Research Article

# A Comparison of Non-Parametric Weighted Linear Models

**Samuel Joel Kamun\***

*Department of Mathematics and Actuarial Sciences, Catholic University of Eastern Africa, Kenya Nairobi, Mombasa, Kenya*

**\*Corresponding Author:** Samuel Joel Kamun, Department of Mathematics and Actuarial Sciences, Catholic University of Eastern Africa, Kenya Nairobi, Mombasa, Kenya.

## Abstract

The analysis of sample-based studies involving sampling designs for small sample sizes is challenging because the sample selection probabilities (as well as the sample weights) are dependent on the response variable and covariates. This research focused on nonparametric weighted linear models in order to find more precise estimators with lower sample bias. The study has used rank-based approaches because they outperform least-squares procedures when the data deviates from normality and/or contains outliers. Weights can be added to these approaches to create weighted strategies (WT). In this paper, we demonstrate how to construct WT estimates using rank-based regression. Rank-based estimators were developed to provide a nonparametric, robust alternative to traditional likelihood or least squares estimators. They are then used to generate estimates with higher relative efficiencies and lower finite small sample bias than the Horvitz-Thompson weighted estimator with unmodified weight. The purpose of our study is to compare estimators using the reciprocal of the sample inclusion probabilities and other weights derived by modifying and rescaling them using relative efficiency, sample bias, and standard error for small sample sizes. The constructed estimates using different modified and rescaled weights are actually the weighted nonparametric estimators. The study compared three new estimators for both the unmodified and modified weights, which were found to have better relative efficiency and smaller finite small sample bias than the estimates from the conventional Horvitz-Thompson weighted estimator.

**Keywords:** Small Samples; Estimators; Relative Efficiency; Sample Bias; Standard Error

## Introduction

The study has focused on using nonparametric linear models with different modified and rescaled weights to estimate the relative efficiency, sample bias, and standard error. They were used in the study to generate estimates with higher relative efficiencies and lower finite small sample bias than the Horvitz-Thompson Weighted Estimator with unmodified weight.

## The model

Suppose the data is produced according to a function:

$$f ( y \mid x; \theta) \, g(x) \qquad \text{---------------- (1.1)}$$

Where y is a response variable which is multivariate and x is a continuous or discrete vector of covariate variables and

$$f ( y \mid x; \theta) \qquad \text{---------------- (1.2)}$$

Is the regression part of the function. The marginal distribution of x is denoted by g(x) which for this study we have used Gaussian density to represent, as shown below

$$k\left(u\right) = \frac{1}{\sqrt{(2\pi)}} e^{\left(\frac{-u^2}{2}\right)} \qquad \text{---------------- (1.3)}$$

Where

$$u = \frac{x_1 - mean(x_1)}{s\tan .dev(x_1)}$$

We describe the conditional distribution of y given $x_1$ as θ. The likelihood is given by

$$\prod f(y|x;\theta)$$

------------(1.4)

As explained, refer to [5,6].

## Rank-regression

The purpose of rank-based regression, like least squares, is to estimate the vector of coefficients, β, of a general linear model of the form:

$$y_i = \alpha + x_i^T \beta + e_i \text{ for I = 1........,n}$$

-------------- (2.1)

$y_i$ is the response variable, $x_i$ is the explanatory variable vector, α is the intercept parameter, and $e_i$ is the error term. We assume that the errors are associated with the probability density function (pdf) f(t). (2.1) is written in matrix notation as follows for convenience.

y = **α**1 + X**β** + e

--------------(2.2)

Therefore y = $[y_1,...,y_n]^T$ is the n x 1 vectors of outcomes, X = $[x_1,...,x_n]^T$ is the n × p design matrix, and e = $[e_1,...,e_n]^T$ is the n × 1 vector of error terms. The model is broad because the sole assumption on the distribution of errors is that it is continuous. Remember that the least squares estimator minimizes the Euclidean distance between y and $\hat{y}_{LS} = X\hat{\beta}_{LS}$. To obtain the R estimate, a new distance measure, refer to [4] based dispersion function is used. The dispersion function refer to [4], is given by

$$D(\beta) = \|y - X\beta\|_{\varphi}$$

---------------- (2.3)

Where $\|.\|_{\varphi}$ is a semi-norm defined as

$$\|u\|_{\varphi} = \sum_{i=1}^{n} a(R(u_i))u_i$$

Where R symbolizes rank, $a(t) = \varphi\left(\frac{t}{n+1}\right)$ and $\varphi$ is a square-integral, non-decreasing score function defined on the interval (0, 1). Assume it is standardized without losing generality, so that $\int \varphi(u)du = 0$ and $\int \varphi^2(u)du = 1$.

The R estimator of β is defined as follows:

$$\hat{\beta}_{\varphi} = Arg\min\|y - X\beta\|_{\varphi}$$

------------ (2.4)

This estimator is exceptionally efficient and resilient in the Y-space.

## Weighted non parametric weight

The process of obtaining the sample weight, which is the reciprocal of the sample inclusion probability, refer to [5,6]. Below is the weight:

$$w_i = \tilde{P}^{-1} = 1/g(x_i) = 1/\exp\left(\frac{(x_i - mean_{x_i})/sd_{x_i}}{\sqrt{\pi_i}}\right)$$

------------- (3.1)

## Weighted conditional non parametric estimator

For the process of obtaining the modified sample weight refer to [5,6]. Below is the modified weight:

$$w_a = \frac{\delta E_S(w_i|y_i,v_i,\theta,\gamma,\beta)}{\delta\beta}$$

------------- (3.2)

## Non-parametric rescaled weight (NPRW(I))

For the process of obtaining the modified sample weight, refer to [5,6]. Below is the modified weight:

$$\tilde{w}_i^m = \left(\frac{w_{reg1}}{E(w_{reg1}|x_i,y_i,z_i)}\right)\left(\frac{E(w_{reg2}|x_i)}{w_{reg2}}\right)\left(\frac{w_i}{E(w_i|x_i)}\right)$$

----------- (3.3)

## Non-parametric rescaled weight (NPRW(II) & NPRW(III))

For the process of obtaining the sample weight, which is the reciprocal of the sample inclusion probability, $w_i$ refer to [5,6]. The Non Parametric Rescaled Weight (NPRW(II)) is obtained following the process given here. The absolute difference between the regressed observed and the observed value of the dependent variable and the conditional expectation of the difference between regressed and observed value of the dependent variable y on the predictor variables $x_1$, $x_1$ and z, y and z, and y, $x_1$ and z is first obtained. Then the product of the reciprocal of the sample inclusion probability and ratios of the absolute difference and conditional expectations is obtained to give the weight. Seen below:

$$\tilde{w}_i^{4m} = w_i\left(\frac{y_{regEst} - y_i}{E(y_{regEst} - y_i|x_i)}\right)\left(\frac{E(y_{reg2Est} - y_i|x_i,z)}{y_{reg2Est} - y_i}\right)$$

$$\left(\frac{y_{reg3Est} - y_i}{E(y_{reg3Est} - y_i|y_i,z)}\right)\left(\frac{E(y_{reg1Est} - y_i|y_i,x_i,z)}{y_{reg1Est} - y_i}\right)$$

--------------(4.7)

Meanwhile, the non-parametric rescaled weight (NPRW-III) is obtained following the process given here. The absolute difference between the regressed observed and observed value of $w_i$, as well as the conditional expectation of the difference between the regressed and observed value of wi on the variables $x_1$ and z, y and z, and y, $x_1$ and z, are obtained first. Then the product of the reciprocal of the sample inclusion probability and the ratios of the absolute difference and conditional expectations is obtained to give the weight. See below

$$w_d = w_i \left( \frac{w_{regEst} - w_i}{\mathrm{E}\left(w_{regEst} - w_i \big| x_i\right)} \right) \left( \frac{\mathrm{E}\left(w_{reg2.Est} - w_i \big| x_1, z\right)}{w_{reg2.Est} - w_i} \right)$$

$$\left( \frac{w_{reg3.Est} - w_i}{\mathrm{E}\left(w_{reg3.Est} - w_i \big| y, z\right)} \right) \left( \frac{\mathrm{E}\left(w_{reg1.Est} - w_i \big| y, x_i, z\right)}{w_{reg1.Est} - w_i} \right)$$

-------------(4.8)

### Non-parametric weight (NPW(I) & NPW(II))

The process of obtaining the sample weight, which is the reciprocal of the sample inclusion probability, is described in [5,6]. The non-parametric rescaled weight (NPW(I)) is calculated using the procedure outlined here. The conditional expectation of the difference between the regressed and observed value of the dependent variable y on the predictor variables $x_1$, $x_2$, x3, and $x_4$ and the sum of the conditional expectation of the difference between the regressed and observed values on $x_1$, $x_2$, x3, and $x_4$ are first obtained. The weight is calculated as the product of the reciprocal of the sample inclusion probability and the ratio of the conditional expectation of the difference between the regressed and observed value of the dependent variable y on the predictor variables x1, x2, x3, and x4, as well as the sum of the conditional expectation of the difference between the regressed and observed values on x1, x2, x3, and x4. See below:

$$w_c = w_i \frac{\mathrm{E}\left(\left(y_{.reg.Est} - y\right)\big| x_{1i}, x_{2i}, x_{3i}, x_{4i}\right)}{\sum_{i=1}^{n} \mathrm{E}\left(\left(y_{.reg.Est} - y\right)\big| x_{1i}, x_{2i}, x_{3i}, x_{4i}\right)}$$

-------------(4.9)

The non-parametric rescaled weight (NPW-II) is obtained following the process given here. The conditional expectation of the difference between the regressed and observed values of the reciprocal of the sample inclusion probability is dependent on y, $x_2$, $x_3$, and $x_4$, and the sum of the conditional expectation of the difference between the regressed and observed values of the reciprocal of the sample inclusion probability on y, $x_2$, $x_3$, and $x_4$ is first obtained. Then the product of the reciprocal of the sample inclusion probability $w_i$ and the ratio of the conditional expectation of the difference between the regressed and observed value of the reciprocal of the sample inclusion probability on y, $x_2$, $x_3$, and $x_4$ and the sum of the conditional expectation of the difference between the regressed and observed values of the reciprocal of the sample inclusion probability on y, $x_2$, $x_3$, and $x_4$ are obtained to give the weight. Seen below:

$$w_e = w_i \frac{\left(w_{reg.Est} - w_i\right)}{\mathrm{E}\left(\left(w_{reg.Est} - w_i\right)\big| y_i, x_{2i}, x_{3i}, x_{4i}\right)}$$

------------- (4.9)

### Weights used for re-weighting estimators

The table below gives weights used to re-weight estimators starting with sample inclusion probability, $\tilde{P}_1$.

| s/n | Estimator | Plan - Type | Weight |
|---|---|---|---|
| | | **Weights used for re-weighting estimators** | |
| 1. | NPWLE | $\tilde{P}^{-1}$ | $w_i = \tilde{P}^{-1} = 1/g(x_i) = 1/\exp\left(\dfrac{\left(x_i - mean_{x_i}\right)\big/ sd_{x_i}}{\sqrt{\pi_i}}\right)$ |
| 2. | WCNPE(I) | $w_a$ | $WCNP(I) = w_a = \dfrac{\delta E_S\left(w_i \mid y_i, v_i, \theta, \gamma, \beta\right)}{\delta\beta_i}$ |
| 3. | WCNPE(II) | $w_b$ | $WCNP(II) = w_b = \dfrac{1}{\dfrac{\delta E_S\left(w_i \mid y_i, v_i, \theta, \gamma, \beta\right)}{\delta\beta_i}}$ |
| 4. | NPRW(I) | $\tilde{w}_i^{3m}$ | $NPRW(1) = \tilde{w}_i^{3m} = \left(\dfrac{w_{reg1}}{E\left(w_{reg1} \mid y_i, x_i, z\right)}\right)\left(\dfrac{E\left(w_{reg2} \mid x_i\right)}{w_{reg2}}\right)\left(\dfrac{w_i}{E\left(w_i \mid x_i\right)}\right)$ |
| 5. | NPRW(II) | $\tilde{w}_i^{4m}$ | $\tilde{w}_i^{4m} = w_i\left(\dfrac{y_{regEst} - y_i}{E\left(y_{regEst} - y_i \mid x_i\right)}\right)\left(\dfrac{E\left(y_{reg2Est} - y_i \mid x_i, z\right)}{y_{reg2Est} - y_i}\right)\left(\dfrac{y_{reg3Est} - y_i}{E\left(y_{reg3Est} - y_i \mid y_i, z\right)}\right)\left(\dfrac{E\left(y_{reg1Est} - y_i \mid y_i, x_i, z\right)}{y_{reg1Est} - y_i}\right)$ |
| 6. | NPW(I) | $=w_c$ | $w_c = w_i \dfrac{E\left(\left(y_{.reg.Est} - y\right) \mid x_{1i}, x_{2i}, x_{3i}, x_{4i}\right)}{\sum_{i=1}^{n} E\left(\left(y_{.reg.Est} - y\right) \mid x_{1i}, x_{2i}, x_{3i}, x_{4i}\right)}$ |
| 7. | NPRW(III) | $= w_d$ | $w_d = w_i\left(\dfrac{w_{regEst} - w_i}{E\left(w_{regEst} - w_i \mid x_i\right)}\right)\left(\dfrac{E\left(w_{reg2.Est} - w_i \mid x_1, z\right)}{w_{reg2.Est} - w_i}\right)\left(\dfrac{w_{reg3.Est} - w_i}{E\left(w_{reg3.Est} - w_i \mid y, z\right)}\right)\left(\dfrac{E\left(w_{reg1.Est} - w_i \mid y, x_i, z\right)}{w_{reg1.Est} - w_i}\right)$ |
| 8. | NPW(II) | $=w_e$ | $w_e = w_i \dfrac{\left(w_{reg.Est} - w_i\right)}{E\left(\left(w_{reg.Est} - w_i\right) \mid y_i, x_{2i}, x_{3i}, x_{4i}\right)}$ |

**Table a**

## Weights used for matching estimators

| s/n | Estimator | Weight-Type | Weighting equation |
|---|---|---|---|
| | | | **Weights used for matching estimators** |
| 1. | NPWLE | $P^{-1}$ | $w_i = 1/g(x_i)$ |
| 2. | WCNPE(I) | $w_a$ | $w_a = WCNP(I)$ |
| 3. | NPRWE(I) | $\tilde{w}_i^{3m}$ | $\tilde{w}_i^{3m} = NPRW(I)$ |
| 4. | NPRWE(II) | $\tilde{w}_i^{4m}$ | $\tilde{w}_i^{4m} = NPRW(II)$ |
| 5. | NPWE(I) | $w_c$ | $w_c = NPW(I)$ |
| 6. | NPRWE(III) | $w_d$ | $w_d = NPRW(III)$ |

**Table b**

### Finite small sample properties of estimators

The first property deals with the mean location of the distribution of the estimator.

Biasedness - The bias of an estimator is defined as:

$$Bias\left(\hat{\theta}\right) = E\left(\hat{\theta}\right) - \theta$$  ------------------ (5.1)

Where θ is an estimator of θ, an unknown population parameter. If E (θ) = θ then the estimator is unbiased. If E (θ) ≠ θ then the estimator has either a positive or negative bias. That is, on average the estimator tends to over (or under) estimate the population parameter.

A second property deals with the variance of the distribution of the estimator. Efficiency is a property usually reserved for unbiased estimators.

Efficiency - Let $\theta_1$ and $\theta_2$ be unbiased estimators of θ with equal sample sizes. Then, $\theta_1$ is a more efficient estimator than $\theta_2$ if

$$\text{var}\left(\hat{\theta}_1\right) < \text{var}\left(\hat{\theta}_2\right)$$  ---------------- (5.2)

refer to [5,6].

### Results

For each sample plan, we ran the simulation 10,000 times. We used software programs built for weighted estimator analysis in R.

| s/n | Estimators | Coeff. Det. | Sample Bias | Standard error | AIC | BIC | rmse | Var. ratio |
|---|---|---|---|---|---|---|---|---|
| 1. | NPWLE (HT) | 0.800345925994816 | 0.0704873959 | 0.1133457038 | 145 | 148 | 61.40643 | 1.000 |
| 2. | WCNPE(I) | 0.999999999999998 | 7.77156e-16 | 1.46134e-15 | -66.9 | 149 | 2.475e-06 | 7.639 |
| 3. | NPRWE(I) | 0.966488126946826 | 0.0109821936 | 0.0206298061 | 108 | 111 | 13.04562 | 3.719 |
| 4. | NPRWE(II) | 0.9999999999999992 | 3.330669e-16 | 3.218755e-16 | -266 | -263 | 2.242e-06 | 2.919 |
| 5. | NPWE(I) | 0.9999999999999603 | 1.154632e-14 | 2.065846e-14 | -267 | 265 | 2.120e-06 | 167.1 |
| 6. | NPRWE(III) | 0.9999999999999603 | 1.187939e-14 | 2.060457e-14 | -267 | -265 | 2.119e-06 | 167.1 |

**Table 1:** Summaries of Estimators Performance for Generated Data, n = 12.

Table 1 shows that the Estimators have greater relative efficiency and coefficients of determination than NPWLE(HT), and so are more efficient for Simulated Data for n = 12. Where NPWLE(HT) is our reference estimator, the Horvitz-Thompson Estimator.

According to the results in table 2, all estimators with a relative efficiency larger than one are more efficient than NPWLE(HT), which is the Horvitz-Thompson Estimator for n = 12.

| s/n | Estimators | Coeff. Det. $R^2$ | Sample Bias | Standard error | AIC | BIC | rmse | Var. ratio |
|---|---|---|---|---|---|---|---|---|
| 1. | NPWLE (HT) | 0.4414891306018555 | 0.3150583870 | 0.1873574593 | -84.6 | 181 | 528.1731 | 1.000 |
| 2. | WCNPE(I) | 0.9999999998264293 | 8.302836e-11 | 9.121266e-11 | 83.8 | 86.2 | 2.056e-05 | $2.1e^{05}$ |
| 3. | NPRWE(I) | 0.9999999999227354 | 3.638156e-11 | 4.524960e-11 | -203 | -201 | 1.372e-05 | $2.1e^{05}$ |
| 4. | NPRWE(II) | 0.9999999999227354 | 3.691025e-11 | 4.379275e-11 | -203 | -201 | 1.372e-05 | $2.1e^{05}$ |
| 5. | NPWE(I) | 0.9999999999999972 | 9.992007e-16 | 1.629883e-15 | -289 | -287 | 2.769e-07 | $1.8e^{04}$ |
| 6. | NPRWE(III) | 0.9999999999999991 | 3.330669e-16 | 4.696557e-16 | -301 | -265 | 1.585e-07 | $1.8e^{04}$ |

**Table 2:** Summaries of Estimators Performance for Real Data, n = 12.

| s/n | Estimators | $\hat{\beta}_0$ | | $\hat{\beta}_1$ | | $\hat{\beta}_2$ | | $\hat{\beta}_3$ | | $\hat{\beta}_4$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Bias | Error | Bias | Error | Bias | Error | Bias | Error | Bias | Error |
| 1 | NPWLE (HT) | 898.741310 | | -3.039630 | | -0.134045 | | 0.224827 | | -0.066847 | |
| | | 53.346 | 441.375 | 0.185 | 2.033 | -0.017 | 0.113 | -0.021 | 0.254 | -0.019 | 0.121 |
| 2 | WCNPE (I) | 370.391641 | | -1.951511 | | -0.032212 | | 0.144849 | | -0.038341 | |
| | | -7.472 | 98.537 | 0.045 | 0.414 | 0.002 | 0.025 | -0.001 | 0.053 | -0.002 | 0.027 |
| 3 | NPRWE (I) | 115.035368 | | -2.484001 | | 0.011200 | | 0.280233 | | 0.018031 | |
| | | 8.8e-07 | 1.8e-05 | -1.2e-09 | 7.4e-08 | -1.9e-10 | 4.5e-09 | 5.6e-10 | 8.0e-09 | -5.4e-10 | 4.9e-09 |
| 4 | NPRWE (II) | 115.035368 | | -2.484001 | | 0.011200 | | 0.280233 | | 0.018031 | |
| | | 8.8e-07 | 1.8e-05 | -1.2e-09 | 7.4e-08 | -1.9e-10 | 4.5e-09 | 5.6e-10 | 8.0e-09 | -5.4e-10 | 4.9e-09 |
| 5 | NPWE (I) | 54.867425 | | -0.189712 | | -0.001530 | | 0.025961 | | -0.013883 | |
| | | -2.2e-06 | 1.9e-05 | -4.5e-09 | 7.0e-08 | 9.6e-10 | 4.8e-09 | -1.1e-10 | 9.7e-09 | 5.3e-10 | 4.8e-09 |
| 6 | NPRWE (III) | 54.867425 | | -0.189712 | | -0.001530 | | 0.025961 | | -0.013883 | |
| | | -2.2e-06 | 1.9e-05 | -4.5e-09 | 7.0e-08 | 9.6e-10 | 4.8e-09 | -1.1e-10 | 9.7e-09 | 5.3e-10 | 4.8e-09 |

**Table 3:** Summary of the Performance of Coefficients of Weighted Estimators based on Bias and Standard Error for Generated Data, n = 12.

The results in table 3 for the generated data show a summary of the performance of the coefficients of weighted estimators based on bias and standard error for n = 12, and the summary shows that the coefficients of the weighted estimators appear to have smaller bias and standard errors than the Horvitz-Thompson Estimator for n = 12.

| s/n | Estimators | $\hat{\beta}_0$ | | $\hat{\beta}_1$ | | $\hat{\beta}_2$ | | $\hat{\beta}_3$ | | $\hat{\beta}_4$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Bias | Error | Bias | Error | Bias | Error | Bias | Error | Bias | Error |
| 1. | NPWLE (HT) | -1.0998e+04 | | 7.7141e+01 | | 1.0998e+00 | | 3.7291e+00 | | -4.3881e-02 | |
| | | 6553.1 | 32876.5 | -21.144 | 97.409 | -0.164 | 5.693 | -6.646 | 41.323 | -0.356 | 1.931 |
| 2. | WCNPE (I) | -62.502529 | | 0.114179 | | 0.005019 | | 0.098543 | | -0.008280 | |
| | | 406.935 | 8169.08 | -0.620 | 12.947 | -0.036 | 0.739 | -0.494 | 10.060 | -0.009 | 0.289 |
| 3. | NPRWE (I) | 96.453305 | | 0.243113 | | 0.000870 | | -0.015891 | | 0.002458 | |
| | | -1.3e-01 | 3.7e+00 | 3.9e-04 | 1.2e-02 | 2.0e-05 | 5.8e-04 | 9.6e-05 | 3.0e-03 | 1.5e-08 | 2.2e-07 |
| 4. | NPRWE (II) | 96.453305 | | 0.243113 | | 0.000870 | | -0.015891 | | 0.002458 | |
| | | -1.3e-01 | 3.7e+00 | 3.9e-04 | 1.2e-02 | 2.0e-05 | 5.8e-04 | 9.6e-05 | 3.0e-03 | 1.5e-08 | 2.2e-07 |
| 5. | NPWE (I) | -9.3860e+01 | | 3.7993e-01 | | 2.8069e-02 | | 5.1453e-02 | | 5.5068e-04 | |
| | | -1.3e-02 | 9.2e-01 | 2.1e-04 | 4.8e-03 | 6.6e-06 | 1.5e-04 | 2.2e-05 | 6.8e-04 | -3.6e-10 | 9.7e-09 |
| 6. | NPRWE (III) | -9.3860e+01 | | 3.7993e-01 | | 2.8069e-02 | | 5.1453e-02 | | 5.5068e-04 | |
| | | -1.3e-02 | 9.2e-01 | 2.1e-04 | 4.8e-03 | 6.6e-06 | 1.5e-04 | 2.2e-05 | 6.8e-04 | -3.0e-10 | 8.2e-09 |

**Table 4:** Summary of the Performance of Coefficients of Weighted Estimators based on Bias and Standard Error for actual Data, n = 12.

Table 4 shows a summary of the performance of the coefficients of weighted estimators based on bias and standard error for n = 12, and the summary shows that the coefficients of the weighted estimators appear to have smaller bias and standard errors than the Horvitz-Thompson estimator for n = 12.

## Conclusion

The reciprocal of the sample inclusion probability was utilized as weights in this work, with the primary goal of developing nonparametric weighted estimators that are more comparatively efficient and have lower sample bias when compared to classical estimators such as the Horvitz-Thompson estimator. The study developed three new estimators: the non-parametric rescaled weighted estimators II and III as well as the non-parametric weighted estimator I.

## Bibliography

1. Hettmansperger T P and J W McKean. "A robust alternative based on ranks to least squares in analyzing linear models". *Technometrics* 19 (1977): 275-284.

2. Hettmansperger T P and J W McKean. "Robust Nonparametric Statistical Methods". Arnold, London (1998).

3. Hollander M and D A Wolfe. "Nonparametric statistical methods". 2nd edition. John Wiley and Sons, New York (1999).

4. Jaeckel L A. "Estimating regression coefficients by minimizing the dispersion of the residuals". *Annals of Mathematical Statistics* 43 (1972): 1449-1458.

5. Kamun, S. J., Simwa, R. and Sewe, S. "On Derivation of the Semi-Parametric Weighted Likelihood Estimator, SPW, and the Weighted Conditional Pseudo Likelihood Estimator, WCPE". *Far East Journal of Theoretical Statistics* © 2021 Pushpa Publishing House, Prayagraj, 62.2 (2021): 81-90.

6. Kamun, S. J., Simwa, R. and Sewe, S. "Comparison of the New Estimators: The Semi-Parametric Likelihood Estimator, SPW, and the Conditional Weighted Pseudo Likelihood Estimator, WPCE". *American Journal of Theoretical and Applied Statistics* 10.4 (2021): 202-207.

7. Kamun, S. J. "On Derivation of Non-Parametric Weighted Linear Models". *IJIRSE* 2.12 (2022): 25-30.