

Estimation of the simple correlation coefficient

GWOWEN SHIEH

National Chiao Tung University, Hsinchu, Taiwan

This article investigates some unfamiliar properties of the Pearson product–moment correlation coefficient for the estimation of simple correlation coefficient. Although Pearson's r is biased, except for limited situations, and the minimum variance unbiased estimator has been proposed in the literature, researchers routinely employ the sample correlation coefficient in their practical applications, because of its simplicity and popularity. In order to support such practice, this study examines the mean squared errors of r and several prominent formulas. The results reveal specific situations in which the sample correlation coefficient performs better than the unbiased and nearly unbiased estimators, facilitating recommendation of r as an effect size index for the strength of linear association between two variables. In addition, related issues of estimating the squared simple correlation coefficient are also considered.

In order to reform statistical practices, Wilkinson and the American Psychological Association Task Force on Statistical Inference (1999), the *Publication Manual of the American Psychological Association* (2001), and the American Educational Research Association Task Force on Reporting of Research Methods (2006) recommended the reporting of effect sizes in all empirical social science research. Accordingly, numerous practical guidelines and suggestions for selecting, calculating, and interpreting effect size indices for various types of statistical analyses have been provided in the literature, such as Alhija and Levy (2009), Breaugh (2003), Durlak (2009), Ferguson (2009), Grissom and Kim (2005), Huberty (2002), Kline (2004), Richardson (1996), Rosenthal, Rosnow, and Rubin (2000), Rosnow and Rosenthal (2003), and Vacha-Haase and Thompson (2004). In particular, Ferguson suggested that effect sizes can be categorized into four general classes: (1) group difference, (2) strength of association, (3) corrected estimates, and (4) risk estimates. Notably, the Pearson product–moment correlation coefficient, or simple correlation coefficient r , is the most commonly used strength-of-association measure in applied research across virtually all disciplines of social sciences. The correlation summarizes the magnitude and direction of linear relationship between two variables. It is generally known that a value of r close to zero suggests that the linear association is weak; however, high correlation does not imply causality.

Although the fundamental results and associated usages of r are described in most introductory textbooks of statistics and quantitative methods, it is not well understood that the underlying probability distribution function of r is complicated in form, under the classical assumption that the two variables follow a **bivariate normal distribution**. The complexity incurs continuous investigations to give various expressions, approximations, and

computing algorithms for examining statistical features of the sample correlation coefficient. For theoretical developments in statistical literature, Johnson, Kotz, and Balakrishnan (1995, chap. 32) and Stuart and Ord (1994, chap. 16) contain comprehensive discussions and technical details. On the other hand, Bobko (2001) and Cohen, Cohen, West, and Aiken (2003) emphasize operational guidelines and practical implications in the behavioral and social sciences.

The purposes of this article are to explicate the intrinsic issues surrounding point estimators of strength of association and to support the use of Pearson's r as an effect size measure in the light of new empirical results based on direct integration and computing techniques. Despite its computational ease and widespread usage, r is not an unbiased estimator of population simple correlation coefficient ρ , except for the special situations of $\rho = -1, 0$, and 1 . It appears that standard textbooks rarely mention this undesirable nature of r , whereas the embedded unbiasedness associated with a sample mean or a sample variance is always emphasized. Notably, Olkin and Pratt (1958) derived the unique minimum variance unbiased estimator of ρ , but unfortunately, the computational complexity of the resulting expression is overwhelming, particularly in the absence of appropriate computer software. This has restricted acceptance of their unbiased formula and may contribute to the continual application of the sample correlation coefficient at the expense of its potentially detrimental consequences. Nonetheless, unbiasedness is certainly not the only criterion of theoretical importance. Another consideration related to the statistical properties of a point estimator deals with the concept of mean squared error (MSE). There is no study to our knowledge that investigates the MSE of r through intensive numerical integration, although some limited simulated results of r 's

G. Shieh, gwshieh@mail.nctu.edu.tw

mean and standard deviation were presented in Zimmerman, Zumbo, and Williams (2003). Therefore, it is vital to provide a unified and rigorous justification for the *MSE* performance of the sample correlation coefficient, along with other prominent estimators of ρ . In this research, the statistical properties of r and competing formulas are examined both numerically and graphically to provide a clear understanding of their advantages and disadvantages in evaluating the extent of the linear relationship between two variables.

As is well known, the squared simple correlation coefficient r^2 can be viewed as a special case of the squared multiple correlation coefficient or coefficient of determination R^2 in the context of multiple linear regression models. In this case, R^2 denotes the percentage of the total variation of the criterion that is accounted for by the relationship with the predictors. The problem of estimating the squared multiple correlation coefficient has been studied by Raju, Bilgic, Edwards, and Fler (1997, 1999), Shieh (2008), and Yin and Fan (2001). Thus, the square of simple correlation coefficient r^2 not only has a distinct interpretation as a percentage measure of variance that is accounted for, but also possesses completely different properties. In contrast to the extensive results related to R^2 , the investigation of r^2 has received little attention, although a notable exception is Wang and Thompson (2007). However, their results are confined to simulated mean biases of some corrected formulas of r^2 . Instead, detailed numerical study is conducted here to assess the exact bias and *MSE* of r^2 and several well-known estimators. The present exposition helps to clarify the unique and contrasting behavior of r and r^2 and to choose an appropriate effect size measure within the framework of correlation analysis.

Estimation of the Simple Correlation Coefficient

Let $(X_1, Y_1), \dots, (X_N, Y_N)$ be independent and identically distributed with bivariate normal distribution with means μ_X, μ_Y , variances σ_X^2, σ_Y^2 , and correlation ρ ($|\rho| < 1$). The sample correlation coefficient r is the sample covariance divided by the product of sample standard deviations

$$r = \frac{S_{XY}}{S_X S_Y},$$

where

$$S_{XY} = \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y}) / (N - 1),$$

$$S_X^2 = \sum_{i=1}^N (X_i - \bar{X})^2 / (N - 1),$$

and

$$S_Y^2 = \sum_{i=1}^N (Y_i - \bar{Y})^2 / (N - 1),$$

with sample means

$$\bar{X} = \sum_{i=1}^N X_i / N$$

and

$$\bar{Y} = \sum_{i=1}^N Y_i / N.$$

The exact density function of r was originally obtained by Fisher (1915), following a geometrical argument. For ease of exposition, the fundamental results associated with r are presented in the Appendix. Accordingly, the probability density function of r given in Equation A1 is extremely complex and does not admit a simplified expression, except in some special cases, and considerable attention has been devoted to the construction of useful approximations. For most practical purposes, inferences about population correlation coefficient ρ are based on the famous Fisher's (1921) Z transformation, which has an approximately normal distribution irrespective of ρ and N :

$$Z = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right) \sim N(\mu_Z, \sigma_Z^2),$$

where

$$\mu_Z = \frac{1}{2} \ln \left(\frac{1+\rho}{1-\rho} \right)$$

and

$$\sigma_Z^2 = \frac{1}{N-3}.$$

Alternatively, exact inferential procedures are available, and interested readers are referred to a recent article by Shieh (2006). Here, we focus on the point estimation problem of ρ under the ultimate notion of choosing a profound correlational effect size measure for the strength of association between the two variables X and Y . It can be seen from Equation A3 that r is a biased estimator of ρ , and the mean and variance of r can be approximated by

$$E[r] \doteq \rho - \frac{\rho(1-\rho^2)}{2(N-1)} \quad \text{and} \quad \text{Var}[r] \doteq \frac{(1-\rho^2)^2}{(N-1)}.$$

It follows that $E[r] < \rho$ or $E[r] > \rho$ if $\rho > 0$ or $\rho < 0$. Hence, on the average, r will underestimate ρ for positive ρ , and it tends to overestimate ρ when ρ is less than 0. In contrast, Olkin and Pratt (1958) have derived the unique minimum variance unbiased estimator (MVUE) $\hat{\rho}_U$ of ρ as given in Equation A2. Although the unbiasedness viewpoint is of theoretical meaning, the computation of $\hat{\rho}_U$ is considerably cumbersome for practical use. Thus, they suggested the approximation

$$\hat{\rho}_{\text{OPA}}(r) = r \left\{ 1 + \frac{1-r^2}{2(N-4)} \right\}.$$

For comparative purposes, three additional different approximations can be obtained from the expansion of $\hat{\rho}_U$:

$$\hat{\rho}_{\text{OP1}}(r) = r \left\{ 1 + \frac{1-r^2}{2(N-2)} \right\},$$

$$\hat{\rho}_{\text{OP2}}(r) = r \left\{ 1 + \frac{1-r^2}{2(N-2)} + \frac{9(1-r^2)^2}{8N(N-2)} \right\},$$

Table 1
Bias for Estimators of ρ With $N = 20$

ρ	r	$\hat{\rho}_{OP1}$	$\hat{\rho}_{OPA}$	$\hat{\rho}_{OP2}$	$\hat{\rho}_{OP5}$	$\hat{\rho}_M$
.00	.000000	.000000	.000000	.000000	.000000	.000000
.05	-.001295	-.000138	.000007	-.000025	-.000001	-.002386
.10	-.002573	-.000272	.000015	-.000049	-.000001	-.004742
.15	-.003816	-.000400	.000027	-.000071	-.000002	-.007038
.20	-.005006	-.000518	.000043	-.000091	-.000002	-.009243
.25	-.006126	-.000623	.000065	-.000108	-.000003	-.011326
.30	-.007156	-.000711	.000095	-.000121	-.000003	-.013254
.35	-.008080	-.000781	.000131	-.000130	-.000003	-.014994
.40	-.008876	-.000830	.000176	-.000134	-.000003	-.016511
.45	-.009526	-.000855	.000228	-.000134	-.000002	-.017768
.50	-.010007	-.000856	.000287	-.000128	-.000002	-.018725
.55	-.010300	-.000832	.000351	-.000119	-.000002	-.019341
.60	-.010379	-.000783	.000417	-.000105	-.000001	-.019569
.65	-.010222	-.000709	.000480	-.000089	-.000001	-.019361
.70	-.009801	-.000614	.000535	-.000070	-.000001	-.018660
.75	-.009090	-.000500	.000574	-.000051	-.000000	-.017407
.80	-.008058	-.000375	.000586	-.000033	-.000000	-.015531
.85	-.006671	-.000247	.000557	-.000018	-.000000	-.012955
.90	-.004894	-.000128	.000468	-.000007	-.000000	-.009584
.95	-.002686	-.000038	.000293	-.000001	-.000000	-.005311
.99	-.000578	-.000002	.000070	.000000	.000000	-.001151

and

$$\hat{\rho}_{OP5}(r) = r \left\{ 1 + \sum_{k=1}^5 \frac{\left[\Gamma\left(\frac{1}{2} + k\right) \right]^2 \Gamma\left(\frac{N-2}{2}\right)}{\left[\Gamma\left(\frac{1}{2}\right) \right]^2 \Gamma\left(\frac{N-2}{2} + k\right)} \cdot \frac{(1-r^2)^k}{k!} \right\}.$$

Moreover, it is noteworthy from the prescribed bivariate normal distribution of X and Y that the maximum likelihood estimators (MLEs) of μ_X , μ_Y , σ_X^2 , σ_Y^2 , and ρ are \bar{X} , \bar{Y} , $\{(N-1)/N\}S_X^2$, $\{(N-1)/N\}S_Y^2$, and r , respectively. Thus, Pearson's r is the joint MLE of ρ . However, this is different from the marginal MLE $\hat{\rho}_M$ based on the probability density function $f(r)$ given in Equation A1. In this case, there is no explicit closed form expression for $\hat{\rho}_M$, although it was shown in Fisher (1915) that

$$\hat{\rho}_M \doteq r \left\{ 1 + \frac{1-r^2}{2N} \right\}.$$

Note that r and $\hat{\rho}_M$ are asymptotically equivalent and yield similar estimation performance in finite samples.

In addition to the unbiasedness consideration, MSE is another useful performance criterion obtained by incorporating the bias (accuracy) and variability (precision) of an estimator. Specifically, the MSE of an estimator $\hat{\rho}$ of ρ is the function

$$MSE(\hat{\rho}, \rho) = E[(\hat{\rho} - \rho)^2] = \{\text{Bias}(\hat{\rho}, \rho)\}^2 + \text{Var}[\hat{\rho}],$$

where $\text{Bias}(\hat{\rho}, \rho) = E[\hat{\rho}] - \rho$. It is possible for a biased estimator $\hat{\rho}$ that a trade-off occurs between bias $\text{Bias}(\hat{\rho}, \rho)$ and variance $\text{Var}[\hat{\rho}]$ in such a way that a larger decrease in variance can be obtained for a small increase in bias, resulting in an improvement in $MSE(\hat{\rho}, \rho)$. This phenomenon is demonstrated in the following numerical investigation of Pearson's r and $\hat{\rho}_M$.

Due to the complexity of the estimation problem, analytical justifications of the theoretical discrepan-

cies of competing estimators are generally not feasible. Thus, a special-purpose computer program has been developed for this study to perform numerical integration with respect to the probability density distribution of r . The exact properties for the estimators of r , $\hat{\rho}_{OP1}$, $\hat{\rho}_{OPA}$, $\hat{\rho}_{OP2}$, $\hat{\rho}_{OP5}$, and $\hat{\rho}_M$ are examined. The computed exact biases for $\rho = .00$ to $.95$, with an increment of $.05$, and $.99$ are presented in Tables 1–3 for $N = 20, 50$, and 100 , respectively. In addition, the corresponding root-mean squared errors ($RMSE = MSE^{1/2}$) are summarized in Tables 4–6 for $N = 20, 50$, and 100 , respectively. Due to the distinct distributional property of r , the corresponding results for negative simple correlation coefficient are not reported here, because the bias associated with $-\rho$ ($\rho > 0$) has a sign opposite to that of ρ —that is, $\text{Bias}(\hat{\rho}, -\rho) = -\text{Bias}(\hat{\rho}, \rho)$. However, the MSE and $RMSE$ are identical for both cases of $-\rho$ and ρ , $MSE(\hat{\rho}, -\rho) = MSE(\hat{\rho}, \rho)$ and $RMSE(\hat{\rho}, -\rho) = RMSE(\hat{\rho}, \rho)$, where $\rho > 0$. For a concise visualization of these results, the exact bias and the $RMSE$ results of r are plotted in Figures 1 and 2, respectively. Overall, the bias and $RMSE$ performance of these estimators improve with increased sample size.

It can be readily seen from Tables 1–3 that both Pearson's r and marginal MLE $\hat{\rho}_M$ underestimate ρ except when $\rho = 0$. However, r performs consistently better than $\hat{\rho}_M$ because $\text{Bias}(\hat{\rho}_M, \rho) < \text{Bias}(r, \rho) < 0$ for $\rho > 0$. As was expected, the other four estimators ($\hat{\rho}_{OP1}$, $\hat{\rho}_{OPA}$, $\hat{\rho}_{OP2}$, $\hat{\rho}_{OP5}$) are nearly unbiased. Since its bias is almost negligible, the five-term approximation $\hat{\rho}_{OP5}$ is basically equivalent to the MVUE $\hat{\rho}_U$. Nonetheless, an appealing feature of the approximate formula $\hat{\rho}_{OPA}$ is that it enjoys both overall accuracy and computational ease.

The computed $RMSE$ results listed in Tables 4–6 reveal complex and unfamiliar relations among the competing formulas. First, the exact MSE performance of the practically unbiased estimator $\hat{\rho}_{OP5}$ and nearly unbiased formula $\hat{\rho}_{OPA}$ cross each other, showing that each

Table 2
Bias for Estimators of ρ With $N = 50$

ρ	r	$\hat{\rho}_{OP1}$	$\hat{\rho}_{OPA}$	$\hat{\rho}_{OP2}$	$\hat{\rho}_{OP5}$	$\hat{\rho}_M$
.00	.000000	.000000	.000000	.000000	.000000	.000000
.05	-.000506	-.000022	-.000001	-.000002	.000000	-.000980
.10	-.001005	-.000044	-.000002	-.000003	.000000	-.001945
.15	-.001490	-.000064	-.000002	-.000005	.000000	-.002883
.20	-.001953	-.000083	-.000002	-.000006	.000000	-.003780
.25	-.002386	-.000099	.000000	-.000007	.000000	-.004621
.30	-.002782	-.000113	.000003	-.000008	.000000	-.005393
.35	-.003135	-.000123	.000008	-.000009	.000000	-.006082
.40	-.003435	-.000129	.000014	-.000009	.000000	-.006672
.45	-.003676	-.000132	.000022	-.000009	.000000	-.007150
.50	-.003850	-.000131	.000031	-.000008	.000000	-.007498
.55	-.003948	-.000126	.000041	-.000007	.000000	-.007702
.60	-.003962	-.000116	.000051	-.000006	.000000	-.007744
.65	-.003884	-.000104	.000061	-.000005	.000000	-.007607
.70	-.003705	-.000088	.000069	-.000004	.000000	-.007274
.75	-.003417	-.000070	.000075	-.000003	.000000	-.006725
.80	-.003011	-.000052	.000077	-.000002	.000000	-.005940
.85	-.002475	-.000033	.000073	-.000001	.000000	-.004899
.90	-.001802	-.000017	.000061	.000000	.000000	-.003578
.95	-.000981	-.000005	.000038	.000000	.000000	-.001954
.99	-.000209	.000000	.000009	.000000	.000000	-.000418

estimator is better only for certain combined configurations of ρ and N . Specifically, when $N = 20$, the order of MSE is

$$MSE(\hat{\rho}_{OPA}, \rho) > MSE(\hat{\rho}_{OP5}, \rho) \text{ for } \rho \leq .80$$

and

$$MSE(\hat{\rho}_{OPA}, \rho) < MSE(\hat{\rho}_{OP5}, \rho) \text{ for } \rho > .80.$$

On the other hand, when $N = 50$ and 100, the resulting behavior is

$$MSE(\hat{\rho}_{OPA}, \rho) \leq MSE(\hat{\rho}_{OP5}, \rho) \text{ for } \rho \leq .15 \text{ and } \rho \geq .80$$

and

$$MSE(\hat{\rho}_{OPA}, \rho) > MSE(\hat{\rho}_{OP5}, \rho) \text{ for } .15 < \rho < .80.$$

Second, it is important to note that the interrelationships between r , $\hat{\rho}_M$, and $\hat{\rho}_{OPA}$ are as follows:

$$MSE(\hat{\rho}_M, \rho) < MSE(r, \rho) < MSE(\hat{\rho}_{OPA}, \rho) \\ \text{for } \rho \leq .50 \text{ when } N = 20, 50, \text{ and } 100,$$

$$MSE(\hat{\rho}_M, \rho) < MSE(r, \rho) < MSE(\hat{\rho}_{OPA}, \rho) \\ \text{for } \rho = .55 \text{ when } N = 20,$$

$$MSE(r, \rho) < MSE(\hat{\rho}_M, \rho) < MSE(\hat{\rho}_{OPA}, \rho) \\ \text{for } \rho = .55 \text{ when } N = 50,$$

$$MSE(r, \rho) < MSE(\hat{\rho}_{OPA}, \rho) < MSE(\hat{\rho}_M, \rho) \\ \text{for } \rho = .55 \text{ when } N = 100,$$

$$MSE(r, \rho) < MSE(\hat{\rho}_{OPA}, \rho) < MSE(\hat{\rho}_M, \rho)$$

Table 3
Bias for Estimators of ρ With $N = 100$

ρ	r	$\hat{\rho}_{OP1}$	$\hat{\rho}_{OPA}$	$\hat{\rho}_{OP2}$	$\hat{\rho}_{OP5}$	$\hat{\rho}_M$
.00	.000000	.000000	.000000	.000000	.000000	.000000
.05	-.000251	-.000006	.000000	.000000	.000000	-.000496
.10	-.000499	-.000011	-.000001	.000000	.000000	-.000984
.15	-.000739	-.000016	-.000001	-.000001	.000000	-.001456
.20	-.000968	-.000021	-.000001	-.000001	.000000	-.001906
.25	-.001182	-.000025	-.000001	-.000001	.000000	-.002327
.30	-.001378	-.000028	.000000	-.000001	.000000	-.002713
.35	-.001551	-.000031	.000001	-.000001	.000000	-.003056
.40	-.001699	-.000032	.000003	-.000001	.000000	-.003348
.45	-.001816	-.000033	.000005	-.000001	.000000	-.003582
.50	-.001900	-.000032	.000007	-.000001	.000000	-.003750
.55	-.001946	-.000031	.000009	-.000001	.000000	-.003844
.60	-.001950	-.000028	.000012	-.000001	.000000	-.003856
.65	-.001909	-.000025	.000014	-.000001	.000000	-.003779
.70	-.001818	-.000021	.000016	.000000	.000000	-.003604
.75	-.001674	-.000017	.000018	.000000	.000000	-.003322
.80	-.001472	-.000012	.000018	.000000	.000000	-.002925
.85	-.001208	-.000008	.000017	.000000	.000000	-.002404
.90	-.000878	-.000004	.000014	.000000	.000000	-.001749
.95	-.000476	-.000001	.000009	.000000	.000000	-.000951
.99	-.000102	.000000	.000002	.000000	.000000	-.000203

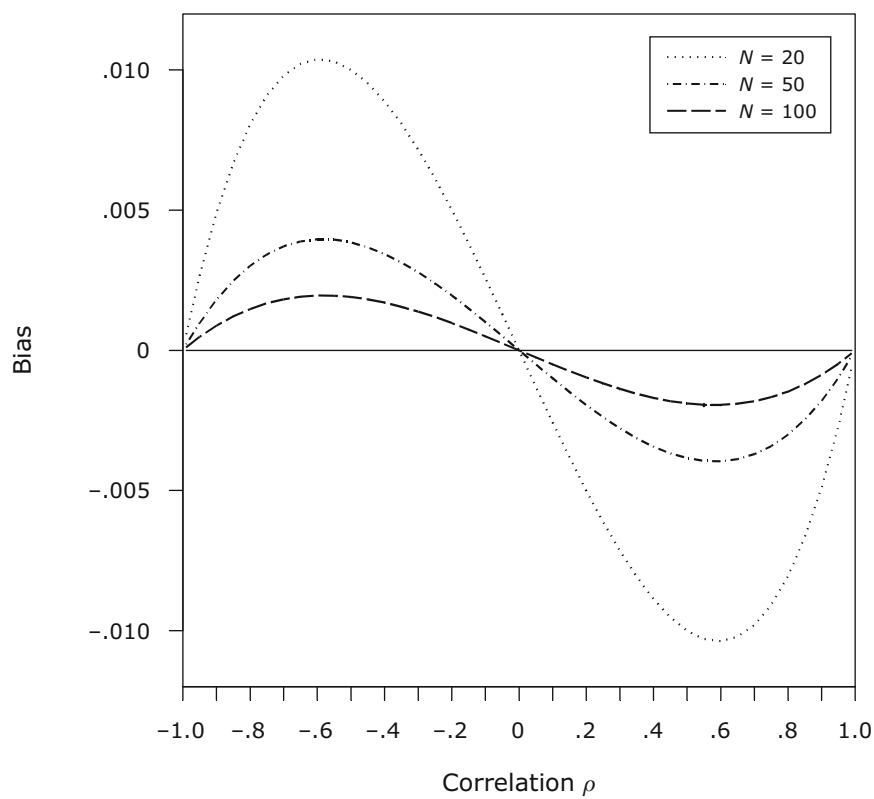
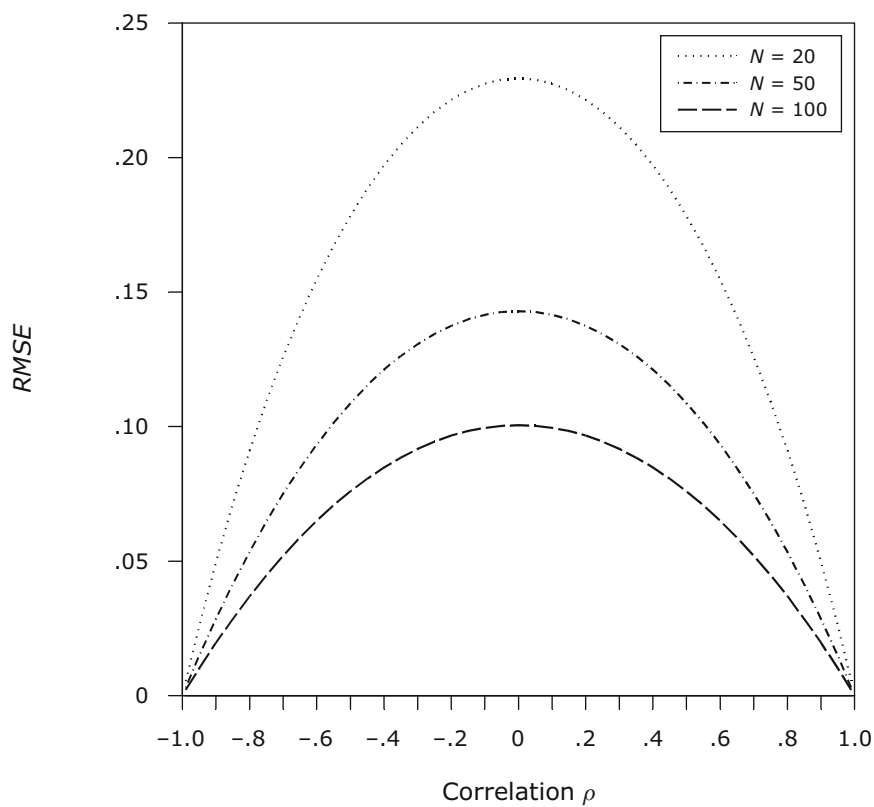
Figure 1. The bias of r .Figure 2. The root-mean squared error (RMSE) of r .

Table 4
Root-Mean Squared Error for Estimators of ρ With $N = 20$

ρ	r	$\hat{\rho}_{OP1}$	$\hat{\rho}_{OPA}$	$\hat{\rho}_{OP2}$	$\hat{\rho}_{OP5}$	$\hat{\rho}_M$
.00	.229416	.234879	.235562	.235414	.235529	.224268
.05	.228920	.234337	.235015	.234865	.234978	.223820
.10	.227431	.232711	.233373	.233219	.233327	.222474
.15	.224944	.229997	.230632	.230473	.230570	.220223
.20	.221449	.226188	.226787	.226620	.226705	.217053
.25	.216933	.221277	.221830	.221653	.221723	.212948
.30	.211379	.215250	.215749	.215564	.215616	.207883
.35	.204764	.208096	.208531	.208338	.208372	.201828
.40	.197064	.199795	.200161	.199963	.199978	.194748
.45	.188245	.190329	.190619	.190421	.190418	.186597
.50	.178271	.179675	.179884	.179691	.179672	.177324
.55	.167097	.167805	.167932	.167751	.167720	.166866
.60	.154672	.154690	.154735	.154575	.154534	.155147
.65	.140935	.140296	.140262	.140132	.140087	.142078
.70	.125816	.124584	.124478	.124389	.124344	.127555
.75	.109230	.107513	.107346	.107307	.107267	.111447
.80	.091078	.089034	.088825	.088842	.088810	.093598
.85	.071242	.069099	.068872	.068943	.068923	.073816
.90	.049578	.047651	.047443	.047552	.047543	.051855
.95	.025906	.024635	.024496	.024600	.024598	.027396
.99	.005373	.005059	.005024	.005057	.005057	.005741

for $\rho = .60$ when $N = 20$,

$$MSE(\hat{\rho}_{OPA}, \rho) < MSE(r, \rho) < MSE(\hat{\rho}_M, \rho)$$

for $\rho = .60$ when $N = 50$ and 100 ,

and

$$MSE(\hat{\rho}_{OPA}, \rho) < MSE(r, \rho) < MSE(\hat{\rho}_M, \rho)$$

for $\rho > .60$ when $N = 20, 50$, and 100 .

The corresponding situations among r , $\hat{\rho}_M$, and $\hat{\rho}_{OP5}$ are identical to those of r , $\hat{\rho}_M$, and $\hat{\rho}_{OPA}$ just described, with the only exception in the case of

$$MSE(\hat{\rho}_{OP5}, \rho) < MSE(r, \rho) < MSE(\hat{\rho}_M, \rho)$$

for $\rho = .60$ when $N = 20$.

Hence, despite the disadvantageous bias in the performance of r and $\hat{\rho}_M$, it is conceivable that they are not dominated by the unbiased or nearly unbiased estimators in terms of MSE . In view of the close behavior and computational requirement between r and $\hat{\rho}_M$, it is worthwhile to consider r , which yields similar results with less computation. Moreover, the prescribed results suggest that $MSE(r, \rho) < MSE(\hat{\rho}_{OPA}, \rho)$ for $\rho \leq .60$, and $MSE(r, \rho) > MSE(\hat{\rho}_{OPA}, \rho)$ for $\rho > .60$. It appears that no absolutely dominant answer is obtained with the exact MSE results, although more information is gathered about Pearson's r with respect to the other prominent estimators. The ultimate implication is that the sample correlation coefficient proves to be computationally and theoretically useful in estimating the strength of association for $|\rho| \leq .60$.

Table 5
Root-Mean Squared Error for Estimators of ρ With $N = 50$

ρ	r	$\hat{\rho}_{OP1}$	$\hat{\rho}_{OPA}$	$\hat{\rho}_{OP2}$	$\hat{\rho}_{OP5}$	$\hat{\rho}_M$
.00	.142857	.144258	.144319	.144317	.144322	.141489
.05	.142520	.143907	.143967	.143966	.143971	.141167
.10	.141508	.142855	.142914	.142911	.142915	.140200
.15	.139820	.141100	.141156	.141152	.141156	.138584
.20	.137453	.138642	.138694	.138688	.138691	.136315
.25	.134402	.135477	.135525	.135516	.135519	.133388
.30	.130664	.131604	.131646	.131635	.131636	.129795
.35	.126231	.127019	.127055	.127041	.127042	.125527
.40	.121097	.121718	.121747	.121731	.121731	.120572
.45	.115253	.115697	.115718	.115700	.115700	.114917
.50	.108688	.108950	.108964	.108945	.108944	.108547
.55	.101391	.101472	.101478	.101460	.101458	.101444
.60	.093349	.093256	.093255	.093238	.093236	.093586
.65	.084547	.084295	.084288	.084274	.084272	.084951
.70	.074969	.074582	.074569	.074559	.074557	.075513
.75	.064597	.064108	.064091	.064086	.064084	.065240
.80	.053409	.052865	.052845	.052846	.052845	.054099
.85	.041382	.040842	.040822	.040828	.040828	.042052
.90	.028491	.028031	.028014	.028023	.028023	.029054
.95	.014708	.014421	.014410	.014418	.014418	.015056
.99	.003017	.002949	.002947	.002949	.002949	.003099

Table 6
Root-Mean Squared Error for Estimators of ρ With $N = 100$

ρ	r	$\hat{\rho}_{OP1}$	$\hat{\rho}_{OPA}$	$\hat{\rho}_{OP2}$	$\hat{\rho}_{OP5}$	$\hat{\rho}_M$
.00	.100504	.101001	.101012	.101012	.101013	.100010
.05	.100260	.100752	.100762	.100763	.100763	.099773
.10	.099527	.100005	.100015	.100015	.100015	.099059
.15	.098305	.098758	.098768	.098768	.098768	.097866
.20	.096594	.097013	.097022	.097021	.097021	.096191
.25	.094390	.094767	.094775	.094774	.094774	.094032
.30	.091694	.092021	.092028	.092026	.092026	.091389
.35	.088501	.088773	.088779	.088776	.088776	.088258
.40	.084810	.085021	.085026	.085023	.085023	.084633
.45	.080618	.080765	.080768	.080765	.080765	.080510
.50	.075921	.076002	.076004	.076001	.076001	.075884
.55	.070714	.070731	.070732	.070728	.070728	.070746
.60	.064993	.064949	.064949	.064946	.064946	.065091
.65	.058754	.058655	.058653	.058651	.058651	.058910
.70	.051990	.051845	.051843	.051841	.051841	.052193
.75	.044695	.044517	.044514	.044513	.044513	.044930
.80	.036862	.036669	.036665	.036665	.036665	.037111
.85	.028485	.028296	.028293	.028294	.028294	.028722
.90	.019555	.019396	.019393	.019395	.019395	.019751
.95	.010063	.009965	.009963	.009965	.009965	.010183
.99	.002059	.002036	.002036	.002036	.002036	.002086

Estimation of the Squared Simple Correlation Coefficient

It is generally known that R^2 is a positively biased estimator of ρ^2 within the multiple regression framework. To correct such overestimation, several modified formulas have been suggested in the literature. Comprehensive discussions and comparisons can be found in the work of Raju et al. (1999), Shieh (2008), and Yin and Fan (2001). Since the squared simple correlation coefficient r^2 can be viewed as a special case of the coefficient of determination R^2 under causal consideration, it is of practical interest to extend the assessment to the exact performance of r^2 as an index of the population coefficient of determination ρ^2 .

As an estimator of ρ^2 , the expected value of r^2 , denoted by $E[r^2]$, is provided in Equation A4. But without a special computing algorithm, it is difficult to conceive the resulting magnitude from the analytical expression, except that $E[r^2] = 1/(N - 1)$ when $\rho^2 = 0$. Moreover, it is natural to assume that $E[r^2] > \rho^2$ under the common conception that $E[R^2] > \rho^2$. It is shown below that although r^2 tends to overestimate ρ^2 , it may be unbiased or negatively biased

for certain values of ρ and N . In this case, the so-called adjusted R^2 formula reduces to

$$\hat{\rho}_E^2(r^2) = 1 - \frac{N-1}{N-2}(1-r^2).$$

Also, the MVUE $\hat{\rho}_U^2$ derived by Olkin and Pratt (1958) is given in Equation A5. It should be noted from Equations A2 and A5 that $\hat{\rho}_U^2$ is not a square function of $\hat{\rho}_U$. Unfortunately, it appears that the desirable property of unbiasedness for $\hat{\rho}_U^2$ is outweighed by its computational complexity, just as $\hat{\rho}_U$ in the estimation of ρ . A useful alternative suggested by Pratt is

$$\hat{\rho}_{PA}^2(r^2) = 1 - \frac{N-3}{N-2}(1-r^2) \left\{ 1 + \frac{2(1-r^2)}{N-3.3} \right\}.$$

Furthermore, simplified approximations of $\hat{\rho}_{OP1}^2$, $\hat{\rho}_{OP2}^2$, and $\hat{\rho}_{OP5}^2$ can be obtained from the expansion of $\hat{\rho}_U^2$ as shown at the bottom of the page.

To delineate the disparate performance by estimators of ρ^2 , the exact bias and MSE of r^2 , $\hat{\rho}_E^2$, $\hat{\rho}_{OP1}^2$, $\hat{\rho}_{PA}^2$, $\hat{\rho}_{OP2}^2$,

$$\hat{\rho}_{OP1}^2(r^2) = 1 - \frac{N-3}{N-2}(1-r^2) \left\{ 1 + \frac{2(1-r^2)}{N} \right\},$$

$$\hat{\rho}_{OP2}^2(r^2) = 1 - \frac{N-3}{N-2}(1-r^2) \left\{ 1 + \frac{2(1-r^2)}{N} + \frac{8(1-r^2)^2}{N(N+2)} \right\},$$

and

$$\hat{\rho}_{OP5}^2(r^2) = 1 - \frac{N-3}{N-2}(1-r^2) \left\{ 1 + \sum_{k=1}^5 \frac{[\Gamma(1+k)]^2 \Gamma\left(\frac{N}{2}\right)}{\Gamma\left(\frac{N}{2} + k\right)} \cdot \frac{(1-r^2)^k}{k!} \right\}.$$

Table 7
Bias for Estimators of ρ^2 With $N = 20$

ρ	r^2	$\hat{\rho}_E^2$	$\hat{\rho}_{OP1}^2$	$\hat{\rho}_{PA}^2$	$\hat{\rho}_{OP2}^2$	$\hat{\rho}_{OP5}^2$
.00	.052632	.000000	.020050	.003212	.005230	.000257
.05	.052275	-.000238	.019926	.003157	.005191	.000254
.10	.051210	-.000945	.019557	.002996	.005073	.000246
.15	.049455	-.002103	.018949	.002733	.004880	.000233
.20	.047037	-.003683	.018117	.002379	.004618	.000215
.25	.043997	-.005642	.017076	.001948	.004295	.000194
.30	.040387	-.007925	.015851	.001456	.003919	.000171
.35	.036273	-.010462	.014468	.000924	.003502	.000145
.40	.031733	-.013171	.012959	.000375	.003057	.000120
.45	.026863	-.015950	.011360	-.000163	.002599	.000095
.50	.021773	-.018684	.009713	-.000664	.002141	.000072
.55	.016592	-.021236	.008061	-.001100	.001699	.000052
.60	.011469	-.023450	.006451	-.001440	.001288	.000035
.65	.006575	-.025144	.004932	-.001660	.000922	.000021
.70	.002108	-.026109	.003552	-.001739	.000611	.000012
.75	-.001704	-.026104	.002357	-.001662	.000365	.000006
.80	-.004597	-.024853	.001390	-.001430	.000187	.000002
.85	-.006266	-.022030	.000678	-.001065	.000075	.000001
.90	-.006351	-.017260	.000234	-.000620	.000019	.000000
.95	-.004432	-.010094	.000034	-.000202	.000002	.000000
.99	-.001115	-.002283	.000000	-.000010	.000000	.000000

and $\hat{\rho}_{OP5}^2$ are computed. With the same settings of ρ and N in the previous examination for simple correlation coefficient, the exact biases are presented in Tables 7–9, and *RMSEs* are summarized in Tables 10–12.

Regarding the accuracy results in Tables 7–9, the biases are smaller for large N with fixed value of ρ . Specifically, the squared simple correlation coefficient has $\text{Bias}(r^2, \rho^2) > 0$ for $0 \leq \rho \leq .70$ and $\text{Bias}(r^2, \rho^2) < 0$ for $\rho \geq .75$. Therefore, r^2 can be overestimated, underestimated, or unbiased. The exact population $\rho^{2*} \in (.70, .75)$ so that $\text{Bias}(r^2, \rho^{2*}) = 0$ can be numerically determined for different sample size N . Also, the adjusted formula $\hat{\rho}_E^2$ is unbiased when $\rho = 0$ and is overadjusted because $\text{Bias}(\hat{\rho}_E^2, \rho^2) < 0$ when $\rho > 0$. The other four estimators are almost unbiased, with the accu-

racy increasing in the order of $\hat{\rho}_{OP1}^2$, $\hat{\rho}_{PA}^2$, $\hat{\rho}_{OP2}^2$, and $\hat{\rho}_{OP5}^2$. Accordingly, it has been reported in Shieh (2008) and Yin and Fan (2001) that $\hat{\rho}_E^2$ is not the most effective estimator in estimating ρ^2 . They recommended $\hat{\rho}_{PA}^2$ for its remarkable simplicity and performance in estimating ρ^2 .

Next, we focus on the *RMSE* results presented in Tables 10–12 for $N = 20, 50$, and 100 , respectively. For the two nearly unbiased estimators $\hat{\rho}_{PA}^2$ and $\hat{\rho}_{OP5}^2$, it can be readily seen that when $N = 20$, $MSE(\hat{\rho}_{PA}^2, \rho^2) < MSE(\hat{\rho}_{OP5}^2, \rho^2)$ for $\rho \leq .70$, and $MSE(\hat{\rho}_{PA}^2, \rho^2) > MSE(\hat{\rho}_{OP5}^2, \rho^2)$ for $\rho > .70$. In the two instances of $N = 50$ and 100 , $MSE(\hat{\rho}_{PA}^2, \rho^2) < MSE(\hat{\rho}_{OP5}^2, \rho^2)$ for $\rho \leq .65$, and $MSE(\hat{\rho}_{PA}^2, \rho^2) > MSE(\hat{\rho}_{OP5}^2, \rho^2)$ for $\rho > .65$. Hence, there is no dominant situation in their *RMSEs*.

Table 8
Bias for Estimators of ρ^2 With $N = 50$

ρ	r^2	$\hat{\rho}_E^2$	$\hat{\rho}_{OP1}^2$	$\hat{\rho}_{PA}^2$	$\hat{\rho}_{OP2}^2$	$\hat{\rho}_{OP5}^2$
.00	.020408	.000000	.003201	.000543	.000362	.000002
.05	.020261	-.000098	.003179	.000533	.000359	.000002
.10	.019823	-.000389	.003113	.000504	.000350	.000002
.15	.019103	-.000864	.003005	.000457	.000334	.000002
.20	.018112	-.001510	.002858	.000394	.000313	.000001
.25	.016871	-.002309	.002676	.000317	.000288	.000001
.30	.015404	-.003234	.002463	.000231	.000259	.000001
.35	.013740	-.004255	.002225	.000139	.000227	.000001
.40	.011916	-.005335	.001969	.000047	.000194	.000001
.45	.009975	-.006432	.001702	-.000042	.000161	.000001
.50	.007963	-.007496	.001431	-.000123	.000128	.000000
.55	.005938	-.008470	.001166	-.000189	.000098	.000000
.60	.003960	-.009291	.000913	-.000238	.000072	.000000
.65	.002099	-.009888	.000682	-.000266	.000049	.000000
.70	.000433	-.010183	.000477	-.000270	.000031	.000000
.75	-.000953	-.010088	.000307	-.000249	.000017	.000000
.80	-.001964	-.009505	.000175	-.000207	.000008	.000000
.85	-.002496	-.008329	.000082	-.000148	.000003	.000000
.90	-.002432	-.006441	.000027	-.000083	.000001	.000000
.95	-.001647	-.003712	.000004	-.000026	.000000	.000000
.99	-.000405	-.000828	.000000	-.000001	.000000	.000000

Table 9
Bias for Estimators of ρ^2 With $N = 100$

ρ	r^2	$\hat{\rho}_E^2$	$\hat{\rho}_{OP1}^2$	$\hat{\rho}_{PA}^2$	$\hat{\rho}_{OP2}^2$	$\hat{\rho}_{OP5}^2$
.00	.010101	.000000	.000800	.000138	.000047	.000000
.05	.010027	-.000049	.000794	.000135	.000046	.000000
.10	.009806	-.000196	.000777	.000128	.000045	.000000
.15	.009442	-.000436	.000749	.000115	.000043	.000000
.20	.008943	-.000762	.000711	.000099	.000040	.000000
.25	.008318	-.001163	.000664	.000079	.000037	.000000
.30	.007581	-.001627	.000609	.000057	.000033	.000000
.35	.006746	-.002139	.000548	.000034	.000028	.000000
.40	.005834	-.002678	.000483	.000010	.000024	.000000
.45	.004865	-.003223	.000416	-.000012	.000020	.000000
.50	.003864	-.003749	.000348	-.000032	.000016	.000000
.55	.002860	-.004228	.000281	-.000048	.000012	.000000
.60	.001884	-.004628	.000219	-.000060	.000009	.000000
.65	.000970	-.004913	.000162	-.000066	.000006	.000000
.70	.000157	-.005045	.000113	-.000067	.000004	.000000
.75	-.000514	-.004983	.000072	-.000061	.000002	.000000
.80	-.000996	-.004680	.000040	-.000050	.000001	.000000
.85	-.001242	-.004086	.000019	-.000035	.000000	.000000
.90	-.001197	-.003148	.000006	-.000019	.000000	.000000
.95	-.000804	-.001807	.000001	-.000006	.000000	.000000
.99	-.000197	-.000402	.000000	.000000	.000000	.000000

For ease of exposition, the following results are summarized for r^2 , $\hat{\rho}_E^2$, and $\hat{\rho}_{PA}^2$ for all three different sample sizes:

$$MSE(\hat{\rho}_E^2, \rho^2) < MSE(\hat{\rho}_{PA}^2, \rho^2) < MSE(r^2, \rho^2)$$

for $\rho \leq .15$,

$$MSE(\hat{\rho}_E^2, \rho^2) < MSE(r^2, \rho^2) < MSE(\hat{\rho}_{PA}^2, \rho^2)$$

for $\rho = .20$ and $.25$,

$$MSE(r^2, \rho^2) < MSE(\hat{\rho}_E^2, \rho^2) < MSE(\hat{\rho}_{PA}^2, \rho^2)$$

for $.30 \leq \rho \leq .65$,

$$MSE(r^2, \rho^2) < MSE(\hat{\rho}_{PA}^2, \rho^2) < MSE(\hat{\rho}_E^2, \rho^2)$$

for $.70 \leq \rho \leq .85$,

$$MSE(\hat{\rho}_{PA}^2, \rho^2) < MSE(r^2, \rho^2) < MSE(\hat{\rho}_E^2, \rho^2)$$

for $\rho > .85$.

The relative performance among r^2 , $\hat{\rho}_E^2$, and $\hat{\rho}_{OP5}^2$ is analogous to the above by replacing $\hat{\rho}_{PA}^2$ with $\hat{\rho}_{OP5}^2$. The only modification is

$$MSE(\hat{\rho}_E^2, \rho^2) < MSE(r^2, \rho^2) < MSE(\hat{\rho}_{OP5}^2, \rho^2)$$

for $\rho = .15$ when $N = 20$.

According to these findings, $\hat{\rho}_E^2$ is advantageous in MSE for small $\rho < .30$, r^2 dominates for $.30 \leq \rho \leq .85$, and $\hat{\rho}_{PA}^2$ performs best for large $\rho > .85$. This information may be useful in selecting an appropriate measure of the propor-

Table 10
Root-Mean Squared Error for Estimators of ρ^2 With $N = 20$

ρ	r^2	$\hat{\rho}_E^2$	$\hat{\rho}_{OP1}^2$	$\hat{\rho}_{PA}^2$	$\hat{\rho}_{OP2}^2$	$\hat{\rho}_{OP5}^2$
.00	.086711	.072739	.078993	.078710	.079300	.080321
.05	.088394	.075241	.081460	.081353	.081944	.083000
.10	.093190	.082189	.088333	.088684	.089278	.090428
.15	.100457	.092322	.098385	.099340	.099935	.101212
.20	.109401	.104326	.110295	.111889	.112478	.113886
.25	.119261	.117143	.122965	.125166	.125733	.127258
.30	.129373	.129978	.135544	.138287	.138809	.140420
.35	.139172	.142212	.147367	.150565	.151013	.152668
.40	.148164	.153333	.157880	.161434	.161778	.163428
.45	.155897	.162879	.166594	.170394	.170604	.172197
.50	.161934	.170406	.173050	.176978	.177027	.178513
.55	.165831	.175456	.176797	.180728	.180597	.181929
.60	.167120	.177544	.177374	.181179	.180859	.181999
.65	.165286	.176135	.174296	.177845	.177344	.178266
.70	.159751	.170621	.167046	.170214	.169560	.170254
.75	.149849	.160303	.155062	.157732	.156980	.157452
.80	.134792	.144353	.137724	.139802	.139033	.139312
.85	.113632	.121772	.114349	.115775	.115097	.115227
.90	.085197	.091326	.084172	.084953	.084480	.084519
.95	.047988	.051437	.046355	.046598	.046410	.046414
.99	.010581	.011338	.009997	.010008	.009997	.009997

Table 11
Root-Mean Squared Error for Estimators of ρ^2 With $N = 50$

ρ	r^2	$\hat{\rho}_E^2$	$\hat{\rho}_{OP1}^2$	$\hat{\rho}_{PA}^2$	$\hat{\rho}_{OP2}^2$	$\hat{\rho}_{OP5}^2$
.00	.034648	.028583	.029659	.029637	.029714	.029750
.05	.036933	.031523	.032662	.032671	.032755	.032794
.10	.042967	.038916	.040218	.040296	.040396	.040444
.15	.051191	.048491	.050000	.050152	.050272	.050330
.20	.060371	.058809	.060518	.060739	.060875	.060943
.25	.069715	.069090	.070955	.071236	.071384	.071459
.30	.078706	.078858	.080811	.081141	.081294	.081373
.35	.086970	.087770	.089725	.090094	.090243	.090324
.40	.094203	.095542	.097403	.097799	.097937	.098016
.45	.100135	.101915	.103581	.103991	.104109	.104183
.50	.104508	.106639	.108006	.108419	.108509	.108576
.55	.107070	.109460	.110435	.110835	.110894	.110951
.60	.107559	.110118	.110620	.110996	.111021	.111067
.65	.105705	.108338	.108312	.108652	.108645	.108680
.70	.101223	.103831	.103258	.103550	.103516	.103540
.75	.093805	.096285	.095195	.095432	.095378	.095393
.80	.083124	.085363	.083854	.084030	.083969	.083977
.85	.068818	.070698	.068957	.069072	.069016	.069019
.90	.050495	.051888	.050215	.050274	.050237	.050237
.95	.027717	.028488	.027331	.027348	.027335	.027335
.99	.005964	.006130	.005837	.005838	.005837	.005837

tion of explained variance when a researcher has some basic conceptual idea about ρ .

Concluding Remarks

This article concerns the use of Pearson's r as a correlational effect size measure. Despite its routine and common application in empirical studies, the fundamental properties of the sample correlation coefficient are often not sufficiently emphasized in applied work. Perhaps the complexity of r 's distributional function contributes to the fact that its estimation behavior has received little attention in standard texts and related research. Contemporary computer capabilities can be used to conduct intensive computation for the exact bias and MSE of r , as well as other notable formulas. The numerical examinations and

graphical displays facilitate the presentation of different aspects of accuracy and precision in estimating population correlation coefficient. Recognition of the different considerations of biasness and MSE helps to clarify the issue of evaluating strength of association and to choose appropriate effect size estimate in correlation analysis. The empirical results recommend the following procedures for the estimation of the simple correlation coefficient and squared simple correlation coefficient. First, the Olkin and Pratt (1958) approximate formula $\hat{\rho}_{OPA}$ is nearly unbiased for estimating ρ and is easier to apply than the unbiased estimator $\hat{\rho}_U$. Under the MSE consideration, $\hat{\rho}_{OPA}$ and Pearson's r have important advantages for different ranges of underlying population correlation coefficient. Second, the formula $\hat{\rho}_{PA}^2$ has desirable overall

Table 12
Root-Mean Squared Error for Estimators of ρ^2 With $N = 100$

ρ	r^2	$\hat{\rho}_E^2$	$\hat{\rho}_{OP1}^2$	$\hat{\rho}_{PA}^2$	$\hat{\rho}_{OP2}^2$	$\hat{\rho}_{OP5}^2$
.00	.017321	.014215	.014491	.014488	.014500	.014502
.05	.019759	.017199	.017521	.017525	.017540	.017543
.10	.025620	.023912	.024337	.024355	.024375	.024379
.15	.032907	.031848	.032388	.032420	.032445	.032450
.20	.040558	.039971	.040616	.040659	.040689	.040695
.25	.048053	.047825	.048551	.048604	.048637	.048644
.30	.055082	.055138	.055913	.055975	.056009	.056017
.35	.061418	.061706	.062492	.062560	.062594	.062602
.40	.066872	.067350	.068106	.068178	.068210	.068217
.45	.071269	.071901	.072584	.072658	.072686	.072692
.50	.074442	.075194	.075763	.075837	.075858	.075864
.55	.076223	.077063	.077481	.077552	.077567	.077572
.60	.076445	.077340	.077576	.077642	.077650	.077654
.65	.074934	.075852	.075886	.075945	.075946	.075949
.70	.071514	.072419	.072247	.072297	.072292	.072294
.75	.066000	.066858	.066493	.066532	.066524	.066525
.80	.058202	.058974	.058455	.058484	.058474	.058474
.85	.047919	.048564	.047962	.047981	.047972	.047972
.90	.034941	.035417	.034841	.034850	.034844	.034844
.95	.019046	.019308	.018913	.018916	.018914	.018914
.99	.004073	.004129	.004030	.004030	.004030	.004030

performance and computational ease for estimating ρ^2 . However, r^2 and $\hat{\rho}_E^2$, which are the simplified version of R^2 and adjusted R^2 , demonstrate their own usefulness in terms of MSE for some subsets of population correlation coefficient. In view of the use of r across a wide variety of disciplines within the social sciences, the updated consideration of its benefits and costs presented here should be essential to researchers for making sound statistical analysis.

AUTHOR NOTE

The author thanks the action editor, Ira Bernstein, and the three anonymous reviewers for their helpful comments. Correspondence concerning this article should be addressed to G. Shieh, Department of Management Science, National Chiao Tung University, 1001 Ta Hsueh Road, Hsinchu, Taiwan 30050 (e-mail: gwshieh@mail.nctu.edu.tw).

REFERENCES

- ALHJA, F. N. A., & LEVY, A. (2009). Effect size reporting practices in published articles. *Educational & Psychological Measurement*, **69**, 245-265.
- AMERICAN EDUCATIONAL RESEARCH ASSOCIATION TASK FORCE ON REPORTING OF RESEARCH METHODS (2006). Standards for reporting on empirical social science research in AERA publications. *Educational Research*, **35**, 33-40.
- AMERICAN PSYCHOLOGICAL ASSOCIATION (2001). *Publication manual of the American Psychological Association* (5th ed.). Washington, DC: Author.
- BOBKO, P. (2001). *Correlation and regression: Applications for industrial organizational psychology and management* (2nd ed.). Thousand Oaks, CA: Sage.
- BREAUGH, J. A. (2003). Effect size estimation: Factors to consider and mistakes to avoid. *Journal of Management*, **29**, 79-97.
- COHEN, J., COHEN, P., WEST, S. G., & AIKEN, L. S. (2003). *Applied multiple regression/correlation analysis for the behavioral sciences* (3rd ed.). Mahwah, NJ: Erlbaum.
- DURLAK, J. A. (2009). How to select, calculate, and interpret effect sizes. *Journal of Pediatric Psychology*, **34**, 917-928.
- FERGUSON, C. J. (2009). An effect size primer: A guide for clinicians and researchers. *Professional Psychology: Research & Practices*, **40**, 532-538.
- FISHER, R. A. (1915). Frequency distribution of the values of the correlation coefficient in samples from an indefinitely large population. *Biometrika*, **10**, 507-521.
- FISHER, R. A. (1921). On the probable error of a coefficient of correlation deduced from a small sample. *Metron*, **1**, 3-32.
- GHOSH, B. K. (1966). Asymptotic expansions for the moments of the distribution of correlation coefficient. *Biometrika*, **53**, 258-262.
- GRISSOM, R. J., & KIM, J. J. (2005). *Effect sizes for research: A broad practical approach*. Mahwah, NJ: Erlbaum.
- HOTELLING, H. (1953). New light on the correlation coefficient and its transforms. *Journal of the Royal Statistical Society B*, **15**, 193-232.
- HUBERTY, C. (2002). A history of effect size indices. *Educational & Psychological Measurement*, **62**, 227-240.
- JOHNSON, N. L., KOTZ, S., & BALAKRISHNAN, N. (1995). *Continuous univariate distributions* (2nd ed., Vol. 2). New York: Wiley.
- KLINE, R. B. (2004). *Beyond significance testing: Reforming data analysis methods in behavioral research*. Washington, DC: American Psychological Association.
- OLKIN, I., & PRATT, J. W. (1958). Unbiased estimation of certain correlation coefficients. *Annals of Mathematical Statistics*, **29**, 201-211.
- RAJU, N. S., BILGIC, R., EDWARDS, J. E., & FLEER, P. F. (1997). Methodology review: Estimation of population validity and cross-validity, and the use of equal weights in prediction. *Applied Psychological Measurement*, **21**, 291-305.
- RAJU, N. S., BILGIC, R., EDWARDS, J. E., & FLEER, P. F. (1999). Accuracy of population validity and cross-validity estimation: An empirical comparison of formula-based, traditional empirical, and equal weights procedures. *Applied Psychological Measurement*, **23**, 99-115.
- RICHARDSON, J. T. E. (1996). Measures of effect size. *Behavior Research Methods, Instruments, & Computers*, **28**, 12-22.
- ROSENTHAL, R., ROSNOW, R. L., & RUBIN, D. B. (2000). *Contrasts and effect size in behavioral research: A correlational approach*. New York: Cambridge University Press.
- ROSNOW, R. L., & ROSENTHAL, R. (2003). Effect sizes for experimenting psychologists. *Canadian Journal of Experimental Psychology*, **57**, 221-237.
- SHIEH, G. (2006). Exact interval estimation, power calculation and sample size determination in normal correlation analysis. *Psychometrika*, **71**, 529-540.
- SHIEH, G. (2008). Improved shrinkage estimation of squared multiple correlation coefficient and squared cross-validity coefficient. *Organizational Research Methods*, **11**, 387-407.
- STUART, A., & ORD, J. K. (1994). *Kendall's advanced theory of statistics* (6th ed., Vol. 1). New York: Halsted Press.
- VACHA-HAASE, T., & THOMPSON, B. (2004). How to estimate and interpret various effect sizes. *Journal of Counseling Psychology*, **51**, 473-481.
- WANG, Z., & THOMPSON, B. (2007). Is the Pearson r^2 biased, and if so, what is the best correction formula? *Journal of Experimental Education*, **75**, 109-125.
- WILKINSON, L., & THE AMERICAN PSYCHOLOGICAL ASSOCIATION TASK FORCE ON STATISTICAL INFERENCE (1999). Statistical methods in psychology journals: Guidelines and explanations. *American Psychologist*, **54**, 594-604.
- YIN, P., & FAN, X. (2001). Estimating R^2 shrinkage in multiple regression: A comparison of different analytical methods. *Journal of Experimental Education*, **69**, 203-224.
- ZIMMERMAN, D. W., ZUMBO, B. D., & WILLIAMS, R. H. (2003). Bias in estimation and hypothesis testing of correlation. *Psicologica*, **24**, 133-158.

APPENDIX

Fundamental Results of Sample Correlation Coefficient

Under the bivariate normal distribution assumption, the probability density function is conveniently expressed in terms of a hypergeometric function by Hotelling (1953):

$$f(r) = \frac{(N-2)(1-\rho^2)^{(N-1)/2} (1-r^2)^{(N-4)/2}}{N^{1/2}(N-2)B\left(\frac{1}{2}, N-\frac{1}{2}\right)(1-\rho r)^{N-3/2}} \cdot F_h\left(\frac{1}{2}, \frac{1}{2}; N-\frac{1}{2}; \frac{1+\rho r}{2}\right), \quad (\text{A1})$$

where $-1 \leq r \leq 1$, $-1 < \rho < 1$, $B(\alpha, \beta)$ is the standard beta function with parameters α and β , $F_h(a, b; c; x)$ is the Gauss hypergeometric function defined as

$$F_h(a, b; c; x) = \sum_{k=0}^{\infty} \frac{\Gamma(a+k)\Gamma(b+k)\Gamma(c)}{\Gamma(a)\Gamma(b)\Gamma(c+k)} \cdot \frac{x^k}{k!},$$

and $\Gamma(\cdot)$ is the gamma function. Moreover, Olkin and Pratt (1958) have shown that the unique minimum variance unbiased estimator $\hat{\rho}_U$ of ρ is of the form

$$\hat{\rho}_U(r) = r \cdot F_h\left(\frac{1}{2}, \frac{1}{2}; \frac{N-2}{2}; 1-r^2\right). \quad (\text{A2})$$

Also, it follows from Ghosh (1966) that the first and second moments of r are

$$E[r] = \frac{2\left[\Gamma\left(\frac{N}{2}\right)\right]^2}{(N-1)\left[\Gamma\left(\frac{N-1}{2}\right)\right]^2} \rho \cdot F_h\left(\frac{1}{2}, \frac{1}{2}; \frac{N+1}{2}; \rho^2\right) \quad (\text{A3})$$

and

$$E[r^2] = 1 - \frac{(N-2)(1-\rho^2)}{(N-1)} \cdot F_h\left(1, 1; \frac{N+1}{2}; \rho^2\right). \quad (\text{A4})$$

Specifically, the exact bias and MSE for an estimator $\hat{\rho} = \hat{\rho}(r)$ of ρ are computed as

$$\text{Bias}(\hat{\rho}, \rho) = \int_{-1}^1 (\hat{\rho} - \rho) f(r) dr \quad \text{and} \quad \text{MSE}(\hat{\rho}, \rho) = \int_{-1}^1 (\hat{\rho} - \rho)^2 f(r) dr,$$

where $f(r)$ is given in Equation A1. Due to the complication, intensive numerical integration using Simpson's rule with respect to the probability density distribution $f(r)$ is conducted to compute the exact values of $\text{Bias}(\hat{\rho}, \rho)$ and $\text{MSE}(\hat{\rho}, \rho)$.

For the estimation of ρ^2 , Olkin and Pratt (1958) also derived the unique minimum variance unbiased estimator $\hat{\rho}_U^2$ of ρ^2 :

$$\hat{\rho}_U^2(r^2) = 1 - \frac{N-3}{N-2} (1-r^2) \cdot F_h\left(1, 1; \frac{N}{2}; 1-r^2\right). \quad (\text{A5})$$

Similarly, the exact bias and MSE for an estimator $\hat{\rho}^2 = \hat{\rho}^2(r^2)$ of ρ^2 can be computed as

$$\text{Bias}(\hat{\rho}^2, \rho^2) = \int_{-1}^1 (\hat{\rho}^2 - \rho^2) f(r) dr \quad \text{and} \quad \text{MSE}(\hat{\rho}^2, \rho^2) = \int_{-1}^1 (\hat{\rho}^2 - \rho^2)^2 f(r) dr,$$

where $f(r)$ is given in Equation A1.