

CHAPTER 2

THEORY

2.1 Outliers

According to Barnett, V. et al.(1994) (access by <http://en.wikipedia.org>), an outlier in the sense of statistics is an observation that is numerically distant from the rest of the data. They can occur by chance in any distribution, but they are often indicative either of measurement error or that the population has a heavy-tailed distribution. In the former case one wishes to discard them or use statistics that are robust to outliers, while in the latter case they indicate that the distribution has high kurtosis and that one should be very cautious in using tool or intuitions that assume a normal distribution. Outliers, being the most extreme observations, will include the sample maximum or sample minimum, or both, depending on whether they are extremely high or low. However, the sample maximum and minimum need not be outliers if they are not unusually far from other observations. This definition is similar to others such as Montgomery, et al. (1982), Moore, et al.(1999), Monhor, D. et al.(2005) and Walfish, S. (2005), etc.

Outliers in this research are considered in the sense of regression outliers. They are the observed data that are distinct from the linear relationship representing most of the data and they can draw a regression line away from the usual data. Nevertheless, they exclude unusual incidents of outliers. Moreover, types of regression outliers (Rousseeuw and Zomeren, 1990: 633) studied in this research are as below:

1) *y-direction* outliers

These are points that are outliers only because they are extreme *y-coordinates*. The extent to which such outliers will affect the parameter is estimated depending on both their *x-coordinate* and the general configuration of other points. Thus, those points could also be a regression outliers or residual outliers.

2) *x-direction outliers*

These are points that deviate only with regard to the *x-coordinates*. Such points can cause some regression estimates to perform poorly. The *x-direction* outliers could also be regression or residual outliers.

3) *xy-direction outliers*

These are points outlying in both *x-* and *y- coordinates*. It may be a regression outliers or residual outliers (Ryan, 1997: 350).

Classical estimate such as the sample mean, variance, covariance and correlation, or the LS fit of a regression model, can be very adversely influenced by outliers, even by a single one, and often fail to provide good fits to the bulk of the data. An alternative approach such as robust approach has been introduced to cope with outliers' problem in order to provide a good fit for the bulk of the data containing outliers, as well as when the data are free of them. Nevertheless, robust approach is quite complicate.

2.2 LS Regression

The second topic to be studied is the regression analysis which is well concluded by Ampanthong (2009: 10-12). Accordingly, regression analysis is a statistical tool for modeling and analyzing several variables which underlie vital assumptions. To estimate the unknown parameters in a regression model is among the most significant objectives of regression analysis. This process is also called “fitting the data to the model”. According to Gauss-Markov theorem, the maximum likelihood estimate of $\underline{\theta} = (\alpha, \beta_1, \beta_2, \dots, \beta_k)$ turns out to be the BLUE of $\underline{\theta}$. One usually tries to estimate the unknown parameters in a regression model from a data set by the LS method to obtain $\hat{\underline{\theta}} = (X'X)^{-1} (X'Y)$. When the method is applied to acquire the estimates $\hat{\theta}_j$ of θ_j , for $j=1, 2, \dots, k$, those so found are called LS estimates of the regression coefficients. The fitted regression equation concluded from the data set is $\hat{y}_i = \hat{\alpha} + x_{i1}\hat{\beta}_1 + x_{i2}\hat{\beta}_2 + \dots + x_{ik}\hat{\beta}_k$ for $i=1, 2, \dots, n$.

This equation is regarded as the estimate of the regression model $\tilde{Y} = \tilde{x}_i\tilde{\theta}$, where $\tilde{x}_i = (x_{i1}, x_{i2}, \dots, x_{ik})$, $i=1, 2, \dots, n$. The residual r_i is defined as the difference between the observed value y_i and the fitted value \hat{y}_i , i.e. $r_i = y_i - \hat{y}_i$. The method of obtaining the LS of $\hat{\theta}_j$, for each $j=1, 2, \dots, p$ is the most popular estimate method. The LS estimator that minimizes the sum of squared residuals, i.e. $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$ will have
$$\text{Min}_{\hat{\theta}} \sum_{i=1}^n r_i^2 = \text{Min}_{\hat{\theta}} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$
. Nevertheless, LS estimator is not suitable when the distribution of their residuals is not normal. In some case, when data have a heavy tail in any direction due to the presence of outliers, LS method of estimate might not be preferable. While in other cases, data may have only a small fraction of outliers, but LS estimate is still not suitable for further analysis. A small fraction of outliers may have a large effect on the LS estimator.

2.3 Two-Limit Tobit Model

The two-limit Tobit model (Tobin, 1958) is the simplest model for censored data. Here, we stimulate the discussion using an example based on Tobin's application of the model. Let a dependent variable be the monthly expenditure on luxury goods of each household and let an independent variable or explanatory variable be such as the monthly income for the corresponding the household's monthly expenditure. The parameters vector θ , which contains the set of population regression parameters related to the variables, need to be estimated. In this example, link variable or Y_i^* might be the capacity of households to spend their income on luxury goods, but this is only realized as actual expenditure, Y_i , if that expenditure exceeds zero. Thus, even if many observations might have value to be zero on the Y_i , they can be considered as having changing values on the link variable Y_i^* . The two-limit Tobit model (Tobin, 1958: 26, Rosett, 1975: 141 and Jöreskog, 2002: 13) can be written as the dependent variable Y_i satisfies

$$Y_i = \begin{cases} L & ; Y_i^* \leq L \\ Y_i^* & ; L < Y_i^* < U \\ U & ; Y_i^* \geq U, \end{cases} \quad (2.1)$$

where Y_i^* , for $i=1,2,\dots,n$, is the link function generated by the linear regression model

$$Y_i^* = \alpha + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i,$$

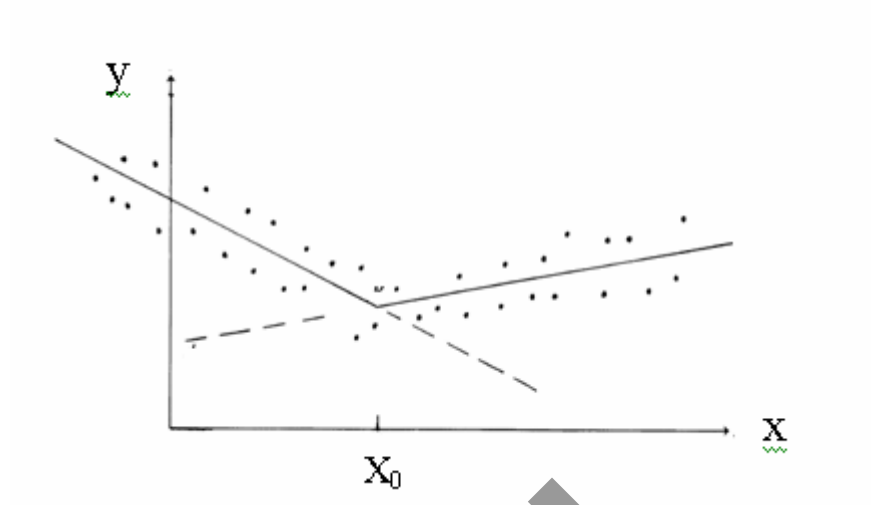
where x_1, x_2, \dots, x_k are regressors, and ε_i 's are the error terms having independent normal distributions with zero mean and constant variance ($\varepsilon_i \sim \text{i.i.d.} N(0, \sigma^2)$) and are independent of x_i . L and U in the model (2.1) are an lower and upper limits, respectively.

The probability density function (p.d.f.) of Y for given values of each L and U is determined by $f_Y(Y_i) = \Phi\left(\frac{L - x_i\theta}{\sigma_i}\right)$ if $Y_i = L_i$, $f_Y(Y_i) = \frac{1}{\sigma_i} \phi\left(\frac{Y_i - x_i\theta}{\sigma_i}\right)$ if $Y_i = Y_i^*$, and $f_Y(Y_i) = 1 - \Phi\left(\frac{U_i - x_i\theta}{\sigma_i}\right)$ if $Y_i = U_i$. Where Φ and ϕ are the cumulative distribution function (c.d.f.) and the p.d.f. of a standard normal distribution, respectively. From the p.d.f. of Y , we then get the log-likelihood function and by the ML fashion, the Tobit estimator was constructed. Actually, LS fashion might be inappropriate in the case that the dependent variable is limited by some desired variable. As the mention of Greene (1981: 505-513) who described that the LS estimator of parameter vector in Tobit model is the bias and also the asymptotic bias of the regression coefficients, this means that the LS based for the limited dependent variable case is inconsistent.

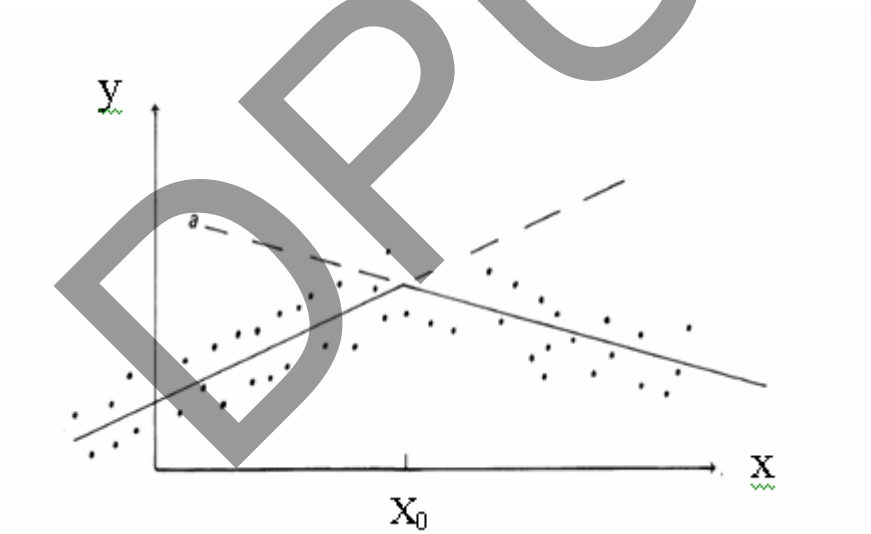
2.4 Piecewise Regression Model

In this research, the structural change in the regression model is taken into account, thus the piecewise regression model (Quandt, 1958: 874, Hudson, 1966: 1097-1129, Goldfeld, Kelejian and Quandt, 1971, Suits, Mason and Chan 1978: 132-133) is considered. Quandt first introduced that economic variables may sometimes be fitted by linear relations with the property that the parameters of the relationship are subject to discontinuous changes. For example, consider the consumption function $Y = \alpha X + \beta$. Aggregate consumption depends upon the level of aggregated income. In addition, it may be hypothesized that consumption depends non-linearly on other factors such as the state of expectations concerning the future of the economy, the volume of installment buying, the level of the interest rate, etc. These other variables may have the effect of altering the parameters of the consumption function in the following fashion: when the critical outside variable i satisfies $i < i^*$ then $Y = \alpha_1 X + \beta_1$ and when $i \geq i^*$ then $Y = \alpha_2 X + \beta_2$, where i^* is the critical level of the outside variable in question. In general, one may not be able to identify the critical outside variable and one may not be able to state at what time the system $Y = \alpha X + \beta$ changes from one regime to the other. In the paper of Quandt, there was indicated an estimating procedure for the switching point under the conditions when it is known that the time period under consideration contains a single switching. Parameters in the piecewise regression model were estimated by the ML method.

Subsequently, Hudson (1966: 1097-1129) discussed a similar estimate problem in which the two regression regimes are required to be intersected, see example on figure 2.1. Parameters estimate based on the LS method and the models were assumed to be joined at the value x_0 .



(a) Model (2.14a).



(b) Model (2.14b).

Figure 2.1 Two possible types of the Piecewise Regression Model

Source: Tishler and Zang, 1981: 117.

Two regression regimes joined at point v (Hudson, 1966) can be represented by

$$Y_i = \begin{cases} \alpha_1 + \beta_1 x_i + \varepsilon_i & ; x_i \leq v \\ \alpha_2 + \beta_2 x_i + \varepsilon_i & ; x_i > v \end{cases} \quad (2.2)$$

where Y_i is a dependent variable, x_i is a corresponding independent variable, and error terms, ε_i 's are assumed to be normally and independently distributed with mean zero and variance σ^2 and are independent of the independent variable. In addition, the model (2.2) is subject to $\alpha_1 + \beta_1 v = \alpha_2 + \beta_2 v$ and can be written as

$$Y_i = \begin{cases} \alpha_1 + \beta_1 x_i + \varepsilon_i & ; x_i \leq v \\ \alpha_1 + \beta_1 x_i + (\beta_2 - \beta_1)(x_i - v) + \varepsilon_i & ; x_i > v. \end{cases} \quad (2.3)$$

Suits et al. (1978: 132-133) extended model (2.3) so that it can be written in the multiple regression model with a dummy variable, D_i , consisting of two independent variables as follows;

$$Y_i = \alpha_1 + \beta_1 x_i + \beta_2 x_i^* + \varepsilon_i, \quad (2.4)$$

where $x_i^* = (x_i - v)D_i$ and $D_i = \begin{cases} 0 & ; x_i \leq v \\ 1 & ; x_i > v \end{cases}$.

2.5 Methods for Finding the Minimum of the Sum of Squares

This research is related to solve nonlinear least square problem. Therefore, there are three methods introduced to find the minimum of sum of squares for that problem. Nonlinear least square problems occur for instance in nonlinear regression namely piecewise and TP regression as considered in this study.

2.5.1 Gauss-Newton Method

The Gauss-Newton approximation (Seber and Wild, 1988: 25) is described as the followings. Suppose $\tilde{\theta}^{(a)}$ is an approximation to the LS estimator $\hat{\theta}$ of a nonlinear model. By the Taylor's expansion, we will get

$$f(X; \tilde{\theta}) \approx f(X; \tilde{\theta}^{(a)}) + F^{(a)} (\tilde{\theta} - \tilde{\theta}^{(a)}), \quad (2.5)$$

where $F^{(a)}$ is $F(X; \tilde{\theta}^{(a)})$. It can be applied to the residual vector, $r(X; \tilde{\theta})$, as

$$\begin{aligned} r(X; \tilde{\theta}) &= Y - f(X; \tilde{\theta}) \\ &= r(X; \tilde{\theta}^{(a)}) - F^{(a)} (\tilde{\theta} - \tilde{\theta}^{(a)}). \end{aligned}$$

From the equation $S(\tilde{\theta}) = \sum_{i=1}^n (Y_i - f(x_i; \tilde{\theta}))^2$, we will get

$$\begin{aligned} S(\tilde{\theta}) &\approx r'(X; \tilde{\theta}^{(a)}) r(X; \tilde{\theta}^{(a)}) - 2r'(X; \tilde{\theta}^{(a)}) F^{(a)} (\tilde{\theta} - \tilde{\theta}^{(a)}) \\ &\quad + (\tilde{\theta} - \tilde{\theta}^{(a)})' F^{(a)'} F^{(a)} (\tilde{\theta} - \tilde{\theta}^{(a)}). \end{aligned} \quad (2.6)$$

Thus, we can conclude that the right hand side of the approximation (2.6) is minimized with respect to $\tilde{\theta}$ when

$$\begin{aligned} \tilde{\theta} - \tilde{\theta}^{(a)} &= \left(F^{(a)'} F^{(a)} \right)^{-1} F^{(a)'} r(X; \tilde{\theta}^{(a)}) \\ &= \tilde{\delta}^{(a)}. \end{aligned} \quad (2.7)$$

This equation gives the approximation of $\underline{\theta}^{(a)}$ then we can say that the next approximation is followed by

$$\underline{\theta}^{(a+1)} = \underline{\theta}^{(a)} + \underline{\delta}^{(a)}. \quad (2.8)$$

The equations (2.7) and (2.8) determine the updating result and the $\hat{\theta}$ can be attained by the equation (2.8). In addition, Seber and Wild (1988) mentioned that the Gauss-Newton algorithm is convergent.

2.5.2 Steepest Descent Method

The steepest descent method is also known as the gradient descent. It is based on the gradient of $\xi' \xi$. Seber and Wild (1988: 594) described the theory of this method as the followings. The steepest descent method is one of iterative processes where an initial guess $\underline{\theta}^{(1)}$ is furnished, from which the algorithm sequentially moves in \Re^p of points $\underline{\theta}^{(2)}, \underline{\theta}^{(3)}, \dots$ which are aimed to converge to a local minimum $\hat{\theta}$. Practically useful is this algorithm method to make sure that $h(\underline{\theta})$, a real-valued function of p parameters vector $\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_p)'$, is reduced at each iteration so that $h(\underline{\theta}^{(a+1)}) < h(\underline{\theta}^{(a)})$. In this research the function $h(\underline{\theta})$ is defined as

$$S(\underline{\theta}) = \sum_{i=1}^n \left(Y_i - f(\underline{x}_i; \underline{\theta}) \right)^2.$$

The approximation of the $(a+1)^{\text{th}}$ iterative is the same as (2.8) nevertheless the updating $\underline{\delta}^{(a)}$, the computation of the a^{th} step, is differently defined as

$$\underline{\delta}^{(a)} = \rho^{(a)} \underline{d}^{(a)}, \quad (2.9)$$

where the vector $\underline{d}^{(a)}$ is called the direction of the step and $\rho^{(a)}$ is the step length. Frequently $\rho^{(a)}$ is selected to approximately minimize $S(\underline{\theta})$ along the line $\underline{\theta}^{(a)} + \rho \underline{d}^{(a)}$, this process known as a line search. If the $\rho^{(a)}$ is the exact minimum at each iteration, then the algorithm is said to have exact line searches; otherwise, it employs approximate line searches. Whenever there exists the convergent of $S(\underline{\theta})$ to the local minimum, this means that the exact line searches are not manipulated and $S(\underline{\theta})$ is sufficiently reduced at each iteration (Gill, 1991: 100). A descent direction $\underline{d}^{(a)}$ can be determined by

$$\underline{g}^{(a)'} \underline{d} = \left. \frac{\partial S(\underline{\theta}^{(a)} + \rho \underline{d})}{\partial \rho} \right|_{\rho=0} < 0 \quad (2.10)$$

By the Taylor expansion, we get

$$S(\underline{\theta}^{(a)} + \rho \underline{d}) \approx S(\underline{\theta}^{(a)}) + \rho \underline{g}^{(a)'} \underline{d} + O(\rho^2). \quad (2.11)$$

Thus, when the approximation (2.11) attains, the decrease in function $S(\underline{\theta}^{(a)})$ can be obtained by the small enough step ρ in the direction \underline{d} . The important theorem is stated in order to know “how descent directions can be calculated?” as followed.

Theorem 1 A direction \underline{d} is a descent direction at parameters vector $\underline{\theta}$ if and only if there exists a positive definite matrix R such that

$$\underline{d} = -R \underline{g}.$$

Proof . It is available on Nonlinear Regression (Seber and Wild, 1988: 595).

Therefore from equations (2.8) and (2.9), we can state that

$$\underline{\theta}^{(a+1)} = \underline{\theta}^{(a)} - \rho^{(a)} R^{(a)} \underline{g}^{(a)}. \quad (2.12)$$

Where the choice of R is $R = I_p$ the descent direction becomes to $\underline{d} = -\underline{g}^{(a)}$ and this is called the steepest descent direction. However, the direction of steepest descent depends entirely upon the scaling of $\underline{\theta}$.

2.5.3 Levenberg-Marquardt Method (Seber and Wild, 1988)

This method is a compromise between the Gauss-Newton and steepest descent methods. As $\underline{d} \rightarrow 0$, the direction approaches Gauss-Newton. As $\underline{d} \rightarrow \infty$, the direction approaches steepest descent. Levenberg-Marquardt method is equivalent to performing a series of ridge regressions and is useful when the parameter estimates are highly correlated or the objective function is not well approximated by a quadratic.

Let a model be fitted into the data, there is likely the function $f(\underline{x}_i; \underline{\theta}^*)$ as expressed in the model $Y_i = f(\underline{x}_i; \underline{\theta}^*) + \varepsilon_i$. The problem is to compute the estimates of parameters which will minimize $S(\underline{\theta}) = \|\underline{Y} - f(\underline{x}_i; \underline{\theta})\|^2$. By utilizing ideas of Levenberg together with Marquardt (1963), thus the Levenberg-Marquardt algorithm adaptively varies the parameter updates between the gradient descent and Gauss-Newton update (Seber and Wild, 1988)

$$\underline{\delta}^{(a)} = - \left(\underline{F}^{(a)'} \underline{F}^{(a)} + \eta^{(a)} \underline{D}^{(a)} \right)^{-1} \underline{F}^{(a)'} \underline{r}(\underline{X}; \underline{\theta}^{(a)}), \quad (2.13)$$

where $\underline{D}^{(a)}$ is a diagonal matrix with positive diagonal element frequently defined to be the same as $\underline{F}^{(a)'} \underline{F}^{(a)}$ and $\eta^{(a)}$ is the a^{th} step direction. When $\underline{D}^{(a)}$ is I_p and $\eta^{(a)} \rightarrow 0$, the direction approaches Gauss-Newton. Whereas, $\eta^{(a)} \rightarrow \infty$, the direction approaches steepest descent. In the case that $\eta^{(a)} > 0$ then $\underline{F}^{(a)'} \underline{F}^{(a)} + \eta^{(a)} \underline{D}^{(a)}$ is positive definite, as $\underline{D}^{(a)}$ is positive definite. Thus the updating function $\underline{\delta}^{(a)}$ as in (2.13) and by the Theorem 1, is a descent direction. And in the case that

$\eta^{(a)} \rightarrow \infty$ then $\delta^{(a)}$ tend to zero. Therefore if we select very large $\eta^{(a)}$ then we can so fast reduce the function $S(\theta) = \|\tilde{Y} - f(\tilde{x}_i; \theta)\|^2$. Nevertheless, for many iterations $\eta^{(a)}$ which are too large, the algorithm adapt with little progress.

Marquardt's finding indicates that the average angle between Gauss-Newton and steepest descent directions is about 90 degree. A choice of initial value of the direction $\eta^{(0)}$ for this research is $\eta^{(0)} = 10^{-3}$ used to start and compute the updating vector $\delta^{(a)}$. If $S(\theta^{(a)} + \delta^{(a)}) < S(\theta^{(a)})$, then η becomes $\frac{\eta}{10}$ for the next iteration. Otherwise $S(\theta^{(a)} + \delta^{(a)}) > S(\theta^{(a)})$, then η is 10η for the next iteration.

2.6 TP Regression Model

There was first interested the combination of two principal ideas, i.e. Tobit and piecewise regression, where each has a different benefit as described before. And there were not any literatures which applied these two ideas to cope with the outliers problem until Mekbunditkul (2010) first introduced the TP (abbreviated from Tobit-piecewise) regression as the derivation of the TP regression model and the log-likelihood function of θ described below:

According to the two-limit Tobit model (2.1) and by assuming that there exists the structural change in regression parameter, the piecewise multiple linear regression (Quandt, 1958: 874, Hudson, 1966: 1097-1129, Goldfeld and Quandt, 1971, Suits et al. 1978: 132-133) is utilized. The link function as mentioned in the model (2.1), Y_i^* can be broken into two regression regimes as

$$Y_i^* = \begin{cases} \alpha_1 + \beta_{11}x_{i1} + \beta_{12}x_{i2} + \dots + \beta_{1k}x_{ik} + \varepsilon_i & ; \text{ if } v_i \leq v, \\ \alpha_2 + \beta_{21}x_{i1} + \beta_{22}x_{i2} + \dots + \beta_{2k}x_{ik} + \varepsilon_i & ; \text{ if } v_i > v, \end{cases} \quad (2.14)$$

where Y_i^* is a dependent variable, x_{ij} represents the i^{th} observation of the j^{th} independent variable, for $j=1, \dots, k$ and $i=1, \dots, n$. In addition, $v_i = x_i \theta$ (Quandt,

1972: 307), where \underline{x}_i is the row vector with k variables of the i^{th} observation and $\underline{\theta}$ is a k -dim vector of unknown parameters. The errors ε_i 's are $N(0, \sigma_i^2)$. Suppose $\underline{x}_0 = (x_{01}, \dots, x_{0k})$ is a vector of regressors at a joined point, i.e., $\underline{x}_0 \underline{\theta} = v$, then, from model (2.5), $\alpha_1 + \sum_{j=1}^k \beta_{1j} x_{0j} = \alpha_2 + \sum_{j=1}^k \beta_{2j} x_{0j}$, i.e., $\alpha_2 = \alpha_1 - \sum_{j=1}^k (\beta_{2j} - \beta_{1j}) x_{0j}$.

For $v_i > v$,

$$Y_i^* = \alpha_1 + \sum_{j=1}^k \beta_{1j} x_{ij} + \sum_{j=1}^k (\beta_{2j} - \beta_{1j}) x_{ij} - \sum_{j=1}^k (\beta_{2j} - \beta_{1j}) x_{0j} + \varepsilon_i.$$

By using a dummy variable, $D_i = \begin{cases} 1 & ; v_i > v, \\ 0 & ; v_i \leq v, \end{cases}$, the model (2.5) can be written in a single equation as (Mekbunditkul, 2010)

$$Y_i^* = \alpha_1 + \sum_{j=1}^k \beta_{1j} x_{ij} + \sum_{j=1}^k \beta_{2j}^* x_{ij} D_i + \beta_3^* D_i + \varepsilon_i, \quad (2.15)$$

where $\beta_2^* = \beta_{2j} - \beta_{1j}$ and $\beta_3^* = -\sum_{j=1}^k (\beta_{2j} - \beta_{1j}) x_{0j}$.

Thus the TP regression model can be written as

$$Y_i = \begin{cases} L_i & ; Y_i^* \leq L_i \\ Y_i^* & ; L_i < Y_i^* < U_i \\ U_i & ; Y_i^* \geq U_i, \end{cases}$$

where $Y_i^* = \alpha_1 + \sum_{j=1}^k \beta_{1j} x_{ij} + \sum_{j=1}^k \beta_{2j}^* x_{ij} D_i + \beta_3^* D_i + \varepsilon_i$. In addition, this model can be

written in the matrix form as

$$\underline{Y} = \underline{X} \underline{\theta} + \underline{\varepsilon}, \quad (2.16)$$

where

$$\begin{aligned}\tilde{Y} &= [\tilde{Y}'_1 \mid \tilde{Y}'_2 \mid \tilde{Y}'_3]'_{n \times 1}, \\ &= \begin{bmatrix} L_{11} & L_{21} & \cdots & L_{n_1 1} & Y_{12} & Y_{22} & \cdots & Y_{n_2 2} & U_{13} & U_{23} & \cdots & U_{n_3 3} \end{bmatrix}',\end{aligned}$$

$\tilde{Y}_2 = \tilde{Y}_2^*$ and $\theta = (\alpha_1, \beta_{11}, \dots, \beta_{1k}, \beta_{21}^*, \dots, \beta_{2k}^*, \beta_3^*)'$ and X is defined as in the equation (2.9). Moreover, the limits L and U are defined by

$$L_{im} = \begin{cases} L_a & ; v_{im} \leq v, \\ L_b & ; v_{im} > v, \end{cases} \text{ and } U_{im} = \begin{cases} U_a & ; v_{im} \leq v, \\ U_b & ; v_{im} > v. \end{cases} \quad (2.17)$$

The vector \tilde{Y}^* can be described as below: Without loss of generality (WLOG), all of the observed data Y_i^* 's as well as x_{i1}, \dots, x_{ik} to which Y_i^* corresponds, for $i=1, \dots, n$ are rearranged. Hence, observation vector \tilde{Y}^* consists of three parts. One of them is observation with values smaller than the lower limit L . The second comprises of values that lie between the limit (L, U) and the third indicates values greater than the upper limit U . To be specific, suppose that the observation in each part are n_1, n_2 and n_3 , respectively. Thus,

$$\begin{aligned}\tilde{Y}^* &= [\tilde{Y}_1^{*'} \mid \tilde{Y}_2^{*'} \mid \tilde{Y}_3^{*'}]'_{n \times 1}, \\ &= \begin{bmatrix} Y_{11}^* & Y_{21}^* & \cdots & Y_{n_1 1}^* & Y_{12}^* & Y_{22}^* & \cdots & Y_{n_2 2}^* & Y_{13}^* & Y_{23}^* & \cdots & Y_{n_3 3}^* \end{bmatrix}'.\end{aligned}$$

The variance-covariance matrix of \tilde{Y}^* is assumed to be

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 & 0 \\ 0 & \Sigma_2 & 0 \\ 0 & 0 & \Sigma_3 \end{bmatrix}, \quad \text{where} \quad \Sigma_m = \text{diag}(\sigma_{1m} \dots \sigma_{n_m m}) \text{ and}$$

$$\sigma_{im}^2 = \begin{cases} \sigma_a^2 & \text{if } v_{im} \leq v, \\ \sigma_b^2 & \text{if } v_{im} > v \end{cases}, m = 1, 2, 3; i = 1, \dots, n_m.$$

Therefore, the inverse matrix of Σ is $\Sigma^{-1} = \begin{bmatrix} \Sigma_1^{-1} & 0 & 0 \\ 0 & \Sigma_2^{-1} & 0 \\ 0 & 0 & \Sigma_3^{-1} \end{bmatrix}$.

Since Σ_m are diagonal matrices, so $\Sigma_m^{-1} = \text{diag}\left(\frac{1}{\sigma_{1m}}, \dots, \frac{1}{\sigma_{n_m m}}\right)$, where $m = 1, 2, 3$.

In addition, it is assumed that $\varepsilon_{im} \sim N(0, \sigma_{im}^2)$.

The matrix of X independent variables corresponding to

$\mathbf{Y} = [\mathbf{L}'_1 \mid \mathbf{Y}'_2 \mid \mathbf{U}'_3]_{n \times 1}$, is

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \mathbf{X}_3 \end{bmatrix} = \begin{bmatrix} 1 & x_{111} & \dots & x_{1k1} & x_{111}^* & \dots & x_{1k1}^* & x'_{11} \\ 1 & x_{211} & \dots & x_{2k1} & x_{211}^* & \dots & x_{2k1}^* & x'_{21} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{n_1 11} & \dots & x_{n_1 k1} & x_{n_1 11}^* & \dots & x_{n_1 k1}^* & x'_{n_1 1} \\ \hline 1 & x_{112} & \dots & x_{1k2} & x_{112}^* & \dots & x_{1k2}^* & x'_{11} \\ 1 & x_{212} & \dots & x_{2k2} & x_{212}^* & \dots & x_{2k2}^* & x'_{22} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{n_2 12} & \dots & x_{n_2 k2} & x_{n_2 12}^* & \dots & x_{n_2 k2}^* & x'_{n_2 2} \\ \hline 1 & x_{113} & \dots & x_{1k3} & x_{113}^* & \dots & x_{1k3}^* & x'_{13} \\ 1 & x_{213} & \dots & x_{2k3} & x_{213}^* & \dots & x_{2k3}^* & x'_{23} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{n_3 13} & \dots & x_{n_3 k3} & x_{n_3 13}^* & \dots & x_{n_3 k3}^* & x'_{n_3 3} \end{bmatrix}_{n \times (2k+2)} \quad (2.9)$$

where $x_{ijm}^* = x_{ijm}D_{im}$, $x'_{im} = D_{im}$ and $D_{im} = \begin{cases} 1 & ; v_{im} > v, \\ 0 & ; v_{im} \leq v, \end{cases}$ where $m=1, 2, 3$;

$j=1, \dots, k$; $i=1, \dots, n_m$, and k is the number of regressor variables. Note that $2k+2$ is less than n .

2.7 TP Estimator

The log-likelihood function of θ in TP regression model was derived by Mekbunditkul (2011) via the p.d.f. of Y_{im} , the im^{th} element of vector \underline{Y} . The p.d.f. of ε_{im} is assumed to be normal with zero mean and σ_{im}^2 variance and it is independent from each other. When given the values of L_{im} and U_{im} , for $i=1, \dots, n_m$ and $m=1, 2$ and 3 , the p.d.f. of \underline{Y} was derived into three parts as the followings.

Part 1. For $Y_{i1} = L_{i1}$, where $i=1, \dots, n_1$:

$$\begin{aligned} f_{Y_{i1}}(L_{i1}) &= P(Y_{i1} = L_{i1}), \\ &= P(Y_{i1}^* \leq L_{i1}), \\ &= P(x_{i1}\theta + \varepsilon_{i1} \leq L_{i1}), \\ &= P\left(\frac{\varepsilon_{i1}}{\sigma_{i1}} \leq \frac{L_{i1} - x_{i1}\theta}{\sigma_{i1}}\right), \\ &= \Phi\left(\frac{L_{i1} - x_{i1}\theta}{\sigma_{i1}}\right). \end{aligned}$$

Part 2. For $L_{i2} < Y_{i2} < U_{i2}$, where $i=1, \dots, n_2$:

$$P(L_{i2} < Y_{i2} < U_{i2}) = P(L_{i2} < x_{i2}\theta + \varepsilon_{i2} \leq y_{i2}^*),$$

$$\begin{aligned}
&= P\left(\frac{L_{i2} - \underline{x}_{i2}\underline{\theta}}{\sigma_{i2}} < \frac{\varepsilon_{i2}}{\sigma_{i2}} \leq \frac{y_{i2}^* - \underline{x}_{i2}\underline{\theta}}{\sigma_{i2}}\right), \\
&= \Phi\left(\frac{y_{i2}^* - \underline{x}_{i2}\underline{\theta}}{\sigma_{i2}}\right) - \Phi\left(\frac{L_{i2} - \underline{x}_{i2}\underline{\theta}}{\sigma_{i2}}\right).
\end{aligned}$$

Hence, $f_{Y_{i2}}(y_{i2}) = \frac{1}{\sigma_{i2}} \phi\left(\frac{y_{i2}^* - \underline{x}_{i2}\underline{\theta}}{\sigma_{i2}}\right) = \frac{1}{\sigma_{i2}} \phi\left(\frac{y_{i2} - \underline{x}_{i2}\underline{\theta}}{\sigma_{i2}}\right)$, where $y_{i2} \in (L_{i2}, U_{i2})$.

Part 3. For $Y_{i3} = U_{i3}$, where $i = 1, \dots, n_3$:

$$\begin{aligned}
P(Y_{i3} = U_{i3}) &= P(\underline{x}_{i3}\underline{\theta} + \varepsilon_{i3} \geq U_{i3}), \\
&= 1 - P\left(\frac{\varepsilon_{i3}}{\sigma_{i3}} < \frac{U_{i3} - \underline{x}_{i3}\underline{\theta}}{\sigma_{i3}}\right), \\
&= 1 - \Phi\left(\frac{U_{i3} - \underline{x}_{i3}\underline{\theta}}{\sigma_{i3}}\right).
\end{aligned}$$

Functions Φ and ϕ are the c.d.f. and p.d.f. of a standard normal distribution, respectively.

Some notations were indicated to be used in the next part as

$$I_L = \{i \mid Y_{i1} = L_{i1}, i=1, \dots, n_1\},$$

$$I_Y = \{i \mid L_{i2} < Y_{i2} < U_{i2}, i = 1, \dots, n_2\}, \text{ and}$$

$$I_U = \{i \mid Y_{i3} = U_{i3}, i = 1, \dots, n_3\}.$$

From the independent property of each element in the vector \underline{Y} , the p.d.f. of \underline{Y} can be expressed as

$$f_{\mathbf{Y}}(\mathbf{y}) = \prod_{i \in I_L} \Phi\left(\frac{L_{i1} - \mathbf{x}_{i1}\boldsymbol{\theta}}{\sigma_{i1}}\right) \cdot \left[\frac{\exp\left\{-\frac{1}{2}(\mathbf{Y}_2 - \mathbf{X}_2\boldsymbol{\theta})' \Sigma_2^{-1}(\mathbf{Y}_2 - \mathbf{X}_2\boldsymbol{\theta})\right\}}{(2\pi)^{n_2/2} |\Sigma_2|^{n_2/2}} \right] \\ \cdot \prod_{i \in I_U} 1 - \Phi\left(\frac{U_{i3} - \mathbf{x}_{i3}\boldsymbol{\theta}}{\sigma_{i3}}\right). \quad (2.19)$$

Thus, the log-likelihood function is thus given by

$$\ln L(\boldsymbol{\theta}; \mathbf{Y}) = \sum_{i \in I_L} \ln \Phi(\lambda_{i1}^-) - \frac{n_2}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_2| - \frac{1}{2} (\mathbf{Y}_2 - \mathbf{X}_2\boldsymbol{\theta})' \Sigma_2^{-1} (\mathbf{Y}_2 - \mathbf{X}_2\boldsymbol{\theta}) \\ + \sum_{i \in I_U} \ln [1 - \Phi(\lambda_{i3})]. \quad (2.20)$$

where $\lambda_{i1}^- = \frac{L_{i1} - \mathbf{x}_{i1}\boldsymbol{\theta}}{\sigma_{i1}}$ and $\lambda_{i3} = \frac{U_{i3} - \mathbf{x}_{i3}\boldsymbol{\theta}}{\sigma_{i3}}$.

The ML estimators of $\boldsymbol{\theta}$ can be obtained straightforwardly from the log-likelihood equation (2.20), which consists of three parts, as

$$\frac{\partial \ln L(\boldsymbol{\theta}; \mathbf{y})}{\partial \boldsymbol{\theta}} = - \sum_{i \in I_L} \left(\frac{\phi(\hat{\lambda}_{i1}^-)}{\Phi(\hat{\lambda}_{i1}^-)} \right) \frac{\mathbf{x}_{i1}'}{\sigma_{i1}} - (\mathbf{X}_2' \Sigma_2^{-1} \mathbf{X}_2) \hat{\boldsymbol{\theta}}_{TP} + (\mathbf{X}_2' \Sigma_2^{-1} \mathbf{Y}_2) \\ + \sum_{i \in I_U} \left(\frac{\phi(\hat{\lambda}_{i3})}{1 - \Phi(\hat{\lambda}_{i3})} \right) \frac{\mathbf{x}_{i3}'}{\sigma_{i3}} = \mathbf{0}_{(2k+1) \times 1}.$$

From Mekbunditkul (2010), $\hat{\boldsymbol{\theta}}_{TP}$ was verified as in the form of

$$\begin{aligned}
\hat{\theta}_{TP} &= \left(X_2' \Sigma_2^{-1} X_2 \right)^{-1} \left[-X_1' \Sigma_1^{-1/2} \left\{ H_1(\hat{\lambda}^-) \right\} + \left(X_2' \Sigma_2^{-1} Y_2 \right) + X_3' \Sigma_3^{-1/2} \left\{ H_3(\hat{\lambda}) \right\} \right] \\
&= \left(X_2' \Sigma_2^{-1} X_2 \right)^{-1} \left(X_2' \Sigma_2^{-1} Y_2 \right) - \left(X_2' \Sigma_2^{-1} X_2 \right)^{-1} \left[X_1' \Sigma_1^{-1/2} \left\{ H_1(\hat{\lambda}^-) \right\} \right] \\
&\quad + \left(X_2' \Sigma_2^{-1} X_2 \right)^{-1} \left[X_3' \Sigma_3^{-1/2} \left\{ H_3(\hat{\lambda}) \right\} \right],
\end{aligned}$$

$$\text{where } H_1(\hat{\lambda}^-) = \left(h(\hat{\lambda}_{11}^-) \quad \dots \quad h(\hat{\lambda}_{n_1 1}^-) \right)' = \left(\frac{\phi(\hat{\lambda}_{11}^-)}{\Phi(\hat{\lambda}_{11}^-)} \quad \dots \quad \frac{\phi(\hat{\lambda}_{n_1 1}^-)}{\Phi(\hat{\lambda}_{n_1 1}^-)} \right)', \quad (2.21)$$

$$\text{and } H_3(\hat{\lambda}) = \left(h(\hat{\lambda}_{13}) \quad \dots \quad h(\hat{\lambda}_{n_3 3}) \right)' = \left(\frac{\phi(\hat{\lambda}_{13})}{1 - \Phi(\hat{\lambda}_{13})} \quad \dots \quad \frac{\phi(\hat{\lambda}_{n_3 3})}{1 - \Phi(\hat{\lambda}_{n_3 3})} \right)', \quad (2.22)$$

where $\hat{\lambda}_{n_1 1}^-$ and $\hat{\lambda}_{i3}$ are estimators of $\lambda_{n_1 1}^-$ and λ_{i3} , respectively.

There exist three parts of TP estimator which the first part is the LS estimator based on n_2 observations where values are not at the limits. The other two parts concern n_1 and n_3 observations whose values are truncated respectively by the lower and upper desired limits.

2.8 Properties of TP Estimator

In Mekbunditkul's dissertation, there were verified some properties of TP estimator in terms of its bias and mean square error (MSE). Two situations were defined: (1) $U \rightarrow \infty$ and L is finite, (2) $L \rightarrow -\infty$ and U is finite. Some vectors/matrices such as \underline{Y} , X , \underline{L} and \underline{U} and Σ were defined as in section (2.4) and the following statements were indicated to refer throughout this section

C1: Assuming that $U \rightarrow \infty$ and L is finite, the response variable \underline{Y} in TP regression

model is as $\underline{Y} = \begin{bmatrix} \underline{L}_1 \\ \underline{Y}_2 \end{bmatrix}_{n \times 1}$, where $n = n_1 + n_2$.

C2: Assuming that $L \rightarrow -\infty$ and U is finite, the response variable \mathbf{Y} in TP regression

model is as $\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_2 \\ \mathbf{U}_3 \end{bmatrix}_{n \times 1}$, where $n = n_2 + n_3$.

The TP estimator, $\hat{\theta}_{TP}$ corresponding respectively to each of C1 and C2 is

$$\hat{\theta}_{TP} = \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{X}_2 \right)^{-1} \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{Y}_2 \right) - \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{X}_2 \right)^{-1} \left[\mathbf{X}_1' \Sigma_1^{-1/2} \left\{ \mathbf{H}_1(\hat{\lambda}) \right\} \right], \quad (2.23)$$

$$\hat{\theta}_{TP} = \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{X}_2 \right)^{-1} \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{Y}_2 \right) + \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{X}_2 \right)^{-1} \left\{ \mathbf{X}_3' \Sigma_3^{-1/2} \left\{ \mathbf{H}_3(\hat{\lambda}) \right\} \right\}, \quad (2.24)$$

where $\mathbf{H}_1(\hat{\lambda})$ and $\mathbf{H}_3(\hat{\lambda})$ are in the forms of equations (2.23) and (2.24).

Theorem 2. Assume C1, $\hat{\theta}_{TP}$, as defined in equation (2.23), is biased where the bounds of the bias are (Mekbunditkul, 2010)

$$\begin{aligned} \text{Bias}_\ell &= -\mathbf{A}_1 (\mathbf{I} + \mathbf{A}_1)^{-1} \hat{\theta} - \left\{ \mathbf{I} - \mathbf{A}_1 (\mathbf{I} + \mathbf{A}_1)^{-1} \right\} \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{X}_2 \right)^{-1} \mathbf{X}_1' \Sigma_1^{-1/2} \left(\mathbf{I} - \Sigma_1^{-1/2} \mathbf{L}_1 \right), \\ \text{Bias}_u &= \left\{ \mathbf{I} - \mathbf{A}_1 (\mathbf{I} + \mathbf{A}_1)^{-1} \right\} \hat{\theta} - \left\{ \mathbf{I} - \mathbf{A}_1 (\mathbf{I} + \mathbf{A}_1)^{-1} \right\} \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{X}_2 \right)^{-1} \\ &\quad \cdot \left(\mathbf{X}_2' \Sigma_2^{-1} \mathbf{L}_2 - \mathbf{X}_1' \Sigma_1^{-1} \mathbf{L}_1 + \mathbf{X}_2' \Sigma_2^{-1/2} \mathbf{I} \right). \end{aligned}$$

Proof. It is available in Mekbunditkul's dissertation.

Theorem 3. Assume C1 holds. The asymptotic variance-covariance matrix of $\hat{\theta}_{TP}$, as defined in (2.23), is (Mekbunditkul, 2010)

$$\begin{aligned} \left(\mathbf{X}' \mathbf{W} \mathbf{X} \right)^{-1} &= \left\{ \left(\mathbf{X}_1' \mathbf{W}_1 \mathbf{X}_1 \right) + \left(\mathbf{X}_2' \mathbf{W}_2 \mathbf{X}_2 \right) \right\}^{-1}, \text{ where } \mathbf{W}_1 = \mathbf{G}_1(\lambda^-) \Sigma_1^{-1} \text{ and} \\ \mathbf{W}_2 &= \Sigma_2^{-1}. \end{aligned}$$

Proof. It is available in Mekbunditkul's dissertation.

Theorem 4. If C2 holds, the estimator $\hat{\theta}_{TP}$, as defined in equation (2.24), is biased where the bounds of the bias are (Mekbunditkul, 2010)

$$\begin{aligned} \text{Bias}_\ell &= \left\{ \mathbf{I} - \mathbf{A}_2 (\mathbf{I} + \mathbf{A}_2)^{-1} \right\} \underline{\boldsymbol{\theta}} - \left\{ \mathbf{I} - \mathbf{A}_2 (\mathbf{I} + \mathbf{A}_2)^{-1} \right\} \\ &\quad \cdot \left\{ \left(\mathbf{X}'_2 \Sigma_2^{-1} \mathbf{X}_2 \right)^{-1} \left(\mathbf{X}'_3 \Sigma_3^{-1} \mathbf{U}_3 - \mathbf{X}'_2 \Sigma_2^{-1} \mathbf{U}_2 - \mathbf{X}'_2 \Sigma_2^{-1} \mathbf{1} \right) \right\}, \\ \text{Bias}_u &= \left\{ \mathbf{I} - \mathbf{A}_2 (\mathbf{I} + \mathbf{A}_2)^{-1} \right\} \left[\left(\mathbf{X}'_2 \Sigma_2^{-1} \mathbf{X}_2 \right)^{-1} \left\{ \mathbf{X}'_3 \Sigma_3^{-1/2} (\mathbf{U}_3 + \mathbf{1}) \right\} \right] - \mathbf{A}_2 (\mathbf{I} + \mathbf{A}_2)^{-1} \underline{\boldsymbol{\theta}}. \end{aligned}$$

Proof. It is available in Mekbunditkul's dissertation.

Theorem 5. Assume C2 holds. The asymptotic variance-covariance matrix of $\hat{\boldsymbol{\theta}}_{\text{TP}}$, as defined in (2.24), is (Mekbunditkul, 2010)

$$\begin{aligned} \left(\mathbf{X}' \mathbf{W} \mathbf{X} \right)^{-1} &= \left\{ \left(\mathbf{X}'_3 \mathbf{W}_3 \mathbf{X}_3 \right) + \left(\mathbf{X}'_2 \mathbf{W}_2 \mathbf{X}_2 \right) \right\}^{-1}, \text{ where } \mathbf{W}_3 = \mathbf{G}_3(\lambda) \Sigma_3^{-1}, \text{ and} \\ \mathbf{W}_2 &= \Sigma_2^{-1} \text{ (Mekbunditkul, 2010).} \end{aligned}$$

Proof. It is available in Mekbunditkul's dissertation.

Theorem 6. Let C1 hold, then the ML estimators of each σ_a^2 and σ_b^2 in a TP regression model are obtained using (Mekbunditkul, 2010)

$$\begin{aligned} \max(c_j^2, f_j) &< \hat{\sigma}_j^2 < \left(d_j + \sqrt{\frac{1}{2} d_j^2 + 2f_j} \right)^2, \quad \text{where } j = a, b, \\ c_j &= 0 \text{ if } d_j - \sqrt{\frac{1}{2} d_j^2 + 2f_j} < 0 \text{ and } c_j = d_j - \sqrt{\frac{1}{2} d_j^2 + 2f_j} \text{ otherwise,} \\ d_j &= \frac{\mathbf{1}'(\hat{\mathbf{Y}}_{j1} - \mathbf{L}_{j1})}{n_{j2}}, \text{ and } f_j = \frac{(\hat{\mathbf{Y}}_{j1} - \mathbf{L}_{j1})'(\hat{\mathbf{Y}}_{j1} - \mathbf{L}_{j1})}{n_{j2}} + \frac{(\mathbf{Y}_{j2} - \hat{\mathbf{Y}}_{j2})'(\mathbf{Y}_{j2} - \hat{\mathbf{Y}}_{j2})}{n_{j2}}. \end{aligned}$$

Proof. It is available in Mekbunditkul's dissertation.

Theorem 7. Let C2 hold, then ML estimators of each σ_a^2 and σ_b^2 in the TP regression model are obtained by (Mekbunditkul, 2010)

$$\max(p_j^2, r_j) < \hat{\sigma}_j^2 < \left(q_j + \sqrt{\frac{1}{2}q_j^2 + 2r_j} \right)^2, \quad \text{where } j = a, b,$$

$$p_j = 0 \text{ if } q_j - \sqrt{\frac{1}{2}q_j^2 + 2r_j} < 0 \text{ and } p_j = q_j - \sqrt{\frac{1}{2}q_j^2 + 2r_j} \text{ otherwise,}$$

$$q_j = \frac{\mathbf{1}'(\mathbf{U}_{j3} - \hat{\mathbf{Y}}_{j3})}{n_{j2}} \text{ and } r_j = \frac{(\mathbf{U}_{j3} - \hat{\mathbf{Y}}_{j3})'(\mathbf{U}_{j3} - \hat{\mathbf{Y}}_{j3})}{n_{j2}} + \frac{(\mathbf{Y}_{j2} - \hat{\mathbf{Y}}_{j2})'(\mathbf{Y}_{j2} - \hat{\mathbf{Y}}_{j2})}{n_{j2}}.$$

Proof. It is available in Mekbunditkul's dissertation.

Whenever the same sample sizes n_1 and n_2 are assumed then an estimate of variance σ_{TP}^2 is as their pooled average of σ_a^2 and σ_b^2 (Snedecor and Cochran, 1989 and Welch, 1947). That is the ML estimator of σ_{TP}^2 is $\hat{\sigma}_{TP}^2 = \frac{\hat{\sigma}_a^2 + \hat{\sigma}_b^2}{2}$, where $\hat{\sigma}_a^2$ and $\hat{\sigma}_b^2$ as shown in Theorems 6 and 7.