

What Does the Correlation Coefficient Really Tell Us About the Individual?

R. C. Gardner and R. W. J. Neufeld

Department of Psychology

University of Western Ontario

ABSTRACT

The Pearson product moment correlation is an index of how well one can predict scores on one variable given scores on another based on a linear regression equation; the higher the correlation the better prediction. This article investigates the nature of this predictability at the individual level for various values of the correlation coefficient. It demonstrates that there is relatively little improvement in predictability of individual status for a large percentage of the population as correlations increase from .10 to .50, and that it is not that much better for correlations as high as .90. A higher correlation does imply better prediction in the general sense that there is a stronger linear relation between the two variables, but substantial qualifications come into play at the individual level of prediction. Regions of predictor-predicted variable distributions allowing greater and lesser predictive confidence specifically at this level are delineated.

Considerable research is concerned with the study of individual differences, often focusing on the correlation between two variables. A significant correlation is considered as an indication of some (linear) consistency in individual differences. But what does it really tell us about individual differences? Cohen (1988) has characterized a correlation of .10 as depicting a small effect, commenting that “many relationships pursued in ‘soft’ behavioral science are of this magnitude” (p.79). He characterizes a correlation of .30 as a medium effect, observing that such values are encountered in behavioral science and that “this degree of relationship would be perceptible to the naked eye of a reasonably sensitive observer” (p. 80). He describes a correlation of .50 as a large effect, remarking that Ghiselli (1964) considered it as the practical upper limit of predictive effectiveness. These values have been described as reflecting 1%, 9%, and 25% of the variance common to two variables respectively.

SIMULATION STUDY

Our interest in the meaning of correlation as it applies to the individual drove us to conduct an empirical study in which we investigated actual data producing a correlation of .50, corresponding to Cohen’s large effect. Specifically, we obtained a sample of 1000 observations from a population correlation of .50 using the Correlated Multivariate Random Normal Scores Generator (Aguinis, 1994). The sample means for the two variables were -.018 and -.006, and the standard deviations were 1.000 and 1.054. The obtained correlation was .509. Figure 1 presents a scatter plot of the empirical data.

Insert Figure 1 about here

These scores were transformed to decile values, and their bivariate frequency distribution

is presented in Table 1 in the form of a 10x10 matrix. If the correlation between the two variables were 0, this would result in a 10x10 matrix with an expected 10 observations in each cell. Thus, given 1000 cases, the probability of being in any given cell is $10/1000 = .01$, and the probability of being in any given decile for one of the variables is $100/1000 = .10$. Correlations greater than 0 will result in some cell frequencies being greater than 10 and others less, but the marginal frequencies would still be 100 and the marginal probabilities would still be .10.

Expressed in this way, it is possible to assess (1) the probability of an observation being in any one cell, or (2) the probability of being in any one Z_y cell or set of Z_y cells given that one is in any one decile of Z_x , or (3) the probability of being in any one cell or set of Z_y cells given that one is in any set of Z_x deciles. Returning to Table 1, consider the following examples:

1. The probability of being in any one cell is simply the frequency of that cell divided by the total N . Thus, the probability of simultaneously being in the 10th decile of Z_x and the 10th decile of Z_y is $36/1000 = .036$, while the probability of being in the 5th decile of Z_x and the 4th decile of Z_y is $14/1000 = .014$.

2. The probability of being in any one cell of Z_y given that one is in any one decile of Z_x is the ratio of the frequency in the particular cell divided by the total number of observations in the Z_x decile. Thus, the probability of being in the 10th decile of Z_y given that one is in the 10th decile of Z_x is $36/100 = .36$. Furthermore, the probability of being in any of a set of Z_y deciles given that one is in one set of Z_x deciles is the sum of the relevant set divided by the number of individuals in the Z_x set. For example, the probability of being in the 8th, 9th, or 10th decile of Z_y given one is in the 10th decile of Z_x is equal to the $(36 + 19 + 14)/100 = .69$.

3. The probability of being in any one cell or set of Z_y cells given that one is in any set of Z_x deciles is equal to the ratio of the frequency in that cell or the sum of frequencies in the set

of cells over the sum of frequencies of the Z_x deciles. For example, the probability of being in the 9th decile cell of Z_y given that one is in the 7th or 8th decile of Z_x is $(10+11)/(100+100) = .105$, or the probability of being in the 7th to 10th deciles of Z_y given that one is in the 8th to 10th deciles of Z_x is $(18+11+13+10+13+12+21+17+36+19+14+8)/(100+100+100) = 192/300 = .64$. In short, one can determine the probability of any given scenario.

This is an empirical example based on 1000 observations drawn from a bivariate normal population with a correlation of .50, and if another sample were drawn the cell frequencies would vary slightly because of sampling fluctuations, resulting in slightly altered probabilities. The same logic can be demonstrated, however, based on theoretical considerations.

THEORETICAL SPECIFICS

Governing Formulae.

The equation for the probability density function of the bivariate normal distribution in standard score form is given by Hayes (1973, p. 659) as:

$$f(Z_X, Z_Y) = \frac{1}{K} e^{-G}$$

where:
$$G = \frac{(Z_X^2 + Z_Y^2 - 2\rho Z_X Z_Y)}{2(1 - \rho^2)}$$

and
$$K = 2\pi\sqrt{(1 - \rho^2)}$$

Examination of the function will reveal that it is symmetrical for any given value of the population correlation, ρ . When ρ is 0, the resulting function is a circle, but for any other value it is an ellipse with a slope given by the sign of ρ , and as the correlation increases in magnitude, the minor axis of the ellipse becomes shorter (i.e., the ellipse becomes thinner). Figure 2 illustrates the parameter space for a positive correlation between two variables, Z_x and Z_y .

 Insert Figure 2 about here

Given the population correlation, ρ , the regression of two sets of standard scores, Z_y on Z_x would have the same slope as that for Z_x on Z_y with a value equal to ρ , and the values of Z_x and Z_y would be distributed within the ellipse. The ellipse expresses the nature of the probability density function of the bivariate distribution and for any given value of ρ would be symmetrical about the center of the ellipse. Figure 1 shows cut-off values for the Z_x and Z_y axes. The area designated A indicates the cases that are above both cut-offs, the area identified as B indicates those that are above the Z_x but below the Z_y cutoffs, and the area C indicates those values that are below the Z_x but above the Z_y cut-offs. Note from the figure that the proportion of cases above the Z_x cut-off that are also above the Z_y cut-off is given by the ratio of $A/(A + B)$, while the proportion of cases that are above the Z_y cut-off that are also above the Z_x cut-off are $A/(A + C)$. In short, given a correlation between Z_x and Z_y , conditional probabilities exist which would increase with increases in the magnitude of their correlation. Thus, the Bayesian expression of the probability of exceeding Z_Y given Z_X is:

$$P(Z_Y|Z_X) = \frac{P(Z_Y)P(Z_X|Z_Y)}{P(Z_X)} = \frac{A}{A+B}$$

Similarly, the probability of exceeding Z_X given Z_Y is:

$$P(Z_X|Z_Y) = \frac{P(Z_X)P(Z_Y|Z_X)}{P(Z_Y)} = \frac{A}{A+C}$$

Note, this relationship does not depend on any cause-effect scenario. It is possible that individual differences in X are responsible for individual differences in Y , or that individual

differences in Y are responsible for individual differences in X , or that individual differences in X and Y are due to some other variable(s). The correlation model assumes simply that there is a linear association between the two variables. In research, a significant correlation implies that there exists in the population an association between the two variables, that the unbiased estimate of the population correlation (ρ) is the value of r obtained in the sample, and that this association has implications for the individuals sampled. The major point is that if there is a significant correlation in the population between the two variables, this means that information about one of the variables has implications about the other variable, regardless of any cause-effect interpretation.

The probability density function presented above can be used to compute the areas in the intersection of the deciles in the parameter space. This is achieved by double integration of $f(Z_X, Z_Y)$ with respect to Z_X and Z_Y over the intervals prescribed by each pair of X and Y deciles. That is:

$$\int_{Z_{xl}}^{Z_{xu}} \int_{Z_{yl}}^{Z_{yu}} f(Z_x, Z_y) dZ_x dZ_y$$

where Z_{xl} , Z_{xu} , Z_{yl} and Z_{yu} are the lower and upper limits of the variables Z_X and Z_Y respectively.

Analytical Results

Tables 2 to 6 present the areas (probabilities) associated with each decile in the bivariate distributions for correlations of .1, .3, .5, .7, and .9. The cell values are given to five decimal places while the marginal values have been rounded to two places. It will be noted that each of the marginal values equal .10, and that the cell values in each row and column add to .10 with rounding, hence the sum of all the cell values equals 1.0. Being areas identified by the density function, each of these values can be viewed as the probability of membership in any given cell,

while the row and column marginals are the probabilities of being in any given decile of either Z_x or Z_y .

Examination of the cell values in each table will also reveal that they are perfectly symmetrical. Although it might be expected that the probabilities would be greatest for deciles in agreement between the two variables (i.e., the diagonal values), this is true only for the values in Table 6 (the correlation of .90). For all other tables, these probabilities are greatest only for deciles 1 and 10; for the correlation of .70, this is also true for deciles 5 and 6 (see Table 5). In all other instances probabilities are larger in some higher deciles in each column for those above the median and lower for those below, and the extent to which this occurs depends on the magnitude of the correlation, reflecting the width of the ellipse described by the sample space.

 Insert Tables 2 to 6 about here

These values can be used to estimate the probability of an individual achieving various decile cut-offs in Y given a specific decile cut-off value in X . Using the values from Table 4 as examples, it can be observed that the probability of being simultaneously in the 10th decile of Z_x and the 10th decile of Z_y is .03218,¹ and the probability of being simultaneously in the 5th decile of Z_x and the 4th decile of Z_y is .01140. Furthermore, the probability of being in the 10th decile of Z_x given that one is in the 10th decile of Z_y is $.03218/.10 = .3218$ and the probability of being in the 8th, 9th or 10th decile of Z_y given they are in the 10th decile of Z_x is $(.01398 + .01912 + .03218)/.10 = .06528/.10 = .6528$. Finally, the probability of being in the 9th decile of Z_y given

¹The examples used here are the same as those used in the empirical example discussed above for a correlation of .50.

that one is in the 7th or 8th decile of Z_x is $(.01223 + .01441)/(1+.1) = .02664/.2 = .1332$, while the probability of being in the 7th to 10th deciles of Z_y given that one is in the 8th to 10th deciles of Z_x is $(.01398 + .01912 + .03218 + .01441 + .01660 + .01912 + .01369 + .01441 + .01398 + .01248 + .01223 + .01043)/(1.0 + 1.0 + 1.0) = .19263/.30 = .6421$. Of course, these are expected values given that they are calculated from values defined by the density function, but it will be noted that they are very similar to those calculated from the empirical sampling distribution of the population correlation of .50 presented in Table 1.

The values in Tables 2 to 6 can be used to estimate the probabilities associated with other scenarios. For example, it can be demonstrated that given a correlation of .10, the probability of one being in the top 50% of one distribution given that they are in the top 50% of the other is .53130, hardly a difference from the value of .50 expected in any event. The corresponding probability for each of the other correlations is .59112 for $\rho = .30$, .66572 for $\rho = .50$, .74598 for $\rho = .70$, and .85570 for $\rho = .90$. Note, therefore that for what has been characterized as a large effect ($r = .50$) in psychological research, roughly two thirds of the sample would be correctly identified as being in the top half of the predicted distribution while one-third would be misidentified. Even for a near perfect correlation ($\rho = .90$), only 86% are correctly placed in the top half of the distribution while 14% are misclassified².

At the individual level, the correlation might be viewed as an index of how consistent an individual's position will be in both variables relevant to the other individuals in the sample, and it is assumed that the degree of consistency will increase with increases in the magnitude of the correlation. Obviously, if the correlation is 1.00, there will be perfect consistency, but what

² These values may appear to be inconsistent with Rosenthal and Rubin's (1982) discussion of the binomial effect size display (BESD). Note, however, that the "r" values to which they refer are phi (ϕ) coefficients so that their medium effect of .30 corresponds to a Pearson correlation of .50 which is a strong effect.

about other values of the correlation? One way of investigating this is to consider the probability of an individual being in the same decile for the predicted variable as the predictor variable as a function of the magnitude of the correlation. Table 7 presents the probability of individuals being in corresponding deciles for each of the five correlations.

 Insert Table 7 about here.

Inspection of Table 7 will reveal some important features. First, note that for each column the values in the top half are mirror images of those in the lower half. The probabilities are highest for the extreme deciles, decrease in magnitude for the more central ones, are symmetrical around the median, and become more variable as the correlation increases. Second, note that the probabilities for each decile tend to increase as the correlations increase; the mean probabilities over all deciles are .11, .14, .17, .23, and .36 for the five correlations as would be expected given that the correlations are increasing. What is probably unexpected is how low these values are. On average, the probability of an individual being in the corresponding deciles is relatively low regardless of the magnitude of the correlation though admittedly it increases with increasing correlation. On closer inspection it will be noted that for correlations of .10, .30, and .50 there is relatively no change in predictability for deciles 2 to 9. For correlations of .70 and .90, there are slight improvements but here again the largest changes are in deciles 1, 2, 9 and 10 (i.e., the extreme 40% of the cases). Stated in another way, there is relatively little predictability for 80% of the cases in correlations ranging from .10 to .50, and 60% for larger correlations.

One might argue that the purpose of correlation is not to assess the predictability of an

individual being in the same decile on two variables, but rather the predictability of an individual scoring as high or higher if above the median or as low or lower if below the median. Table 8 presents the resulting probabilities for these scenarios.

Insert Table 8 about here

As was the case with Table 7, the values in the top part of each column are mirror images of those in the lower part. In Table 8, however, the values are lowest on the extreme ends and increase as they approach the middle, and the variability of the values in each column decreases as the correlations increase. Note too that like Table 7, the probabilities for each row tend to increase as the correlation increases; the mean probabilities over all rows are .39, .42, .46, .50, and .59. That is, on average, for correlations less than .70, less than 50% of the cases would be correctly identified as doing as well or more extremely on Z_Y given their decile position on Z_X . For individuals in deciles 5 or 6, there is virtually no difference in predictability, and relatively not very much more for deciles 4 to 7. That is, as in Table 7, much of the improved predictability as the correlations increase involves the extremes of the distributions, but even here the hit rate exceeds 50% only when the correlation is .90.

CONCLUDING COMMENTS

Clearly, as the magnitude of the correlation increases the degree of linear relationship between two variables increases. Nonetheless, when it comes to predictability at the individual level relatively little is gained. If the magnitude of the correlation is considered to imply a degree of predictability about the individual obtaining similar, or similar and more extreme scores, on both variables then improvement of predictability holds primarily at the ends of the

distributions. The view that a higher correlation implies better prediction is true only in the general sense that there is a stronger linear relationship between the two variables, but is much less so when predictions about an individual are addressed.

(Dated: March 30, 2012)

References

- Aguinis, H. (1994). A quickbasic program for generating correlated multivariate random normal scores. *Educational and Psychological Measurement, 54*, 687-689.
- Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences (Second Edition)*. Hillsdale, NJ: Lawrence Erlbaum.
- Ghiselli, E. E. (1964). Dr. Ghiselli comments on Dr. Tupe's note. *Personnel Psychology, 17*, 61-63.
- Hayes, W. L. (1973). *Statistics for the Social Sciences (Second Edition)* New York: Holt, Rinehart and Winston.
- Rosenthal, R. & Rubin, D.B. (1982). A simple, general-purpose display of magnitude of experimental effect. *Journal of Educational Psychology, 74*, 166-169.

Author Notes

This research was supported by a grant from the Social Sciences and Humanities Research Council of Canada.

Figure 1: Scatter plot of the Empirical Data

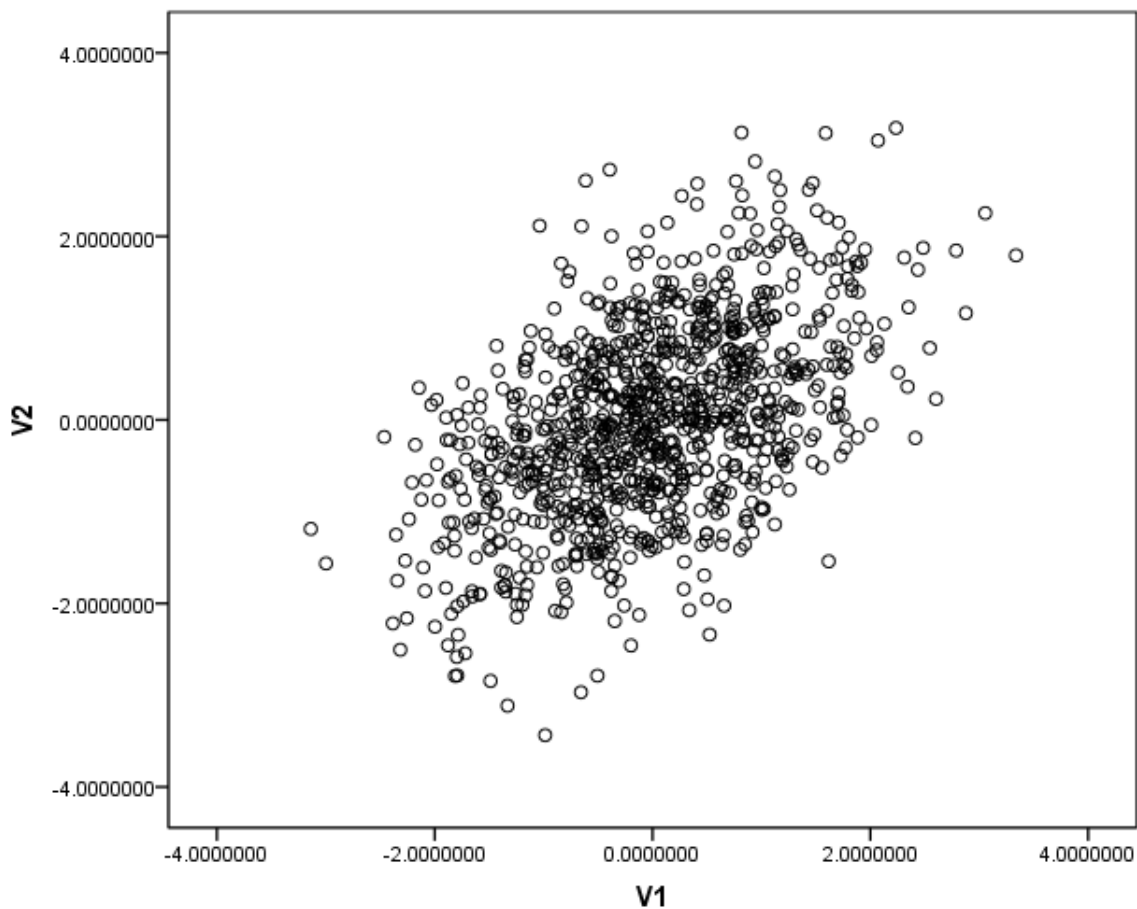
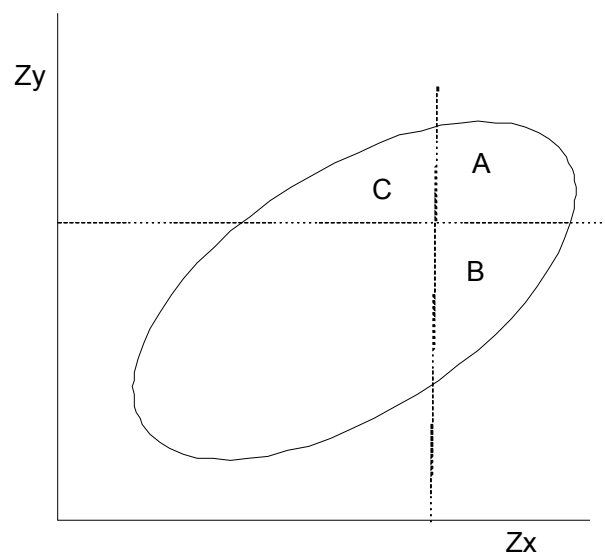


Figure 2: Parameter Space of a Bivariate Distribution



Y Variable	X Variable										
Decile	1	2	3	4	5	6	7	8	9	10	Row N
10	1	0	0	5	7	9	11	18	13	36	100
9	1	9	7	10	10	11	10	11	12	19	100
8	4	5	7	3	7	13	13	13	21	14	100
7	6	10	7	9	8	16	9	10	17	8	100
6	3	9	12	16	9	10	14	8	11	8	100
5	5	8	10	10	12	17	12	9	12	5	100
4	14	11	13	9	14	5	10	12	9	3	100
3	14	18	9	13	11	6	10	11	2	6	100
2	18	13	19	15	11	7	7	6	3	1	100
1	34	17	16	10	11	6	4	2	0	0	100
Column N	100	100	100	100	100	100	100	100	100	100	1000

Table 1. Sample Space from a Population Correlation of $\rho = .50$.

Zy	Zx										
Decile	1	2	3	4	5	6	7	8	9	10	Total
10	.00719	.00821	.00886	.00922	.00955	.00998	.01055	.01124	.01186	.01332	.10
9	.00821	.00894	.00941	.00959	.00975	.01001	.01039	.01084	.01114	.01186	.10
8	.00886	.00941	.00976	.00984	.00992	.01009	.01037	.01070	.01084	.01124	.10
7	.00922	.00959	.00984	.00984	.00984	.00994	.01014	.01037	.01039	.01055	.10
6	.00955	.00975	.00992	.00984	.00977	.00981	.00994	.01009	.01001	.00998	.10
5	.00998	.01001	.01009	.00994	.00981	.00977	.00984	.00992	.00975	.00955	.10
4	.01055	.01039	.01037	.01014	.00994	.00984	.00984	.00984	.00959	.00922	.10
3	.01124	.01084	.01070	.01037	.01009	.00992	.00984	.00976	.00941	.00886	.10
2	.01186	.01114	.01084	.01039	.01001	.00975	.00959	.00941	.00894	.00821	.10
1	.01332	.01186	.01124	.01055	.00998	.00955	.00922	.00886	.00821	.00719	.10
Total	.10	.10	.10	.10	.10	.10	.10	.10	.10	.10	1.00

Table 2. Proportions of the area for Deciles in the Population Space for a correlation of .10.

Zy	Zx										
Decile	1	2	3	4	5	6	7	8	9	10	Total
10	.00302	.00477	.00608	.00717	.00825	.00950	.01106	.01305	.01554	.02151	.10
9	.00477	.00661	.00780	.00864	.00940	.01024	.01123	.01239	.01352	.01554	.10
8	.00608	.00780	.00881	.00942	.00984	.01050	.01117	.01189	.01239	.01305	.10
7	.00717	.00864	.00942	.00979	.01008	.01040	.01079	.01117	.01123	.01106	.10
6	.00825	.00940	.00994	.01008	.01014	.01024	.01040	.01050	.01024	.00950	.10
5	.00950	.01024	.01050	.01040	.01024	.01014	.01008	.00994	.00940	.00825	.10
4	.01106	.01123	.01117	.01079	.01040	.01008	.00979	.00942	.00864	.00717	.10
3	.01305	.01239	.01189	.01117	.01050	.00994	.00942	.00881	.00780	.00608	.10
2	.01554	.01352	.01239	.01123	.01024	.00940	.00864	.00780	.00661	.00477	.10
1	.02151	.01554	.01305	.01106	.00950	.00825	.00717	.00608	.00477	.00302	.10
Total	.10	.10	.10	.10	.10	.10	.10	.10	.10	.10	1.00

Table 3. Proportions of the area for Deciles in the Population Space for a Correlation of .30.

Zy	Zx										
Decile	1	2	3	4	5	6	7	8	9	10	Total
10	.00075	.00190	.00314	.00447	.00599	.00789	.01043	.01398	.01912	.03218	.10
9	.00190	.00395	.00568	.00721	.00871	.01034	.01223	.01441	.01660	.01912	.10
8	.00314	.00568	.00750	.00887	.01006	.01125	.01248	.01369	.01441	.01398	.10
7	.00447	.00721	.00887	.00992	.01071	.01140	.01204	.01248	.01223	.01043	.10
6	.00599	.00871	.01006	.01071	.01106	.01129	.01140	.01125	.01034	.00789	.10
5	.00789	.01034	.01125	.01140	.01129	.01106	.01071	.01006	.00871	.00599	.10
4	.01043	.01223	.01248	.01204	.01140	.01071	.00992	.00887	.00721	.00447	.10
3	.01398	.01441	.01369	.01248	.01125	.01006	.00887	.00750	.00568	.00314	.10
2	.01912	.01660	.01441	.01223	.01034	.00871	.00721	.00568	.00395	.00190	.10
1	.03218	.01912	.01398	.01043	.00789	.00599	.00447	.00314	.00190	.00075	.10
Total	.10	.10	.10	.10	.10	.10	.10	.10	.10	.10	1.00

Table 4. Proportions of the area for Deciles in the Population Space for a Correlation of .50.

Zy	Zx										
Decile	1	2	3	4	5	6	7	8	9	10	Total
10	.00004	.00026	.00074	.00153	.00277	.00472	.00786	.01313	.02227	.04652	.10
9	.00026	.00120	.00267	.00455	.00687	.00974	.01331	.01759	.02173	.02227	.10
8	.00074	.00267	.00502	.00745	.00993	.01246	.01499	.01712	.01759	.01313	.10
7	.00153	.00455	.00745	.00989	.01191	.01355	.01472	.01499	.01331	.00786	.10
6	.00277	.00687	.00993	.01191	.01310	.01366	.01355	.01246	.00974	.00472	.10
5	.00472	.00974	.01246	.01355	.01366	.01310	.01191	.00993	.00687	.00277	.10
4	.00786	.01331	.01499	.01472	.01355	.01191	.00989	.00745	.00455	.00153	.10
3	.01313	.01759	.01712	.01499	.01246	.00993	.00745	.00502	.00267	.00074	.10
2	.02227	.02173	.01759	.01331	.00974	.00687	.00455	.00267	.00120	.00026	.10
1	.04652	.02227	.01313	.00786	.00472	.00277	.00153	.00074	.00026	.00004	.10
Total	.10	.10	.10	.10	.10	.10	.10	.10	.10	.10	1.00

Table 5. Proportions of the area for Deciles in the Population Space for a Correlation of .70.

Zy	Zx										
Decile	1	2	3	4	5	6	7	8	9	10	Total
10	.00000	.00000	.00000	.00001	.00008	.00040	.00170	.00653	.02248	.06870	.10
9	.00000	.00000	.00005	.00035	.00144	.00447	.01135	.02375	.03628	.02248	.10
8	.00000	.00005	.00050	.00214	.00596	.01263	.02141	.02809	.02375	.00653	.10
7	.00001	.00035	.00214	.00638	.01283	.01972	.02387	.02141	.01135	.00170	.10
6	.00008	.00144	.00596	.01283	.01915	.02203	.01972	.01263	.00447	.00040	.10
5	.00040	.00447	.01263	.01972	.02203	.01915	.01283	.00596	.00144	.00008	.10
4	.00170	.01135	.02141	.02387	.01972	.01283	.00638	.00214	.00035	.00001	.10
3	.00653	.02375	.02809	.02141	.01263	.00596	.00214	.00050	.00005	.00000	.10
2	.02248	.03628	.02375	.01135	.00447	.00144	.00035	.00005	.00000	.00000	.10
1	.06870	.02248	.00653	.00170	.00040	.00008	.00001	.00000	.00000	.00000	.10
Total	.10	.10	.10	.10	.10	.10	.10	.10	.10	.10	1.00

Table 6. Proportions of the area for Deciles in the Population Space for a Correlation of .90.

Decile	$\rho = .10$	$\rho = .30$	$\rho = .50$	$\rho = .70$	$\rho = .90$
10	.13	.22	.32	.47	.69
9	.11	.14	.17	.22	.36
8	.11	.12	.14	.17	.28
7	.10	.11	.12	.15	.24
6	.10	.10	.11	.14	.22
5	.10	.10	.11	.14	.22
4	.10	.11	.12	.15	.24
3	.11	.12	.14	.17	.28
2	.11	.14	.17	.22	.36
1	.13	.22	.32	.47	.69

Table 7. Probability of Co-occurrence in the Deciles.

Scenario	$\rho = .10$	$\rho = .30$	$\rho = .50$	$\rho = .70$	$\rho = .90$
$P((9-10) 9)$.23	.29	.36	.44	.59
$P((8-10) 8)$.33	.37	.42	.48	.58
$P((7-10) 7)$.41	.44	.47	.51	.58
$P((6-10) 6)$.50	.51	.52	.54	.59
$P((5-10) 5)$.59	.58	.58	.58	.61
$P((1-5) 5)$.50	.51	.52	.54	.59
$P((1-4) 4)$.41	.44	.47	.51	.58
$P((1-3) 3)$.33	.37	.42	.48	.58
$P((1-2) 2)$.23	.29	.36	.44	.59

Table 8. Probability of Attaining Initial or More Extreme Decile in Z_Y given Z_X .