

Simulation Based Comparison among Fifteen Estimators of Correlation Coefficient

Md. Ashad Alam

Department of Statistics

Hajee Mohammad Danesh Science and Technology University

Dinajpur, Bangladesh.

E-mail: kakon_bud@yahoo.com

Mohammed Nasser

Department of Statistics

Rajshahi University

Rajshahi, Bangladesh

E-mail: mnasser.ru@gmail.com

[Received April 4, 2006; Revised November 10, 2008; Accepted November 13, 2008]

Abstract

In this paper, fifteen estimators of correlation coefficient available in the literature have been compared through simulation. We have considered simulating sampling distribution at bivariate standard normal and its five contaminated forms with three different correlation coefficients and three different sample sizes. We also have considered a real population of size 1491. The estimators have been compared with regard to bias, standard error, mean square error and length of 90% percentile interval. The class of robust estimators, especially the normal score estimator has been found superior to the class of non-robust estimators over all. To measure correlation coefficient we have recommend using for normal score estimator especially when sample is large and contamination is probable.

Keywords and Phrases: Correlation Coefficient, Robust Estimator of Correlation Coefficient and Monte Carlo Simulation.

AMS Classification: 62H20, 62F35, 65C05.

1 Introduction

Francis Galton (1822-1911), an English anthropologist and eugenicist, is generally regarded as the founder of correlation analysis. On February 9, 1877, Galton presented a lecture at the Royal Institution of Great Britain entitled “Typical Laws of Heredity in Man” which introduced the concepts of regression (termed “reversion”) and correlation. Galton’s work greatly influenced the career of Karl Pearson (1857-1936), who systematized the application of correlation and developed present day version of product moment correlation coefficient r in 1896. The first published work on rank correlation appeared in 1904 in a psychological study of intelligence by Charles E. Spearman (1863-1945). Gini (1914) developed modified footrule correlation from Spearman correlation coefficient. The next rank correlation was discovered by Kendall (1938). During fifties and sixties of the last century attention was drawn to the fact that inference procedures based on the use of the product-moment correlation r are heavily dependent on the assumption of bivariate normality and very much sensitive to outliers. Various researchers (Blomqvist, 1950; Sheppard, 1899; Kendall, 1970; Fieller and Pearson, 1961; Fisher and Yates, 1963; Gnanadesikan and Kettenring, 1972 etc.) proposed different robust estimators and developed various techniques to assess effect of outliers on them. Basing on the principle of maximum deviations Gideon et al. (1987) developed a new rank correlation coefficient resistant to outliers. Rodgers and Nicewander (1988) provided thirteen ways to look at the product moment correlation coefficient. Gideon (1998) forwarded a generalized interpretation of Pearson’s r , which is important as well as different from those found in Rodgers and Nicewander (1988) as well as formulated three new estimators of correlation coefficient. All these fifteen estimators have been yet to be compared simultaneously. In this article we have considered all the fifteen estimators of correlation coefficient and used Monte Carlo simulation to compare them.

Section 2 describes our methods and models that are very much similar to those in Devlin et al. (1975). Section 3 presents all fifteen estimators of correlation coefficient briefly. Section 4 discusses the results of our study and conclusion is in the last section 5.

2 Methodology and Materials

In this work we have mainly used Monte Carlo simulation. We simulate sampling distribution at bivariate standard normal and its five contaminated forms with three different correlation coefficients (0, 0.5, and 0.9) and three different sample sizes (10, 20, and 50), and a real population of size 1491. The estimators have been compared with regard to bias, standard error, MSE and length of 90% percentile interval. The ‘null’ model is the bivariate normal and there are five ‘non null’ bivariate distributions to represent the situation with outliers: Laplace; contaminated normal; Cauchy; Slash and artificial outlier in minor axis. We have generated 10000 samples in each case

and calculated values of each estimators for each sample. Hence altogether we have simulated total 550000 samples. Devlin et al. (1975) considered only 500 values for each situation and did not use Slash model that is more diverse than Cauchy. Our target population are as follows:

- (a) Standard Bivariate Normal, $F=BVSN(0, \Gamma)$, Γ is a bivariate correlation matrix with off-diagonal element ρ (0, 0.5 and 0.9).
- (b) The distribution with contaminated normal, $0.9 N(0, \Gamma) + 0.1N(0, 9\Gamma)$.
- (c) The distribution with outliers along the minor axes.
- (d) The distribution with 90% bivariate normal and 10% bivariate Cauchy.
- (e) The distribution with 90% bivariate normal and 10% bivariate Laplace.
- (f) The distribution with 50% bivariate normal and 50% bivariate Slash.

We have also considered a real set of health data of 1491 Japanese adult male students from various districts of Japan as population. Four head measurements - head length, head breadth, head height and head circumference were taken by one observer collected by Fumio Ohtsuki[see Hossain et al. (2005)]. We have considered two variables head length and head circumference which possess correlation coefficient 0.7052.

3 Estimators of Correlation Coefficient

3.1 Definition of Pearson correlation coefficient

Correlation methods for determining the strength of the linear relationship between two or more variables are among the most widely applied statistical techniques. Theoretically, the concept of correlation has been a starting point of a building block in the development of a number of areas of statistical research. For discussions of correlation in situations involving more than two variables, the reader should consult the articles on general topics, such as canonical analysis, factor analysis, path analysis, and time series analysis. The correlation between variables X and Y is defined as

$$\rho = corr(XY) = \frac{Cov(X, Y)}{[var(X)var(Y)]^{1/2}} \quad (1)$$

where

$$Cov(X, Y) = E[(X - E(X))(Y - E(Y))]$$

and for a pair of sample $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, the sample correlation is given by

$$r_1 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y} \quad (2)$$

where

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}; \quad \bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$

and

$$s_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}; \quad s_y = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}$$

3.2 Definition of other estimators of correlation coefficient

We have considered fifteen estimators of correlation coefficients available in the literature. Among them four (i-iv) estimators of correlation coefficient (CC) use original data directly, four (v-viii) take account of rank value of original data and the rest are robust estimators according to Devlin et al. (1975).

- i) The Pearson correlation coefficient, r_1 (Pearson, 1896);
- ii) An absolute value CC, r_2 (Gideon, 1998);
- iii) An absolute vale from median CC, r_3 (Gideon, 1998);
- iv) A median-type CC, r_4 (Gideon, 1998);
- v) Spearman's CC, r_5 (Spearman, 1904);
- vi) Spearman's Modified Footrule CC, r_6 (Gini, 1914);
- vii) Kendall's CC, r_7 (Kandall, 1938);
- viii) The Greatest Deviation CC, r_8 (Gideon et al. 1987);
- ix) The quadrant Estimate CC estimate, r_9 (Blomqvist, 1950; Sheppard, 1899);
- x) A transformation of Kendall's CC estimate, r_{10} (Kendall, 1970);
- xi) The Normal Scores CC estimate, r_{11} (Fielier et al. 1957; Fielier and Pearson, 1961);
- xii) The Sum and differences of the standardized observed values CC estimate, r_{12} (Gnanadesikan and Kettenring, 1972);

- xiii) A Bivariate trimming CC estimate, r_{13} (Gnanadesika and Kettenring, 1972);
- xiv) A Bivariate Winsorizing CC estimate, r_{14} (Devlin et al. 1975);
- xv) A CC estimate by trimming with respect to the principal components, r_{15} (Devlin et al. 1975);

The Pearson CC and rank corelation estimators are well known. Gideon (1998) forwarded the other three estimators of CC using original data but he did not compare them with their competitors. The robust estimators were already discussed in the work of Devlin et al. (1975). Here we have briefly defined robust estimators r_{11} , r_{12} , r_{13} , r_{14} , and r_{15} of CC as we have adopted their definitions with slight change. We have also discussed newly proposed Gideon's three estimators mentioned earlier.

Definition 1. The Normal Scores CC estimate, r_{11} . Suppose that there are n pairs of associated ranking $i = u_1, u_2, \dots, u_n$ and $p_i = v_i, v_2 \dots v_n$ where the integers $u_i (i = 1, 2, \dots, n)$ may be taken in ascending order $1, 2, 3, \dots, n$ and v'_i s are a permutation of these integers.

Let $\xi(i|n)$ be a so-called normal order statistic i.e. the expected value of the ith largest standardized deviates in a sample of n observations from a normal population. Following suggestions by Fisher and Yates (1938), Fieller et al. (1957) forwarded a measure of rank correlation obtained from the product moment correlation coefficient of these scores namely.

$$r_{11} = \frac{\sum_{i=1}^n \xi(i|n) \xi(v_i|n)}{\sum_{i=1}^n \xi^2(i|n)} \quad (3)$$

Convenient tables of the individual $\xi(i|n)$ as well as of $\sum_i \xi^2(i|n)$ are given, for example, in Fisher and Yates (1938, Tables XX and XXI). The tables need to be extended for large n 's.

Definition 2. The sum and differences of the standardized observed values CC estimate, r_{12} .

Let $(x_i, y_i), i = 1, 2, \dots, n$ be the bivariate data set. Suppose that there are robust estimator (median) m_x and m_y of x and y variable respectively. Also we have median absolute deviation(MAD), $MAD_x = median|x_i - median(x)|$ and $MAD_y = median|y_i - median(y)|$ as the robust estimators of scale for variables X and Y respectively. Let us consider the robust standardization of the observations x and y to

$$\tilde{x}_i = \frac{x_i - m_x}{MAD_x}$$

$$\tilde{y}_i = \frac{y_i - m_x}{MAD_y}$$

Then the sum and difference are obtained as

$$z_{i1} = \tilde{x}_i + \tilde{y}_i$$

$$z_{i2} = \tilde{x}_i - \tilde{y}_i$$

Now we have to calculate robust variance estimates, MAD_{z_1} and MAD_{z_2} of z_1 and z_2 . By using of this robust variance estimates we obtained

$$r_{12} = \frac{MAD_{z_1} - MAD_{z_2}}{MAD_{z_1} + MAD_{z_2}} \quad (4)$$

In practice, MAD generally involves a multiplicative factors for making it a consistent estimator of the variance in question. The factor is typically a function of both the presumed distribution of the data and the nature of MAD.

Definition 3. A bivariate trimmed CC estimate, r_{13} .

The steps for computing Bivariate trimming CC estimate are as follow:

we

Step 1 start with an initial robust estimate of mean vector and covariance vector. Let

us consider an initial robust estimate of the mean vector $m^* = \begin{bmatrix} m_x \\ m_y \end{bmatrix}$

and covariance matrix, $V^* = \begin{bmatrix} MAD_x & MAD_{xy} \\ MAD_{yx} & MAD_y \end{bmatrix}$

Step 2 then use elliptical metric (the square Mahalanobis distance) which is defined as $d_i^2 = (x_i - m^*)' V^{*-1} (x_i - m^*)$. and temporarily set aside the $[\alpha n]$ points as bivariate trimmed.

Step 3 update m^* and V^* by replacing current values with the sample mean and covariance matrix of the untrimmed points.

Step 4 repeat step (2) and (3) until reasonable convergence is achieved.

Step 5 at last define a bivariate trimmed CC estimate as

$$r_{13} = v_{xy}^* (v_x^* v_y^*)^{-\frac{1}{2}} \quad (5)$$

using the elements of V^* .

Definition 4. A bivariate winsorized CC estimate, r_{14} .

The difference is that the orientation information in the $[\alpha n]$ most extreme points is preserved and exploited by moving them towards the current m^* until their distances are equal to that for the most distant unadjusted point. The adjusted and unadjusted points are then all used to form an updated m^* and V^* .

$$r_{14} = v_{xy}^* (v_x^* v_y^*)^{-\frac{1}{2}} \quad (6)$$

using the elements of V^* .

Definition 5. A CC estimate by trimming with respect to the principal components, r_{15} .

The beginning point in deducing internal structure by principal components is the set of n p -dimensional observations, the columns of $p \times n$ matrix X , which are considered for purposes of the analysis as an unstructured multivariate sample. The usual sample co-variance matrix S or correlation matrix R may then be computed. In the case of S , the linear principal components transformation of the data is given by

$$Z = L(X - \bar{X}) \quad (7)$$

where the p rows of the orthogonal matrix L are eigenvectors of S customarily chosen to correspond to its eigenvalues in descending order of magnitude, and \bar{X} is a $p \times n$ matrix all of whose columns are sample mean vector \bar{x} . Each row, $l'_i (i = 1, 2, \dots, p)$ of L provides a principal component coordinate and each row of z gives the deviations of the projections of the original sample from the projection of the sample centroid \bar{x} onto a specific principal components coordinate.

Thus, for instance, with $p=2$ we can compute a CC estimate by trimming with respect to the principal components r_{15} . The steps are as follows:

Step 1 at first estimate sample covariance matrix from the original data, $S^* = \begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix}$

Step 2 then find a matrix of eigenvector of S and arrange in descending order of magnitude, $L^* = \begin{bmatrix} l_{11} & l_{12} \\ l_{21} & l_{22} \end{bmatrix}$

Step 3 calculate principal components as (7)

Step 4 at last compute univariate trimmed data corresponding to each PC. Hence find revised covariance matrix, V^* from original untrimmed data corresponding to trimmed data Z and we get

$$r_{15} = v_{12}^* (v_{11}^* v_{22}^*)^{-\frac{1}{2}} \quad (8)$$

Definition 6. An absolute value CC, r_2 .

Let us consider (x_i^*, y_i^*) , $i = 1, 2, \dots, n$ be the n pair of observations deviation from their mean i.e $x_i^* = x_i - \bar{x}$ and $y_i^* = y_i - \bar{y}$. The sum of absolute values about the mean is denoted by $SA_x = \sum |x_i - \bar{x}|$ and $SA_y = \sum |y_i - \bar{y}|$ for x and y variables respectively. Thus the absolute value CC is defined as

$$r_2 = \left\{ \sum \left(\frac{x_i^*}{SA_x} + \frac{y_i^*}{SA_y} \right)^2 - \sum \left(\frac{x_i^*}{SA_x} - \frac{y_i^*}{SA_y} \right)^2 \right\} / 2 \quad (9)$$

where the denominator is 2 because $\sum \left| \frac{x_i^*}{SA_x} \right| + \sum \left| \frac{y_i^*}{SA_y} \right| = 2$. It is noticeable that that the same heuristic motivation for Pearson's r holds for absolute value CC, r_2 .

Definition 7. An absolute value from median CC, r_3 .

Because the median is the value of "a" that minimizes $\sum |x_i - a|$, r_2 in equation (9) could be modified as follows to give another correlation coefficient.

$$r_3 = \frac{1}{2} \left(\sum \left| \frac{x_i - m_x}{SAM_x} + \frac{y_i - m_y}{SAM_y} \right| - \sum \left| \frac{x_i - m_x}{SAM_x} - \frac{y_i - m_y}{SAM_y} \right| \right) \quad (10)$$

where m_x, m_y are sample medians, and $SAM_x = \sum |x_i - m_x|$, $SAM_y = \sum |y_i - m_y|$. Unlike r_2 this correlation, r_3 has not been to be bounded between -1 and +1.

Definition 8. A median-type CC, r_4 .

A median-type correlation coefficient is defined as

$$r_4 = \frac{1}{2} \left(\text{med} \left| \frac{x_i - \text{med}(x)}{MAD_x} + \frac{y_i - \text{med}(y)}{MAD_y} \right| - \text{med} \left| \frac{x_i - \text{med}(x)}{MAD_x} - \frac{y_i - \text{med}(y)}{MAD_y} \right| \right) \quad (11)$$

It is evidently not true that $|r_4| \leq 1$, let $x_i^* = \frac{x_i - \text{med}(x)}{MAD_x}$ and similarly for y_i^* . Now $\text{med}|x_i^*| = \text{med}|y_i^*| = 1$. The proof that $|r_4| \leq 1$, breaks down because the median

of the sum of two sets of nonnegative numbers is not always less than the sum of the medians. Simulation studies of r_{mad} show it to behave very much like other correlation coefficients even with the anomaly of being greater than one.

3.3 Computation of all estimators of CC

In this subsection we have drawn a sample of size 10 (Table 1) from the real population and have contaminated each with one outlier and calculated all estimator of correlation coefficient. Using data in Table 1 we have Pearson Correlation Coefficient, $r_1 = +0.7738$. We have replaced x_{10th} observation by 1850. Hence we have compared different estimators by using data in Table 1. These results are given in Table 2 for original data (OD) and contaminated data (CD).

Table 1: Generated data of real population

| | | | | | | | | | | |
|-----------------------|------|------|------|------|------|------|------|------|------|------|
| Head length(x) | 184 | 177 | 182 | 175 | 175 | 174 | 185 | 174 | 181 | 185 |
| Head circumference(y) | 55.2 | 54.6 | 54.2 | 55.0 | 54.2 | 53.4 | 56.9 | 54.2 | 55.1 | 56.2 |

Table 2: Calculation result of original and contaminant data

| | | | | | | | | |
|-----------|--------|----------|----------|----------|----------|----------|----------|-------|
| Nonrobust | r_1 | r_2 | r_3 | r_4 | r_5 | r_6 | r_7 | r_8 |
| OD | 0.7738 | 0.6026 | 0.5769 | 0.5556 | 0.2606 | 0.20000 | 0.1556 | 0.0 |
| CD | 0.4465 | 0.3001 | 0.1921 | 0.500 | 0.2606 | 0.20000 | 0.1556 | 0.0 |
| Robust | r_9 | r_{10} | r_{11} | r_{12} | r_{13} | r_{14} | r_{15} | |
| OD | 0.809 | 0.2419 | 0.613 | 0.3659 | 0.7224 | 0.6369 | 0.8435 | |
| CD | 0.8090 | 0.2419 | 0.6130 | 0.3176 | 0.7197 | 0.7806 | 0.5695 | |

4 Results

In this section we have considered simulation results of three different sample sizes (10, 20, and 50) with three values of ρ (0, 0.5 and 0.9). At first we have regarded all model individually. Individual results are not given in the table for the sake of space but major findings are reported below. In order to show the result simply we have divided all model into two categories ('null' and 'nonnull' models). The 'null' sampling distribution is the bivariate normal and there are five 'nonnull' bivariate distributions to represent the situation with outliers: contaminated normal, artificial outlier in minor axis, Cauchy, Laplace, and Slash. To construct simple comparison for 'nonnull' models, we have calculated averages of bias, standard error and mean square error of five models for each situation. The results have been represented in different tables and boxplots. We have repeated the same work with real world population.

4.1 Summary of results

Summary of results is as follows:

- Table 3 presents the name of the estimators that give the minimum value of different criteria described below the table. We see that the estimators r_{11} , r_{13} , and r_{15} have given 53 times, 24 times, and 18 times the minimum value respectively. On other hand, only 12 times the estimator r_1 has given the minimum value. That is, most of the situation the robust estimators have given the minimum value.
- In terms of standard error the performances of nonrobust estimators are very poor. The nonrobust estimator r_1 has given the minimum standard error only 6 times (null cases) out of 54 times (6 models, 3 correlation coefficient, ρ and 3 sizes) but the robust estimators r_{11} , r_{15} and r_{13} , have given 42, 5 and 1 times the minimum standard error respectively. For simplification when we consider ‘null’ and ‘nonnull’ model, we see that the estimator r_{11} has given 14 times out of 18 the minimum standard error where as r_1 offers us the minimum value only one time. With respect to other criteria related to standard error such as minimum of maximum of standard error, minimum of minimum of standard error etc the robust estimators, r_{11} has shown the best performance. In every respect of standard error, which we consider, undoubtedly r_{11} is the winner.
- We see that robust estimators has given the smallest bias 38 times out of 54. More over, only the estimator r_{13} has given 27 times minimum bias but r_1 only 10 times, and r_{11} 3 times. When we look into the Table no 3, we find that with respect to different criteria related to bias r_{13} is undoubtedly the best performer for nonnull models, r_{11} is better than r_1 , but for null models r_1 is the best. If we consider bias and standard error jointly we find r_{11} as the champion estimator. The normal score estimator, r_{11} have given alone 29 times the minimum mean square error out of 54 times and the robust estimators, in total, 45 times. On the other hand, Pearson estimator has given only 3 times the minimum mean square error. When we examine the aggregate performances, we see that the estimator r_1 has given three times minimum MSE, all for null model. On the contrary, the robust estimators r_{11} have given, in total, 9 times the minimum MSE, 7 times for ‘nonnull’ models.
- We have seen that robust estimators have given 48 times out of 54 the minimum length of 90% percentile interval and the nonrobust estimator r_1 has given only 6 times the minimum length. Only the estimator, r_{11} has given 39 times the minimum length of 90% percentile interval and the estimator, r_{15} 9 times. From Table 3 we find that out of 18 times, 11 times r_{11} has performed best and r_{15} for the rest of the cases.

- We have taken random sample of size 10, 20 and 50 from the Japanese data. the Table 4 offers us bias, standard error and MSE of each of the estimators for each sample. We have observed that the robust estimators have given minimum values 7 times out of 9 times. The estimator, r_{11} has given the minimum value of standard error all times (3), the minimum value of MSE for two times out of three and the rest has been provided by r_1 . We can infer that r_{11} has become the winner for the set of real data also.

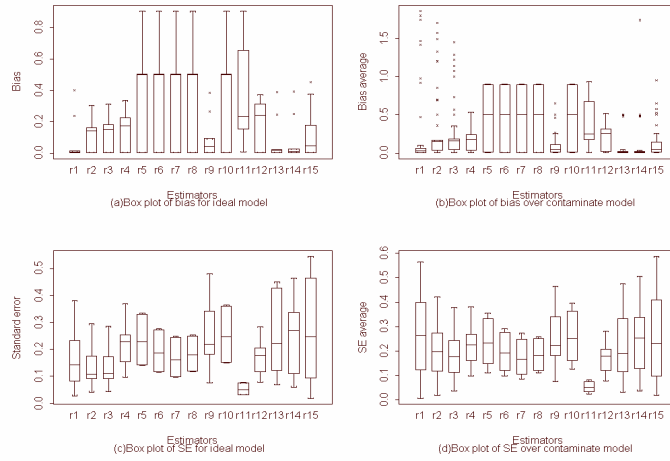


Figure 1: Boxplot of all estimates for simulation bias and standard error

The Fig 1(a) and 1(b) present the variation of bias in ‘null’ and ‘nonnull’ models and the Fig 1(c) and 1(d), the variation of standard error of different estimators. We see that the variation of robust estimator is, in general, smaller than that of nonrobust estimators. In case of standard the normal score estimator has given the minimum value in both ‘null’ and ‘nonnull’ models.

Table 3: Estimate with minimum bias, standard error, and length (w.r.t different criteria) for different models

| Criteria n | $\rho = 0$ | | | $\rho = .5$ | | | $\rho = .9$ | | |
|------------|------------|----------|----------|-------------|----------|----------|-------------|----------|----------|
| | 10 | 20 | 50 | 10 | 20 | 50 | 10 | 20 | 50 |
| Min(Bi) | r_2 | r_3 | r_7 | r_{15} | r_1 | r_1 | r_{15} | r_1 | r_{12} |
| Min(Ba) | r_{13} | r_4 | r_{13} | r_{13} | r_{13} | r_{13} | r_{11} | r_{13} | r_{15} |
| Min(Bm) | r_{13} | r_4 | r_{13} | r_{11} | r_{13} | r_{14} | r_{13} | r_{13} | r_{13} |
| Min(Bmi) | r_{14} | r_6 | r_1 | r_{13} | r_{13} | r_9 | r_{13} | r_{13} | r_{13} |
| Min(IMVB) | r_{13} | r_{11} | r_{13} | r_{11} | r_{13} | r_{14} | r_{13} | r_{13} | r_{13} |
| Min(Sei) | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_1 | r_{15} |
| Min(Sea) | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{15} | r_{15} |
| Min(Sem) | r_{11} | r_1 | r_1 | r_{11} | r_{11} | r_{11} | r_{11} | r_{15} | r_{15} |
| Min(Semi) | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_1 | r_1 |
| Min(IMVS) | r_{11} | r_4 | r_4 | r_{11} | r_9 | r_{13} | r_{11} | r_4 | r_{15} |
| Min(MSEi) | r_7 | r_7 | r_7 | r_{11} | r_{10} | r_1 | r_1 | r_{11} | r_1 |
| Min(MSEa) | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{15} | r_{15} |
| Min(Li) | r_{15} | r_{11} | r_{11} | r_{11} | r_{15} | r_{11} | r_{15} | r_{15} | r_{15} |
| Min(La) | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{11} | r_{15} | r_{15} |

(Bi=Bias of ideal model, Ba=Bias averaged over contaminated model, Bm = Bias maximum over contaminated model, Bmi = Bias minimum over contaminated model, IMVB = Inter Model Variation of Bias over contaminated model, Sei= Standard error of ideal model, Sea= Standard error of averaged over contaminated model, Sem = Standard error maximum over contaminated model, Semi = Standard error minimum over contaminated model, IMVS = Inter Model Variation of Standard error over contaminated model, MSEi = Means square error of ideal model, MSEa = Means square error of averaged over contaminated model, Li = Length of ideal model, La = Length of averaged over contaminated model).

Table 4: Bias, standard error and mean square error using real data ($\rho = .7052$)

| Estimates | Bias | | | S.E | | | M.S.E | | |
|-----------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| | n=10 | n=20 | n=50 | n=10 | n=20 | n=50 | n=10 | n=20 | n=50 |
| r_1 | 0.0136 | 0.0037 | 0.0028 | 0.1899 | 0.1273 | 0.0788 | 0.0362 | 0.0162 | 0.0062 |
| r_2 | 0.1669 | 0.1616 | 0.1622 | 0.1768 | 0.1183 | 0.0735 | 0.0591 | 0.0401 | 0.0317 |
| r_3 | 0.1857 | 0.1716 | 0.1667 | 0.1743 | 0.1179 | 0.0734 | 0.0649 | 0.0433 | 0.0332 |
| r_4 | 0.2355 | 0.2064 | 0.1786 | 0.3021 | 0.2085 | 0.1321 | 0.1467 | 0.0861 | 0.0493 |
| r_5 | 0.6720 | 0.6789 | 0.6789 | 0.3314 | 0.2275 | 0.1428 | 0.5613 | 0.5126 | 0.4813 |
| r_6 | 0.6793 | 0.6845 | 0.6840 | 0.2741 | 0.1867 | 0.1162 | 0.5366 | 0.5033 | 0.4813 |
| r_7 | 0.6743 | 0.6790 | 0.6791 | 0.2478 | 0.1604 | 0.0975 | 0.5161 | 0.4868 | 0.4706 |
| r_8 | 0.6846 | 0.6871 | 0.6856 | 0.2510 | 0.1805 | 0.1170 | 0.5317 | 0.5047 | 0.4838 |
| r_9 | 0.0424 | 0.0175 | 0.0434 | 0.3203 | 0.2168 | 0.1328 | 0.1044 | 0.0473 | 0.0195 |
| r_{10} | 0.6599 | 0.6653 | 0.6646 | 0.3621 | 0.2441 | 0.1513 | 0.5666 | 0.5022 | 0.4646 |
| r_{11} | 0.1995 | 0.0383 | 0.0534 | 0.0321 | 0.0756 | 0.0498 | 0.0408 | 0.0072 | 0.0053 |
| r_{12} | 0.3334 | 0.3137 | 0.3000 | 0.2493 | 0.1635 | 0.0987 | 0.1733 | 0.1251 | 0.0997 |
| r_{13} | 0.0201 | 0.0029 | 0.0049 | 0.2590 | 0.1655 | 0.1008 | 0.0675 | 0.0274 | 0.0102 |
| r_{14} | 0.0234 | 0.0037 | 0.0022 | 0.2745 | 0.1318 | 0.0804 | 0.0759 | 0.0174 | 0.0065 |
| r_{15} | 0.0601 | 0.0974 | 0.1174 | 0.2562 | 0.1479 | 0.0796 | 0.0692 | 0.0314 | 0.0201 |

5 Conclusion

So far we know, nobody has yet considered all the fifteen estimators at a time like us. Researchers (Devlin et al., 1975; Rodgers and Nicewander, 1988 and Gideon, 1998) considered only a group of estimators for comparison. For instance, Devlin et al. (1975) proposed graphical methods that can detect observations that may unduly influence the sample correlation coefficient and developed robust estimator of correlation coefficient. But they compared only seven estimators of correlation coefficient. In this study, we have considered fifteen estimators of correlation coefficients available in the literature. Over all, we have found that the class of robust estimators have performed better than the class of nonrobust estimators. Especially the normal score estimator, r_{11} has given the best performance among all estimators.

References

- [1] Blomqvist, N. (1950). On a Measure of Dependence Between Two Random Variables, *Annals of Mathematical Statistics*, **21**, 593-600.
- [2] Devlin, J. S, Gnanadesikan, R. and Kettenring, R. J. (1975). Robust estimation and outlier detection with correlation coefficients, *Biometrika*. **72**, 531-545.
- [3] Duncan, G. T. and Layard, M. W. J. (1973). A Monte-Carlo study of asymptotically robust tests for correlation coefficients, *Biometrika*, **60**, 551-559.
- [4] Fieller, E. C., Hartley, H. O. and Pearson, E. S. (1957). Test for rank correlation coefficients I, *Biometrika* **44**, 470-481.
- [5] Fieller, and Pearson, E. S. (1961). *Test for rank correlation coefficients II*, *Biometrika* **48**, 29-40.
- [6] Fisher, R. A. Yates, F. (1938). *Statistical tables for biological, agricultural and medical research*, 3rd ed. London: Oliver & Boyd, 2627.
- [7] Gideon, R. A. and Hollister, R. A. (1987). A Rank Correlation Coefficient Resistant to Outliers, *Journal of the American Statistical Assoc*, **82**, 656-666.
- [8] Gideon, R. A. (1998). A Generalized Interpretation of Pearson's r , <http://www.math.umt.edu/gideon>.
- [9] Gini, C. (1914). 'L' Ammontare della Composizioni della Ricchezza della Nazioni, Bocca, Torino.
- [10] Gnanadesikan, R. and Kettenring, J. R. (1972). Robust estimates, residuals, and outlier detection with multiresponse data, *Biometrics*, **28**, 81-124.

- [11] Hossain, M. G., Lestrel, P. E. and Ohtsuki, F. (2005). Secular changes in head dimensions of Japanese adult male students over eight decades, *Journal of Comparative Human Biology*, **55**, 239-250.
- [12] Kendall, M. G. (1970). *Rank Correlation Methods*. 4th ed. Chaise Griffin, London.
- [13] Kendall, M. G. (1938). A New Measure of Rank Correlation, *Biometrika*, **30**, 81-93.
- [14] Pearson, K. (1896). Mathematical Contributions to the Theory of Evolution, III, Regression, Heredity and Panmixia, *Philosophical Transactions of the Royal Society of London*, **187**, 253-318.
- [15] Rodgers, J. L. and Nicewander, W. A. (1988). Thirteen Ways to Look at the Correlation Coefficient, *The American Statistician*, **42**, **1**, 59-66.
- [16] Spearman, C. (1904). The Proof and Measurement of Association Between Two Things, *American Journal of Psychiatry* **15**, 72-101.
- [17] Sheppard, W. F. (1899). On the Application of the theory of error to cases of normal Distribution and Normal correlation, *Phil. Trans. R. Soc. A*. **192**, 101-167.