

基于 ARIMA 模型上的 CPI 的建模及预测

黄梓鹏

华南农业大学数学与信息学院, 广东 广州 510640

DOI:10.61369/ASDS.2025040002

摘要 : 本文基于 ARIMA 模型对我国消费者价格指数 (CPI) 进行建模与预测, 结合最小二乘法和贝叶斯估计方法对比分析模型性能。通过差分处理非平稳时间序列, 选择 ARIMA(3,1,0) 模型进行参数估计, 利用 AIC 准则和残差检验验证模型有效性。实证分析显示, 最小二乘法与贝叶斯方法均能较好拟合 CPI 趋势, 预测值与实际值偏差在可控范围内。贝叶斯方法通过引入先验分布和 MCMC 抽样, 增强了参数不确定性建模能力, 尤其适用于小样本或高波动场景。研究结果表明, 两种方法在 CPI 预测中均有效, 在实证例子中, 贝叶斯推断在融合先验信息与动态更新后验分布方面更具优势, 预测值也能为宏观经济政策评估和通胀调控提供了理论支持。

关键词 : ARIMA 模型; CPI 预测; 贝叶斯估计; 最小二乘法

Modeling and Forecasting CPI Based on the ARIMA Model

Huang Zipeng

College of Mathematics and Information, South China Agricultural University, Guangzhou, Guangdong 510640

Abstract : This paper models and forecasts China's Consumer Price Index (CPI) using the ARIMA model, comparing the performance of the least squares method and Bayesian estimation. To handle non-stationary time series, differencing is applied, and the ARIMA(3,1,0) model is selected based on parameter estimation. The model's validity is confirmed through AIC criterion and residual diagnostics. Empirical analysis shows that both the least squares and Bayesian methods effectively capture the CPI trend, with prediction errors within a controllable range. The Bayesian approach, incorporating prior distributions and MCMC sampling, enhances the ability to model parameter uncertainty, making it especially suitable for small samples or highly volatile scenarios. The results demonstrate that both methods are effective for CPI forecasting, while Bayesian inference shows advantages in integrating prior knowledge and dynamically updating posterior distributions. The predicted values provide theoretical support for macroeconomic policy assessment and inflation control.

Keywords : ARIMA; CPI forecasting; Bayesian estimation; LSM

引言

消费者价格指数 (CPI, Consumer Price Index) 是衡量居民消费商品及服务价格水平变动的关键指标, 其波动受到多重因素影响, 并与货币政策、产业结构、国际环境等密切相关, 它通常用来衡量一篮子消费品和服务的价格变动, 反映通货膨胀或通货紧缩的水平, 因此预测 CPI 还可以帮助评估已实施的经济政策的效果。例如, 通过比较政策实施前后的 CPI 预测值与实际值, 可以判断政策是否达到了预期的调控目标, 为未来政策调整提供依据。

消费者价格指数 (CPI) 的预测常采用 ARIMA (自回归积分移动平均) 模型, 该模型通过分析时间序列的趋势和季节性进行预测。分析数据的自相关性和移动平均特性来预测的特点, 适用于平稳或可差分平稳的时间序列数据, 国内外就有很多通过 ARIMA 模型预测的案例, 如基于 ARIMA 模型, 马瑶等 (2021) 对中、美、德三国 2000-2020 年月度 CPI 数据建模分析, 验证其预测能力, 发现政府疫情干预对 CPI 影响显著^[1-2]。

在 ARIMA (Autoregressive Integrated Moving Average model) 模型中, 最小二乘法主要用于参数估计阶段。当构建 ARIMA 的自回归和移动平均部分时, 可通过最小二乘法来估计模型参数, 通过最小化预测值与实际观测值的残差平方和来优化模型拟合。接着, 通过 AIC、BIC 准则定阶后, 最后用最小二乘法于参数估计^[3-5]。

贝叶斯推断是一种基于贝叶斯定理的统计方法, 通过结合先验知识与新数据动态更新概率分布, 形成后验概率。其核心思想是先将先

验分布（初始假设）与似然估计（观测数据）结合，利用贝叶斯公式计算后验分布（更新后的概率）^[6]。贝叶斯推断将概率视为对事件发生的主观信心而非长期频率，尤其在小数据场景下能保留不确定性，因此可通过马尔可夫链蒙特卡罗（MCMC）等方法进行参数推断，引入先验知识，相较于最小二乘法的点估计，可增强对不确定性的建模能力，更适合处理小样本或高不确定性场景^[7]，但在可找到的案例中使用较少。

一、模型及分布类型介绍

（一）模型介绍

1. ARIMA 模型

ARIMA(p,d,q) 称为差分自回归移动平均模型，AR (Autoregressive) 是自回归，p 为自回归项；MA (Moving Average) 为移动平均，q 为移动平均项数；I (Integrated) 指的是差分整合，d 为时间序列成为平稳时所做的差分次数 ARIMA(p,d,q) 模型模型具有如下的架构^[8]：

$$\begin{cases} \Phi(B)\nabla^d x_t = \Theta(B)\varepsilon_t \\ E(\varepsilon_t) = 0, \text{Var}(\varepsilon_t) = \sigma_\varepsilon^2, E(\varepsilon_t \varepsilon_s) = 0, s \neq t \\ E_s \varepsilon_t = 0, \forall s < t \end{cases}, \quad (1.1)$$

式中：

$\Phi(B) = 1 - \phi_1 B - \dots - \phi_p B^p$ ，为平稳可逆 ARMA(p,q) 模型的自回归系数多项式；

$\Theta(B) = 1 - \theta_1 B - \dots - \theta_q B^q$ ，为平稳可逆 ARMA(p,q) 模型的移动平滑系数多项式；

式 (1.1) 可以简记为：

$$\nabla^d x_t = \frac{\Theta(B)}{\Phi(B)} \varepsilon_t \quad (1.2)$$

其中 B 为延迟算子 $Bx_t = x_{t-1}$ ，即，为零均值白噪声序列。

显然，ARIMA 模型的实质就是差分运算与 ARMA 模型的组合。所以我们在处理非平稳的时间序列时，只要对其进行差分运算就可得到一平稳的时间序列。

此外，可以看出由于 $\{\varepsilon_t\}$ 为零均值白噪声序列，它要求影响每个观察数据的因素大致一样，本文选用 ARIMA 模型对 CPI 进行研究的原因也是基于 CPI 具有这个性质^[6]。

2. 最小二乘估计

最小二乘估计的原理是使残差平方和达到最小的那组参数值即为最小二乘估计值。由上文可知，ARIMA 模型是针对非平稳时间序列的模型，为了使其变成平稳一般使用差分方法，然后对平稳的时间序列使用 ARMA 模型进行建模，所以实际上它们的原理是相同的，为了简化说明，使 n 阶差分 $x_t = (1 - B)^n y_t$ ，然后直接套用 ARMA 模型进行说明：

对于 ARMA(p,q)，记

$$\tilde{\beta} = (\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q)^T \quad (1.3)$$

$$F_t(\tilde{\beta}) = \phi_1 x_{t-1} + \dots + \phi_p x_{t-p} - \theta_1 \varepsilon_{t-1} - \dots - \theta_q \varepsilon_{t-q} \quad (1.4)$$

残差项为：

$$\varepsilon_t = x_t - F_t(\tilde{\beta}) \quad (1.5)$$

残差平方和为：

$$\begin{aligned} Q(\tilde{\beta}) &= \sum_{t=1}^n \varepsilon_t^2 \\ &= \sum_{t=1}^n (x_t - \phi_1 x_{t-1} - \dots - \phi_p x_{t-p} - \theta_1 \varepsilon_{t-1} - \dots - \theta_q \varepsilon_{t-q})^2 \end{aligned} \quad (1.6)$$

当残差平方和达到最小时，所得的参数即为 $\tilde{\beta}$ 的最小二乘估计，即：

$$\begin{aligned} Q(\tilde{\beta}) &= \min Q(\tilde{\beta}) \\ &= \min \sum_{t=1}^n (x_t - \phi_1 x_{t-1} - \dots - \phi_p x_{t-p} - \theta_1 \varepsilon_{t-1} - \dots - \theta_q \varepsilon_{t-q})^2 \end{aligned} \quad (1.7)$$

由此看见最小二乘估计充分应用了每一个观察值所提供的信息，因而其估计精度高，计算复杂度不高，是人们使用较为普遍的方法^[8-9]。

3. 贝叶斯估计方法

本节为方便描述，直接用 ARMA 模型进行说明。设模型的参数向量为 $X = (x_1, x_2, \dots, x_t)$ ， $\eta = (\phi, \theta)$ ，其中 X 为观察序列， $\phi = (\phi_1, \dots, \phi_p)$ 是自回归参数向量， $\theta = (\theta_1, \dots, \theta_q)$ 是移动平均参数向量。

贝叶斯理论是利用与未知变量有关的统计数据来获得未知变量的条件分布函数，即未知变量的后验概率密度函数，基本公式为：

$$h(\eta | X) = \frac{\pi(\eta)f(X|\eta)}{g(X)} = \frac{f(X|\eta)\pi(\eta)}{\int f(X|\eta)\pi(\eta)d\eta} \quad (1.8)$$

其中， $h(\eta | X)$ 为后验概率密度函数， $f(X|\eta)$ 为条件密度函数， $g(X)$ 为先验分布函数， $\pi(\eta)$ 为 η 的概率密度函数。

而且 η 的贝叶斯估计为

$$\hat{\eta} = E[h(\eta | X)] \quad (1.9)$$

显然，即使所有的条件都给出，但是要计算式 (2.9) 存在困难，因此人们引入了马尔科夫蒙特卡罗法^[10] (Markov Chain Monte Carlo)。

MCMC 方法的基本思想是：建立一个马尔科夫链对未知变量 U_t 进行模拟，当链达到稳态分布时即得所求的后验分布。随机点 U_t 来自于分布 $\pi(U)$ ，由不同的抽样方法得到了不同的 MCMC 方法，如 Metropolis-Hastings 方法、Gibbs 抽样方法以及各种复合方法^[11]。

Gibbs 抽样是最简单也是应用最广泛的一种抽样方法。在上述假设条件下，首先给定各个参数的初始值 $\hat{\eta}^{(0)} = (\phi^{(0)}, \theta^{(0)})$ ，然后从上面的分析得到的各个参数 ϕ, θ 的条件后验密度中循环抽取的 n 次，Gibbs 抽样的第一次迭代如下：

$$\begin{aligned}
\tilde{\eta}_1^{(1)} &\sim f(\tilde{\eta}_1/\tilde{\eta}_2^{(0)}, \dots, \tilde{\eta}_s^{(0)}) \\
\tilde{\eta}_2^{(1)} &\sim f(\tilde{\eta}_2/\tilde{\eta}_1^{(1)}, \tilde{\eta}_3^{(0)}, \dots, \tilde{\eta}_s^{(0)}) \\
&\vdots \\
\tilde{\eta}_k^{(1)} &\sim f(\tilde{\eta}_k/\tilde{\eta}_1^{(1)}, \dots, \tilde{\eta}_{k-1}^{(1)}, \tilde{\eta}_{k+1}^{(0)}, \dots, \tilde{\eta}_s^{(0)}) \\
&\vdots \\
\eta_s^{(1)} &\sim f(\eta_k/\eta_1^{(1)}, \eta_2^{(1)}, \dots, \eta_{s-1}^{(1)})
\end{aligned} \tag{1.10}$$

其中 \sim 表示左边从右边抽取。以上完成了一次 Gibbs 迭代过程，即完成了由 $\eta^{(0)} = (\phi^{(0)}, \theta^{(0)})$ 到 $\eta^{(1)} = (\phi^{(1)}, \theta^{(1)})$ 经过 n 次迭代，则可以得到各参数的 n 次抽样值。当马尔可夫链在循环迭代 $m(m < n)$ 次后收敛时，由蒙特卡罗积分公式可以得到各个参数的后验均值和方差分别为

$$\begin{cases} E(\tilde{\eta}_k) \approx \frac{1}{n-m} \sum_{t=m+1}^n \tilde{\eta}_k^{(t)}, \\ Var(\tilde{\eta}_k) \approx \frac{1}{n-m} \sum_{t=m+1}^n (\tilde{\eta}_k^{(t)})^2 - \left(\frac{1}{n-m} \sum_{t=m+1}^n \tilde{\eta}_k^{(t)} \right)^2 \end{cases} \tag{1.11}$$

其中 $\tilde{\eta}_k$ 表示参数向量 $\tilde{\eta} = (\phi, \theta)$ 中的任意参数。

二、实证分析

1. 建模前数据处理

CPI 数据的收集比较容易，政府部门每隔一段固定的时间都会发布该阶段的 CPI 数据，以供社会参考。

下面是我国 2021 年 6 月到 2024 年 6 月的 CPI 数据：

表 2.1 2021 年到 2024 年 1-6 月我国 CPI 指数

	1月	2月	3月	4月	5月	6月
2021						101.1
2022	101.5	100.7	100.8	101	101.3	100.9
2023	99.2	99.7	99.5	99.8	100	100.1
2024	100.7	100.1	100.1	100.3	100.3	100.2

表 2.2 2021 年到 2024 年 7-12 月我国 CPI 指数

	7月	8月	9月	10月	11月	12月
2021	101	101.1	100.7	100.9	101.5	102.3
2022	100.4	99.8	99.7	100.2	100.9	101.5
2023	99.7	100	100.1	99.8	99.5	99.7
2024	100.5					

首先，先得到数据的序列图以及自相关系数图（图 2.1），明显地序列并不平稳，所以我们进行一、二阶差分运算，得到新的序列图以及自相关系数图。

由时序图 2.2、2.3 可以看出差分后的图像趋向于平稳，再由自相关系数图得自相关系数基本在两倍的标准差内波动，进一步确定了差分后的序列的平稳性。

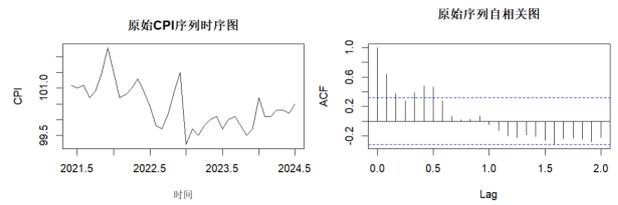


图 2.1 CPI 数据时序图、自相关系数图

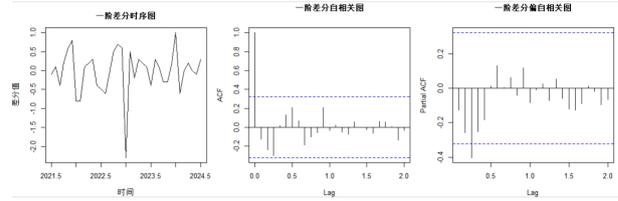


图 2.2 CPI 数据一阶差分时序图、自相关系数、偏自相关系数图

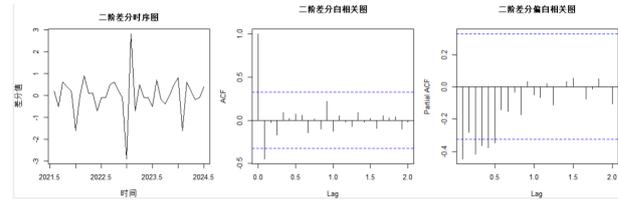


图 2.3 CPI 数据二阶差分时序图、自相关系数、偏自相关系数图

2. 建立模型并预测

由图 2.2 和 2.3 看一、二阶差分后的序列图明显是平稳的，计算一、二阶差分序列 $\{v_{1t}\}$ 、 $\{v_{2t}\}$ 方差分别为 0.3392、0.7833，后者方差偏大，选用一阶差分。由图 2.1 可看出数据无明显的季节效应，选用 ARIMA(3,1,0) 和 ARIMA(1,1,1) 模型。

然后，分别对数据进行建模，分析模型参数结果：

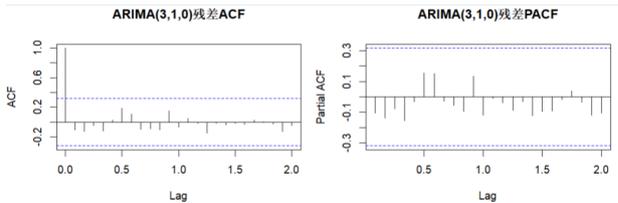


图 2.4 ARIMA(3,1,0) 模型残差检验图、自相关系数残差图

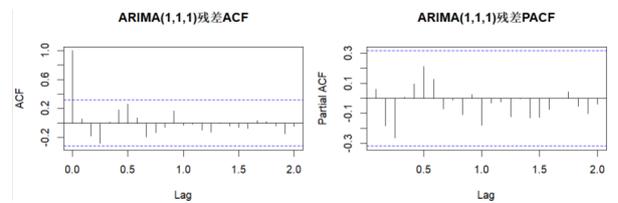


图 2.5 ARIMA(1,1,1) 模型残差检验图、自相关系数残差图

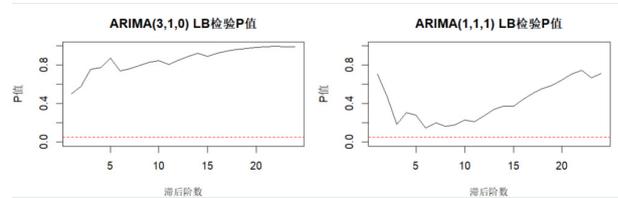


图 2.6 模型 LB 统计量的 P 值检验图

由 AIC 准则，数值越小模型拟合越好，以及残差检验图 P 值越大越拒绝模型无效的假设，最终选用 ARIMA(3,1,0)。根据最小二乘法运算，所以模型为：

$$x_t = 0.7411x_{t-1} - 0.0473x_{t-2} - 0.0730x_{t-3} + 0.3792x_{t-4} + \varepsilon_t \quad (2.1)$$



图2.7 预测与实际相比较（红色线为实际数据的走向）

可见，拟合结果不完美，但仍然在上下限两倍的标注差中，在可控制的范围内^[12]。

3. 使用贝叶斯统计分析

基于前文对两个 ARIMA 模型的分析，加上对后面的结果有可比性，我们仍然 ARIMA(3,1,0) 模型。对于以贝叶斯统计分析和贝叶斯统计推断为基础的数学模型构建过程为：将先验信息和样本数据信息通过贝叶斯定理，建立后验理论模型并进行参数估计，最后进行模型检验。得到如表 3.3 结果，得模型为：

$$x_t = 0.7450x_{t-1} - 0.0706x_{t-2} - 0.0782x_{t-3} + 0.4037x_{t-4} + \varepsilon_t \quad (2.2)$$

表 2.3 30000 次 Gibbs 抽样迭代的参数后验估计统计量

参数	均值	标注差	标准误	2.5%	97.5%
ϕ_1	-0.10283	0.2660	0.00096	-0.7479	0.6053
ϕ_2	-0.17603	0.2623	0.00094	-0.8333	0.4703
ϕ_3	-0.23108	0.2632	0.00094	-0.8995	0.4640

下面是各参数的 Gibbs 抽样过程图：

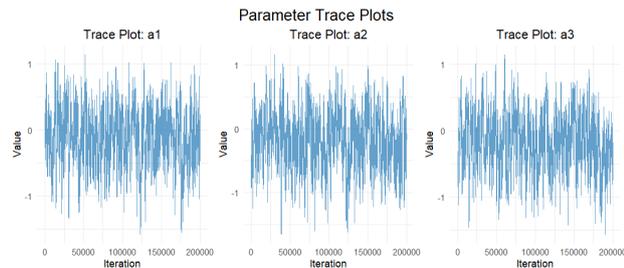


图2.8 各参数的 Gibbs 抽样过程图

参考文献

- [1] 邓迎春. 经济时间序列 ARMA 模型的贝叶斯分析及其应用. 湖南大学学院硕士毕业论文, 2006, 6.
- [2] 马瑶, 孙中玉, 邹益玲. 基于 ARIMA 模型对中, 美, 德三国 CPI 的分析与预测 [J]. 经济研究导刊, 2021(8):5.
- [3] 方珂昊, 赵凌. 基于偏最小二乘方法的 ARIMA 模型在股票指数预测中的应用 [J]. 四川文理学院学报, 2018, 28(5):5. DOI:CNKI:SUN:DXSZ.0.2018-05-002.
- [4] 洪京一. 基于 ARIMA 模型的上海市居民消费价格指数实证分析 [J]. 中小企业管理与科技, 2021.
- [5] 肖曼君, 夏荣尧. 中国的通货膨胀预测基于 ARIMA 模型的实证分析. 上海金融, 2008, (8): 38-42.
- [6] 刘红梅. ARIMA 模型在股票价格预测中的应用. 广西轻工业, 2008, (115): 92-93.
- [7] 刘乐平. 贝叶斯计量经济学从先验到结论. 中国经济学年会, 2006, (3): 53-56.
- [8] 王燕. 应用时间序列分析. 北京: 中国人民大学出版社, 2000. 140-174.
- [9] 孙荣恒. 应用数理统计. 第二版. 北京: 科学出版社, 2002. 67-85.
- [10] 吴海霞, 刘路锋. 蒙特卡罗方法在实际问题中的应用. 太原师范学院学报, 2009, (1): 5-8.
- [11] W N Venables, B D Ripley Springer. Modern Applied Statistics with S Fourth edition. New Jersey: U S Patent, 2002. 1-147.
- [12] Frank R Kleibergen, Henk Hoek. Bayesian Analysis of ARMA Models. Tinbergen Institute Discussion Paper, 2000, (3): 23-29.

下面是个参数的后验分布概率密度估计图：

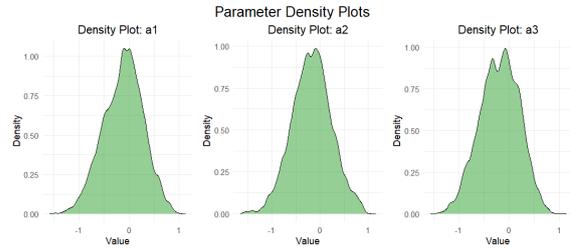


图 2.9 各参数后验概率密度图

然后，我们再利用所得到的模型进行往后6阶的预测，预测图与实际图对比：

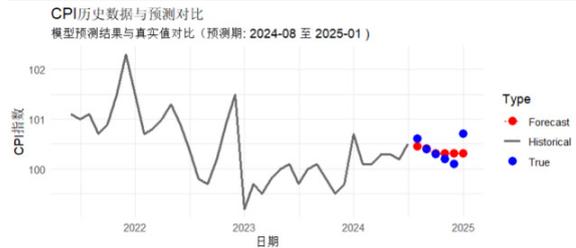


图 2.10 预测与实际相比较（红色线为实际数据的走向）

可见预测效果。根据误差指标，在这里贝叶斯方法的误差更小，两种方法都得到的大致真实的走向，这也说明了两种方法都是有效的。

表 2.4 误差指标

指标	MAE	MSE	RMSE
最小二乘法	0.1610	0.0400	0.2000
贝叶斯方法	0.1474	0.0394	0.1986

三、结论

近年我国经济发展仍将面临复杂的国际和国内环境，但是总的来说经济回升，基础需进一步夯实。对比历史数据和现实情况可以发现，2023以来，CPI的涨幅处于低位，物价整体平稳，但低于通常认为经济运行健康状态下的CPI涨幅2%-3%，物价上涨动力不足，而真实的2024年下半年的数据显示比估计略高，这说明CPI有所回升，但是经济恢复向好态势仍有待进一步巩固。