# SEMIPARAMETRIC MODELING
# OF CONSUMER PRICE INDEX

**Budi Lestari**

Jurusan Matematika, Fakultas MIPA, Universitas Jember
Jl. Kalimantan 37 Jember 68121, Indonesia

**Abstract.** Many classical data, for example, exchange rate, stock price, and consumer price index (CPI) data cannot be analyzed under independent observation assumption. In addition, some time series data cannot be modeled well into a fully linear model, for instance, CPI, price of raw materials for some certain industries and price of some industrial products data in which monetary crisis of Indonesia in 1998 has caused a dramatic effect on the time series of CPI, price of raw materials and industrial products. A semiparametric model is a mixture model between parametric and nonparametric models. If we apply it to time series data, we will obtain a semiparametric time series model known as partly linear autoregressive model:

$$y_t = \beta \, y_{t-1} + g(y_{t-2},...,y_{t-p}) + \varepsilon_t \ \text{ for } \ t \ge p+1.$$

Here $\beta$ is an unknown parameter to be estimated, $g(.)$ is an unknown function in $R^{p-1}, \varepsilon_t$ are i.i.d. random errors with $E(\varepsilon_1) = 0$ and $E(\varepsilon_1^2) < \infty$, and $\varepsilon_t$ are independent of $y_s$ for all $s = 1,2,..., p$ and $t \ge p+1$. Based on the model above, we investigate a model for general consumer price index (GCPI) of Jember data recorded monthly from January 1998 to December 2002 by Statistic Center Bureau of Jember.

**Keywords:** Least Square Estimator, Kernel Estimator, Semiparametric Model, Consumer Price Index

**MSC 2020** : 62F99

## 1. Introduction

In last twenty years, many authors have shown interest in semiparametric regression models, especially, partly linear regression models, $y_j = x_j'\beta + g(t_j) + \varepsilon_j$ for $j = 1,2,...$; e.g. Ansley and Wecker (1983), Heckman (1986), Rice (1986), Chen (1988), Chen and Shiau (1991), and Andrews (1991) investigated asymptotic behaviors by using spline smoothing techniques and Robinson (1988), Speackman (1988), Gao (1992), and Gao, et al (1994) also obtained some asymptotic results by using kernel or nearest neighbor smoothing under the case where the $(x_j, t_j)$ are i.i.d. random variables or fixed design points, and Wahba (1990), Budiantara (1999) and Hardle (2000) studied partly linear model for independent observations case. However, many classical data, for example, exchange rate data, stock price data, and consumer price index data, cannot be analyzed under independent observation assumption. In addition, some time series data cannot be

modeled as a fully linear model, for instance, consumer price index (CPI) data recorded monthly from January 1998 to December 2002 by Statistic Center Bureau of Jember in which monetary crisis of Indonesia in 1998 has caused a dramatic effect on the time series of CPI. Lestari (2001), Lestari (2003), Lestari (2004), and Lestari (2005) studied estimation, asymptotic normality and iterated logarithm distribution of partly linear autoregressive estimators; and applied partly linear autoregressive model to the consumer price index of Jember data. Beside that, latest development in the semiparametric regression has given a strong foundation to semiparametric time series analysis.

The important thing is how we analyze and apply semiparametric regression model to time series data. The resulting model is known as partly linear autoregressive model given by :

$$y_t = \beta \, y_{t-1} + g(y_{t-2}, ...., y_{t-p}) + \varepsilon_t, \quad t \geq p+1 \tag{1}$$

where $\beta$ is an unknown parameter to be estimated, $g(.)$ is an unknown function in $R^{p-1}$, $\varepsilon_t$ are i.i.d. random errors with $E(\varepsilon_1) = 0$ and $E(\varepsilon_1^2) < \infty$, and $\varepsilon_t$ are independent of $y_s$ for all $s = 1,2,...,p$ .

In this paper, we investigate a semiparametric time series model for the general consumer price index (GCPI) of Jember data which is known as partly linear autoregressive model.

## 2. Methodology

The model (1) contains unknown parameter $\beta$ and unknown function $g(.)$ to be estimated.  Least square method is used to obtain estimator $\hat{\beta}_T$ of $\beta$. The estimate of $g(.)$ can be obtained by using kernel estimator approach. In association with the GCPI data, of course, Firstly we identify the parametric and nonparametric components of data. Secondly, we determine an optimum bandwidth for kernel estimator approach with generalized cross validation (GCV) criterion and Gaussian kernel function. Finally, we estimate the model of GCPI by investigating and comparing their mean square errors (MSE), errors and plot for three models approach, i.e., linear autoregressive model, ARIMA model, and partly linear autoregressive model. For estimating and plotting, we use MINITAB, S-PLUS, and MATLAB programs.

## 3. Results

### 3. 1.  Estimation of Function $g(.)$ and Parameter $\beta$

Assume that $\{y_t, t = p+1, p+2,...., T\}$ satisfy the partly linear autoregressive model :

$$y_t = \beta y_{t-1} + g(y_{t-2},...., y_{t-p}) + \varepsilon_t, \ t \geq \text{p+1}. \tag{2}$$

Then,  we have :

$$
\begin{aligned}
g(y_{t-2},..., y_{t-p}) &= E\big((y_t - \beta y_{t-1}) \mid y_{t-2},..., y_{t-p}\big) \\
&= E(y_t \mid y_{t-2},..., y_{t-p}) - \beta E(y_{t-1} \mid y_{t-2},..., y_{t-p}) \\
&= g_1(y_{t-2},...., y_{t-p}) - \beta g_2(y_{t-2},...., y_{t-p}) \ \text{(say)}
\end{aligned}
\tag{3}
$$

Hence, the natural estimates of $g_1(.)$, $g_2(.)$ and $g(.)$ of (3) can be obtained by using kernel estimator approach defined by :

$$\hat{g}_{1T}(x_t) = \sum_{s=p+1}^{T} W_{Ts}(x_t) y_s, \tag{4}$$

$$\hat{g}_{2T}(x_t) = \sum_{s=p+1}^{T} W_{Ts}(x_t) y_{s-1}, \tag{5}$$

and

$$\hat{g}_T(x_t) = \hat{g}_{1T}(x_t) - \beta \hat{g}_{2T}(x_t) \tag{6}$$

where $x_t = (y_{t-2},..., y_{t-p})$ and $W_{ts}(.)$ are kernel weight function Nadaraya-Watson given by :

$$W_{Ts}(x) = K\big((x-x_s)/h\big) \Big/ \sum_{t=p+1}^{T} K\big((x-x_t)/h\big) \tag{7}$$

where $K : R^{p-1} \to R$ is a function and $\{h = h_T; T > p\}$ is a sequence of nonnegative real numbers.

Now, based on the model $y_t = \beta y_{t-1} + \hat{g}_T(y_{t-2},...., y_{t-p}) + \varepsilon_t$, $t = p+1,...,T$, we have the least square estimator $\hat{\beta}_T$ of $\beta$ :

$$\hat{\beta}_T = \frac{\displaystyle\sum_{t=p+1}^{T}\big(y_{t-1} - \hat{g}_{2T}(y_{t-2},..., y_{t-p})\big)\big(y_t - \hat{g}_{1T}(y_{t-2},..., y_{t-p})\big)}{\displaystyle\sum_{t=p+1}^{T}\big(y_{t-1} - \hat{g}_{2T}(y_{t-2},..., y_{t-p})\big)^2} \tag{8}$$

### 3.2. Investigation of GCPI Model.

Based on the model (1), we investigate a model for general consumer price index (GCPI) data recorded monthly from January 1998 to December 2002 by Statistic Center Bureau of Jember. The GCPI data and their series plot are shown in Table 1 and Figure 1,

respectively.

Table 1.  General consumer price index  from january 1998 to december 2001

| Year | GCPI | Year | GCPI | Year | GCPI | Year | GCPI | Year | GCPI |
|------|------|------|------|------|------|------|------|------|------|
| 1998 | 183.69 | 1999 | 187.8 | 2000 | 196.99 | 2001 | 198 | 2002 | 218.92 |
|      | 185.86 |      | 188.01 |      | 197.12 |      | 199.87 |      | 217.1 |
|      | 186.65 |      | 190.45 |      | 197.23 |      | 197.28 |      | 210.27 |
|      | 192.55 |      | 193.1 |      | 197.29 |      | 199.67 |      | 218.07 |
|      | 195.14 |      | 194.25 |      | 197.37 |      | 198.83 |      | 219.48 |
|      | 194.77 |      | 192.35 |      | 197.4 |      | 199.54 |      | 219.96 |
|      | 193.35 |      | 193.5 |      | 197.5 |      | 199.79 |      | 218.88 |
|      | 192.75 |      | 195.2 |      | 197.73 |      | 199.92 |      | 216.89 |
|      | 193.1 |      | 196.37 |      | 197.76 |      | 200.21 |      | 210.09 |
|      | 189.92 |      | 195.88 |      | 197.84 |      | 218.09 |      | 216.25 |
|      | 188.1 |      | 198.22 |      | 197.9 |      | 219.28 |      | 215.59 |
|      | 186.69 |      | 196.41 |      | 197.95 |      | 220 |      | 219.06 |

**Source :** *Statistic Center Bureau of Jember*

### 3.2.1.  Autoregressive Model Approach

In this approach, we consider a linear model as follows :
$$y_t = \beta_0 + \beta_1 y_{t-1} + \beta_2 y_{t-2} + ... + \beta_p y_{t-p} + \varepsilon_t . \tag{9}$$

The linear regression analysis result on the data gives an estimated model :
$$\hat{y}_t = 8.81 + 0.923 y_{t-1} + 0.036 y_{t-2} \tag{10}$$

Validation result of $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\beta}_2$ shows that $\hat{\beta}_1$ is significant, and $\hat{\beta}_0$ and $\hat{\beta}_2$ are not significant. Plot between $y_t$ and $\hat{y}_t$ is shown in Figure 2.
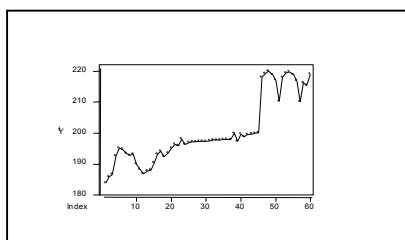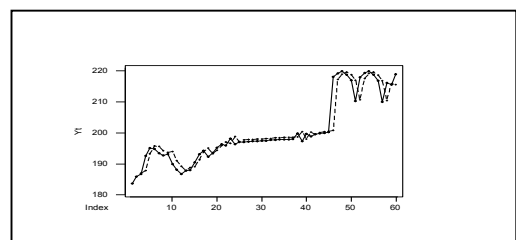


Figure 1.  Series Plot of GCPI Data

Figure 2.  Plot Between $y_t$ ( —— ) and
$$\hat{y}_t = 8.81 + 0.923 y_{t-1} + 0.036 y_{t-2} \ (\text{---})$$

In addition, plotting of residuals as shown model (10) shows that there are a right skew, outliers and indicates no normal distribution. Next, by getting out variables with no significant coefficients from the model, we obtain :
$$y_t = 1.00287 y_{t-1} + \varepsilon_t \tag{11}$$

Plotting of residuals model (11) shows that there are a right skew, outliers and indicates no normal distribution. Plot between $y_t$ and $\hat{y}_t = 1.00287 y_{t-1}$ is shown in Figure 3. Residuals plots for models (10) and (11) are shown in Figure 4 and Figure 5, respectively.
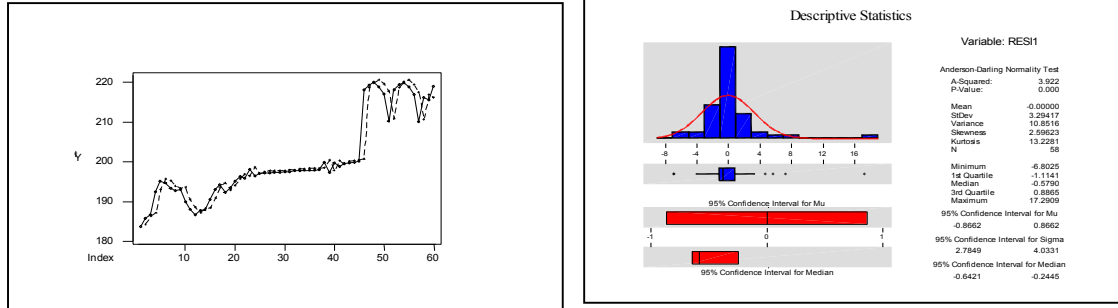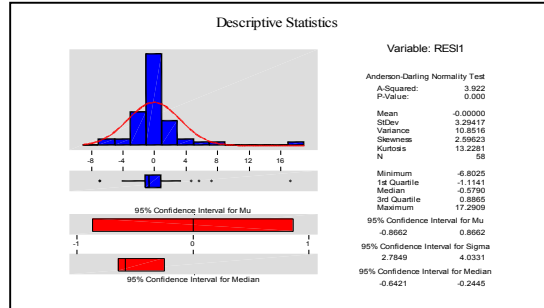


Figure 3. Plot between $y_t$ and $\hat{y}_t = 1.00287 y_{t-1}$   Figure 4.  Residuals plot for Model (10)

### 3.2.2.  ARIMA Model Approach

Figure 1 shows a series with trend. It means that the series has not been stationary yet. Next, by using difference $d = 1$, i.e., $w_t = y_t - y_{t-1}$, we have series plot of $w_t$ as shown in Figure 6. Figure 6 shows that data has been stationary. Based on the autocorrelation function (ACF) and partial autocorrelation function (PACF) plots in which they cut off after lag 5 and lag 6, we have the estimated model of **ARIMA(1,1,0)(1,0,0)[6]**
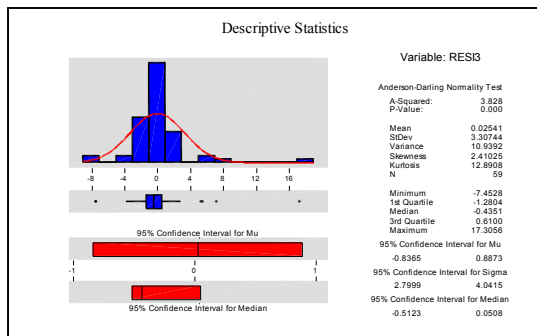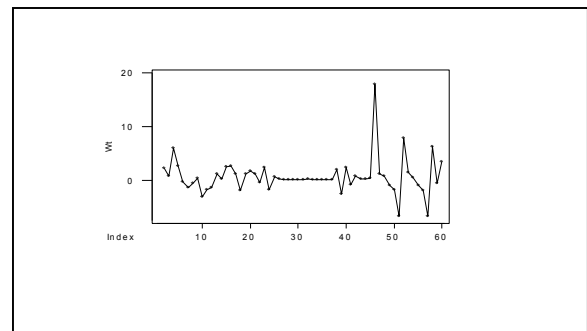


Figure 5.  Residuals plot for model (11)   Figure 6. Series plot with difference $d = 1$

Although diagnostic checking result shows that residuals has been white noise but estimation of parameters result shows that SAR parameter is significant different from zero and AR parameter is not significant. Plot of residuals normality shows that residuals are not normal and have outliers as shown in Figure 7.

### 3.2.3.  Semiparametric Model Approach

By plotting between $y_t$ and $y_{t-1}$; and between $y_t$ and $y_{t-2}$ we knew that plotting result between $y_t$ and $y_{t-1}$ was relatively linear (Figure 8). On the other hand, plotting result between $y_t$ and $y_{t-2}$ was not linear and had no certain pattern (Figure 9). Therefore, $y_{t-1}$ and $y_{t-2}$ were as parametric and nonparametric components, respectively.
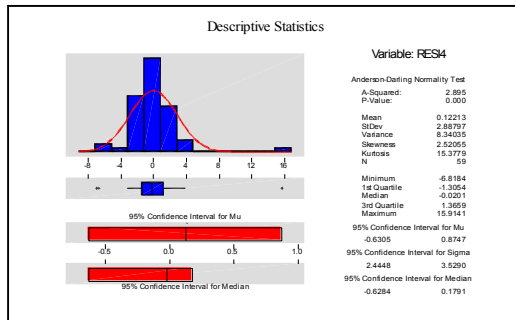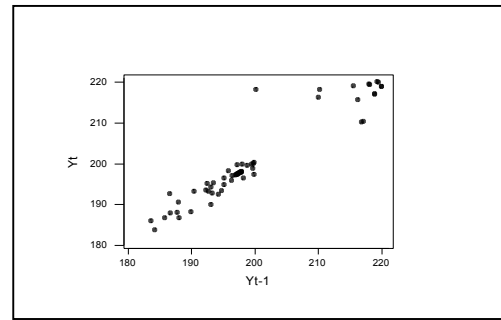
Figure 7. Residuals plot for ARIMA Model



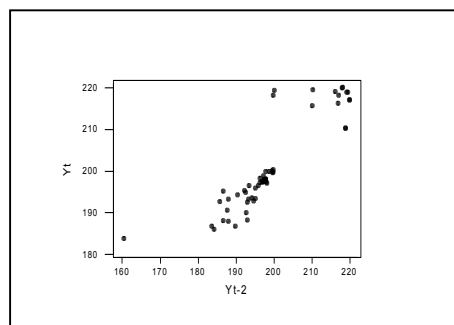Figure 8. Plot between $y_t$ and $y_{t-1}$



Figure 9. Plot between $y_t$ and $y_{t-2}$

It is also supported by linear regression plot results between $y_t$ and $y_{t-1}$ with $R^2 = 90.6$ %, and between $y_t$ and $y_{t-2}$ with $R^2 = 78.6$ % as shown in Figure 10 and Figure 11, respectively.
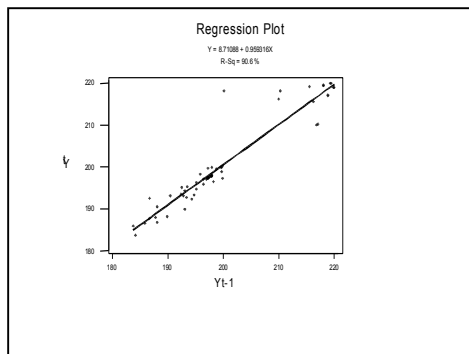




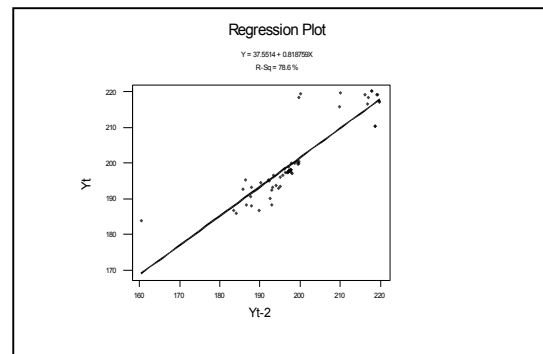Figure 10. Regression plot between $y_t$ and $y_{t-1}$    Figure 11. Regression plot between

$$y_t \text{ and } y_{t-2}$$

After determining both parametric and nonparametric components, we find the optimum bandwidth for Gaussian kernel estimator approach with the generalized cross validation (GCV) criterion. So that, by using programs Matlab and Minitab, we obtained some bandwidth and GCV values as shown in Table 2 and Figure 12, respectively. Table 2 and

Figure 12 show that the optimum bandwidth is 0.017, because GCV minimum, i.e. 1.8447, is obtained at bandwidth (h) = 0.017.

Next, based on that optimum bandwidth we estimated $\hat{g}_T(y_{t-2})$, $\hat{\beta}_T$, $\hat{y}_t$ and $\varepsilon_t$ values by using program Matlab. We got $\hat{\beta}_T = 0.9192$; and $\hat{y}_t$, $\hat{g}_T(y_{t-2})$ and $\varepsilon_t$ as shown in Table 3. It shows that errors are relatively small or leads to zero for every t. So, we obtained the mean square error (MSE), i.e., 0.007083. The relationship of $\hat{y}_t$, $y_{t-1}$ and $y_{t-2}$ formed a surface as shown in Figure 13.
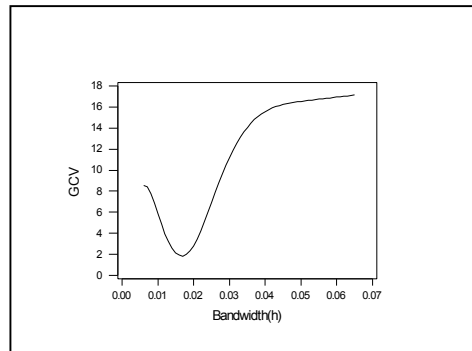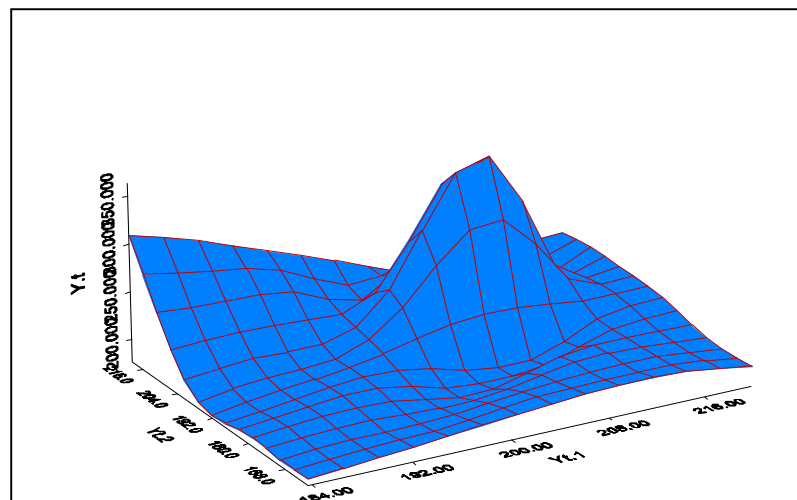


Figure 12. GCV versus bandwidth plot



Figure 13. Surface performed by $\hat{y}_t$, $y_{t-1}$ and $y_{t-2}$

Next, based on that optimum bandwidth we estimated $\hat{g}_T(y_{t-2})$, $\hat{\beta}_T$, $\hat{y}_t$ and $\varepsilon_t$ values by using program Matlab. We got $\hat{\beta}_T = 0.9192$; and $\hat{y}_t$, $\hat{g}_T(y_{t-2})$ and $\varepsilon_t$ as shown in Table 3. It shows that errors are relatively small or leads to zero for every t. So, we obtained the mean square error (MSE), i.e., 0.007083. The relationship of $\hat{y}_t$, $y_{t-1}$ and $y_{t-2}$ formed a surface as shown in Figure 13.

Table 2.  Some GCV values for some bandwidths

| Bandwidth | GCV | Bandwidth | GCV |
|---|---|---|---|
| 0.01 | 5.83 | 0.017 | 1.8447 |
| 0.011 | 4.8474 | 0.018 | 1.9758 |
| 0.012 | 3.9612 | 0.019 | 2.3006 |
| 0.013 | 3.2064 | 0.02 | 2.8035 |
| 0.014 | 2.6038 | 0.021 | 3.4589 |
| 0.015 | 2.1676 | 0.022 | 4.2347 |
| 0.016 | 1.9107 | 0.023 | 5.0963 |

Table 3.  Estimated values of $\hat{y}_t$, $\hat{g}_T(y_{t-2})$ and $\varepsilon_t$

| No | $\hat{y}_t$ | $\hat{g}_T(y_{t-2})$ | $\varepsilon_t$ | No | $\hat{y}_t$ | $\hat{g}_T(y_{t-2})$ | $\varepsilon_t$ |
|---|---|---|---|---|---|---|---|
| 1 | 183.685 | 14.3408 | 0 | 31 | 197.519 | 16.0686 | -0.024 |
| 2 | 185.855 | 17.0072 | 0 | 32 | 197.701 | 16.1586 | 0.024 |
| 3 | 186.645 | 15.8024 | 0 | 33 | 197.755 | 16.0012 | 0 |
| 4 | 192.545 | 20.9762 | 0 | 34 | 197.832 | 16.0513 | 0.0024 |
| 5 | 194.972 | 17.9796 | 0.1632 | 35 | 197.897 | 16.0425 | -0.0024 |
| 6 | 194.765 | 15.392 | 0 | 36 | 197.945 | 16.035 | 0 |
| 7 | 193.347 | 14.3142 | -0.0021 | 37 | 198.018 | 16.0627 | -0.0238 |
| 8 | 192.745 | 15.0174 | 0 | 38 | 199.786 | 17.784 | 0.079 |
| 9 | 193.095 | 15.919 | 0 | 39 | 197.331 | 13.6103 | -0.0563 |
| 10 | 189.915 | 12.4172 | 0 | 40 | 199.87 | 18.5302 | -0.2053 |
| 11 | 188.23 | 13.6553 | -0.1349 | 41 | 199.14 | 15.6035 | -0.3156 |
| 12 | 186.685 | 13.7834 | 0 | 42 | 199.535 | 16.7701 | 0 |
| 13 | 187.795 | 16.1895 | 0 | 43 | 199.785 | 16.3674 | 0 |
| 14 | 188.168 | 15.5424 | -0.1632 | 44 | 199.915 | 16.2676 | 0 |
| 15 | 190.445 | 17.6261 | 0 | 45 | 200.205 | 16.4381 | 0 |
| 16 | 193.095 | 18.0332 | 0 | 46 | 217.879 | 33.8462 | 0.2053 |
| 17 | 194.245 | 16.7472 | 0 | 47 | 219.274 | 18.8057 | 0 |
| 18 | 192.21 | 13.6553 | 0.1349 | 48 | 219.919 | 18.3572 | 0.0747 |
| 19 | 193.495 | 16.6867 | 0 | 49 | 218.914 | 16.69 | 0 |
| 20 | 195.195 | 17.3295 | 0 | 50 | 217.084 | 15.8526 | 0.0102 |
| 21 | 196.365 | 16.9369 | 0 | 51 | 210.265 | 10.7066 | -0.0008 |
| 22 | 195.873 | 15.3693 | 0.0021 | 52 | 218.064 | 24.7841 | 0 |
| 23 | 198.089 | 18.0365 | 0.1252 | 53 | 219.474 | 19.0241 | 0 |
| 24 | 196.405 | 14.2008 | 0 | 54 | 220.029 | 18.2827 | -0.0747 |
| 25 | 196.985 | 16.4446 | 0 | 55 | 218.874 | 16.6868 | 0 |
| 26 | 197.24 | 16.1667 | -0.1252 | 56 | 216.894 | 15.6998 | -0.0102 |
| 27 | 197.225 | 16.0319 | 0 | 57 | 210.083 | 10.7181 | 0.0008 |
| 28 | 197.285 | 15.9908 | 0 | 58 | 216.244 | 23.1296 | 0 |
| 29 | 197.355 | 16.0061 | 0.0096 | 59 | 215.584 | 16.8071 | 0 |
| 30 | 197.082 | 15.6599 | 0.3122 | 60 | 219.054 | 20.8838 | 0 |

### 3.3.  Comparing of Three Models Approach

Based on the discussion above, we can compare mean square errors (MSE), Errors ($\varepsilon_t$) and plots for  these three models approach, i.e., autoregressive, ARIMA and semiparametric (partly linear autoregressive) models as follows:

(i).  MSE of autoregressive, ARIMA, and semiparametric models are 11, 8.502 and 0.007083, respectively. We can see that MSE of semiparametric model approach is smaller than both autoregressive and ARIMA models approach.

(ii).  Errors ($\varepsilon_t$) of  both autoregressive and ARIMA models approach are not normal, but errors ($\varepsilon_t$) of semiparametric model approach is normal as shown in Figure 14. However, these three models approach satisfy white noise condition.

(iii).Plot between estimated values ($\hat{y}_t$) and observation values ($y_t$) for autoregressive, ARIMA and semiparametric models approach are given by Figure 15, Figure 16, and Figure 17, respectively. We can see that plot result of semiparametric model approach gives estimated values ($\hat{y}_t$)  exactly close to observation values ($y_t$). It is different from both autoregressive and ARIMA models approach.
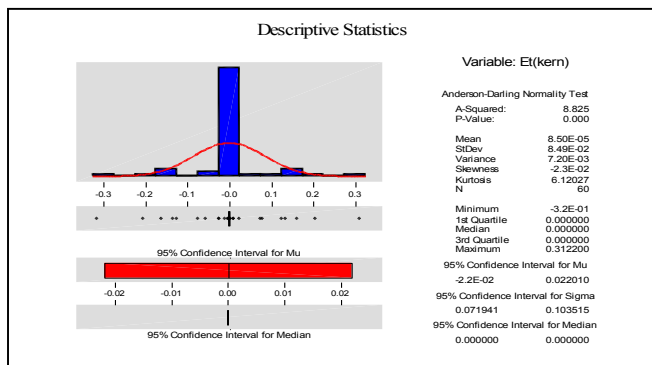


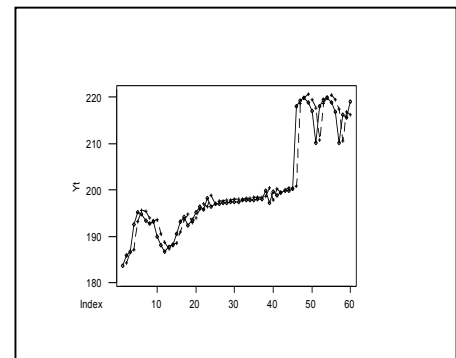Figure 14. Residuals plot for semiparametric Model.

Figure 15. Series plot between $y_t$ and $\hat{y}_t$ for autoregressive model
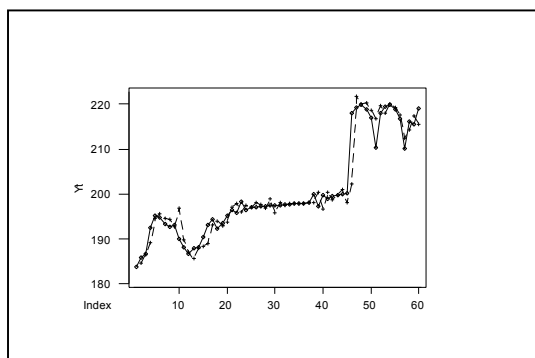


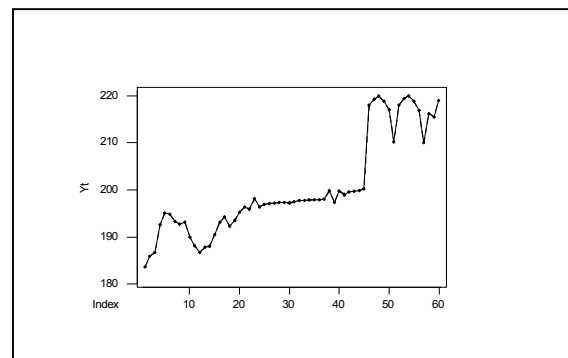Figure 16. Series plot between $y_t$ and $\hat{y}_t$ for ARIMA model.

Figure 17.  Series plot between $y_t$ and $\hat{y}_t$ for semiparametric model.

## 4. Conclusion

Based on the discussion above, we known that using of semiparametric model approach for estimating a model of general consumer price index of Jember is better than both autoregressive and ARIMA models approach. We obtained an estimated model of general consumer price index of Jember as a semiparametric time series model which is called as partly linear autoregressive model as follows :

$$\hat{y}_t = 0.9192 y_{t-1} + \hat{g}_T(y_{t-2}) \tag{12}$$

Where $\quad \hat{g}_T(y_{t-2}) = \hat{g}_{1T}(y_{t-2}) - 0.9192 \hat{g}_{2T}(y_{t-2})$,

$$\hat{g}_{1T}(y_{t-2}) = \sum_{s=3}^{T} W_{Ts}(y_{t-2}) y_s,$$

$$\hat{g}_{2T}(y_{t-2}) = \sum_{s=3}^{T} W_{Ts}(y_{t-2}) y_{s-1},$$

$$W_{Ts}(y_{t-2}) = K\left((y_{t-2} - y_s)h^{-1}\right) \Bigg/ \sum_{j=3}^{T} K\left((y_{t-2} - y_j)h^{-1}\right) \quad \text{and}$$

$$K(y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2}, \quad -\infty < y < \infty.$$

## References

[1]   Andrew, D.W.K. (1991) Asymptotic Theory of Estimates for Nonparametric and Semiparametric Regression Models, *Econometrica*, 59, 307-345.

[2]   Ansley, C.F. and Wecker, W.E. (1983) "Extentions and Examples of The Signal Extration Approach to Regression", In : *Applied Time Series Analysis of Economic Data*, 181-192.

[3]   Budiantara, I. N. (1999) *Estimator Spline Dalam Regresi Nonparametrik dan Semiparametrik*, Doctor Thesis, Gadjah Mada University, Yogyakarta.

[4]   Chen, H. (1988) Convergence Rates for Parametric Components in a Partly Linear Model, *Ann. Statist.*, 16, 136-146.

[5]   Chen, H. and Shiau, J. G. (1991) A Two-stage Spline Smoothing Method for Partially Linear Models, *J. Statist. Plann. Inference*, 25, 187-201.

[6]   Gao, J. T. (1992) *Theory of Large Sample in Semiparametric Regression*, Ph.D. Thesis, Graduate School at China University of Science and Technology, The People Republic of China.

[7]     Gao, J.T., Hong, S.Y., Liang, H., and Shi, P.D. (1994) Survey of Chinese Work on Semiparametric Regression Models, *Chin. Appl. Probab. and Statist.*, 1, 96-104.

[8]     Gao, J. T. (1995) Asymptotic Properties of Some Estimators for Partly Linear Stationary Autoregressive Models, *J. Commun. Statist. Theory Meth.*, 14(8), 2011-2026.

[9]     Hardle, W. (2000) *Partially Linear Models*, Springer-Verlag, London.

[10]   Heckman, H. (1986) Spline Smoothing in Partly Linear Models, *J. Roy. Statist. Soc.*, Ser. B48, 244-248.

[11]   Lestari, B. (2001a) *Kenormalan Asimtotik Dan Distribusi Iterasi Logaritma Estimator Autoregressive Linier Parsial*, Statistical Theory Paper Presented on Statistics National Seminar V at ITS, Surabaya, 1-10.

[12]   Lestari, B. (2001b) *Pemodelan Autoregressive Linier Parsial Untuk Data Indeks Harga Konsumen Kabupaten Jember*, Statistical Applied Paper Presented on Statistics National Seminar V at ITS, Surabaya, 1-13.

[13]   Lestari, B. (2003a) *Estimasi Dan Sifat Asimtotik Estimator Autoregressive Linier Parsial*, Paper of  National Seminar on District Conference IX of Indonesian Mathematician Association at Sebelas Maret University, Surakarta, 1-13.

[14]   Lestari, B. (2003b) *Estimasi Dan Distribusi Asimtotik Estimator AR Linier Parsial*, Paper of Mathematics and Statistics National Seminar VI at  ITS, Surabaya, 1-10.

[15]   Lestari, B. (2004a) *Asymptotic Normality of Partly Linear Autoregressive Estimators*, Paper of  National Seminar on National Basic Sciences Meeting I at Brawijaya University, Malang, 1-13.

[16]   Lestari, B. (2004b) *Semiparametric Time Series Modeling of CPI Data*, Statistics Paper Proposed to Jurnal Teknik Industri dan Informasi (JTII), University of Surabaya (UBAYA), 1-10.

[17]   Lestari, B. (2005) Estimation and Asymptotic Normality of Partly Linear Autoregressive Estimators for GCPI Data, *Jurnal Ilmu Dasar FMIPA* Vol. 6, No.1, University of Jember.

[18]   Rice, J. (1986) Convergence Rates for Partially Splined Models, *Statist. And Probab. Lett.*, 4, 203-208.

[19]   Robinson, P. M. (1988) Rott-N-Consistent Semiparametric Regression, *Econometrica*, 56, 931-954.

[20]    Speckman, P. (1988) Kernel Smoothing in Partial Linear Models, *J. Roy. Statist. Soc.*, Ser. B50, 413-436.

[21]    Wahba, G. (1990) *Spline Models for Observational Data*, SIAM Publication.