

# DISTRIBUTION OF SAMPLE CORRELATION COEFFICIENTS

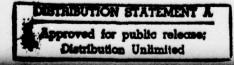
BY

KHURSHEED ALAM

NSS September, 1977

and the second second

TECHNICAL REPORT No. 266





# DISTRIBUTION OF SAMPLE CORRELATION COEFFICIENTS

Khursheed Alam Clemson University

# ABSTRACT

Let  $(Y, X_1, \ldots, X_K)$  be a random vector distributed according to a multivariate normal distribution, where  $X_1, \ldots, X_K$  are considered as predictor variables and Y is the predictand. Let  $r_i$  and  $R_i$  denote the population and sample correlation coefficients, respectively, between Y and  $X_i$ . The population correlation coefficient  $r_i$  is a measure of the predictive power of  $X_i$ . The author has derived the joint distribution of  $R_1, \ldots, R_K$ and its asymptotic property. The given result is useful in the problem of selecting the most important predictor variable, corresponding to the largest absolute value of  $r_i$ .

Key Words: Multivariate Normal Distribution; Correlation Coefficients; Predictor Variables.

AMS Classification: Primary 62E15

Secondary 62E20

*The author's work Research under Contract	was supported by N00014-75-0451.	the Office	White Section Of Naval Buff Section UNANINOUNCED JUSTIFICATION	800
			BY DISTRIBUTION/AVAILABILITY CODE Dist	

#### 1. INTROUCTION

The problem of selecting a variable or several variables from a set of predictor variables  $\{X_i\}$  occurs frequently in the design of experiments. The correlation between a predictor variable  $X_i$  and the predictand Y measures the "leverage" of  $X_i$ upon Y. If  $X_i$  and Y are jointly distributed according to the standard bivariate normal distribution with correlation coefficient  $r_i$  then the conditional distribution of Y given  $X_i$  is normal  $N(r_iX_i, 1-r_i^2)$ . The larger the absolute value of  $r_i$ , the smaller is the variance of the conditional distribution, and therefore higher is the predictive power of  $X_i$ . Thus the predictor variable corresponding to the largest value of  $r_i^2$  may be considered as the most important (best) predictor variable.

Let the random vector  $(Y, X_1, \ldots, X_K)$  be distributed according to a multivariate normal distribution. Suppose that a sample of n observations is taken from the given distribution. Let  $r_i$  and  $R_i$  denote the population and sample correlation coefficients between Y and  $X_i$ , respectively. Let  $r_i^* = r_i (1-r_i^2)^{-\frac{1}{2}}$ and  $R_i^* = R_i (1-R_i^2)^{-\frac{1}{2}}$ . In this paper we derive the asymptotic distribution of  $\underline{B}^* = (R_1^*, \ldots, R_K^*)$ . The given result is useful in the problem of selecting the best predictor variable. The selection problem has been considered recently by Ramberg (1976). Rizvi and Solomon (1976) and Alam, Rizvi and Solomon (1976) have considered the m of selecting from  $p \ge 2$  given multivariate populations, the form in the largest multiple correlation between a single variate, classified as the predictand, and the remaining variates.

-1-

# 2. ASYMPTOTIC DISTRIBUTION OF R\*

First we prove a lemma which will be used in the proof of the main result. Let  $(Z_{1t}, \ldots, Z_{Kt})$  denote the t-th observation in a sample of n observations from a K-variate normal distribution  $N(0,\Omega)$ , and let  $U_i = \sum_{t=1}^{n} Z_{it}^2$  $U = (U_1, \ldots, U_K)'$ . Let  $\Omega = (\omega_{ij})$ ,  $e = (\omega_{11}, \ldots, \omega_{KK})'$  and  $\tilde{\Omega} = (\omega_{ij}^2)$ .

Lemma 2.1. The asymptotic distribution of U is multivariate normal N(ne,  $2n\hat{\Omega}$ ).

Proof: Let  $\theta_1, \ldots, \theta_K$  be K imaginary numbers and let D be a diagonal matrix whose ith diagonal element is  $\theta_i$ . Let  $\lambda_1, \ldots, \lambda_K$  denote the characteristic roots of  $\Omega^{\frac{1}{2}} D \Omega^{\frac{1}{2}}$ , where  $\Omega^{\frac{1}{2}}$ denotes a symmetric square root of  $\Omega$ . We have

$$\sum_{i=1}^{K} \lambda_{i} = \text{trace } \Omega^{\frac{1}{2}} D \Omega^{\frac{1}{2}}$$
$$= \text{trace } D \Omega$$
$$= \theta' \theta$$
$$\sum_{i=1}^{K} \lambda_{i}^{2} = \text{trace } (\Omega^{\frac{1}{2}} D \Omega^{\frac{1}{2}})^{2}$$
$$= \text{trace } (D \Omega)^{2}$$
$$= \theta' \Omega \theta$$

where  $\theta = (\theta_1, \dots, \theta_K)'$ .

The characteristic function of the distribution of U is given by

$$E \exp(\theta'U) = |I-2\Omega^{\frac{1}{2}}D\Omega^{\frac{1}{2}}|^{-n/2}$$

where I denotes an identity matrix. The characteristic function of the normalized distribution of U is given by

$$E \exp\left(\frac{1}{\sqrt{2n}} \frac{\theta}{\theta} \left( \underbrace{\mathbf{U}}_{-\mathbf{n}\underline{e}} \right) \right) = |\mathbf{I} - \sqrt{\frac{2}{n}} \widehat{\Omega}^{\frac{1}{2}} D \widehat{\Omega}^{\frac{1}{2}}|^{-n/2} \exp\left(-\sqrt{\frac{2}{n}} \frac{\theta}{\theta} \cdot \underline{e}\right)$$
$$= \Pi_{i=1}^{K} (1 - \sqrt{\frac{2}{n}} \lambda_{i})^{-n/2} \exp\left(-\sqrt{\frac{2}{n}} \frac{\theta}{\theta} \cdot \underline{e}\right)$$
$$= \exp\left(\sqrt{\frac{2}{n}} \left(\sum_{i=1}^{K} \lambda_{i} - \frac{\theta}{\theta} \cdot \underline{e}\right) + \frac{1}{2} \sum_{i=1}^{K} \lambda_{i}^{2}\right) (1 + O(n^{-\frac{1}{2}}))$$
$$= \exp\left(\frac{1}{2} \frac{\theta}{\theta} \cdot \widehat{\Omega} \frac{\theta}{\theta}\right) (1 + O(n^{-\frac{1}{2}})). \qquad (2.1)$$

The lemma follows since the limiting value of the right hand side of (2.1) as  $n \rightarrow \infty$  is equal to the characteristic function of the multivariate normal distribution N(0, $\hat{\Omega}$ ).

Now we consider the distribution of  $\mathbb{R}^*$ . Without loss of generality we can assume that the variables  $Y, X_1, \ldots, X_K$  are standarized, that is, they are distributed with mean 0 and variance 1. Let  $r_{ij}$  denote the correlation coefficient between  $X_i$  and  $X_j$  and let  $\tilde{\Sigma} = (r_{ij})$  and  $\hat{\Sigma} = (r_{ij}^2)$ . Let  $(Y_t, X_{1t}, \ldots, X_{Kt})$  denote the t-th observation in the sample, and let  $\overline{X}_i = \frac{1}{n} \sum_{t=1}^n X_{it}, \ \overline{Y} = \frac{1}{n} \sum_{t=1}^n Y_t, \ S^2 = \sum_{t=1}^n (Y_t - \overline{Y})^2$ 

$$V_{i} = (\sum_{t=1}^{n} (Y_{t} - \overline{Y}) X_{it}) / S$$
 (2.2)

$$W_{i} = \sum_{t=1}^{n} (x_{it} - \overline{x})^{2} - v_{i}^{2}. \qquad (2.3)$$

From the theory of linear regression analysis it is seen that  $W_i^d(1-r_i^2)\chi_{n-2}^2$ , chi-square with n-2 degrees of freedom, independent of  $V_i$  and  $Y = (Y_1, \ldots, Y_n)'$  and  $V_i^d N(r_i S, 1-r_i^2)$  and

-3-

 $cov(V_i, V_j) = r_{ij} - r_i r_j$ , conditionally given Y. Let  $\lambda_{ij} = r_{ij} - r_i r_j$ and  $\Omega = (\lambda_{ij})$ . It is also seen that  $W_i$  can be represented as the sum of squares of (n-2) orthogonal linear functions of  $X_{i1}, \dots, X_{in}$ . That is

$$W_i \stackrel{d}{=} \sum_{t=1}^{n-2} z_{it}^2$$
 (2.4)

Where  $Z_t = (Z_{1t}, \dots, Z_{kt})'$  are identically and independently distributed as N(0,  $\Omega$ ), independent of V<sub>1</sub>,...,V<sub>K</sub> and S.

Let  $\underline{T} = (\underline{T}_1, \dots, \underline{T}_K)'$  be a random vector distributed as  $N(O, \Omega)$ , independent of S and  $W = (W_1, \dots, W_K)'$ . Then

$$R_{i}^{*} = V_{i} (W_{i})^{-\frac{1}{2}}$$
  
$$d_{x} (T_{i} + r_{i}S) W_{i}^{-\frac{1}{2}}.$$
 (2.5)

Therefore

<u>Theorem 2.1</u>. The joint distribution of the sample correlation coefficients between the predictand and the predictor variables of a multivariate normal distribution is given by (2.5), where  $T \stackrel{d}{\underset{\sim}{\sim}} N(0,\Omega)$ ,  $S^2 \stackrel{d}{\underset{\approx}{\sim}} \chi^2_{n-1}$ , the distribution of W is given by (2.4). Moreover, (S,T,W) are jointly independent.

For large n, W is asymptotically distributed as N((n-2)f, 2(n-2) $\hat{\Omega}$ ) by Lemma 2.1, where  $f = (1-r_1^2, \ldots, 1-r_K^2)'$  and  $\hat{\Omega} = (\lambda_{ij}^2)$ . Therefore

<u>Corollary 2.1</u>. The asymptotic distribution of R\* is given by (2.5), where  $T \stackrel{d}{\approx} N(0,\Omega)$ ,  $s^2 \stackrel{d}{\approx} \chi^2_{n-1}$ ,  $W \stackrel{d}{\approx} N((n-2)f, 2(n-2)\Omega)$  and (S,T,W) are jointly independent.

Let  $\gamma_{ij} = (r_{ij} - r_i r_j)^2 (1 - r_i^2)^{-1} (1 - r_j^2)^{-1}$  and  $\Gamma = (\gamma_{ij})$ .

From (2.5) we have for large n

$$\sqrt{n} (R_{i}^{*} - r_{i}^{*}) = T_{i} + \sqrt{n} r_{i}^{*} ((\frac{(1 - r_{i}^{2})s^{2}}{W_{i}})^{\frac{1}{2}} - 1) + O_{p}(n^{-\frac{1}{2}})$$
$$= T_{i} + \frac{r_{i}^{*}}{\sqrt{2}} (A - B_{i}) + O_{p}(n^{-\frac{1}{2}})$$
(2.6)

Where  $A_{-N}^{d}(0,1)$ ,  $B = (B_1, \ldots, B_K) \stackrel{d}{=} N(0,\Gamma)$ . Moreover, (T,A,B)are jointly independent. Therefore  $\sqrt{n}(R^*-r^*)$  is asymptotically distributed as  $N(0, \Omega + \frac{r_1^{*2}}{2}(\Gamma + E))$  where  $E = (e_{ij})$ ,  $e_{ij} = 1$ . Let

$$= \Omega + \frac{r_{i}^{*2}}{2}(\Gamma + E).$$

The elements of C are given by

$$C_{ii} = 1 - r_i^2 + r_i^{*2}$$

$$C_{ij} = r_{ij} - r_i r_j + r_i^{*} r_j^{*} \left[1 + \frac{(r_{ij} - r_i r_j)^2}{(1 - r_i^2)(1 - r_i^2)}\right]. \quad (2.7)$$

Therefore

<u>Corollary 2.2</u>. For large n,  $\sqrt{n}(R^*-r^*)$  is asymptotically distributed as N(O,C), where C is given by (2.7).

It is interesting to consider the following special cases: (1)  $r_i = 0$ ,  $r_{ij} = 0$ ,  $i \neq j$  for all i and j, that is, the variables Y,  $X_1, \ldots, X_K$  are jointly independent. We have C = Iand  $\sqrt{n} (R^* - r^*) \stackrel{d}{\sim} N(0, I)$ , asymptotically. (2)  $r_i = 0$ ,  $r_{ij} = \rho$ ,  $i \neq j$  for all i and j, that is the predictor variables  $X_1, \ldots, X_K$ are equi-correlated and independent of Y. We have  $C_{ii} = 1$ ,  $C_{ij} = \rho$ ,  $i \neq j$ . (3)  $r_i = \rho$ ,  $r_{ij} = 0$ ,  $i \neq j$  for all i and j, that is the predictor variables are jointly independent and equi-correlated with Y. We have

-5-

$$c_{ii} = 1 - \rho^{2} + \frac{\rho^{2}}{1 - \rho^{2}}$$

$$c_{ij} = \frac{\rho^{2}}{2(1 - \rho^{2})} (1 + \frac{\rho^{4}}{(1 - \rho^{2})^{2}}) - \rho^{2}.$$

Now we consider the problem of selecting the best predictor variable. A standard procedure is to select the variable from the predictor variables corresponding to the largest value of the squared correlation coefficients  $R_1^2, \ldots, R_K^2$  or equivalently  $R_1^{\star 2}, \ldots, R_K^{\star 2}$ . By Corollary 2.2 the probability of a correct selection can be derived for large n from the multivariate normal distribution function. In the special case (1) we have that

n max  $(R_1^{\pm^2}, \ldots, R_K^{\pm^2})$  is distributed as the largest order statistic in a sample of K observations from  $\chi_1^2$ -chi-square with l degree of freedom. This result can be used also to test the significance of the correlation between the selected predictor variable and the predictand. Similar results are obtained for the cases (2) and (3).

-6-

# References

- [1] Alam, K., Rizvi, M. H. and Solomon, H. (1976). Selection of largest multiple correlation coefficients: exact sample size case. <u>Ann. Statist</u>. (4) 614-620.
- [2] Gupta, S. S. (1963). Probability integrals of multivariate normal and multivariate t. <u>Ann. Math. Statist</u>. (34) 792-828.
- [3] Ramberg, J. S. (1977). Selecting the best Predictor variate. Unpublished Manuscript, University of Iowa.
- [4] Rizvi, M. H. and Solomon, H. (1973). Selection of largest multiple correlation coefficients: asymptotic case. J. Amer. Statist. Assoc. (68) 184-188.

ECURITY CLASSIFICATION OF THIS PAGE (When Date Entered 1 NO 6	TR-266
REPORT DOCUMENTATION PAGE	READ INSTRUCTIONS
	BEFORE COMPLETING FORM
N-88	
TITLE (and Subtitio)	S. TYPE OF REPORT & PERIOD COVERED.
Distribution of Sample Correlation	
Coefficients.	Ortechnical rept.
	T. CEALODHING ONC. ASPART HUMBER
AUTHOR(e)	8. CONTRACT OR GRANT NUMBER(*)
Khursheed/Alam	
(1)	N66914-75-C-9451
PERFORMING ORGANIZATION NAME AND ADDRESS	10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
Clemson University V Dept. of Mathematical Sciences	
Clemson, South Carolina 29631	NR 042-271
1. CONTROLLING OFFICE NAME AND ADDRESS	12. REPORT DATE
Office of Naval Research	(11)/Sep 77 /
Code 436 Arlington, Va. 22217	13. NUNDER OF PLES
ATTINGTON, VA. 22217 14. MONITORING AGENCY NAME & ADDRESS(II different from Controlling Office)	15. SECURITY CLASS. (of this report)
	G11-
	Unclassified (2) 11 p
	15. DECLASSIFICATION / DOINGRADING
7. DISTRIBUTION STATEMENT (of the obstract entered in Block 20, il different h	
S. SUPPLEMENTARY NOTES	
9. KEY WORDS (Continue on reverse elde II necessary and identify by block number Multivariate normal distribution; corre selection of variables.	
selection of variables. Sub <sup>1</sup> Sub <sup>1</sup> Sub <sup>1</sup> Sub <sup>1</sup> Let (Y,X <sub>1</sub> ,,X <sub>n</sub> ) be a random vect to a multivariate normal distribution, sample correlation coefficient between has derived the joint distribution of F its asymptotic property.	lation coefficient; $7 \pm 83$ sub- for distributed according and let R, denote the Y and X The author
<ul> <li>S. KEY WORDS (Continue on reverse side if necessary and identify by block number Selection of variables.</li> <li>Multivariate normal distribution; correspondent of variables.</li> <li>AUSTRACT (Continue on reverse side if necessary and identify by block number Let (Y, X,, X, ) be a random vect to a multivariate normal distribution, sample correlation coefficient between thas derived the joint distribution of Fits asymptotic property.</li> </ul>	lation coefficient; $7 \pm 83$ for distributed according and let R <sub>i</sub> denote the Y and X <sub>i</sub> . The author $1, \dots, R_{KA}^{i}$ and has given