

经典 SIR 模型辨识和参数估计问题*

何艳辉, 唐三一

(陕西师范大学 数学与信息科学学院, 西安 710062)

(本刊编委唐三一来稿)

摘要: 可辨识与否是模型的基本特性,也是研究模型参数估计的基础.在传染病动力学中,SIR 模型仍是最常用的模型.该文研究了如何用高阶导数法与多观测点法判定 SIR 模型的可辨识性.研究表明针对 SIR 模型的可辨识技巧中多观测点法优于高阶导数法.不仅从理论上判定了 SIR 模型的可辨识性,而且结合流感疫情数据通过参数估计进一步验证了 SIR 模型的可辨识性.文中发展的技巧和方法有望推广到其他类型的传染病模型辨识和参数确定上.

关键词: 可辨识性; 高阶导数法; 多观测点法; 辨识函数; 参数估计

中图分类号: O213 **文献标志码:** A

DOI: 10.3879/j.issn.1000-0887.2013.03.005

引 言

日益完善的各种疾病观测数据为辨识性分析提供了保障,也为针对传染病动力学模型的辨识性研究提供了数据支撑.模型可辨识性的直观理解即在给定可观测值与初值时,模型中所有参数是不是有且仅有唯一估计值.对于模型可辨识性严格的数学定义,可参见文献[1].

为了使模型具有实际应用,理论上应该先从模型可辨识性入手.如果是可辨识的,才能论及模型参数估计与相应的统计推断.然而实际操作中更多的是直接估计模型参数并做统计推断,此时一旦有参数不可辨识,则所做的统计推断是不合理的.因为如果模型中有参数不可辨识,那么它的估计值就不唯一,即该参数不收敛.本文运用高阶导数法^[2-3]与多观测点法^[4],通过构建适当的辨识函数,在先从代数的角度给出了 SIR(susceptible-infectious-recovered)模型的可辨识性条件.然后结合英国流感疫情数据^[5]计算得到所需的辨识等式判定了 SIR 模型是可辨识的.最后,采用非线性最小二乘法(NLS)来估计 SIR 模型中可辨识的参数,进一步验证了其可辨识性.需要说明,本文均用圆点数或圆括号里的数表示各阶导数.

1 SIR 模型可辨识性

假设:1) 一类传染病在某一地区中传播,且所研究地区的种群的总数 N 是常数;2) S 类、 I 类、 R 类种群都具有出生能力且新出生的个体都属于易感者类;3) 每一类都具有自然出生率

* 收稿日期: 2013-01-15

基金项目: 国家自然科学基金资助项目(11171199)

作者简介: 何艳辉(1987—),女,江西人,硕士(E-mail:yanhuiqun1988@163.com);

唐三一(1970—),男,教授,博士(通讯作者. E-mail: sytang@snnu.edu.cn).

与死亡率.

模型所考虑地区中的种群数量 N 分成 3 类: 易感者类 (susceptible, S 类)、染病者类 (infective, I 类) 和移除者类 (removed, R 类), 即有 $N = S + I + R$, 则该 SIR 模型为

$$\begin{cases} \dot{S} = bN - \beta SI - dS, \\ \dot{I} = \beta SI - \gamma I - dI, \\ \dot{R} = \gamma I - dR, \end{cases} \quad (1)$$

其中, 自然出生率与死亡率分别用 b, d 表示, 单位时间内一个病人传染易感者的概率用 β 表示, 用 γ 表示单位时间内从染病者中的移出率, 并用 $\theta = (b, \beta, d, \gamma)$ 表示该模型参数向量. 在实际生活中, 我们容易收集到病人数 I 的数据, 而易感者 S 和移除者 R 的数据难以获得, 所以 I 是观测状态变量, S 和 R 是未观测 (或潜在) 状态变量.

代数上系统 (1) 的可辨识性定义是指对于某个时间 t^* , 某个正整数 k , 存在一个函数 $\Phi: R^4 \times R^{4(k+1)} \rightarrow R^4$ 满足以下条件:

$$\begin{cases} \det \frac{\partial \Phi}{\partial \theta} \neq 0, \\ \Phi(\theta, I, \dot{I}, \dots, I^{(k)}) = 0, \end{cases} \quad (2)$$

其中, $\dot{I}, \dots, I^{(k)}$ 是 I 在 $[0, t^*]$ 上关于时间 t 的各阶导数, 且分别称等式 (2) 为辨识方程, 函数 $\Phi(\cdot)$ 为辨识函数.

由上述定义可知, 要分析模型 (1) 的可辨识性需先构建适当的辨识函数. 而且, 为了消去其未观测 (或潜在) 状态变量 (即 S, R), 应对其观测状态变量 (即 I) 求高阶导数. 由模型 (1) 前两个等式易得 I 的二阶导数 (记为 \ddot{I}) 为

$$\begin{aligned} \ddot{I} &= \beta \dot{S} \dot{I} + \beta \dot{S} I - (\gamma + d) \dot{I} = \\ &= \beta \dot{S} \dot{I} + \beta (bN - \beta SI - dS) I - (\gamma + d) \dot{I} = \\ &= \frac{1}{I} (\dot{I} + \gamma I + dI) (\dot{I} - \beta I^2 - dI) + \beta bNI - (\gamma + d) \dot{I} = \\ &= \frac{\dot{I}^2}{I} + \beta bNI - \beta I \dot{I} - \beta (\gamma + d) I^2 - d \dot{I} - d (\gamma + d) I. \end{aligned} \quad (3)$$

显然上式已经不含未观测状态变量 S 与 R 了, 否则还需求 I 的更高阶导数. 为了方便, 我们记 $c = \beta bN, \eta = \gamma + d$, 易知 $\theta = (b, \beta, d, \gamma)$ 与 $\theta^* = (c, \beta, \eta, d)$ 之间存在一个一一映射, 则等式 (3) 右端可表示为

$$f(t, \theta^*, I, \dot{I}) = \frac{\dot{I}^2}{I} + cI - \beta I \dot{I} - \beta \eta I^2 - d \dot{I} - d \eta I. \quad (4)$$

因此, 我们可以得到

$$\begin{cases} \frac{\partial f}{\partial c} = I, \\ \frac{\partial f}{\partial \beta} = -I \dot{I} - I^2 \eta, \\ \frac{\partial f}{\partial d} = -\dot{I} - I \eta, \\ \frac{\partial f}{\partial \eta} = -\beta I^2 - dI. \end{cases}$$

由上述过程可知, 要判断 θ^* 或 θ 中 4 个参数的可辨识性, 我们需要得到 4 个可辨识的等式, 即辨识性的关键是如何构建辨识函数. 接下来, 论文将介绍两种如何通过等式 (3) 和 (4) 来

构建辨识函数的方法.

1.1 高阶导数法 (higher-order derivative method, HODM)

该方法的基本思想是通过不断求观测状态变量的更高阶导数,直至消去所有的未观测状态变量,从而得到一系列关于模型参数的辨识函数.根据参考文献中对该方法过程的描述^[2-3],我们可以构建出模型(1)对应的辨识函数:

$$\Phi_0 = (\ddot{I} - f, I^{(3)} - \dot{f}, \dots, I^{(5)} - f^{(3)})' = 0. \quad (5)$$

如果 $\det(\partial\Phi_0/\partial\theta^*) \neq 0$,则由隐函数定理易知 $\theta^* = (c, \beta, \eta, d)$ 可辨识,即如果

$$\text{rank} \left(\frac{\partial\Phi_0}{\partial\theta^*} \right) = \text{rank} \begin{pmatrix} \frac{\partial f}{\partial c} & \frac{\partial f}{\partial \beta} & \frac{\partial f}{\partial d} & \frac{\partial f}{\partial \eta} \\ \frac{\partial f^{(1)}}{\partial c} & \frac{\partial f^{(1)}}{\partial \beta} & \frac{\partial f^{(1)}}{\partial d} & \frac{\partial f^{(1)}}{\partial \eta} \\ \frac{\partial f^{(2)}}{\partial c} & \frac{\partial f^{(2)}}{\partial \beta} & \frac{\partial f^{(2)}}{\partial d} & \frac{\partial f^{(2)}}{\partial \eta} \\ \frac{\partial f^{(3)}}{\partial c} & \frac{\partial f^{(3)}}{\partial \beta} & \frac{\partial f^{(3)}}{\partial d} & \frac{\partial f^{(3)}}{\partial \eta} \end{pmatrix} = 4, \quad (6)$$

则 θ^* 就是可辨识的.易知辨识函数(5)中含有 I 的 5 阶导数,则计算时至少需要 6 个 I 的观测数据.而且需要指出,当 θ^* 为高维参数时,该方法的高阶导数很复杂且式(6)中仍含参数,使得判定模型参数的辨识性很困难,比如通过式(4)计算得到式(6)中的项 $\partial f^{(3)}/\beta$ 可以表示为

$$\frac{\partial f^{(3)}}{\beta} = -2(\dot{I})^2 - 2\ddot{I}I^{(3)} - 3\dot{I}I^{(3)} - II^{(4)} - 2\eta(3\ddot{I} + II^{(3)}),$$

其他项由计算公式可以类似得到.因此,下面将提供另一种计算量相对较少的方法.

1.2 多观测点法 (multiple time points method, MTPM)

假设已有 (I, \dot{I}, \ddot{I}) 在时刻 t_1, \dots, t_4 的值,并用 $(I_i, \dot{I}_i, \ddot{I}_i)$ 表示 (I, \dot{I}, \ddot{I}) 在 $t = t_i, i = 1, \dots, 4$ 的值.令

$$f_1 = f(t_1, \theta^*, I_1, \dot{I}_1), \dots, f_4 = f(t_4, \theta^*, I_4, \dot{I}_4),$$

则由式(3)和式(4)可得

$$\Phi_1 = (\ddot{I}_1 - f_1, \dots, \ddot{I}_4 - f_4)' = 0. \quad (7)$$

如果 $\det(\partial\Phi_1/\partial\theta^*) \neq 0$,则由隐函数定理易知 $\theta^* = (c, \beta, \eta, d)$ 可辨识,该条件等价于

$$\text{rank} \left(\frac{\partial\Phi_1}{\partial\theta^*} \right) = \text{rank} \begin{pmatrix} \frac{\partial f_1}{\partial c} & \frac{\partial f_1}{\partial \beta} & \frac{\partial f_1}{\partial d} & \frac{\partial f_1}{\partial \eta} \\ \frac{\partial f_2}{\partial c} & \frac{\partial f_2}{\partial \beta} & \frac{\partial f_2}{\partial d} & \frac{\partial f_2}{\partial \eta} \\ \frac{\partial f_3}{\partial c} & \frac{\partial f_3}{\partial \beta} & \frac{\partial f_3}{\partial d} & \frac{\partial f_3}{\partial \eta} \\ \frac{\partial f_4}{\partial c} & \frac{\partial f_4}{\partial \beta} & \frac{\partial f_4}{\partial d} & \frac{\partial f_4}{\partial \eta} \end{pmatrix} = 4. \quad (8)$$

同理,通过式(4)计算得到式(8)中的项 $\partial f_3/\beta$ 可以表示为

$$\frac{\partial f_3}{\beta} = -I_3\dot{I}_3 - I_3^2\eta,$$

其他项由计算公式可以类似得到.注意到该方法计算 \ddot{I} 至少需要 3 个 I 的观测数据,那么得到模型(1)对应的辨识函数至少需要 6 个 I 的观测数据,与 HODM 结论一致.

2 SIR 模型可辨识性应用

本节将用英国传染病监测中心在 1978 年发布的有关流感病人的统计数据^[5]来验证 SIR 模型的可辨识性. 该数据统计了两周内每天英格兰北部一所男孩寄宿学校流感爆发和流行情况, 总共有 763 个学生, 从发生一个染病者开始统计染病学生的人数, 总共统计了 15 天, 每一天统计得到的染病人数分别为

$$1, 3, 7, 25, 72, 222, 282, 256, 233, 189, 123, 70, 25, 11, 4.$$

由于该数据收集时间很短, 学生没有因病死亡, 即人口总数是常数, 所以我们进一步假设 $b = d = 0$, 从而得到相应的简化后 SIR 模型为

$$\begin{cases} \dot{S} = -\beta SI, \\ \dot{I} = \beta SI - \gamma I, \\ \dot{R} = \gamma I. \end{cases} \quad (9)$$

由本文第 1 节内容可知, 要分析模型(9)的辨识性, 首先需构建相应的辨识函数. 由模型(9)中前两个等式易得 I 的二阶导数为

$$\begin{aligned} \ddot{I} &= \beta S \dot{I} + \beta \dot{S} I - \gamma \dot{I} = \\ &= \beta S \dot{I} + \beta(-\beta SI)I - \gamma \dot{I} = \\ &= \frac{1}{I}(\dot{I} + \gamma I)\dot{I} - \beta I^2 \frac{1}{I}(\dot{I} + \gamma I) - \gamma \dot{I} = \\ &= \frac{\dot{I}^2}{I} - \beta I \dot{I} - \beta \gamma I^2. \end{aligned} \quad (10)$$

记 $\mu = \beta\gamma$, 易知 $\omega = (\beta, \gamma)$ 与 $\omega^* = (\beta, \mu)$ 之间存在一个一一映射, 则等式右端可表示为

$$\tilde{f}(t, \omega^*, I, \dot{I}) = \frac{\dot{I}^2}{I} - \beta I \dot{I} - \mu I^2. \quad (11)$$

因此, 我们可以得到

$$\begin{cases} \frac{\partial \tilde{f}}{\partial \beta} = -I \dot{I}, \\ \frac{\partial \tilde{f}}{\partial \mu} = -I^2. \end{cases}$$

2.1 HODM 判定模型(9)可辨识

由 HODM 可以构建出模型(9)对应的辨识函数为 $\tilde{\Phi}_0 = (\ddot{I} - \tilde{f}, I^{(3)} - \tilde{f}^{(1)})' = 0$. 如果 $\det(\partial \tilde{\Phi}_0 / \partial \omega^*) = I^2(I^2 - I\dot{I}) \neq 0$, 则由隐函数定理易知 ω^* 可辨识, 即等价于

$$\text{rank} \left(\frac{\partial \tilde{\Phi}_0}{\partial \omega^*} \right) = \text{rank} \begin{pmatrix} \frac{\partial \tilde{f}}{\partial \beta} & \frac{\partial \tilde{f}}{\partial \mu} \\ \frac{\partial \tilde{f}^{(1)}}{\partial \beta} & \frac{\partial \tilde{f}^{(1)}}{\partial \mu} \end{pmatrix} = \text{rank} \begin{pmatrix} -I \dot{I} & -I^2 \\ -I^2 - I \dot{I} & -2I \dot{I} \end{pmatrix} = 2.$$

又根据参考文献[6]的研究上式可化为

$$\text{rank} \left(\frac{\partial \tilde{\Phi}_0}{\partial \omega^*} \right) = \text{rank} \left(\begin{pmatrix} I & 0 \\ \dot{I} & I \end{pmatrix} \begin{pmatrix} \dot{I} & I \\ -\ddot{I} & \dot{I} \end{pmatrix} \right) = \text{rank} \begin{pmatrix} \dot{I} & I \\ -\ddot{I} & \dot{I} \end{pmatrix} = 2. \quad (12)$$

从流感爆发后的观测数据可知, 病人数量 I 和其变化速率即 \dot{I} 在不断变化, 所以 \dot{I} 和 \ddot{I} 都不为 0, 显然式(12)成立, 因此称模型(9)或 ω^* 可辨识.

2.2 MTPM 判定模型(9)可辨识

同样地,我们可由 MTPM 构建出模型(9)对应的辨识函数 $\tilde{\Phi}_1 = (\dot{I}_1 - \tilde{f}_1, \dot{I}_2 - \tilde{f}_2)' = 0$. 如果 $\det(\partial \tilde{\Phi}_1 / \partial \omega^*) = I_1 \dot{I}_1 I_2^2 - I_1^2 I_2 \dot{I}_2 \neq 0$, 则由隐函数定理易知 ω^* 可辨识, 即

$$\text{rank} \left(\frac{\partial \tilde{\Phi}_1}{\partial \omega^*} \right) = \text{rank} \begin{pmatrix} \frac{\partial \tilde{f}_1}{\partial \beta} & \frac{\partial \tilde{f}_1}{\partial \mu} \\ \frac{\partial \tilde{f}_2}{\partial \beta} & \frac{\partial \tilde{f}_2}{\partial \mu} \end{pmatrix} = \text{rank} \begin{pmatrix} -I_1 \dot{I}_1 & -I_1^2 \\ -I_2 \dot{I}_2 & -I_2^2 \end{pmatrix} = 2.$$

同理上式可化为

$$\begin{aligned} \text{rank} \left(\frac{\partial \tilde{\Phi}_1}{\partial \omega^*} \right) &= \text{rank} \begin{pmatrix} I_1(\dot{I}_1 - I_1) & I_1^2 \\ I_2(\dot{I}_2 - I_2) & I_2^2 \end{pmatrix} = \text{rank} \left(\begin{pmatrix} I_1 & 0 \\ 0 & I_2 \end{pmatrix} \begin{pmatrix} \dot{I}_1 - I_1 & I_1 \\ \dot{I}_2 - I_2 & I_2 \end{pmatrix} \right) = \\ &= \text{rank} \begin{pmatrix} \dot{I}_1 - I_1 & I_1 \\ \dot{I}_2 - I_2 & I_2 \end{pmatrix} = \text{rank} \begin{pmatrix} \dot{I}_1 & I_1 \\ \dot{I}_2 & I_2 \end{pmatrix} = 2. \end{aligned} \quad (13)$$

同理可知 I_1, I_2 和 \dot{I}_1, \dot{I}_2 不都相等, 即有式(13)成立, 因此称模型(9)或 ω^* 可辨识.

2.3 参数估计

在 2.1 节与 2.2 节我们分别用了两种可辨识性技巧来研究模型(9)的可辨识性, 以及在观测数据零误差时判定可辨识性所需的最少样本数据点. 然而实际中不可能达到零误差的理想状态, 因此, 在实际判定模型可辨识性时, 需要更多的样本数据点. 本小节通过 Monte Carlo (MC) 随机模拟得到了足够多的样本观测值, 然后采用最常用的 NLS 来估计参数, 结果验证了模型(9)是可辨识的.

对于仅有观测变量 I 的模型(9), 其观测模型即

$$I(t_i) = I(t_i; \omega) + \varepsilon_i, \quad i = 1, 2, \dots, n, \quad (14)$$

其中, $I(t_i; \omega)$ 表示含未知参数 ω 的模型(9)在 t_i 处的解, $I(t_i)$ 为对应的观测值, 随机误差 ε_i 通常设为相互独立同服从 0 均值, σ^2 方差的 Gauss 分布, 即 $\varepsilon_i \sim N(0, \sigma^2)$. 此时, 估计参数用 NLS 与极大似然法 (MLE) 效果是等价的. 易知, NLS 最小化的目标函数为

$$I_{\text{SSR}} = \sum_{i=1}^n [I(t_i) - I(t_i; \omega)]^2, \quad (15)$$

其中, I_{SSR} 表示残差平方和.

接下来, 用 MC 方法在 Matlab 中模拟了模型(9)的观测值 $I(t)$. 由之前给出的实际数据^[5] 可给定初值 $(S_0, I_0, R_0) = (762, 1, 0)$, 又根据已有研究结果^[5] 选取初值参数值即 $(\beta, \gamma) = (0.0022, 0.4529)$, 再通过 MC 随机模拟得到 I 在时间点 (单位: d) $t = (0, 1, 2, \dots, 14)$ 处的值 $I(t; \omega)$, 然后在等式(14)的基础上对 $I(t; \omega)$ 加入 Gauss 随机误差 ε 就得到了样本观测值 $I(t)$.

最后我们采用 NLS 方法利用模拟得到的观测值 $I(t)$ 重新估计模型(9)中的参数.

表 1 $\sigma = 0.01$ 时, 参数估计结果

Table 1 Estimation result for $\sigma = 0.01$

parameters	est. value	std. error	true value
β	0.0022	0.0008-5E	0.0022
γ	0.4527	0.3519-5E	0.4529

选取 $\sigma = 0.01$ 时, 每个时间点模拟 100 次得到的结果见表 1 与图 1. 作为比较, 取 $\sigma = 100$, 在每个时间点模拟 1000 次得到相应结果见表 2 与图 2.

表 2 $\sigma = 100$ 时, 参数估计结果
Table 2 Estimation result for $\sigma = 100$

parameters	est. value	std. error	true value
β	0.002 2	0.000 2	0.002 2
γ	0.456 6	0.054 6	0.452 9

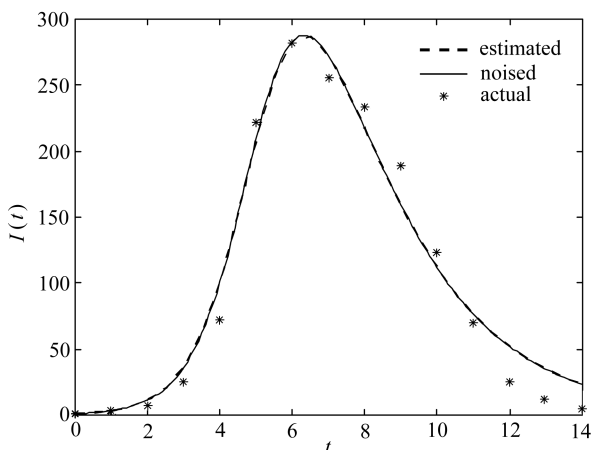


图 1 $\sigma = 0.01$ 时, 加入 Gauss 随机误差的观测值、估计的参数对应的解曲线与实际观测数据对比图

Fig. 1 Initial values S_0 , I_0 and R_0 are known and $\sigma = 0.01$ (The solid line represents the solution with noise, the dashed line represents the solution of estimated parameters without noise, and the stars represent the actual observed data of the infectious)

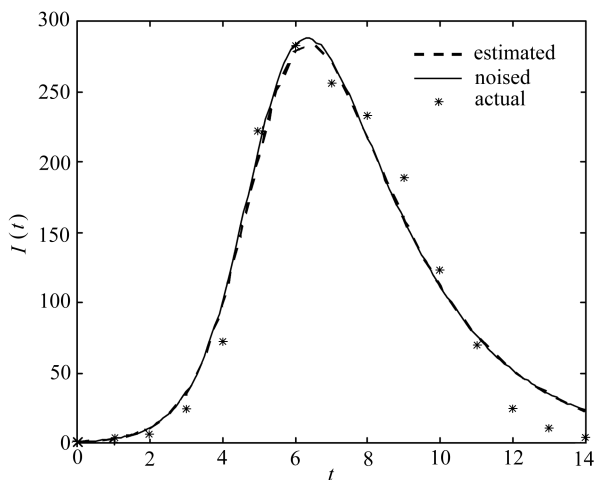


图 2 $\sigma = 100$ 时, 加入 Gauss 随机误差的观测值、估计的参数对应的解曲线与实际观测数据对比图

Fig. 2 Initial values S_0 , I_0 and R_0 are known and $\sigma = 100$ (The solid line represents the solution with noise, the dashed line represents the solution of estimated parameters without noise, and the stars represent the actual observed data of the infectious)

表 1 与表 2 中信息显示, 参数估计值与数据真实值非常接近, 即使观测数据数据噪音很强, 即 σ 值很大时, 参数估计值同样非常接近参数真实值。从图 1 与图 2 可以看出, 估计曲线与实际观测数据也拟合得非常好。因此, 我们在数值上验证了模型(9)是可辨识的。

3 结果与讨论

论文重点探讨了 SIR 模型的辨识性, 通过包括 HODM 与 MTPM 在内的方法来构建辨识函

数,从数学上分析了其辨识性.结果表明两种方法得出的结论一致,即该模型是可辨识的,且它们在分析时所需的最少数据点相同.研究表明针对 SIR 模型的可辨识技巧中 MTPM 优于 HODM.结合文献[5]中给出的实际数据,在 Matlab 中用我们采用 MC 方法进行数值研究,并通过 NLS 估计模型参数,所得结论进一步验证了 SIR 模型的可辨识性.需要指出:当模型的参数向量维数很高时,通过 MTPM 或 HODM 来构建辨识函数变得很复杂,判定也相当困难,建议寻求更优的方法,我们将在以后的工作中做进一步的研究与讨论.

参考文献(References):

- [1] Conte G, Moog C H, Perdon A M. *Nonlinear Control Systems: an Algebraic Setting*[M]. London: Springer, 1999.
- [2] Xia X, Moog C H. Identifiability of nonlinear systems with application to HIV/AIDS models [J]. *IEEE Transactions on Automatic Control*, 2003, **48**(2): 330-336.
- [3] Jeffrey A M, Xia X. *Identifiability of HIV/AIDS Models(Chaptern)*. In: *Deterministic and Stochastic Models of AIDS Epidemics and HIV Infections With Intervention*[M]. Singapore: World Scientific Publishing, 2005.
- [4] Wu H, Zhu H, Miao H, Perelson A S. Parameter identifiability and estimation of HIV/AIDS dynamic models[J]. *Bulletin of Mathematical Biology*, 2008, **70**(3): 785-799.
- [5] 肖燕妮,周义仓,唐三一.生物数学原理[M].西安:西安交通大学出版社,2012.(XIAO Yan-ni, ZHOU Yi-cang, TANG San-yi. *The Principle of Biomathematics*[M]. Xi'an: Xi'an Jiaotong University Press, 2012(in Chinese))
- [6] Roger A H, Charles R J. 矩阵分析[M].杨奇译.机械工业出版社,1985.(Roger A H, Charles R J. *Matrix Analysis*[M]. YANG Qi, Transl. Mechanical Industry Press, 1985. (in Chinese))

Identification and Parameter Estimation for Classical SIR Model

HE Yan-hui, TANG San-yi

(School of Mathematics and Information Science, Shaanxi Normal University,
Xi'an 710062, P. R. China)

Abstract: Whether a model can be identified is a basic characteristic of the model before studying parameter estimation. Until recently, the classical susceptible-infectious-recovered (SIR) model is still one of the most commonly used models. In present work the algebraic identifiability of the SIR model by using high-order derivative method (HODM) and multiple time points method (MTPM) was studied. The results indicate that the SIR model can be identified if only the infectious was reported, and MTPM is much better than HODM. Using the data of the flu, the least square method was adopted to estimate the parameters of the SIR model. The result further confirmed that the SIR model was identifiable. The methods developed here could be applied to investigate other type models and left those for future studies.

Key words: identifiability; HODM; MTPM; identification function; parameter estimation