模型。但由于自己判断主观性比较强，所以用 auto.arima 函数进行重新定阶，之后采取 AIC 准则定阶的方法，比较选出 AIC 最小的模型。图定阶法构建的 ARIMA(2,2,4)模型的 AIC 值为 737.48，而 ARIMA(0,2,4) auto.arima 定阶的模型 AIC 数值为 733.98，更小一点，因此应该对原序列—全国财政收入序列构建 ARIMA(0,2,4)模型。最终原始序列输出形式为：

$$\nabla^2 x_t = \varepsilon_t + 0.0113\varepsilon_{t-1} - 0.3918\varepsilon_{t-2} - 0.0522\varepsilon_{t-3} + 0.7936\varepsilon_{t-4}$$

3.4. ARIMA 模型检验

对选取的全国财政数据建立 ARIMA(2,2,4)模型后, 需要对其显著性进行检验, 本文采取的是 `ts.diag` 进行检验, 如图 6 所示。

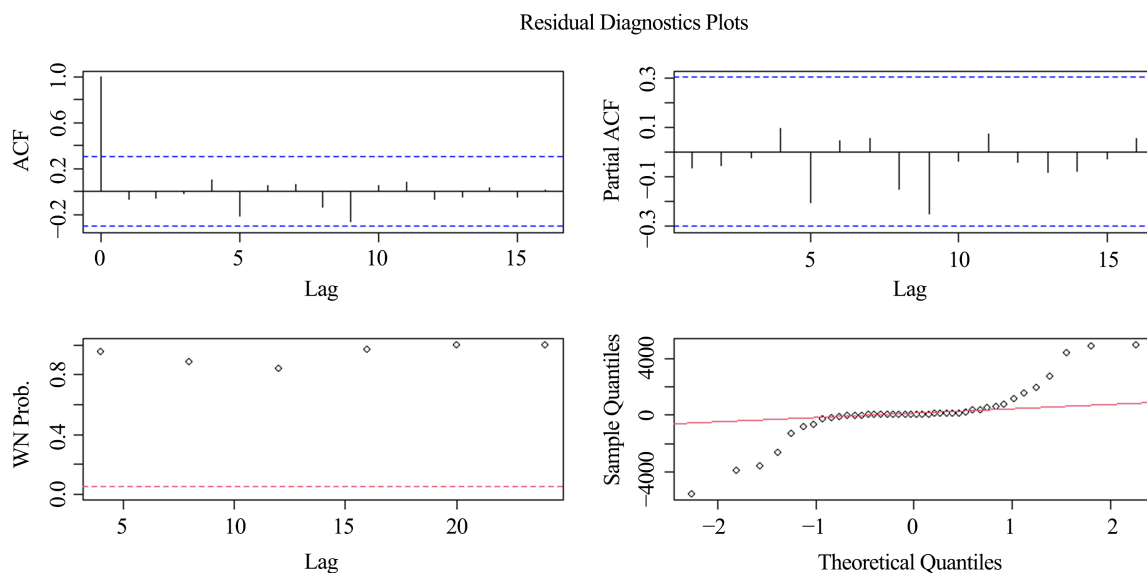


Figure 6. Significance test of ARIMA(0,2,4) model for state revenue from 1978 to 2019

图 6. 1978~2019 年国家财政收入 ARIMA(0,2,4)模型显著性检验

从上图可知考察残差序列白噪声检验结果(图 6), 可以看出各阶延迟下的白噪声检验统计量的 P 值都显著大于 0.05, 可以认为拟合模型的残差序列属于白噪声序列, 拟合模型显著成立。

3.5. ARIMA 模型预测及分析

用构建的 ARIMA(0,2,4)模型预测 2020 至 2024 年的全国财政收入数据, 其时序图如下图 7:

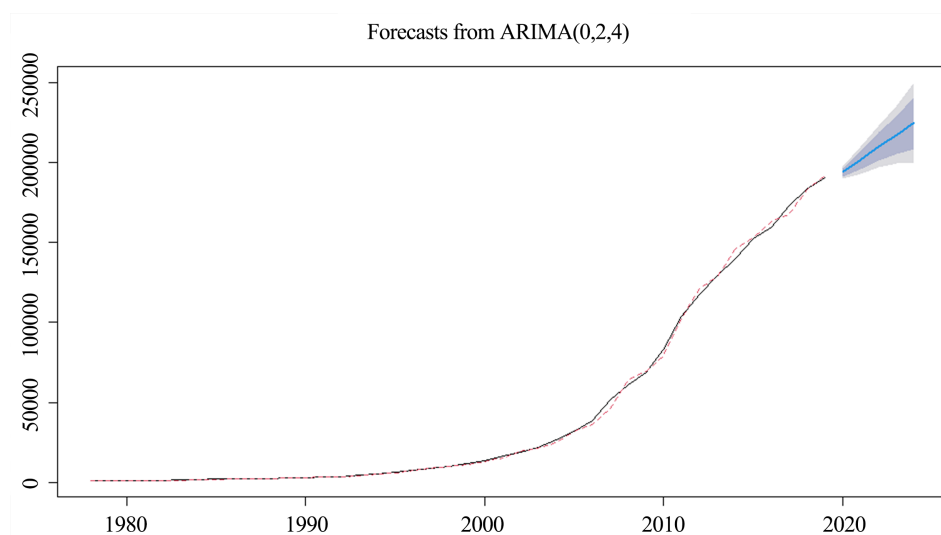


Figure 7. ARIMA model projections of state revenues from 2020 to 2024

图 7. ARIMA 模型对国家财政收入 2020~2024 年的预测图

用 2020 至 2021 年真实的国家财政收入数据对比，其结果见下表 2:

Table 2. 2020~2021 national fiscal revenue ARIMA forecasts vs. real values

表 2. 2020~2021 年全国财政收入 ARIMA 预测值与真实值对比表

年份	真实值	预测值	相对误差
2020	182913.88	193856.020	-10942.140
2021	202538.88	201736.984	801.896

预测图来看，其模型拟合度还是不错的，图中实线为观察值，虚线为模拟拟合值，阴影部分实线为预测值，深色阴影为序列 80% 置信区间，浅色阴影为序列 95% 置信区间。

从数据对比看 2020 年预测效果较差，但 2021 年预测效果较好，可能原因是因为 2020 年疫情大面积爆发，对财政收入影响较大。

4. 基于 ARIMAX 模型的全国财政收入的预测与分析

4.1. 选择数据

本文选取的是 1978~2021 年的国家财政收入和国民生产总值的数据，数据见下表 3:

Table 3. State revenues and gross domestic product from 1978 to 2021 (100 million yuan)

表 3. 1978~2021 年的国家财政收入与国内生产总值(单位: 亿元)

年份	国家财政入	GDP	年份	国家财政入	GDP
1978	1132.26	3678.70	2000	13395.23	100280.14
1979	1146.38	4100.45	2001	16386.04	110863.12
1980	1159.93	4587.58	2002	18903.64	121717.42
1981	1175.79	4935.83	2003	21715.25	137422.03
1982	1212.33	5373.35	2004	26396.47	161840.16
1983	1366.95	6020.92	2005	31649.29	187318.90
1984	1642.86	7278.50	2006	38760.2	219438.47
1985	2004.82	9098.95	2007	51321.78	270092.32
1986	2122.01	10376.15	2008	61330.35	319244.6128
1987	2199.35	12174.59	2009	68518.3	348517.7437
1988	2357.24	15180.39	2010	83101.51	412119.2558
1989	2664.9	17179.74	2011	103874.43	487940.1805
1990	2937.1	18872.87	2012	117253.52	538579.9535
1991	3149.48	22005.63	2013	129209.64	592963.2295
1992	3483.37	27194.53	2014	140370.03	643563.1045
1993	4348.95	35673.23	2015	152269.23	688858.218
1994	5218.1	48637.45	2016	159604.97	746395.0595
1995	6242.2	61339.89	2017	172592.77	832035.9486
1996	7407.99	71813.63	2018	183359.84	919281.1291
1997	8651.14	79715.04	2019	190390.08	986515.2023
1998	9875.95	85195.51	2020	182913.88	1013567
1999	11444.08	90564.38	2021	202538.88	1143669.7

数据同样来源于锐思数据库,本节应用多元动态回归模型——协整模型,对全部数据选取 1978~2019 年的国家财政收入数据和国内生产总值进行建模,其中国内生产总值(GDP)作为输入序列,全国财政收入作为响应序列。做出二者的时序图如图 8:

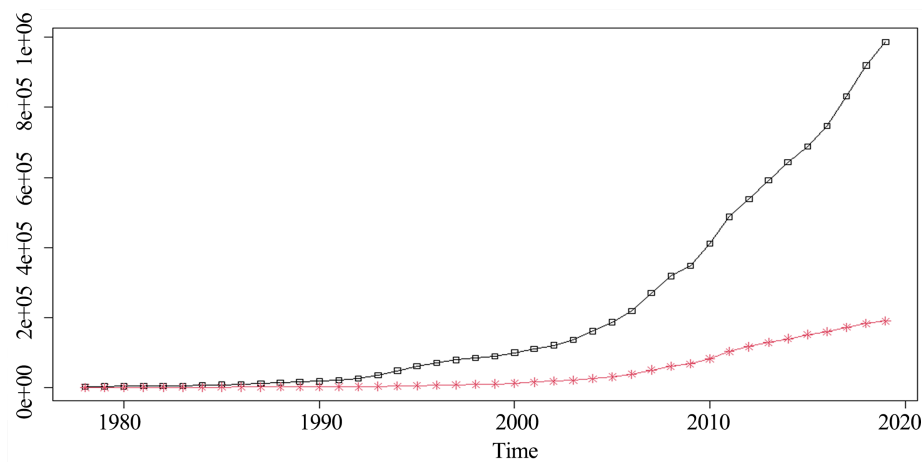


Figure 8. Time series of state revenues and GDP from 1978 to 2019
图 8. 1978~2019 年国家财政收入和国内生产总值的时序图

从两个时间序列的时序图可以看出 1978~2019 年的国家财政收入和国内生产总值数据呈现出较稳定的长期均衡关系,当国内生产总值增长时,全国财政收入也会相应增长,因而可以考虑对其建立多元动态回归模型,即 ARIMAX 模型。

4.2. 数据的检验

I. 首先对数据进行同阶单整检验,同阶单整是协整检验的前提,由于 GDP 数据与之前国家财政收入的数据一样具有长期趋势,因此我们同样对 GDP 数据做一阶差分和二阶差分,时序图如图 9 和图 10。之后进行平稳性检验和白噪音检验,结果显示二阶差分后的国民生产总值序列是平稳非白噪声时间序列。因此国家财政收入序列和国民生产总值序列均在二阶差分后平稳,两个序列满足同阶单整的条件,即都是二阶单整的。

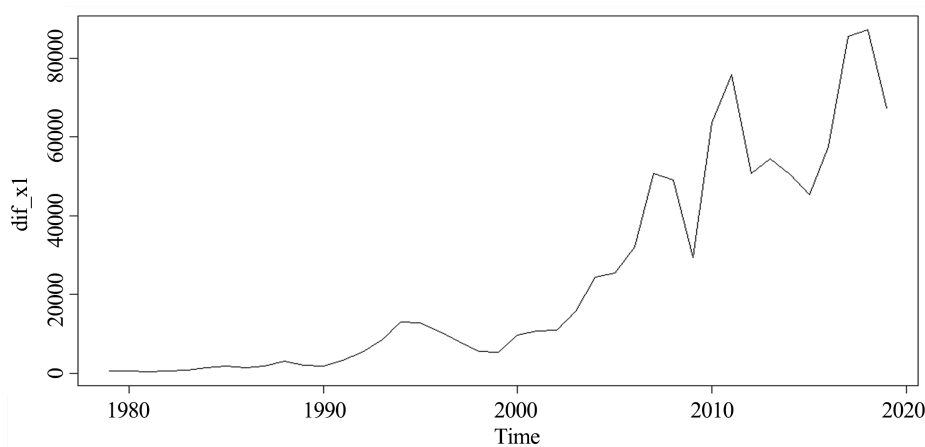


Figure 9. GDP data first-order post-differential time-series plot from 1978 to 2019
图 9. 1978~2019 GDP 数据一阶差分后时序图

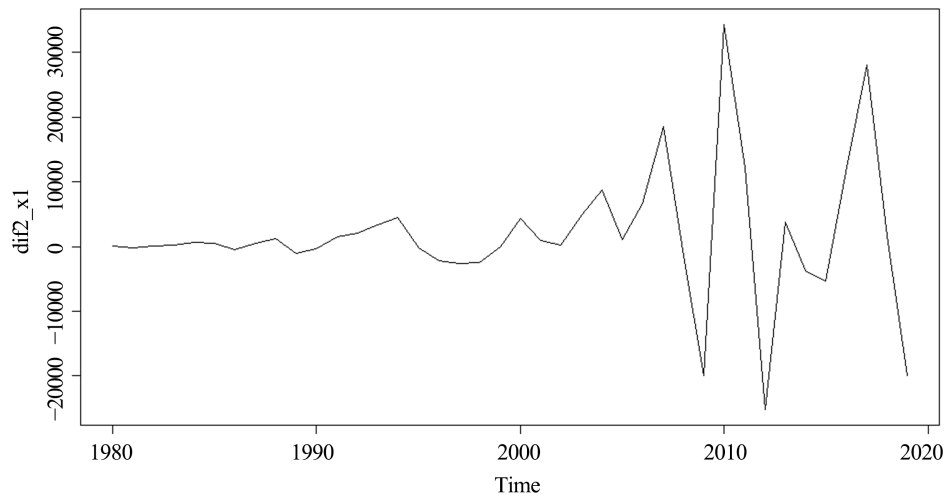


Figure 10. GDP data second-order post-differential time series plot from 1978 to 2019
图 10. 1978~2019 GDP 数据二阶差分后时序图

II. 再进行协整检验，首先绘制互相关图，结果如图 11 所示，拟合回归模型如图 12 所示，可以发现国家财政收入序列与国家生产总值序列存在着很强的相关性，在延迟阶数为零时，两者的协相关系数最大，即当期 GDP 对国家财政收入影响达到最大，因此，在构建归国家财政收入与国民生产总值回归模型时，自变量使用的是 GDP 当期序列同期。

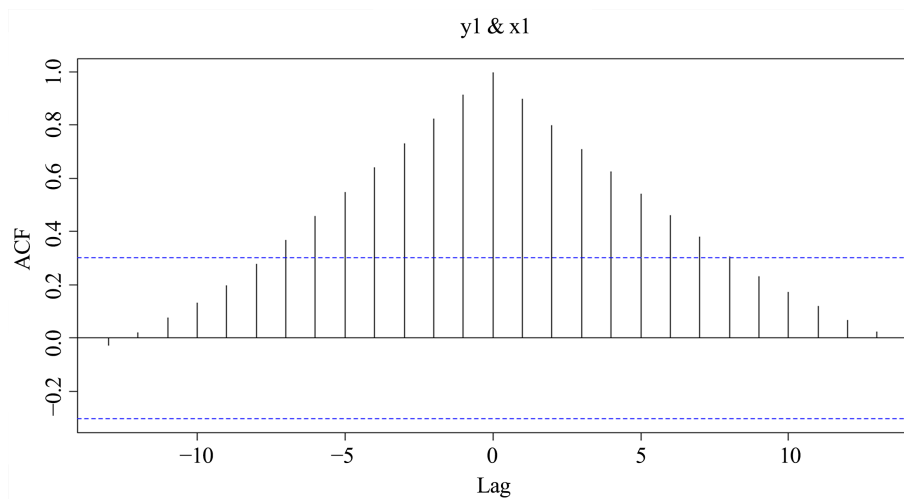


Figure 11. Correlation of state revenues and gross national product from 1978 to 2019
图 11. 1978~2019 年国家财政收入和国民生产总值的互相关图

```
call:
arima(x = y1, xreg = x1, include.mean = F)

coefficients:
      x1
    0.2049
s.e.  0.0023

sigma^2 estimated as 27917174: log likelihood = -419.64, aic = 843.27
```

Figure 12. Fitted regression model for state revenues and GDP from 1978 to 2019
图 12. 1978~2019 年国家财政收入和国民生产总值的拟合回归模型

之后对回归残差序列进行平稳性检验，回归残差序列 ADF 检验可以观察到，类型 1 中延迟 1, 2, 3 阶的检验结果的 p 值都是显著小于 0.05，所以可以认为回归残差序列是平稳的，也就是说国家财政收入序列和国民生产总值序列之间存在协整关系。因此可以对国家财政收入序列和国民生产总值序列建立多元动态回归模型而不用担心虚假回归的问题。

III. 之后对残差序列进行白噪声检验，LB 检验显示 p 值显著小于 0.05，回归残差序列不是白噪声序列还需进一步提取残差序列中蕴含的信息。

III. 拟合协整动态回归模型即 ARIMAX 模型，首先绘制残差序列自相关图和偏自相关图，如图 13 和图 14 所示，残差自回归系数拖尾，偏自回归系数 2 阶结尾，对残差序列拟合 AR(2)模型，拟合协整动态回归模型如下：

$$y_t = 0.1956x_t + \frac{\varepsilon_t}{1 - 1.6995B + 0.8011B^2}$$

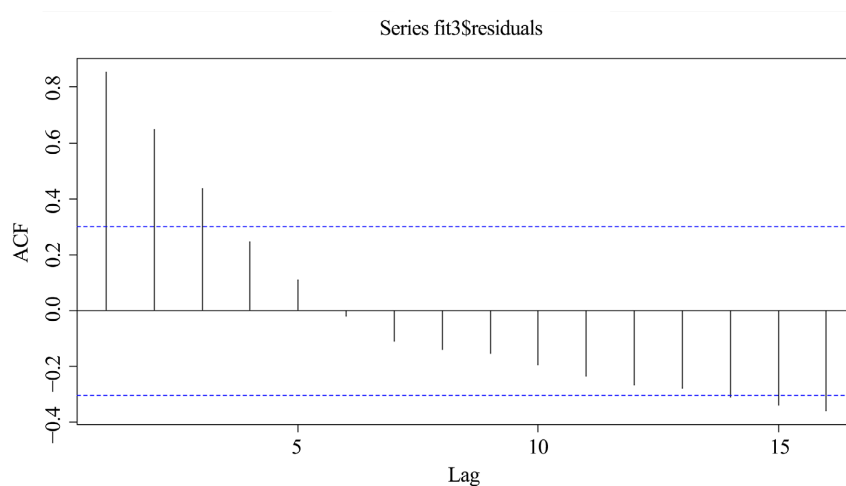


Figure 13. Autocorrelation plot of regression residual series for state revenue and GNP from 1978 to 2019

图 13. 1978~2019 年国家财政收入和国民生产总值的回归残差序列自相关图

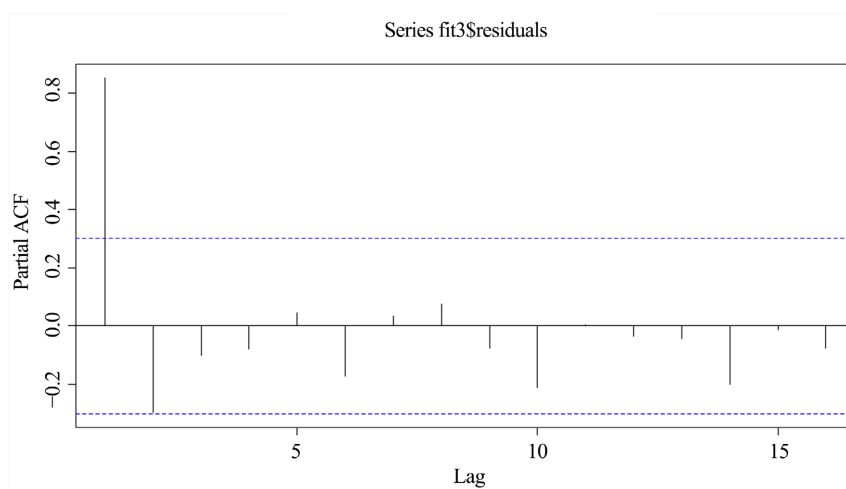


Figure 14. Partial autocorrelation of the regression residual series for state revenue and GDP from 1978 to 2019

图 14. 1978~2019 年国家财政收入和国民生产总值的回归残差序列偏自相关图

IV. 模型显著性检验，如图 15 可以看出残差序列为白噪声序列，说明拟合协整动态回归模型显著成立，我们可以利用这个模型进行预测分析。

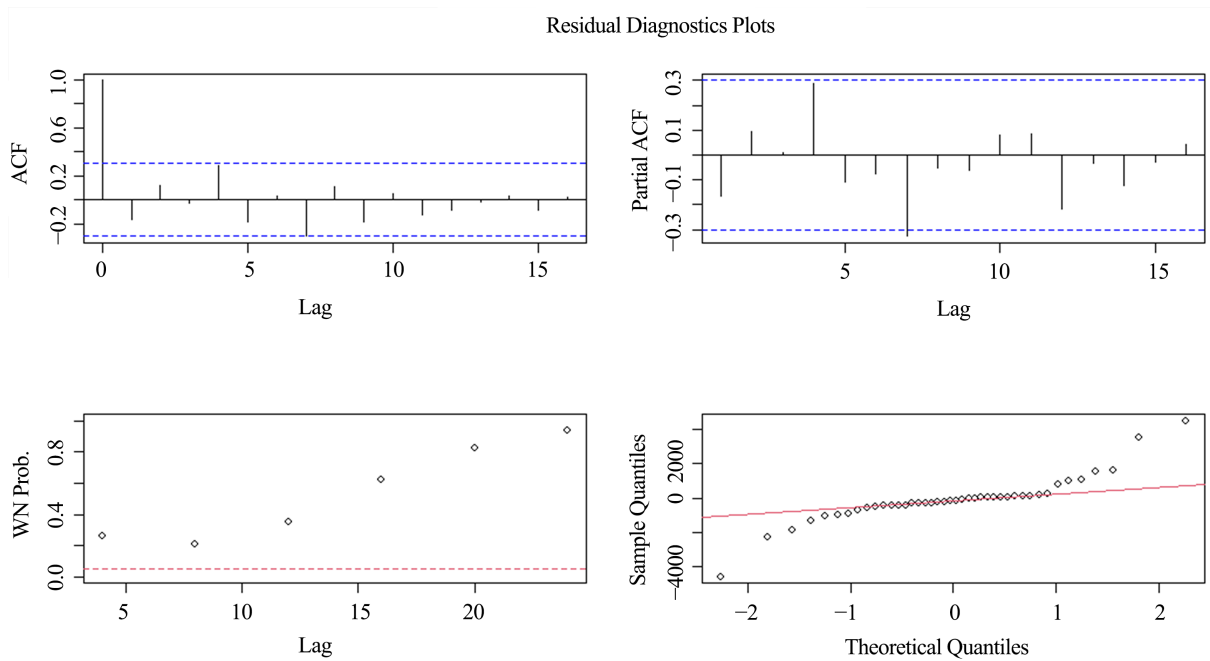


Figure 15. Significance test of the fitted cointegration dynamic regression model for state revenue and GDP from 1978 to 2019
图 15. 1978~2019 年国家财政收入和国民生产总值的拟合协整动态回归模型显著性检验

V. 序列预测，使用回归模型预测国家财政收入时，需要先获得 GDP 输入序列的未来预测值，可以基于单变量预测，在带入回归模型里，就可以获得响应序列的预测值，预测数值如表 4，预测效果图如图 16。

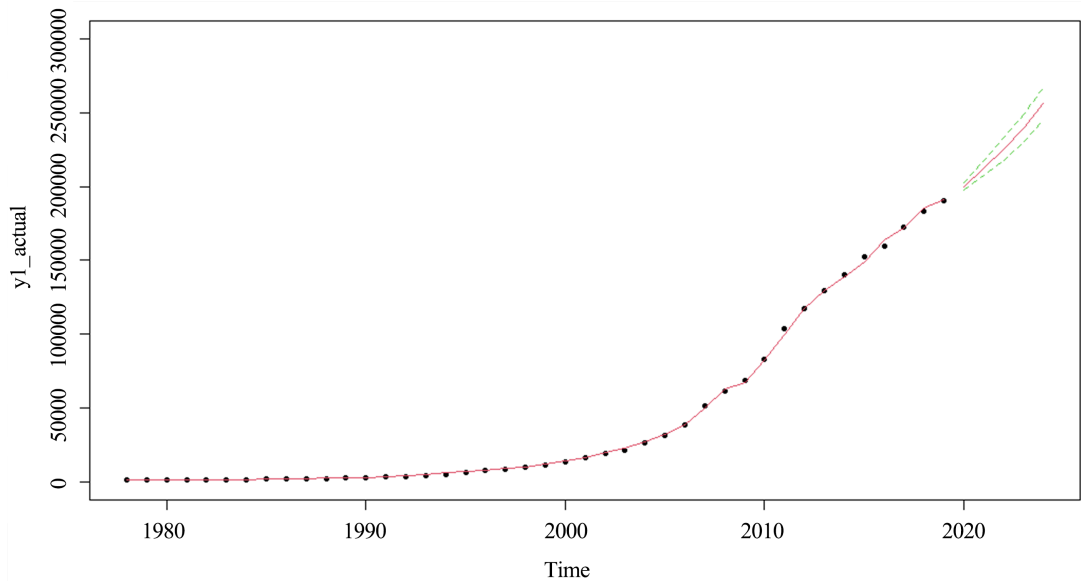


Figure 16. Effect of state revenue projections from 1978 to 2019
图 16. 1978~2019 年国家财政收入预测效果图

Table 4. Projections of state revenues from 2020 to 2024**表 4.** 2020~2024 年国家财政收入预测值

年份	预测值	L95	U95
2020	199965.8	197319.5	202612.1
2021	212561.4	207343.1	217779.7
2022	225491.6	217892.7	233090.5
2023	239905.0	230354.6	249455.4
2024	256091.8	245117.2	267066.5

5. 总结

本文构建 ARIMA 和 ARIMAX 模型, 探究国家财政收入与 GDP 之间的规律, 进一步对比近两年预测值和真实值的误差, 精准预测未来几年的国家财政收入。从 ARIMA 模型预测值与 ARIMAX 模型预测值对比看, 两者预测值相差不多, 从 2020 和 2021 真实值对比看, 发现 ARIMA 模型预测误差更小。这有可能是由于没有进行误差修正模型的原因, 误差修正模型是一个负反馈机制。以及预测这两年由于疫情影响较大, 对比不出真实的预测拟合效果, 这些都是本文后续需要进一步改进的地方。但 ARIMA 模型适用于短期预测, 长期误差会越来越大, 从预测效果图对比看, ARIMAX 模型上下浮动差是比 ARIMA 模型要小的。

参考文献

- [1] 王燕. 应用时间序列分析[M]. 北京: 中国人民大学出版社, 2005.
- [2] 吴笑晗, 王美桃. 国外财政收入预测实践与经验借鉴[J]. 公共财政研究, 2016(3): 34-44.
- [3] 毛琴, 李明江, 刘彦. 基于逐步回归法的国家财政收入数据回归模型分析[J]. 电子技术与软件工程, 2013(19): 227-228.
- [4] 姜昕, 赵洋. 中国财政收入影响因素的实证研究[J]. 中国乡镇企业会计, 2019(8): 105-106.
- [5] 郑鹏辉, 单锐, 陈静. 时间序列分析在我国财政收入预测中的应用[J]. 重庆文理学院学报(自然科学版), 2008, 27(2): 15-18.
- [6] 李伟. 神经网络在财政数据中的应用[D]: [硕士学位论文]. 长春: 吉林大学, 2011.
- [7] 连强. 基于多因素灰色模型的河南省财政收入预测[J]. 中国市场, 2020(22): 34-36.
- [8] 赵海华. 基于灰色 RBF 神经网络的多因素财政收入预测模型[J]. 统计与决策, 2016(13): 79-81.
- [9] 刘茂茹, 王丰效. 基于 GM(1,1)-神经网络模型的安徽省财政收入预测[J]. 高师理科学刊, 2022, 42(6): 17-22.