

非最小相位线性非高斯序列的替代数据检验*

刘耀宗 温熙森 胡芑庆

(国防科技大学机电工程与自动化学院,长沙 410073)

(2000 年 11 月 17 日收到)

替代数据法作为检验时间序列非线性和混沌的统计方法获得了广泛应用. 常用的替代数据法的零假设为“原序列来自(经过单调静态非线性变换的)平稳线性高斯随机过程”. 拒绝此假设,并不能说明序列必然来自确定性的非线性动力系统,非最小相位的线性非高斯序列也会导致基于相位随机化的替代数据检验拒绝此假设.

关键词:替代数据,时间序列分析,非线性,非最小相位

PACC:0545

1 引 言

自然界中广泛存在各种各样形态复杂的时间序列,如太阳黑子序列、癫痫病人脑电图序列、股票价格序列等. 长期以来,人们一直用线性随机过程的理论与方法来处理这些看来杂乱无章的序列. 近年来,非线性动力学特别是混沌理论的发展为貌似随机的复杂序列提供了新的解释和分析方法. 从混沌理论出发,研究人员设计了多种判断序列是否来自混沌动力系统的算法,如最大 Lyapunov 指数估计、关联维估计、非线性预测误差等^[3]. 但是,随着研究的深入,人们发现对于来自实验中的含噪声、有限长的时间序列,这些方法的结果并不十分可靠^[1,2],而且受到噪声水平、参数选取、序列长度等诸多因素的影响. 近年来,基于统计假设检验的替代数据法^[3]获得了广泛的应用.

常用的替代数据法^[3]假设原序列来自平稳的线性高斯随机过程,或来自经过单调的静态非线性变换的线性高斯随机过程. 在保证满足此假设的条件下,随机选择其他特性,生成一组替代数据序列. 分别计算原数据和替代数据的某非线性统计量,如果原数据与替代数据的统计量显著不同,则拒绝零假设,认为序列来自确定性非线性动力系统. 尽管替代数据检验的思想在更早的文献中已经应用,文献^[3]提出了“替代数据”(surrogate data)一词并使这种方法获得了更广泛的应用. 随着研究的进一步

深入,人们发现替代数据法存在一些缺陷. 文献^[3]的作者随后指出^[1],对于具有长相干时间(coherence time)的序列,替代数据检验常常拒绝零假设. 文献^[2]指出循环平稳的线性随机序列也会导致替代数据检验得出错误结论. 文献^[4]指出替代数据生成算法 amplitude adjusted fourier transform(缩写为 AAFT)^[3]可能导致替代序列与原序列的功率谱差异,并提出了相应的叠代算法——IAAFT(iterated AAFT). 文献^[5]研究了序列来自静态非单调非线性变换的线性高斯过程时替代数据的生成算法.

上述文献的研究结果说明了基于相位随机化的替代数据检验虽然是一种严格的统计检验方法,但由于其零假设是“序列来自平稳的线性高斯随机过程”,非线性动力系统并不是其唯一的备选假设,所以,拒绝此假设并不能说明序列一定来自非线性动力系统^[6,7]. 在研究这一问题时,作者发现,对非最小相位的线性非高斯序列,同样会导致基于相位随机化的替代数据检验得出错误的结果. 下文首先介绍替代数据法的具体步骤,然后详细介绍非最小相位的线性非高斯序列的替代数据检验结果,并给出相应结论.

2 替代数据法

对于一个动力学系统,可以用如下的微分方程组表示:

$$\frac{dX}{dt} = F(X, t, \mu) + \eta, \quad (1)$$

*国家自然科学基金(批准号:59775025)资助的课题.

$$V_t = G(X_t) + \varepsilon, \quad (2)$$

其中 X 为 m 维状态变量, μ 为系统参数, η 为系统动态噪声, V 为 k 维观测向量, ε 为观测噪声. 如果函数 F 关于状态变量 X 为线性的, 则称该系统为线性系统, 否则称为非线性系统. 如果系统动态噪声 $\eta = 0$, 则系统为确定性系统, 否则为随机系统. 由于在工程实际中干扰和噪声普遍存在, 严格地讲, 实际系统都可以称为随机系统. 当函数 F 和 G 中有一个是非线性时, 观测序列 V 就是非线性的. 称函数 F 为非线性而函数 G 为线性时的观测序列 V 为“动力学非线性序列(dynamical nonlinearity)”, 函数 F 为线性而函数 G 为非线性时的观测序列 V 为“观测非线性序列(measurement nonlinearity)”. 替代数据检验的目标就是在仅仅已知观测序列 V 时, 判定函数 F 是否为非线性的.

替代数据法的实质是一种统计假设检验, 它包含四个方面的内容: 零假设(null hypothesis), 替代数据生成算法、检验统计量(discriminating statistic), 统计检验方法.

零假设是对数据特征的一种假设性猜测, 它是产生替代数据的依据. 零假设限定了替代数据的特性, 要求替代数据必须保留原数据的某些特征, 如均值、方差等, 其他特征则可随机选取. 常用的替代数据法的零假设为原序列来自平稳的线性高斯随机过程. 设 $\{x(t)\}$ 为一个观测到的标量序列, $\{x(t)\}$ 是平稳的, 且对每个 $t \in Z$, 有

$$x(t) + \varphi_1 x(t-1) + \dots + \varphi_p x(t-p) = \varepsilon(t) + \theta_1 \varepsilon(t-1) + \dots + \theta_q \varepsilon(t-q), \quad (3)$$

其中 $\{\varepsilon(t)\}$ 为均值为零、方差为 σ^2 的高斯白噪声, 记为 $\{\varepsilon(t)\} \sim N(0, \sigma^2)$. 实质上 $\{x(t)\}$ 是一个高斯白噪声激励的 ARMA(p, q) 模型的输出序列. 另一常用的零假设为原序列来自经过单调静态非线性变换的线性高斯随机过程, 这一假设适用于非高斯分布的序列.

替代数据生成算法与零假设密切相关, 不同的零假设对应不同的生成算法. 替代数据法把原序列当成满足零假设的某随机过程的一个实现, 生成算法研究如何生成该过程的其他实现, 即替代数据. 替代数据生成算法有两种实现方式, 即经典实现(typical realization)和约束实现(constrained realization)^[8]. 对应“平稳线性高斯随机过程”零假设的约束实现生成算法是随机化相位的傅里叶变换(PRFT), 对应“单调静态非线性变换线性高斯过

程”的是 AAF^T^[3].

检验统计量是定量表征时间序列某个特性的特征量, 它是判别原数据与零假设是否符合的依据. 如果原数据的检验统计量与满足零假设的替代数据的检验统计量不相符, 则认为零假设不成立, 从而说明原数据具有零假设以外的特征. 在某些零假设下某些检验统计量的理论分布可以通过解析推导得到, 这是传统的非线性检验方法的基础. 在替代数据法中, 检验统计量的近似分布是通过 Monte-Carlo 仿真方法得到的. 在满足零假设(保持与原数据有相同的特征如均值、方差、功率谱等)的前提下随机选择其他特性生成一组替代数据, 分别计算这些替代数据的检验统计量, 用它们的分布来逼近原数据的检验统计量在零假设下的分布. 这一方法虽然计算量较大, 但解决了解析推导分布形式困难甚至不可能的问题, 因而带来了零假设和检验统计量选择的灵活性. 特别地, 零假设和检验统计量的选择可以是完全独立的. 替代数据法实际上是现代统计学中“bootstrap”方法的基本应用. 文献[8]指出, 如果替代数据是用约束实现的方法生成的, 则理论上检验统计量可以是任何非线性测度. 常用的检验统计量有: 关联维、最大 Lyapunov 指数、一步预测误差等.

替代数据法中常用的统计检验方法有两种. 一种方法(σ 检验)是假设满足零假设的序列的检验统计量服从正态分布, 求出替代数据的统计量的均值 \bar{Q}_s 和方差 σ_s^2 , 计算 σ 检验量 $\sigma = |Q_0 - \bar{Q}_s| / \sigma_s$, 通过正态分布的参数检验进行. 另一种方法(rank-order 检验)是产生 $1/\alpha - 1$ 组(单边检验)或 $2/\alpha - 1$ 组(双边检验)替代数据, 计算原数据和替代数据的检验统计量, 然后排序, 如果原序列的检验统计量为最小或最大, 则拒绝零假设, 检验的显著水平为 $1 - \alpha$.

3 时间序列的线性性、最小相位性、高斯性

(3) 式可以看作是以 $\{\varepsilon(t)\}$ 为输入, $\{x(t)\}$ 为输出的线性系统. 如果该系统的所有零点位于单位圆内, 则称其为最小相位系统, 否则, 称为非最小相位系统. 如果输入信号 $\{\varepsilon(t)\}$ 为白噪声, 则(3)式称为 ARMA(p, q) 模型, 它表示一个线性随机过程. 若 $\{\varepsilon(t)\}$ 为高斯分布的白噪声, 则称此 ARMA(p, q) 过程为线性高斯过程, 相应的输出序列 $\{x(t)\}$ 称为线性高斯序列; 若 $\{\varepsilon(t)\}$ 为非

高斯分布的白噪声,则对应地称 ARMA(p, q)为线性非高斯过程, $\{x(t)\}$ 为线性非高斯序列。

由以上概念可知,常用的替代数据法的零假设为线性高斯过程,并没有涵盖线性随机过程的全部,所以,拒绝此假设并不足以说明序列来自非线性动力系统。但是,在实际算法构造中,替代数据生成算法 PRFT 并不需要序列必须是高斯分布的,AAFT 算法也可以生成非高斯分布时间序列的替代数据,这些算法是否也适用于线性非高斯序列呢?换言之,人们常常通过功率谱研究线性随机过程,PRFT 和 AAFT 算法保证了序列的功率谱不变,是否适用于所有线性随机过程呢?本文将通过数值计算的方法研究这一问题。

4 仿真计算

本节将使用替代数据法对已知动力学特性的线性非高斯随机序列进行非线性检验,以考察该方法的实用性。文献[9]对替代数据法常用的检验统计量进行了对比研究,结果表明,非线性一步预测误差对大部分序列具有较好的区分能力,所以在本文的后续研究中将使用它作为检验统计量。该统计量的计算方法如下:

$$t^{PE}(m, \tau, K) = \left(\frac{1}{N-1} \sum_{i=1}^N [x(i+1) - \mu(i+1)]^2 \right)^{1/2}, \quad (4)$$

$$\mu(i+1) = \sum_{k=1}^K \omega_k x(j_k + 1), \quad (5)$$

$$\omega_k = \frac{\|X(i) - NN(j_k)\|^{-2}}{\sum_{l=1}^K \|X(i) - NN(j_l)\|^{-2}}, \quad (6)$$

其中 $\{x(i)\}$ 为观测序列, N 为序列长度, m 为嵌入维数, τ 为时延, K 为用于预测的邻点个数, $NN(j_k)$ 为当前点 $X(i)$ 的第 k 个邻近点, 距离为欧氏距离。(5)式的实质是用当前点邻点的一步演化的加权平均作为当前点的一步预测结果。

为了避免序列首尾跳变的影响^[3,10],在生成替代数据之前对原序列进行首尾匹配^[10]。分别使用 PRFT^[3]和 IAAFT^[4]方法生成 49 组替代数据。统计检验采用常用的 sigma 方法。

算例 1 考虑如下 ARMA 模型:

$$\begin{aligned} x(t) - 0.8x(t-1) + 0.65x(t-2) \\ = \epsilon(t) - 2\epsilon(t-1), \end{aligned} \quad (7)$$

其零极点分布如图 1(a)所示,显然该模型是非最小

相位的,即有零点位于单位圆外。 $\{\epsilon(t)\}$ 取均值为 0, 方差为 1 的均匀分布白噪声序列。图 1(b)为一组原序列和 49 组对应的替代数据的一步预测误差棒线图,其中第一条对应原序列,其他对应经过首位匹配的 PRFT 替代数据序列,图 1(b)中可见,原序列的预测误差最小。按照 rank-order 法检验,在 98% 的显著水平下可拒绝零假设。取 $\{\epsilon(t)\}$ 的多个实现来获得该线性非高斯随机序列的多个实现,其 sigma 检验结果如图 1(c)所示。从图 1(c)中可以看出,虽然 sigma 值并不很大,但在 95% 的显著水平下,所有实现都得出拒绝零假设的结论。

考虑下式的 ARMA 模型:

$$\begin{aligned} x(t) - 0.8x(t-1) + 0.65x(t-2) \\ = \epsilon(t) - 0.5\epsilon(t-1), \end{aligned} \quad (8)$$

其零极点分布如图 2(a)所示, $\{\epsilon(t)\}$ 取均值为 0, 方差为 2 的均匀分布白噪声序列。该模型是与(7)式相应的最小相位系统,即两系统具有相同的极点,但(7)式单位圆外的零点在(8)式中为其共轭倒数所替代。两系统的输出序列具有几乎完全相同的功率谱,如图 2(c)所示。图中为了避免两曲线重叠在一起,把非最小相位序列的功率谱上移,图中虚线为其零点。图 2(b)为与图 1(b)相应的最小相位序列预测误差图。图 1(c)同时示出两序列多次实现的 sigma 统计量,显然,最小相位序列的替代数据检验得出了正确的结论,而非最小相位序列的检验结果拒绝了零假设,经常被解释为非线性而得出了错误的结论。

算例 2 考虑如下 ARMA 模型:

$$\begin{aligned} x(t) - 0.7x(t-1) + 0.1x(t-2) \\ = \epsilon(t) + 0.6\epsilon(t-1) + 0.18\epsilon(t-2), \end{aligned} \quad (9)$$

$$\begin{aligned} x(t) - 0.7x(t-1) + 0.1x(t-2) \\ = \epsilon(t) + 3.333\epsilon(t-1) + 5.556\epsilon(t-2). \end{aligned} \quad (10)$$

(9)和(10)式也是一组对应的线性非最小相位和最小相位系统,其对应的零极点分布如图 3(a)和(b)所示。激励白噪声 $\{\epsilon(t)\}$ 取为指数分布,均值为 0, 方差为 1。文献[4]指出,AAFT 算法^[3]虽然保证了替代数据与原序列具有相同的时域分布,却容易导致自相关函数出现差异,并提出了 IAAFT 算法来解决这一问题,但计算量较大。算例 1 的结果表明,对非最小相位的线性非高斯过程,基于 PRFT 算法的替代数据检验会导致拒绝零假设的结果。本算例将采用 IAAFT 算法生成替代数据,以考察此

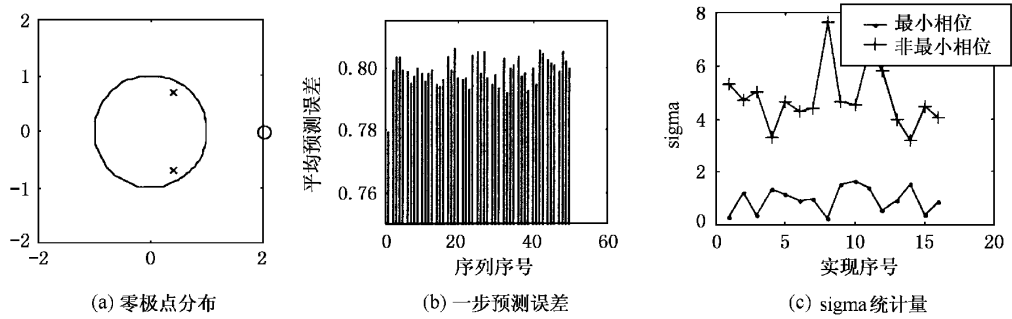


图 1

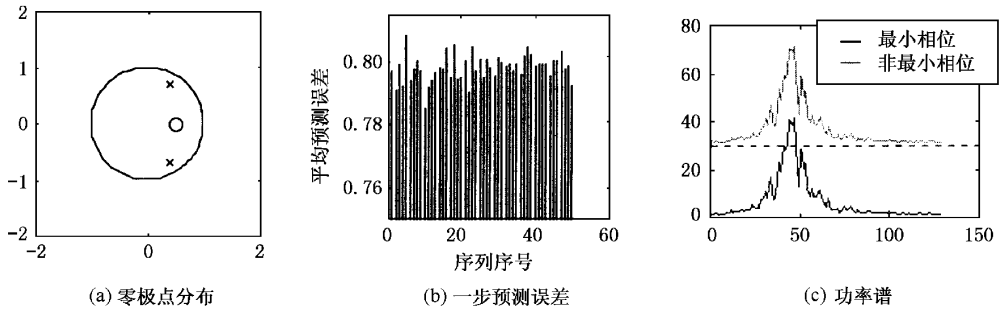


图 2

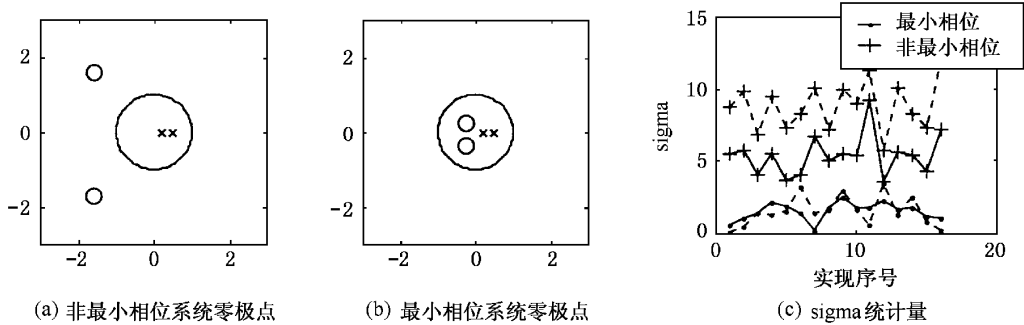


图 3

算法能否适用于非最小相位的线性非高斯序列. 图 3(c)中实线为指数分布白噪声激励式(9)(10)所示 ARMA(2,2)模型输出序列的多次实现的替代数据检验结果,由图可见,对于非最小相位的线性非高斯序列,基于 IAAFT 算法的替代数据检验同样得出了拒绝零假设的结论. 图 3(c)中虚线为该算例使用 PRFT 算法生成的替代数据的检验结果,对比虚线和实线可以发现,由于 IAAFT 算法保证了替代数据序列和原序列的自相关函数更好地匹配,所以检验结果的 sigma 值比 PRFT 替代数据小一些,但对应于非最小相位序列仍能保证在 95% 的显著水平拒绝零假设.

算例 1 和算例 2 分别使用了不同的非最小相位系统、不同的激励白噪声分布和不同的替代数据生

成算法,都得到了拒绝零假设的结果. 如果用高斯分布的白噪声作为激励,无论线性系统是否为最小相位的,也无论使用 PRFT 还是 AAFT,IAAFT 算法,替代数据检验都能得到正确的检验结论. 这一方面说明算例中拒绝零假设的结果不是个别系统的特例,也不是某种分布的白噪声或某种替代数据生成算法导致的,而是由于产生原序列的线性随机过程的非最小相位性和激励白噪声的非高斯性共同作用的结果. 另一方面,也说明根据常用的替代数据检验得出拒绝零假设的结论证明原序列来自非线性动力系统是不合适的,非最小相位的线性非高斯过程也会导致此检验拒绝零假设. 除了上文给出的算例外,作者还用其他的线性系统、其他分布的非高斯白噪声激励和其他的检验统计量进行了大量的仿真

计算,都得到了类似的结果.

5 结 论

本文在详细分析常用的替代数据检验法的基础上,指出由于其零假设为“原序列来自线性高斯过程或经过单调的静态非线性变换的线性高斯过程”,所

以检验得出拒绝零假设的结论时并不能说明原序列来自非线性动力系统,并通过算例说明,非最小相位的线性非高斯过程也会导致常用的基于相位随机化的替代数据检验拒绝零假设.如何解决非最小相位线性非高斯序列的非线性检验问题,作者将在另文中说明.

-
- [1] J. Theiler, P. S. Linsay, D. M. Rubin, *Time Series Prediction: Forecasting the Future and Understanding the Past* (Addison-Wesley, Reading Mass., 1993), p. 429.
- [2] J. Timmer, *Phys. Rev.*, **E58**(1998) 5153.
- [3] J. Theiler, S. Eubank, A. Longtin *et al.*, *Physica*, **D58**(1992), 77
- [4] T. Schreiber, A. Schmitz, *Phys. Rev. Lett.*, **77**(1996), 635.
- [5] D. Kugiumtzis, *Phys. Rev.*, **E62**(2000), 25.
- [6] M. Palus, D. Novotna, *Neural Network World*, **45**(1997), 451.
- [7] J. Timmer, *Phys. Rev. Lett.*, **85**(2000), 2647.
- [8] J. Theiler, D. Prichard, *Physica*, **D94**(1996), 221.
- [9] T. Schreiber, A. Schmitz, *Phys. Rev.*, **E55**(1997), 5443.
- [10] C.J. Stam, J. P. M. Pijn, W. S. Pritchard, *Physica*, **D112**(1998), 361.

SURROGATE DATA TEST FOR THE LINEAR NON-GAUSSIAN TIME SERIES WITH NON-MINIMUM PHASE*

LIU YAO-ZONG WEN XI-SEN HU NIAO-QING

(*Institute of Mechatronics Engineering and Automation, National University of Defense and Technology, Changsha 410073, China*)

(Received 17 November 2000)

ABSTRACT

Surrogate data testing is a popular method to detect nonlinearity and chaos in time series and has been vastly used in many applications with erratic time series. The explicit null hypothesis often used is that the time series is generated from a linear, stochastic, Gaussian stationary process, including a possible invertible nonlinear static observation function. It is pointed out that the rejection of such a hypothesis may not only result from an underlying nonlinear or even chaotic system, but also from, e. g., a linear, stochastic, non-Gaussian and non-minimum phase sequence. We investigate the power of the test against non-minimum phase sequence.

Keywords: surrogate data, time series analysis, nonlinearity, non-minimum phase

PACC: 0545

* Project supported by the National Natural Science Foundation of China (Grant No. 59775025).