

鲁棒 (Robust) 谱估计算法研究

曹长修 夏建昆

(重庆大学自动化系)

摘要

本文在估计理论、系统辨识, B.D. Martin D. Thomson 工作的基础上, 提出一种改进的鲁棒谱估计算法。该方法易于直接用于在线鲁棒谱估计。文中给出仿真实例。

一、引言

通常所谓的鲁棒估计, 是指对数据中的异常小变化所得估计也呈小变化(即估计对异常小变化不敏感), 称此估计为鲁棒估计。

考虑具有有限方差的一类广义平稳随机过程 $\{x_t\}$, $t = 0, 1, 2, \dots$, 若其统计特征值

$$Ex_t \equiv \mu,$$

$$\text{cov}(x_t, x_{t-l}) = c(l), \quad l = 0, \pm 2, \pm 4, \dots$$

已知, 则其频域描述为

$$S_x(f) = \sum_{l=-\infty}^{\infty} c(l) e^{-i 2 \pi f l},$$

式中 $S_x(f)$ 为 x_t 的功率谱密度。目前工程上普遍采用的谱估计方法是以 FFT 技术为基础, 加上各种窗孔技术的平滑周期图法, 以及象征着现代谱估计的参数模型法(如 AR、ARMA 谱等)。实践和理论表明^[1], 它们都不是鲁棒估计, 数据中出现的异常小变化足以使估计失败。本文的目的就在于获得 x_t 的鲁棒谱估计 $\hat{S}_x(f)$ 。

二、时序模型

现代研究平稳时间序列的方法是, 把问题进一步简化为具有有理谱的一类正态平稳时间序列来研究, 在许多情况下这种做法是一种合理的简化。但是众所周知, 观测记录均在不同程度上受到污染, 即存在“异常小变化”^[5], 称为野值(Outliers)(关于野值的严格统计学定义详见文献[1])。因此, 由于野值的存在, 观测序列就不可能满足正态假设, 而是一种带很长尾巴的近似正态分布, 称为受污染的正态分布或拖尾正态分布(Contaminated-normal distribution or heavy-tailed normal distribution), 可用如下分布形式描述^[4]:

$$CN(t; \gamma, \sigma_0^2, \sigma^2) = (1-\gamma)N(t; 0, \sigma_0^2) + \gamma N(t; 0, \sigma^2),$$

式中 $N(x; \mu, \sigma^2)$ 是均值为 μ , 方差为 σ^2 的正态分布。

通常, 时序中只要有一、两点较大的野值存在, 其方差估计、参数估计等就会产生较大的误差。在常规谱估计中, 如果对谱密度低幅部分的谱峰信号感兴趣, 则野值的大小不必大于观测值, 而只需大于过程的新息序列就足以使这些谱峰信号消失(见后面实例)。

由于野值产生的机制复杂, 出现的时刻任意, 很不容易获得其特别的统计知识。统计资料表明, 在通常获得的数据记录中至少包含有 5—10% 左右的野值。因此, 野值的一般性统计知识为: 相对于整个观测序列, 野值以小概率出现。换句话说, 在大部分观测时间内, 观测序列仍可以认为具有正态平稳特征。

本文所引用的时间序列模型是一种不需野值特别知识, 而又能获得对野值不敏感的时序估计模型, 称为 AO 模型(Additive outliers)^{[1], [4]}

$$\begin{cases} x_t = \sum_{i=1}^p \phi_i x_{t-i} + \varepsilon_t, \\ y_t = x_t + v_t, \end{cases} \quad (1)$$

式中 $\{x_t\}$ 是所考虑的平稳过程, 通常称为核过程, $\varepsilon_t \sim N(0, \sigma_\varepsilon^2)$, y_t 为观测数据, v_t 为野值, 描述观测中受污染的成分, 即

$$v_t \sim CND(v_t; \gamma, 0, \sigma^2) = (1-\gamma)N(0, 0) + \gamma N(0, \sigma^2),$$

式中 $0.01 \leq \gamma \leq 0.25$ 。式(1)表明, 观测序列中 $100(1-\gamma)\%$ 的记录为 x_t 过程, 仅 $100\gamma\%$ 的记录有野值污染。

三、鲁棒谱估计

B. D. Martin, D. Thomson 提出了如下鲁棒谱估计表达式(基本思想源于预白化思想):

$$\hat{S}_x(f) = \overline{S}_x(f) / |\hat{D}(f)|^2, \quad (2)$$

式中 $\hat{D}(f) = 1 - \sum_{l=1}^p \hat{\phi}_l e^{i 2 \pi f l}$, $f = k/n$, $k = 0, 1, \dots, [n/2]$ 。 $\hat{D}(f)$ 称为预白化滤波因子, $\hat{\phi}_l$ 为滤波器参数的鲁棒估计, $\overline{S}_x(f)$ 为预白化滤波输出(即残差序列的平滑周期图估计)。可见, 估计 $\hat{S}_x(f)$ 的鲁棒性取决于参数估计 $\hat{\phi}_l$ 的鲁棒性以及预白化滤波器的鲁棒性。

鲁棒谱估计 $\hat{S}_x(f)$ 的计算归结为下面三个部分:

- 1° 参数 ϕ_l ($l = 1, 2, \dots, p$) 的鲁棒估计 ($2 \leq p \leq 6$);
- 2° 观测序列 $\{y_t\}$ 的鲁棒滤波器设计及残差序列 $\{r_t\}$ 的计算;
- 3° $\{r_t\}$ 序列的平滑周期图估计。

本文综合现代控制理论中的状态空间法, 根据系统辨识理论及鲁棒估计的思想, 导出了参数 ϕ_l 鲁棒估计的递推公式; 改进了鲁棒滤波器中的鲁棒代价函数与权函数; 采用常规算法计算出残差序列的平滑周期图估计, 通过简单的换算式(2)最后可得到鲁棒谱估计 $\hat{S}_x(f)$ 。下面仅给出计算步骤 1°、2° 的公式及其有关说明。

1. 鲁棒滤波器

由AO模型(式(1)), 令

$$\mathbf{x}_t^T = [x_t, x_{t-1}, \dots, x_{t-p+1}]_{1 \times p},$$

$$u_t^T = [e_t, 0, \dots, 0]_{1 \times p},$$

$$\Phi = \begin{pmatrix} \phi_1 & \phi_2 & \cdots & \phi_{p-1} & \phi_p \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{pmatrix}_{p \times p},$$

$$H = [1, 0, \dots, 0]_{1 \times p},$$

则AO模型可用下列状态方程表示

$$\begin{cases} \dot{\mathbf{x}}_t = \Phi \mathbf{x}_{t-1} + u_t, \\ y_t = H \mathbf{x}_t + v_t, \end{cases} \quad (3)$$

式中 \mathbf{x}_t^T 的第1列元素是 x_t , 记成 $x_t = (\mathbf{x}_t)_1$.

由 Masreliz 定理^[3], 得鲁棒滤波公式为

$$\begin{cases} \hat{\mathbf{x}}_t = \Phi \hat{\mathbf{x}}_{t-1} + (m_t/s_t^2) s_t \phi(r_t/s_t), \\ M_t = \Phi P_{t-1} \Phi^T + Q, \\ P_t = M_t - w(r_t/s_t) m_t m_t^T / s_t^2, \end{cases} \quad (4)$$

式中 $r_t = y_t - \hat{y}_t^{t-1}$ 为预报残差;

$\hat{y}_t^{t-1} = (\Phi \hat{\mathbf{x}}_{t-1})_1$ 为预报值;

$M_t = [m_{11t}, m_{21t}, \dots, m_{p1t}]_{p \times p}$ 为预报矩阵, $m_t = m_{11t}$;

$m_{11t}^T = [m_{11t}, m_{21t}, \dots, m_{p1t}]_{1 \times p}$;

$s^2 = m_{11t}$ 为预报残差尺度;

P_t 为 $p \times p$ 阶滤波误差协方差阵;

$$Q = \begin{bmatrix} q_{11} & 0 \\ 0 & 0 \end{bmatrix}_{p \times p}, \quad q_{11} = \sigma_e^2,$$

$\hat{\phi}_i$ ($i = 1, \dots, p$) 及 σ_e^2 为参数的鲁棒估计;

$\phi(\cdot)$ 为鲁棒代价函数;

及

$w(\cdot)$ 为鲁棒权函数。

本文给出如下 $\phi(\cdot)$ 、 $w(\cdot)$ 函数的修改形式:

$$\phi_{IK}(t) = \begin{cases} t & |t| \leq a \\ b\text{sign}(t) & a < |t| \leq b \\ \text{sign}(t) \frac{b}{c-b}(c - \text{sign}(t)t) & b < |t| \leq c \\ 0 & |t| > c \end{cases}$$

$$w_{BS}(t) = \begin{cases} [1 - (t/b)^2][1 - \phi(t/b)^2] & |t| < b/2 \\ 0 & |t| \geq b/2 \end{cases}$$

因此, 当 $v_t = 0$ (即设有野值污染) 时, \hat{y}_{t-1}^{t-1} 为 x_t 的线性预测值, 而当 $v_t \neq 0$ 时, \hat{y}_{t-1}^{t-1} 为 x_t 的非线性预测值(内插值), 剔除了野值的影响.

上述鲁棒滤波公式由 Martin, Thomson 导出, 采用的是 Hampel 提出的鲁棒代价函数^[4](有界连续). 笔者将其它一些鲁棒函数与其作了分析比较(由于篇幅有限这里省略), 在此基础上提出了一种改进形式的鲁棒代价函数.

2. 鲁棒参数估计

鲁棒滤波器中参数 ϕ 、 σ_e^2 的估计值 $\hat{\phi}$ 、 $\hat{\sigma}_e^2$, 可由下列 GM 鲁棒估计方程^[3]获得

$$\left\{ \begin{array}{l} \sum_{i=2}^n w(z_i) z_i \phi \left(\frac{y_i - z_i^T \hat{\phi}}{\hat{s}_e} \right) = 0, \\ \sum_{i=2}^n w(z_i) \left\{ \phi^2 \left(\frac{y_i - z_i^T \hat{\phi}}{\hat{s}_e} \right) - B \right\} = 0, \end{array} \right. \quad (5)$$

式中 $z_i^T = [y_{i-1}, y_{i-2}, \dots, y_{i-p}]_{1 \times p}$, 当 $i \leq 0$ 时 $y_i = 0$; B 为补偿系数, $w(z_i)$ 为权函数, 以及 $\phi(\cdot)$ 为鲁棒代价函数. Martin-Thomson 采用加权迭代公式得出 $\hat{\phi}$, $\hat{\sigma}_e^2$.

由 AO 模型对观测值 y_t 进行预白化运算, 得

$$\eta_t = y_t - \phi_1 y_{t-1} - \dots - \phi_p y_{t-p} = \varepsilon_t + v_t - \phi_1 v_{t-1} - \dots - \phi_p v_{t-p}. \quad (6)$$

当 $v_t = 0$ 时

$$\eta_t \sim N(\eta_t; 0, \sigma_e^2),$$

当 $v_t \neq 0$ 时

$$\eta_t \sim CN(\eta_t; \gamma, \sigma_e^2, \sigma^2) = (1-\gamma)N(\eta_t; 0, \sigma_e^2) + \gamma N(\eta_t; 0, \sigma^2),$$

即为受污染的正态分布.

于是, 我们所讨论的常参数模型可用状态空间法描述如下:

$$\begin{cases} \hat{\Theta}_k = \Theta_{k-1} \\ y_k = z_k^T \Theta_k + \eta_k \end{cases} \quad (1 < k \leq n), \quad (7)$$

式中

$$\Theta_k^T = [\phi_1(k), \phi_2(k), \dots, \phi_p(k)]_{1 \times p},$$

$$z_k^T = [y_{k-1}, y_{k-2}, \dots, y_{k-p}]_{1 \times p}.$$

由系统辨识理论与卡尔曼滤波理论，可导出参数辨识的鲁棒估计递推公式如下：

$$\begin{cases} \hat{\Theta}_k = \hat{\Theta}_{k-1} + K_k \hat{s}_\epsilon \phi(r_k / \hat{s}_\epsilon), \\ K_k = P_{k-1} z_k' (z_k'^T P_{k-1} z_k' + \lambda)^{-1}, \\ P_k = P_{k-1} - w(r_k / \hat{s}_\epsilon) K_k z_k'^T P_{k-1} / \lambda, \end{cases}$$

式中 $r_k = y_k - z_k'^T \Theta_{k-1}$ 为估计误差；

$z_k'^T = z_k^T \bar{W}_k$ 为加权观测阵；

$\bar{W}_k = \text{diag}(w_{k-1}, \dots, w_{k-p})$ ；

$$w_i = w(r_i / \hat{s}_\epsilon), \quad i = k-1, k-2, \dots, k-p$$

λ —— 遗忘因子，取 0—1 之间的值；

\hat{s}_ϵ —— 估计残差尺度，由下列辅助方程得到：

$$\hat{s}_\epsilon^{k+1} = \left\{ \left[\sum_{k=p+1}^n w_k (r_k \cdot w_k)^2 \right] / [B(n-2p)] \right\}^{1/2}.$$

采用 Tukey 提出的下列鲁棒代价函数及鲁棒权函数：

$$\phi_{BS}(t) = \begin{cases} t[1 - (t/cb)^2]^2 & |t| < cb \\ 0 & |t| \geq cb \end{cases}$$

$$w(t) = \begin{cases} \phi'_{BS}(t) & |t| < cb/2 \\ 0 & |t| \geq cb/2 \end{cases}$$

初值 $\Theta_0, P_0, \hat{s}_\epsilon^0$ 由 RLS 算法得到， $\hat{s}_\epsilon^0 = \hat{\sigma}_\epsilon^0$ ，可取 $\hat{\sigma}_\epsilon^0 = \text{Median}(r_0 / 0.6745)$ 。

此递推估计公式与 Martin - Thomson 采用迭代公式^[4]比较，有如下特点：

1) 计算工作量、存贮空间减少，提高了算法的可行性，此为递推算法之特点。

2) 引入遗忘因子 λ ，使算法可对慢时变参数进行实时估计。

3) 无需特别构造出 $w(z_i)$ 权函数（式（5）中）直接按 $\phi(r_t/s)$ 对野值的诊断结果，由鲁棒权函数 $w(r_t/s)$ 构成对角权阵 \bar{W}_k 对观测阵加权，即可使观测阵 z_k^T 有界（限制野值的影响），从而极大地简化了计算方法，提高了计算效率。

4) 鲁棒代价函数采用 ϕ_{BS} 而不会产生寄生根，而且 $r_t = y_t - z_k'^T \Theta_{k-1}$ 中 观测序列

为加权后的 \mathbf{z}_k^T , 故提高了 $\phi_{BS}(r_t/s_\epsilon)$ 对野值的分辨率, 相应的权函数采用 ϕ_{BS} 的导函数, 从而使估计为 M 意义下^[5]的最佳估计, 即估计具有更高的可靠性.

3. 几点说明

1) 一般地, 观测方程应为

$$y_t = \mu + x_t + v_t,$$

式中 μ 为位置参数 (Location Parameter), 这时预白化滤波分为两步:

a) 估计位置参数 μ : 当 $v_t \equiv 0$ 时, 估计 $\hat{\mu}$ 为样本均值; 当 $v_t \neq 0$ 时, 估计 $\hat{\mu}$ 为鲁棒位置估计, 参见文献[2].

b) 将观测序列坐标中心平移 $\hat{\mu}$, 即

$$\tilde{y}_t = y_t - \hat{\mu},$$

得

$$\tilde{y} = x_t + v_t.$$

然后, 对平移了的观测值 \tilde{y}_t 进行预白化滤波. 若剔除野值 v_t 后的观测值记为 $\tilde{c}y_t$, 则最后可以得到剔除野值后的正态过程观测序列

$$\tilde{c}y_t = \tilde{c}y + \hat{\mu}.$$

2) AO 模型(式(1))中阶 p 的选择问题, 由式(2)鲁棒估计表达式以及式(3)预白化滤波器可以看出, 阶 p 决定了什么数量级的野值被检测出来并除掉(野值出现的地方做插值运算). 通常, $2 \leq p \leq 6$ 就满足要求了. 如果单纯为 AR(p) 过程的参数鲁棒估计, 那么阶 p 的估计也是其中重要部分^[3].

本文实例中阶的确定是这样进行的: 从 $p=1$ 开始, 逐次增加 1, 当 $\hat{s}_\epsilon(p+1) \approx \hat{s}_\epsilon(p)$ 时, 这时阶 p 即为选定的阶.

四、仿 真 实 例

模拟的随机过程 x_t 由下列各式组成:

$$x_t = 1.33u_t + w_t + z_t,$$

式中

$$u_t = 0.975u_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N(0, 1)$$

$$w_t = 0.95w_{t-1} - 0.90w_{t-2} + \eta_t, \quad \eta_t \sim N(0, 1)$$

$$z_t = 0.33z_{t-1} - 0.90z_{t-2} + \zeta_t, \quad \zeta_t \sim N(0, 1)$$

x_t 为正态平稳随机序列, 模拟的观测序列由下式给出:

$$y_t = x_t + v_t, \quad 1 \leq t \leq 520.$$

共取 520 个样本点, 其中 $v_t \sim CND(0.1, 0, 10) = 0.9N(0, 0) + 0.1N(0, 10)$. 这就是说, 在 90% 的观测值中野值 v_t 为零, 只在 10% 的观测值中存在野值污染.

这里, 我们已知 u_t 、 w_t 、 z_t 的谱分布情况. u_t 在低频部分有较大的能量分布; w_t 带宽为 0.1, 峰位为 0.17Hz 的窄带信号; z_t 带宽为 0.1, 峰位为 0.22Hz 的窄带信号. 从复合信号 x_t 的构成中可以看出, w_t 、 z_t 信号相对地比 u_t 信号弱得多. 假设我们只是

对 w_t, z_t 信号感兴趣, 而 u_t 仅作为一种环境噪声考虑, 那么可以看出, 当观测值中不存在野值污染 ($v_t \equiv 0$) 时, 常规谱估计(基于正态假设)和式(2)谱估计均能正确地检测出 w_t, u_t 信号; 而当观测值中存在野值污染 ($v_t \neq 0$) 时, 常规谱估计失败, 而式(2)谱估计低幅部分两个窄带谱峰信号仍然清晰可见, 即式(2)谱估计对野值不敏感, 估计是鲁棒估计。

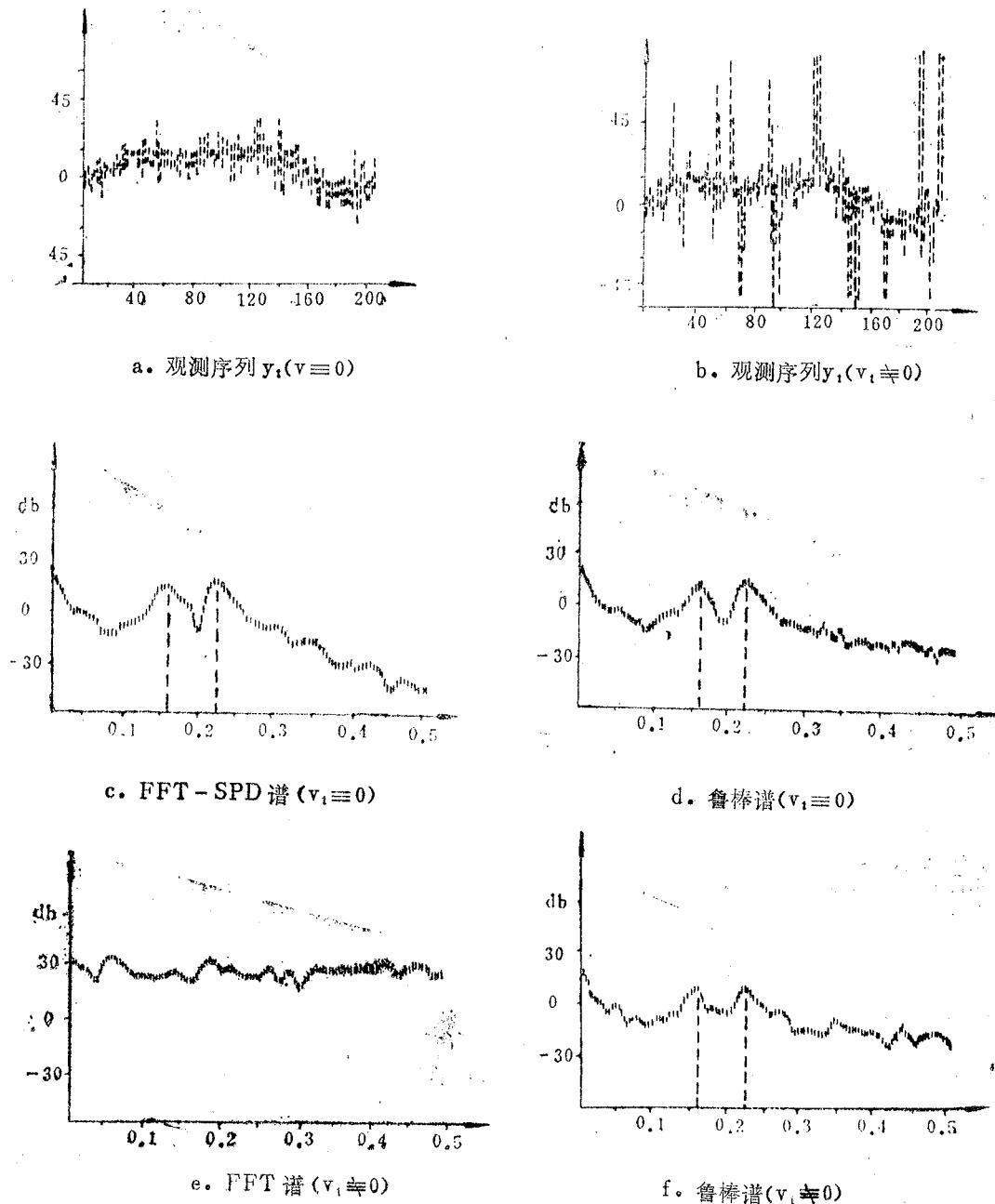


图 a~f 观测 y_t 及各种谱估计

图a, b 分别给出了野值 $v \equiv 0$ 时及野值 $v \neq 0$ 时观测序列 y_t 的样本段(200点). 图c, d 分别给出没有野值污染时常规谱估计(FFT-PSD)与鲁棒谱估计的结果. 图e, f 分别给出在野值污染下, 常规谱估计与鲁棒谱估计结果. 显然, 式(2)谱估计是鲁棒谱估计.

参 考 文 献

- [1] Martin, R. D., Robust Resistant Spectral Analysis. In: D. R. Brillinger and P. R. Krishnaiah eds., Handbook of Statistics, 3, (1983), 185-219.
- [2] Martin, R. D. and Zeh. J. E., Generalized M-estimates for Autoregressions, Including Small sample Efficiency Robustness, Technical Report #214, Dept. of Electrical Engineering, University of Washington, (1978), 5-6.
- [3] Martin, R. D., Robust Method for Time Series, The 2nd Applied Time Series Symp., Tulsa OK, March 3:5, (1980), 1-76.
- [4] Martin, R. D. and Thomson, D., Robust resistant Spectrum Estimation, Proceedings of the IEEE, 70:9, (Sept. 1982), 1097-1115.
- [5] Ершов, А. А.,参数估计的坚韧(Robust)方法, 应用数学与计算数学, 6, (1981), 29-48.

Study of Robust Spectrum Estimation

Cao Changxiu, Xia Jiankun

(Department of Automation, Chongqing University)

Abstract

In this paper, a improved robust spectrum estimation algorithm is presented, which is based on the estimation theory, systems identification techniques and works of Martin and Thomson. This method can be easily and directly used for the on-line robust estimation of spectrum. Simulation results are presented.