# A COMPARATIVE STUDY OF THE EFFICIENCY PROPERTIES OF IMPROVED ESTIMATORS IN THE LINEAR REGRESSION MODEL

**Neeraj Rani[1], Daljeet Kaur[2]**
[1,2]Guru Kashi University, Talwandi Sabo

## ABSTRACT

*The emergence of various assessors of boundaries of direct relapse models, particularly when applied to genuine conditions, can be followed to the non-legitimacy of the suppositions under which the model is generated. Regression analysis can be used to produce predictions in any case. Because regression analysis typically involves non-experimental data, related variables are frequently included in the analysis. Multicollinearity in relapse models happens when at least two indicator factors are related with each other. Because of this issue, the worth of the least squares registered relapse coefficients can become restrictive on the connected indicator factors in the model. As a result, this study took into account the concept of an upgraded estimate using Principal Component Regression, as well as the theoretical features of the suggested estimator.*

*Keywords: - Regression, Linear, Model, Variable, Predictor.*

## I. INTRODUCTION

Linear regression has a substantial body of research in regression models due to its simplicity, ease of analysis, and well-developed inferential techniques. In a range of domains, including social, behavioural, medical, management, and other applied sciences, the linear regression model has proven to be effective. In these domains, it is regarded as one of the most important devices. The simplest regression model is the two- variable straight relapse model, which shows the direct connection between the two factors. Linearity alludes to the linearity of the boundaries to be assessed in this situation. A more wide procedure is the different direct relapse model, which expects that the reaction variable is a straight capacity of the model boundaries and that the model include numerous independent variables. A multiple regression model can be used to investigate the impact of numerous independent factors on response at the same time.

For estimating the parameters in regression models, methods such as the Least Squares technique, the Maximum Likelihood Estimation strategy, the Minimum Chi Squares strategy, and others are accessible. These techniques vary regarding registering effortlessness, presence of a shut structure arrangement, power, and hypothetical suppositions important to help the ideal measurable properties. The Least Squares Method

is the most notable and generally utilized. In the late eighteenth and mid nineteenth century, Legendre and Gauss presented this method independently. The objective to observe boundary gauges by picking the relapse line that is nearest to all information focuses inspires the least squares technique. The best-fitting line for the noticed information is gotten by diminishing the amount of the squared deviations from every information highlight the line. The least squares way to deal with assessing enjoys the benefit of utilizing the example information and giving unprejudiced evaluations, as well as being easy to apply and low in processing cost. In inferential strategies, the least squares method has generally been advocated by two presumptions: it creates most extreme probability assessments of obscure relapse coefficients and all straight fair-minded estimators, and it has the least variation about the regression line.

## II. USE OF LINEAR REGRESSION MODEL FOR SOLVING RELATIONSHIPS

The most frequently involved measurable methodology for settling utilitarian association issues between factors is the direct relapse model. It assists with associating noticed upsides of at least one autonomous factors with perceptions of a reliant variable y, X1, X2,..., Xp. While endeavoring to make sense of the reliant variable, anticipating the upsides of the reliant variable is basic. The straight relapse model is additionally predicated on a bunch of basic suppositions. Among these presumptions, regressors are non-stochastic (fixed in continued inspecting) and autonomous. Additionally, the mistake terms ought to be free, have a steady change, and be unaffected by the regressors. The Ordinary Least Square (OLS) assessor is utilized when all of the assumptions of a standard linear regression model are met:

$$\hat{\beta} = (X^1 X)^{-1} X^1 Y$$

III. Some of the optimal or optimum qualities of an estimator are known to be linearity, unbiasedness, and efficiency. The Best Linear Unbiased Estimator was created by combining these (BLUE). These presumptions, notwithstanding, are not generally met, in actuality, circumstances. Therefore, various techniques for assessing model boundaries have been created.

IV. The presumption of non-stochastic regressors isn't generally fulfilled, particularly in business, financial aspects, and sociologies, in light of the fact that their regressors are much of the time produced by stochastic cycles outside their ability to control. Many creators have talked about conditions and cases in which this supposition that is abused, as well as the ramifications for utilizing the OLS assessor to gauge model boundaries. Regardless of whether the regressors are stochastic and autonomous of the mistake terms, the OLS assessor, regardless of whether it isn't BLUE, is as yet impartial and has the most reduced difference. They additionally referenced that assuming the mistake terms are believed to be typical, traditional hypothesis testing can still be used. However, changes must be made to the confidence intervals produced for each sample and the test's power.

**V.** When the assumption of independent regressors is broken, multicollinearity occurs. Because the regression model's core assumption has been broken when regressors are heavily correlated, the interpretation offered to the relapse results may as of now not be legitimate. Albeit the relapse coefficients assessed by the OLS assessor are as yet impartial as long as multicollinearity is more than a little flawed, the relapse coefficients might have critical testing blunders, harming both deduction and guaging. Various strategies for determining model parameters have been established when a data collection exhibits multicollinearity. Two of these estimators are Ridge Regression and Principal Component Regression Estimator.

## VI. THE ESTIMATORS AND THEIR PROPERTIES

### The concept behind regression estimation

At the point when the assistant variable x is straightly connected with y yet doesn't go through the beginning, a direct relapse assessor is sufficient. This isn't to recommend that the relapse gauge can't be utilized when the block is close to nothing. In such cases, the two assessments, relapse and proportion, might be practically indistinguishable, and you can pick which to utilize.

Moreover, assuming that various helper factors have a straight relationship with y, numerous relapse assessments might be suitable.

The straight connection among y and realized x-values can be utilized to appraise the mean and complete of y-values, addressed as and.

Let's start with a simple example:

$$\hat{y} = a + bx ,$$

which is our basic regression equation.

Then,

$$b = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \quad \text{and}$$

$$a = \bar{y} - b\bar{x}$$

Then to estimate the mean for *y*, substitute as follows:

$$x = \mu_x, a = \bar{y} - b\bar{x}, then$$
$$\hat{\mu}_L = (\bar{y} - b\bar{x}) + b\mu_x$$
$$\hat{\mu}_L = \bar{y} + b(\mu x - \bar{x}), \hat{\mu}_L = a + b\mu_x$$

At the point when the helper variable x is directly connected with y yet doesn't go through the beginning, a straight relapse assessor is satisfactory. This isn't to recommend that the relapse gauge can't be utilized when the block is close to nothing. In such cases, the two assessments, relapse and proportion, might be almost identical, and you can choose which to use.

Furthermore, if numerous helper factors have a straight relationship with y, different relapse appraisals might be fitting.

The straight connection among y and realized x-values can be utilized to assess the mean and complete of y-values, addressed as and.

$$\hat{V}ar(\hat{\mu}_L) = \frac{N-n}{N \times n} \cdot \frac{\sum_{i=1}^{n}(y_i - a - bx_i)^2}{n-2}$$
$$= \frac{N-n}{N \times n} \cdot MSE$$

where *MSE* is the *MSE* of the linear regression model of *y* on *x*.

Therefore, an approximate (1-α)100% CI for μ is:

$$\hat{\mu}^L \pm t_{n-2,\alpha/2}\sqrt{\hat{V}ar(\hat{\mu}_L)}$$

It follows that:

$$\hat{\tau}_L = N \cdot \hat{\mu}_L = N\bar{y} + b(\tau_x - N\bar{x})$$

$$\hat{V}ar(\hat{\tau}_L) = N^2\hat{V}ar(\hat{\mu}_L)$$
$$= \frac{N \times (N-n)}{n} \cdot MSE$$

And, an approximate (1-α)100% CI for τ is:

$$\hat{\tau}_L \pm t_{n-2,\alpha/2}\sqrt{\hat{V}ar(\hat{\tau}_L)}$$

## Properties of Estimators

The conventional least squares gauge of β is gotten by applying the least squares guideline.

b = 〚(X'X)〛 ^(- 1) X'y

which is notable to be the unprejudiced assessor with difference - covariance grid given by

V (b) = σ^2 〚(X'X)〛 ^(- 1)

Linearity and absence of prejudice are overlooked A class of assessors proposed by James and Stein (1961) is better than the standard conventional least squares given by (4.1) assuming the coefficient vector is bigger than two aspects. Considering this, how about we center around the accompanying Stein rule assessor:

## VII. VARIOUS COMPARISON CRITERIA FOR PERFORMANCE PROPERTIES OF ESTIMATORS

A variety of estimators for estimating regression coefficients can be found in the literature. The Mean Squared Error (MSE) criterion is one of the most extensively used approaches for comparing the performance of competing estimators. The mean squared blunder, which was designed via Carl Friedrich Gauss in his work on measurements, can be utilized to evaluate an assessor and the genuine worth of a boundary. It is feasible to determine the error's second instant using MSE, which combines variance and bias into a single metric. The MSE matrix of estimator p is shown in multiple linear regression. By

$$M(\hat{\beta}) \; = \; E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)']$$

The problem with mean squared error is that it gives extreme values too much weight. Because terms are squared, large errors are more severely weighted than smaller ones.

In a decision-theoretic framework, L (,) is a loss function that gives us a measure of how much we lose if we use the estimator when the true parameter is. The average or expected value of the loss function must be used when calculating the "risk function," and this value is provided by

$$R(\hat{\beta}, \beta) \; = \; E[L(\hat{\beta}, \beta)]$$

The overall standard is to pick a gauge that limits risk for each worth of. It's called an allowable gauge. We must ensure that the estimators' properties are not changed by the loss function in order for the decision rule to be resilient to it. The problem is that estimating the loss function for a specific problem is tough, thus it's all too easy to fall back on traditional loss functions.

## VIII. COMPARISION OF ESTIMATORS ON THE BASIS OF RISK PERFORMANCE

The risk performance of the estimator () is compared to that of other weighted estimators in this section. Changing the weights w1 and w2 can assist assess the effectiveness of

different estimators when dealing with different loss structures. Compare the performance of () with the weighted Stein-rule estimator _s. Dube and Manocha (2006) calculated the risk of the weighted Stein-rule estimator and the conditions under which it outperforms the ordinary least squares estimator using a non-normal error distribution. Let's use their method to find the risk difference between the weighted Stein-rule estimator and that of (), which is provided by).

$$R(\widetilde{\beta}_s) - R(\widetilde{\beta}) = -\frac{2k \, \gamma_1 \sigma^3}{\beta'X'X\beta}(1-w_1)(1-w_2)\beta'X'(1*\overline{P}_X)e$$

$$+ \sigma^4 \frac{k^*(1-w_1)}{\beta'X'X\beta}\left[(k+k_1)(1-w_1)\{\gamma_2 \mathrm{tr}(\overline{P}_X * \overline{P}_X)\right.$$

$$+ (n-p)(n-p+2)\}$$

$$\left. -2(1-w_2)\{ \gamma_2 \mathrm{tr}\, M_2(1*\overline{P}_X)+(n-p)(p-2)\}\right]$$

where k* = (k - k]) where k is any non-negative scalar that describes the Stein rule estimator.

We can see from (4.1) that the leading term defines the difference in associated risks of $\widetilde{\beta}$ and $\widetilde{\beta}_s$ for skewed distributions of disturbances. The superiority of $\widetilde{\beta}$ over $\widetilde{\beta}_s$ under a balanced loss function will be ensured if this term is positive.

The degree of efficiency is given by the term of order $O(\sigma^4)$ for symmetric distributions $\gamma_1 = 0$. Therefore, when $\gamma_1 = 0$ reduces to

where

$$q = \gamma_2 \, tr \, (\overline{P_X} * \overline{P_X}) + (n - p)(n - p + 2)$$

$$g = \gamma_2 \, tr \, [M_2(I * \overline{P_X})] + (n - p)(p - 2)$$

The criterion of dominance of p over Ps may be easily proven for symmetric leptokurtic and symmetric platykurtic disturbance distributions (4.2). Based on Vinod and Ullah (1978), we propose the following notations for this purpose:

$$\theta = \frac{\gamma_2}{n - p + \gamma_2}, G = (X'X)'1 \, X'(I * Px)X, \qquad \phi = \frac{trG}{p}$$

(4.2) can be recast using the amended notations as

$$R(\widetilde{\beta}_s) - R(\widetilde{\beta}) = \frac{\sigma^4(1-w_1)}{\beta'X'X\beta} k^* \left[(k+k_1)(1-w_1)q\right.$$

$$\left. -2(1-w_2)(n-p+\gamma_2)\{p(1-\theta\Phi)-2(1-\theta\frac{\beta'X'XG\beta}{\beta'X'X\beta})\}\right]$$

It can be shown from the updated notations that

$0 < \theta < 1$ whenever $\gamma_2 > 0$

$\theta = 0$ whenever $\gamma_2 = 0$

$\theta < 0$ whenever $\gamma_2 < 0$

assuming symmetrical leptokurtic distribution of disturbances, we can see that i.e. $\gamma_1 = 0$ and $\gamma_2 > 0$, $\tilde{\beta}$ dominates $\tilde{\beta}_s$ so long as k* > 0 and

$$k + k_1 > 2 \frac{(1-w_2)(n-p+\gamma_2)}{(1-w_1)q} \{(1-\theta\Phi)p - 2(1-\theta\frac{\beta'X'GX\beta}{\beta'X'X\beta})\}$$

As (4.4) It can't be applied in real-life circumstances because it involves unknown parameters. We observe that

$$\max \frac{\beta'X'XG\beta}{\beta'X'X\beta} = \eta_P$$

And $\min \frac{\beta'X'XG\beta}{\beta'X'X\beta} = \eta_l$

The smallest and greatest eigen values of the matrix G are $\eta_1$ and $\eta_2$. As a result, 4.4 can be rewritten as

$$k > \frac{(1-w_2)(n-p+\gamma_2)}{(1-w_1)q}\{(1-\theta\Phi)p - 2(1-\theta\eta_l)\} \quad ; \quad p > \frac{2(1-\theta\eta_l)}{(1-\theta\Phi)}$$

whenever k > k₁

Similarly for symmetric platykurtic distribution of disturbances i.e. $\gamma_1 = 0$ and $\gamma_2 < 0$, Rao (2002) showed that $n - p > 2$ which implies $2 + \gamma_2 \geq 0$. As a result, the dominance of

$\tilde{\beta}$ over $\tilde{\beta}_s$(4.4) holds so long as

$$k > \frac{(1-w_2)(2+\gamma_2)}{(1-w_1)q}\{p(1-\theta\Phi) - 2(1-\theta\eta_l)\} \quad ; \quad p > \frac{2(1-\theta\eta_l)}{(1-\theta\Phi)}$$

For normally distributed errors i.e. $\gamma_1 = 0$ and $\gamma_2 = 0$ the above dominance conditions reduces to

$$k > \frac{(1-w_2)(p-2)}{(1-w_1)(n-p+2)} " p > 2,$$

Assuming the choice of $k_1 = \frac{1}{n-p}$ and $\frac{p}{n-p}$ in the (4.5) - (4.7) Weighted FMMSE estimator and AFFMSE estimator can be used to obtain the dominance condition over $\tilde{\beta}_s$.

## IX. CONCLUSION

Small disturbance asymptotic approximations have been used to study a family of better estimators and their generalisation at the point when mistakes in the straight relapse model are not typical all of the time. They were contrasted with standard least squares and Stein-rule assessors for execution properties under quadratic mistake and adjusted misfortune capacities. An examination of the exhibition of numerous assessors has been made.

Assuming at least two of the free factors in a relapse model are intercorrelated or subject to one another, the issue of multicollinearity arises. One of the most obvious outcomes is that a lot of the model's variables are related to or dependent on one another. In this setting, non-conventional estimate methodologies have been developed because traditional estimators are ineffective. When there is a complex relationship between the variables, regression is a great method to apply in exploratory research.

## REFERENCES: -

[1] Suhail, Muhammad & Chand, Sohail& Babar, Iqra. (2021). On some new ridge m-estimators for linear regression models under various error distributions. Pakistan Journal of Statistics. 37. 369-391.

[2] Ayinde, K., Lukman, A. F., Alabi, O. O. and Bello, H. A. (2020). A new approach of principal component regression estimator with applications to collinear data. International Journal of Engineering Research and Technology, 13(7), 1616-1622.

[3] Ahmed, Syed. (2014). Estimation Strategies in Multiple Regression Models. 10.1007/978-3-319-03149-1_4.

[4] Ibikunle, Ogunyinka&Sodipo, Ademola. (2013). Efficiency of Ratio and Regression Estimators Using Double Sampling. Journal of Natural Sciences Research, IISTE. 3. 2224-3186.

[5] K. Ayinde, E. Apata and O. Alaba (2012) "Estimators of Linear Regression Model and Prediction under Some Assumptions Violation," *Open Journal of Statistics*, Vol. 2 No. 5, 2012, pp. 534-546. doi: 10.4236/ojs.2012.25069.

[6] Hoque, Zahirul. (2012). Improved estimation for dynamic linear regression model. New Developments in Applied Statistics. 17. 319-330.

[7] Bhushan, Shashi & Pandey, Anshula& Singh, R.. (2009). IMPROVED CLASSES OF REGRESSION TYPE ESTIMATORS. International Journal of Agricultural and Statistics Sciences. 5. 73-84.

[8] Alpuim, Teresa & El-Shaarawi, Abdel. (2008). On the efficiency of regression analysis with AR(p) errors. Journal of Applied Statistics. 35. 717-737. 10.1080/02664760600679775.

[9]  Maroco, João. (2007). Consistency and Efficiency of Ordinary Least Squares, Maximum Likelihood, and Three Type II Linear Regression Models. Methodology: European Journal of Research Methods for the Behavioral and Social Sciences. 3. 81-88. 10.1027/1614-2241.3.2.81.

[10]  Kibria, B M Golam. (2006). Applications of Some Improved Estimators in Linear Regression. Journal of Modern Applied Statistical Methods. 5. 367-380. 10.22237/jmasm/1162354200.

[11]  Toutenburg, Helge & Srivastava, Viren &Schaffrin, Burkhard &Heumann, Christian. (2003). Efficiency properties of weighted mixed regression estimation. Metron - International Journal of Statistics. LXI. 91-103.

[12]  Vasconcellos, Klaus & Cordeiro, Gauss & Barroso, Lúcia. (2000). Improved estimation for robust econometric regression models. Brazilian Journal of Probability and Statistics. 14. 141-157.