# Gompertz: A Scilab Program for Estimating Gompertz Curve Using Gauss-Newton Method of Least Squares

**Surajit Ghosh Dastidar**

ICFAI University

### Abstract

A computer program for estimating Gompertz curve using Gauss-Newton method of least squares is described in detail. It is based on the estimation technique proposed in Reddy (1985). The program is developed using Scilab (version 3.1.1), a freely available scientific software package that can be downloaded from http://www.scilab.org/. Data is to be fed into the program from an external disk file which should be in Microsoft **Excel** format. The output will contain sample size, tolerance limit, a list of initial as well as the final estimate of the parameters, standard errors, value of Gauss-Normal equations namely $GN_1$ $GN_2$ and $GN_3$, No. of iterations, variance($\sigma^2$), Durbin-Watson statistic, goodness of fit measures such as $R^2$, $D$ value, covariance matrix and residuals. It also displays a graphical output of the estimated curve vis a vis the observed curve. It is an improved version of the program proposed in Dastidar (2005)

*Keywords*: Scilab, Gompertz, statistical software, growth modeling.

## 1. Introduction

Despite the diversity of available statistical software packages such as SAS (SAS Institute Inc 2005), SPSS (SPSS Inc 2003), PASS (NCSS Statistical Software 2005) there is a shortage of stand-alone programs that allows the user to estimate the parameters of the Gompertz curve using Gauss-Newton method of least squares. SPSS 12.0 (SPSS Inc 2003) provides SPSS Regression Models as add-on enhancements to the full SPSS (SPSS Inc 2003) Base System which has a general non-linear regression. SAS 9.1 (SAS Institute Inc 2005) however does provide a programming interface to develop a program for Gompertz growth models using modified Gauss-Newton method. PASS (NCSS Statistical Software 2005) provides users with a modified Gompertz growth model but the method for estimating the parameters of Gompertz curve cannot be specified by the user.

This paper documents a computer program for estimating the parameters of the Gompertz curve using Gauss-Newton method of least squares. It is an improved version of the program proposed in Dastidar (2005). The program is based on the estimation technique proposed in Reddy (1985). In this program

- The sample size need not be a multiple of three.

- The user can specify the range of the data for calculating the initial estimates of $A_0$, $B_0$ and $C_0$ but the range should be a multiple of three.

In this program the critical input data to apply the computer code are:

1. Sample Size

2. Tolerance Limit

3. Input Data File

4. First Observation No.

5. Last Observation No.

It should be noted that (Last Observation No.-First Observation No.) should be a multiple of three.

The output will contain the following:

1. Sample Size

2. Tolerance Limit

3. Initial estimate of the parameters

4. Final estimate of the parameters

5. Standard Errors

6. No. of iterations

7. Value of Gauss-Normal equations namely $GN_1$, $GN_2$ and $GN_3$

8. Variance ($\sigma^2$)

9. $R^2$

10. $D$ where $D = \sum_{t=1}^{n}(y_t - \hat{y}_t)^2 / \sum_{t=1}^{n}(y_t - \bar{y}_t)^2$

11. Covariance Matrix

12. Residuals

13. Graphical display of the estimated Gompertz curve vis a vis the observed data.

The program is written in Scilab (version 3.1.1) (INRIA-ENPC 2005) and was run on a desktop PC with a Pentium IV processor and 256 KB of RAM. Scilab is a scientific software package to develop programs that requires numerical computations. It provides a powerful open computing environment for engineering and scientific applications. It can be downloaded freely from http://www.scilab.org/

## 2. A brief survey of well known statistical software

SPSS 12.0 (SPSS Inc 2003) provides general non-linear regression in the add-on enhancement (SPSS Regression Models) to the full SPSS Base System. However, SAS 9.1 (SAS Institute Inc 2005) does provide a programming interface to develop a program using NLIN procedure. PASS (NCSS Statistical Software 2005) software uses the following modified Gompertz growth model in its package. But the embedded procedure for calculating the estimated parameters cannot be specified by the user.

$$Y = A\exp(-\exp(-B(X - C)))\qquad(1)$$

This model is mentioned in Seber and Wild (1989)

## 3. Program description

The program proceeds as follows:

Before executing the program the following points needs to be checked.

1. The Scilab (version 3.1.1) software should be installed in the computer.

2. A new directory named Scilab should be created in C:\

3. The necessary input file should be created in Microsoft **Excel** format and the data should be entered in the format as specified.(Figure 3) It should be noted that separate excel files need to be created for separate variables.

4. All required input files should be saved in C:\Scilab directory.

5. The executable file 'Gompertz.sce' also needs to be stored in C:\Scilab directory.

6. The program is executed using the following command in the Scilab command prompt. exec('C:\Scilab\Gompertz.sce');

The critical user input to apply the computer program are sample size and tolerance limit. After these two inputs are entered from the keyboard, the program asks for the input file (stored in C:\Scilab directory). After the necessary data file has been chosen, it displays the input data and then the program output. There is only one output format listing sample size, tolerance limit, the names of the parameters, initial estimates, final estimates, standard errors, value of Gauss-Normal equations namely $GN_1$, $GN_2$ and $GN_3$, No. of iterations, variance ($\sigma^2$), Durbin-Watson statistic (DW), goodness of fit measures such as $R^2$, $D$ value where $D = \sum_{t=1}^{n}((y_t - \hat{y}_t)^2 / \sum_{t=1}^{n}((y_t - \bar{y}_t)^2$ ,covariance matrix and residuals. It also displays a graphical output of the estimated curve vis a vis the observed curve.(Figure 1, Figure 2)

## 4. Notation and theory

Suppose we have observed time series data $y_t$ for $t = 1, 2, \ldots, n$

$$y_t = a \, b^{c^t} \, u_t \tag{2}$$

Where $a$, $b$ and $c$ are positive constants and $\log u_t$'s are independent and identically distributed normal random variates with mean 0 and variance $\sigma^2$. Using logarithmic transformation we have

$$
\begin{aligned}
\log y_t &= \log a + c^t \log b + \log u_t \qquad \text{or,} \\
Y_t &= A + C^t B + U_t
\end{aligned}
\tag{3}
$$

where $Y_t = \log y_t$, $C = c$, $B = \log b$ and $U_t = \log u_t$.

From the least squares technique that is by taking partial derivatives of $\sum_{t=1}^{n} U_t^2$ with respect to $A$, $B$ and $C$ and equating each one of them to zero we obtain

$$
\begin{aligned}
GN_1 &= \sum_{t=1}^{n} (Y_t - A - BC^t) = 0 \\
GN_2 &= \sum_{t=1}^{n} (Y_t - A - BC^t)C^t = 0 \\
GN_3 &= \sum_{t=1}^{n} (Y_t - A - BC^t)B^t C^{t-1} = 0
\end{aligned}
\tag{4}
$$

The estimates of $A$, $B$ and $C$ are obtained by solving the above system of Gauss-Normal Equations 4. Also, it is easy to note that the least squares residual sum of squares is given by

$$[\sum_{t=1}^{n}(Y_t - \bar{Y})^2] = [\sum_{t=1}^{n}(Y_t - \bar{Y})^2(C^t - C)^2]/[\sum_{t=1}^{n}(C^t - \bar{C})^2] \tag{5}$$

where $\bar{Y} = (1/n)\sum_{t=1}^{n} Y_t$ and $\bar{C} = (1/n)\sum_{t=1}^{n} C^t$

Unlike the least squares estimators in the linear regression, the estimates of $A$, $B$ and $C$ may not be unique.

Further, it may be noted here that it is not possible to obtain analytical expressions for $C$ which satisfies Equations 4 and also minimizes Equation 5 (Reddy 1977). However, given such a value of $C$, the estimates of $A$ and $B$ are given by

$$
\begin{aligned}
\hat{A} &= \bar{Y} - \hat{B}\bar{C} \\
\hat{B} &= \sum_{t=1}^{n}(Y_t - \bar{Y})(C_t - C)/\sum_{t=1}^{n}(C^t - \bar{C})^2
\end{aligned}
\tag{6} \tag{7}
$$

Gauss-Newton method of least squares is applied in the algorithm to obtain the estimators of $A$, $B$ and $C$ and their standard errors apart from the usual goodness of fit measures such as $R^2$ and $D$ where

$$D = \sum_{t=1}^{n}(y_t - \hat{y}_t)^2 / \sum_{t=1}^{n}(y_t - \bar{y}_t)^2 \tag{8}$$

where $\log \hat{y}_t = \log \hat{a} + \hat{c}^t . \log \hat{b}$. This Gauss-Newton method consists of taking linear expansion of $Y_t = f_t(A, B, C) = A + B.C^t$ around $A_0$, $B_0$ and $C_0$ and retaining the first degree terms and then using ordinary least squares method to obtain $A$, $B$ and $C$. The initial values of $A_0$, $B_0$ and $C_0$ are given as follows:

$$C_0 = (D_2/D_1)^{1/r} \tag{9}$$
$$B_0 = [(1 - C_0)/C_0].[D_1^3/(D_1 - D_2)^2] \tag{10}$$
$$A_0 = (1/3r)[(S_1 + S_2 + S_3) - [(D_1^2 + D_1 D_2 + D_2^2)/(D_1 - D_2)]] \tag{11}$$

Where

$$S_1 = \sum_{t=1}^{r} Y_t,$$

$$S_2 = \sum_{t=r+1}^{2r} Y_t,$$

$$S_3 = \sum_{t=2r+1}^{3r} Y_t,$$

$$D_1 = S_1 - S_2,$$

$$D_2 = S_2 - S_3.$$

And without the loss of generality it is assumed that the sample size for calculating the initial estimates $A_0$, $B_0$ and $C_0$ is a multiple of three(n=3r). The above analytical expressions for $A_0$, $B_0$ and $C_0$ are obtained by assuming the deterministic relation.

$$Y_t = A + BC^t \tag{12}$$

For $t = 1, 2, \ldots, n$ and solving for $A$, $B$ and $C$ using essentially three sets of equations. The estimation procedure of $A_0$, $B_0$ and $C_0$ is termed as 'three point' method and this is a modified version of the method given in Chakravarti and Laha (1967).

Using vector notation let

$$\theta' = [\theta_1, \theta_2, \theta_3] = [A, B, C], \tag{13}$$
$$Y' = (Y_1, Y_2, \ldots, Y_n), \tag{14}$$
$$f'(\theta) = (f_1(\theta_1), f_2(\theta_2), \ldots, f_n(\theta_n)). \tag{15}$$

Let $F$ be an $n \times 3$ matrix of partial derivatives where $F = \delta f_i(\theta_i)/\delta(\theta_j)$; $i = 1, 2, \ldots, n$ and $j = 1, 2, 3$ and let $F_0$ be the value of $F$ evaluated at the initial values $\theta_0 = (\theta_{10}, \theta_{20}, \theta_{30})$. From the Gauss-Newton method we get the estimate of

$$(\theta - \theta_0) = (F_0' F_0)^{-1} F_0' (Y - f(\theta_0)) \tag{16}$$

Or the estimate of $\theta$ is given by

$$\hat{\theta}_0 = \theta_0 + (F_0' F_0)^{-1} F_0' (Y - f(\theta_0)) \tag{17}$$

This procedure is repeated with the new values as the starting values and the procedure is terminated if the successive values of $(\theta - \theta_0)$ or $((\theta - \theta_0)/\theta_0)$ are very small and less than a pre-specified limit.

Since the residuals $U_i$ are independent and identically distributed normal variates with mean 0 and variance $(\sigma^2)$ and if $\hat{\theta}$ is the final estimate of $\theta$, then

$$\hat{\sigma}^2 = (1/n)\sum_{t=1}^{n}(Y_i - f_i(\hat{\theta}))^2,$$

$$R^2 = 1 - \left(n\hat{\sigma}^2/\sum_{t=1}^{n}(Y_i - \bar{Y})^2\right).$$

And $\hat{\theta}$ is approximately normally distributed with mean $\theta$ and covariance matrix $\sigma^2\,[F(\theta)'F(\theta)]^{-1}$. The estimate of the covariance matrix is given by

$$\hat{\sigma}^2[F(\hat{\theta})'F(\hat{\theta})]^{-1} \tag{18}$$

The estimates of $GN_1$, $GN_2$ and $GN_3$ are obtained by substituting the final estimates of $A$, $B$ and $C$ in Equations 4.

To fit Gompertz curve for the given data, using Gauss-Newton method of least squares the following procedure is adopted in the computer program:

1. the estimates of $A$, $B$ and $C$ obtained from the 'three point' method are used as the initial values of $\theta_0' = (\theta_{10}, \theta_{20}, \theta_{30})$

2. final convergence of $\theta$ occurs at the $r$th iteration if

$$(\theta_{io}^{(r+1)}) - \theta_{io}^r)/\theta_{io}^r \leq \text{tolerance limit (user input)} \tag{19}$$

For all values of i, where $\theta_{io}^{(r+1)}$ denotes the values of i obtained at the (r+1) th iteration. The tolerance limit can be specified by the user.

The advantage of Gauss-Newton method is that it is possible to conclude statistically that

1. $y$ is increasing at an increasing rate, i.e., there is acceleration in $y$ if $\log b > 0$ and $c > 1$

2. $y$ is increasing at a decreasing rate, i.e., there is deceleration in $y$ if $\log b < 0$ and $c < 1$

3. there is neither acceleration nor deceleration in $y$ if the 95% confidence interval of $\log b$ contains 0 and the 95% confidence interval of $c$ contains 1.

# 5. Restriction

All input parameters are checked for nullity and a fault message is returned if there is an illegal entry. Added to that the following are checked:

1. Sample size should be greater than five

2. Tolerance limit should be greater than 0

3. Input file name cannot be left blank

# 6. Sample input/output cases

The output of the computer program is in tabular form. It will contain a list of sample size, tolerance limit, initial as well as the final estimate of the parameters, standard errors, value of Gauss-Normal equations namely $GN_1$, $GN_2$ and $GN_3$, No. of iterations, variance, Durbin-Watson statistic ($DW$), goodness of fit measures such as $R^2$, $D$, covariance matrix and residuals. It also displays a graphical output of the estimated curve (Figure 1, Figure 2). The program has been analyzed using , data on Tiwari's estimates of Gross Domestic Product (TGDP), Industrial Production (TIP), see Table 1. The sample size taken was 15 and tolerance limit = 0.005. All the 15 observation were taken for calculating the initial estimates $A_0$, $B_0$ and $C_0$.

The input file will be an **Excel** format as shown in Figure 3. The output file will be in text format. The name of the executable file is `Gompertz.sce` which contains the code for estimating Gomepertz curve.
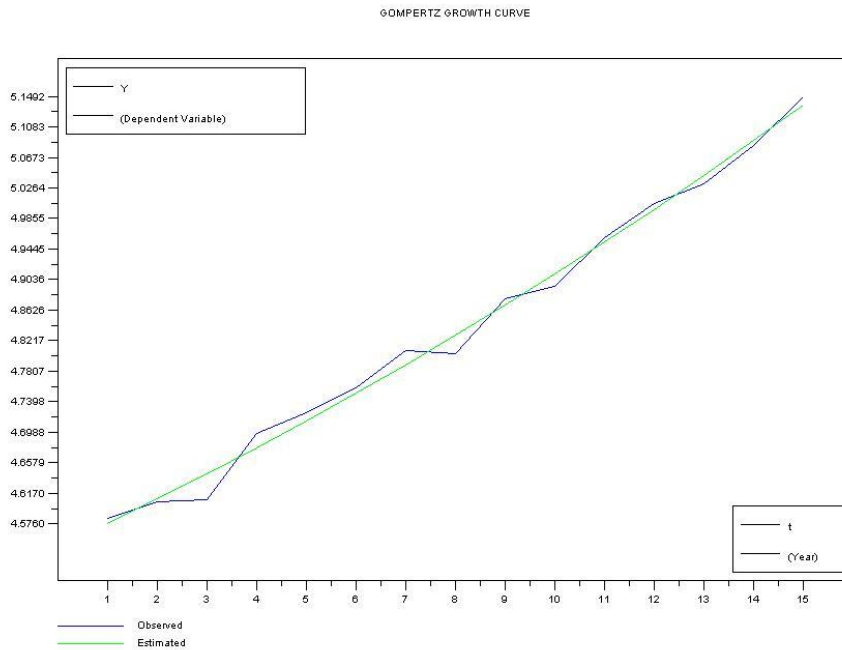


Figure 1: Plot of TGDP against year

The TGDP series is depicted in Figure 1, the output from Scilab is given below.

```
REPORT SHOWING RESULTS
----------------------
Sample Size = 15                    Tolerance Limit = 0.005000
```

```
PARAMETER INITIAL ESTIMATES   FINAL ESTIMATES   STD. ERRORS      DW
-------  ------------------    ---------------   -----------    ------


A        3.583742               3.448600          0.571081     2.578518
B        0.961720               1.095294          0.560680
C        1.032668               1.029317          0.012232


 No. of Iterations: = 3

GN1 = 0.0000000

GN2 = 0.0000000

GN3 = 0.0000001

Sigma_Hat_Square= 0.000232      R_Square= 0.992328      D=0.007672

 COVARIANCE MATRIX
 ------------------------
0.326134        -0.320167        0.006975
-0.320167        0.314363        -0.006851
0.006975        -0.006851        0.000150

 SHOWING RESIDUALS

Y = 0.007635
Y = -0.003788
Y = -0.034316
Y = 0.019377
Y = 0.010319
Y = 0.007159
Y = 0.019571
Y = -0.023907
Y = 0.008351
Y = -0.016757
Y = 0.006385
Y = 0.007299
Y = -0.011531
Y = -0.006885
Y = 0.011088
```
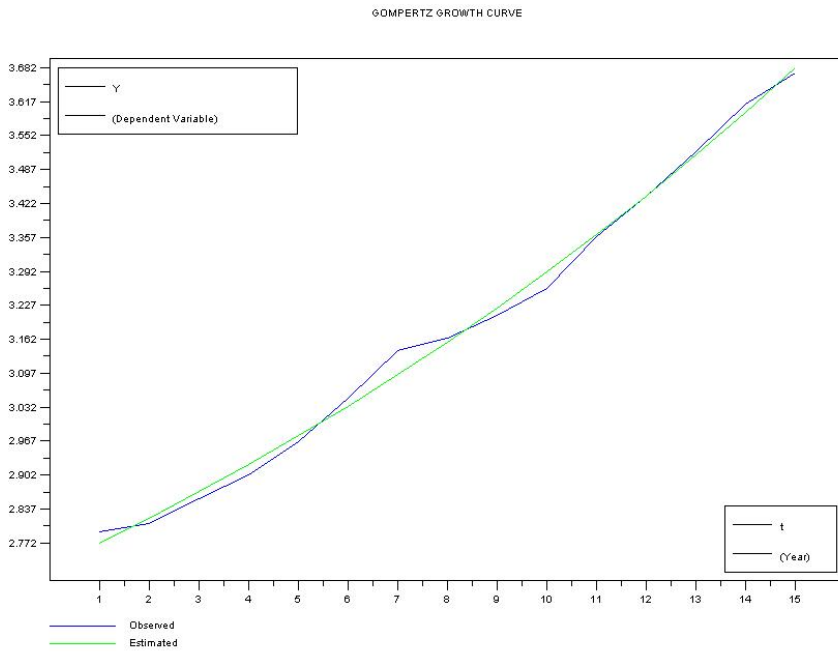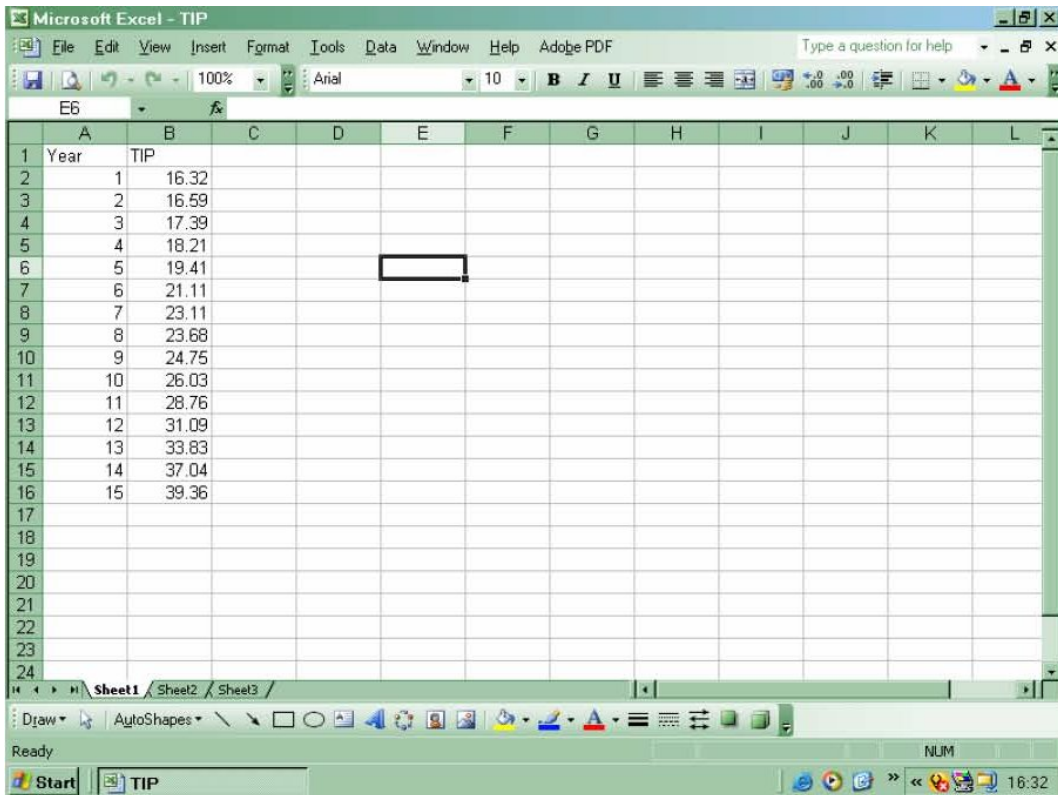
Figure 2: Plot of TIP against year



Figure 3: Input **Excel** data file for TIP

The TIP series is depicted in Figure 1, the corresponding input **Excel** data file can be seen in Figure 3 and the output from Scilab is given below.

```
REPORT SHOWING RESULTS
----------------------
Sample Size = 15                     Tolerance Limit = 0.005000

PARAMETER INITIAL ESTIMATES    FINAL ESTIMATES    STD. ERRORS       DW
-------- ----------------      ---------------    -----------     ------


A        1.270229                1.726868          0.293523      1.250812
B        1.436541                0.999303          0.281053
C        1.035029                1.045774          0.009470

 No. of Iterations: = 3

GN1 = -0.0000005

GN2 = -0.0000009

GN3 = -0.0000106

Sigma_Hat_Square= 0.000349       R_Square= 0.995587        D=0.004413

 COVARIANCE MATRIX
 ------------------------
0.086156          -0.082460        0.002769
-0.082460          0.078991        -0.002656
0.002769          -0.002656        0.000090

 SHOWING RESIDUALS

Y = 0.020478
Y = -0.010949
Y = -0.013880
Y = -0.020121
Y = -0.011014
Y = 0.015730
Y = 0.046414
Y = 0.008207
Y = -0.013036
Y = -0.031044
Y = -0.002874
Y = 0.000186
Y = 0.006381
Y = 0.015183
Y = -0.009662
```

| Year | TGDP | TIP |
|---|---|---|
| 1951-52 | 100.01 | 16.59 |
| 1952-53 | 100.36 | 17.39 |
| 1953-54 | 109.67 | 18.21 |
| 1954-55 | 112.67 | 19.41 |
| 1955-56 | 116.56 | 21.11 |
| 1956-57 | 122.61 | 23.11 |
| 1957-58 | 122.10 | 23.68 |
| 1958-59 | 131.31 | 24.75 |
| 1959-60 | 133.50 | 26.03 |
| 1960-61 | 142.61 | 28.76 |
| 1961-62 | 149.18 | 31.09 |
| 1962-63 | 153.20 | 33.83 |
| 1963-64 | 161.28 | 37.04 |
| 1964-65 | 172.30 | 39.36 |

Table 1: Data on Tiwari's estimates of Gross Domestic Product (TGDP), Industrial Production (TIP), given in crores (10 million) of Indian Rupees in 1960/61 prices.

# Acknowledgments

# References

Chakravarti IM, Laha RG (1967). *Handbook of Methods of Applied Statistics.* John Wiley and Sons, New York.

Dastidar SG (2005). "**Gompertz** (Version 1.00): A Computer Program for Estimating Gompertz Curve Using Gauss-Newton Method of Least Squares." *The ICFAI Journal of Systems Management*, **III**(4).

INRIA-ENPC (2005). *Scilab (Version 3.1.1).* URL http://www.scilab.org/.

NCSS Statistical Software (2005). *Power Analysis Software (PASS).* URL http://www.ncss.com/.

Reddy VN (1977). "Statistical Fitting of Growth Curves with illustrations from Data on Indian Economy." *Technical report*, Center for Management and Development Studies, IIM, Kolkata.

Reddy VN (1985). "On Estimating Gompertz Curve." Indian Institute of Management(IIM), Kolkata, Unpublished manuscript.

SAS Institute Inc (2005). *SAS/STAT® User's Guide, Version 9.1.3.* Cary, NC.

Seber GAF, Wild CJ (1989). *Nonlinear Regression.* John Wiley and Sons, New York.

SPSS Inc (2003). *SPSS Regression Models 12.0*. URL http://www.spss.com/.

**Affiliation:**

Surajit Ghosh Dastidar
ICFAI Institute for Management Teachers (IIMT)
3rd Floor, Astral Heights, 6-3-352/2 & 3
Road No.1, Banjara Hills
Hyderabad - 500 034, India
E-mail: sghoshdastidar@gmail.com