

ON THE ASYMPTOTIC DISTRIBUTION OF GOWER'S m^2 GOODNESS-OF-FIT CRITERION IN A PARTICULAR CASE

A. W. DAVIS

(Received Nov. 25, 1976)

1. Summary

The asymptotic behaviour of Gower's m^2 goodness-of-fit criterion is considered in the case of two independent sets of Mahalanobis's D^2 distances, obtained from the same p -variate normal population (equal numbers of samples, equal sample sizes). The asymptotic distribution is approximated by a multiple of chi-square, and some Monte Carlo results are presented to illustrate the approach to this distribution.

2. Introduction

The m^2 (originally R^2) statistic was introduced by Gower ([2], [3]) to compare sets of distances constructed between multivariate samples or populations. For example, we may wish to compare (i) distances derived by different approaches (e.g. Mahalanobis's D^2 and Pearson's coefficient of racial likeness) from the same observations on the same samples, (ii) distances derived by the same or different methods from different subsets of observations on the same samples, or (iii) distances derived from different samples from the same populations.

Given k populations and two sets of distances between them, (d_{ij}) and (d_{ij}^*) ($i, j=1, \dots, k$) say, Gower suggested (a) applying principal coordinate analysis for example, (Gower [1]), to map the distances onto two sets of geometric points (P_i) , (P_i^*) ($i=1, \dots, k$) in Euclidean p -space, in such a way that $d_{ij}/(d_{ij}^*)$ is the Euclidean distance between P_i and P_j (P_i^* and P_j^*); then (b) moving the P_i^* relative to the P_i in p -space by means of translations, rotations, reflections and scalings until the "residual" sum of squares

$$(1) \quad m^2 = \sum_{i=1}^k \delta^2(P_i, P_i^*)$$

Key words: Multivariate analysis; Gower's goodness-of-fit criterion; Asymptotic distribution; Non-central Wishart distribution; Canonical analysis.

is minimum, where δ denotes Euclidean distance.

Gower [2] showed that the required translation consists of shifting the two sets of points to a common centroid, which we shall take to be the origin. If X and Y denote the $k \times p$ matrices whose i th rows are the resulting coordinates of P_i and P_i^* respectively, then the minimum value of (1) following suitable rotation and reflection is

$$(2) \quad m^2 = \text{trace} [X'X + Y'Y - 2(X'YY'X)^{1/2}] .$$

The reader is referred to Gower [3] for the generalization to more than two sets of distances, with allowance for scaling (Generalized Procrustes Analysis).

It was pointed out by Gower [2] that the distributional properties of m^2 in "null" situations are fundamental for any statistical inference based on the above approach. As a starting point for the investigation of these he suggested a particular case of (iii), namely, that in which Mahalanobis's D^2 distances are constructed for two independent sets of samples from the same p -variate normal populations $N(\mu_i, \Sigma)$ ($i=1, \dots, k$), with mean vectors μ_i and common covariance matrix Σ . The rows of X and Y are then the canonical mean vectors (centroid at the origin) arising from separate canonical variate analyses of the two sets of samples. In the present note we shall further assume that the individual samples have equal size n , and show that, as $n \rightarrow \infty$, nm^2 is asymptotically distributed as a central positive definite quadratic form in normal variables. The complexity of this result suggests that the exact distribution of m^2 may be exceedingly difficult to obtain. The asymptotic mean and variance are also derived, and used to obtain an asymptotic approximation $nm^2 \sim c\chi_a^2$.

3. Asymptotic distribution

We shall refer to the two sets of samples as the x -set and y -set respectively, and introduce the following notation for the x -set:

\bar{x}_i = mean vector of the x -sample (size n) from $N(\mu_i, \Sigma)$, ($i=1, \dots, k$)

$\bar{x} = k^{-1} \sum_{i=1}^k \bar{x}_i$ = grand mean vector,

$\mathcal{X}(k \times p) = (\bar{x}_1 - \bar{x}, \dots, \bar{x}_k - \bar{x})'$

S_x = pooled within-samples covariance matrix on $\nu = k(n-1)$ degrees of freedom,

with a corresponding notation for the y -set. We assume $k \geq p$. Now let

$$z'_i = \sqrt{n}(\bar{x}'_i, \bar{y}'_i), \quad (i=1, \dots, k),$$

$$\bar{z} = k^{-1} \sum_{i=1}^k z_i,$$

$$\mathcal{Z}(k \times 2p) = (z_1 - \bar{z}, \dots, z_k - \bar{z})'$$

$$B(2p \times 2p) = \mathcal{Z}'\mathcal{Z} = n \begin{bmatrix} \mathcal{X}'\mathcal{X} & \mathcal{X}'q_j \\ q_j'\mathcal{X} & q_j'q_j \end{bmatrix} = \begin{bmatrix} B_{xx} & B_{xy} \\ B'_{xy} & B_{yy} \end{bmatrix},$$

where B_{xx} , B_{yy} are the between-samples sums of squares and products matrices for the x - and y -sets, respectively.

In carrying out a canonical variate analysis for the x -set, say ([8], Chapter 7), an orthogonal matrix H_x is found such that

$$S_x^{-1/2} B_{xx} S_x^{-1/2} = H_x \Lambda_x H_x',$$

where Λ_x is the diagonal matrix of eigenvalues of the left-hand side matrix. Defining

$$X(k \times p) = \mathcal{X} S_x^{-1/2} H_x,$$

the rows of X are seen to be the canonical mean vectors of the x -samples, with centroid at the origin as required. Similarly, we construct

$$Y(k \times p) = q_j S_y^{-1/2} H_y$$

for the y -set, and Gower's m^2 statistic (equation (2)) takes the form

$$(3) \quad m^2 = n^{-1} \text{trace} [B_{xx} S_x^{-1} + B_{yy} S_y^{-1} - 2 \{S_y^{-1/2} B'_{xy} S_y^{-1} B_{xy} S_y^{-1/2}\}^{1/2}].$$

(Note that if H is orthogonal and A is positive definite symmetric, then $\text{trace}(H A H')^{1/2} = \text{trace} H A^{1/2} H' = \text{trace} A^{1/2}$.)

To discuss the asymptotic properties of this quantity for large n , we first note that since the z_i are independent $2p$ -variate normal vectors with means $\sqrt{n}(\mu'_i, \mu'_i)'$ and covariance matrix

$$(4) \quad \begin{bmatrix} \Sigma & O_p \\ O_p & \Sigma \end{bmatrix},$$

where O_p is the $p \times p$ zero matrix, the matrix B has the non-central Wishart distribution [6] with $q = k - 1$ degrees of freedom, population covariance matrix (4), and matrix of non-centrality parameters $n\Omega/2$, where

$$(5) \quad \begin{aligned} \Omega &= \begin{bmatrix} \Theta & \Theta \\ \Theta & \Theta \end{bmatrix}, \\ \Theta &= \Sigma^{-1/2} \left\{ \sum_{i=1}^k (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})' \right\} \Sigma^{-1/2}, \\ \bar{\mu} &= k^{-1} \sum_{i=1}^k \mu_i. \end{aligned}$$

Now let K be an orthogonal matrix reducing Θ to diagonal form,

$$(6) \quad \Theta = K \Lambda K', \quad \Lambda = \text{diag}(\theta_1, \dots, \theta_p).$$

It will be assumed that Θ is positive definite, that is, all $\theta_i > 0$. If we transform to new variables

$$x^* = K \Sigma^{-1/2} x, \quad y^* = K \Sigma^{-1/2} y,$$

then m^2 is given by (3) in terms of the corresponding "starred" quantities; νS_x^* and νS_y^* have central Wishart distributions with ν degrees of freedom and unit population covariance matrices I_p , and B^* has a non-central Wishart distribution with q degrees of freedom, population covariance matrix I_{2p} , and non-centrality matrix $n\Omega^*/2$, where

$$(7) \quad \Omega^* = \begin{bmatrix} \Lambda & \Lambda \\ \Lambda & \Lambda \end{bmatrix} = (\omega_{ij}).$$

Now

$$(8) \quad E(B^*) = n\Omega^* + qI_{2p}$$

(see for example [4]), so that we may write

$$(9) \quad n^{-1}B^* = \Omega^* + (q/n)I_{2p} + U,$$

where

$$(10) \quad U = \begin{bmatrix} U_{xx} & U_{xy} \\ U'_{xy} & U_{yy} \end{bmatrix} = (U_{ij})$$

has zero expectation. Also let

$$(11) \quad S_x^* = I_p + V_x, \quad (V_x = (V_{ij}^x))$$

with a similar notation for S_y^* . Then U , V_x and V_y are $O(n^{-1/2})$ (equations (A7), (A8) in the Appendix), and to order n^{-1} (Appendix (a))

$$(12) \quad \begin{aligned} m^2 \sim & 2pq/n + \sum_{i=1}^p (U_{ii} + U_{i+p, i+p} - 2U_{i, i+p}) \\ & - \sum_{i=1}^p \sum_{j=1}^p (U_{ij} V_{ij}^x + U_{i+p, j+p} V_{ij}^y) \\ & + \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p (\theta_i + \theta_j)^{-1} \{ -(U_{i, j+p} U_{j, i+p})^2 \\ & + 4U_{i, j+p} (\theta_j V_{ij}^x + \theta_i V_{ij}^y) + \theta_i \theta_j (V_{ij}^x - V_{ij}^y)^2 \}. \end{aligned}$$

It follows from (9) that the first two terms sum to trace $(\mathcal{X}^* - \mathcal{Y}^*)(\mathcal{X}^* - \mathcal{Y}^*)$, which is the trace of a central Wishart matrix with q degrees of freedom and population covariance matrix $(2/n)I_p$, and is

thus distributed as $(2/n)\chi_{pq}^2$. Considering the remaining terms, we note that the elements of the non-centrality matrix $n\mathbf{Q}^*/2 \rightarrow \infty$ as $n \rightarrow \infty$, so that the variates $\sqrt{n}U_{ij}$ are asymptotically jointly normal [4] with zero means, and covariances given by equation (A7). The $\sqrt{n}V_{ij}^x$ and $\sqrt{n}V_{ij}^y$ also approach joint normality, with zero means and covariances given by (A8), and hence nm^2 is asymptotically distributed as a central quadratic form in normal variables ([5], Chapter 29). Since nm^2 is non-negative, this form is necessarily positive-definite, although it seems difficult to show this explicitly. Generally such distributions are well approximated by multiples of chi-square, $c\chi_\alpha^2$, on non-integer degrees of freedom, and in order to evaluate c and α it is sufficient to have the asymptotic mean μ'_1 and variance μ_2 of nm^2 ; these are derived in Appendix (b):

$$(13) \quad \mu'_1 = 2 \left[p \left(q - \frac{1}{2}(p-1) \right) + k^{-1} \sum_{1 \leq i \leq j \leq p} \frac{\theta_i \theta_j}{\theta_i + \theta_j} \right],$$

$$(14) \quad \mu_2 = 8 \left[p \left(q - \frac{1}{2}(p-1) \right) + k^{-1}(p+1) \sum_{i=1}^p \theta_i \right. \\ \left. - k^{-1} \sum_{1 \leq i \leq j \leq p} \frac{\theta_i^2 + \theta_j^2}{\theta_i + \theta_j} + k^{-2} \sum_{1 \leq i \leq j \leq p} \frac{\theta_i^2 \theta_j^2}{(\theta_i + \theta_j)^2} \right].$$

Then

$$(15) \quad c = \mu_2 / 2\mu'_1, \quad \alpha = 2(\mu'_1)^2 / \mu_2.$$

4. Some Monte Carlo results

Ten 5-variate population mean vectors were constructed by sampling uniformly over a hypersphere of radius 6 (Table 2), and 200 x -sets and y -sets of samples (sizes 50, 100, 200) were generated with population covariance matrices $\mathbf{\Sigma} = \mathbf{I}_5$. The matrix $\mathbf{\Theta}$ thus reduces to the sum of squares and products matrix of the mean vectors,

$$(16) \quad \mathbf{\Theta} = \begin{bmatrix} 72.96 & 13.30 & 6.29 & -3.55 & 6.56 \\ * & 46.40 & 11.71 & -23.88 & -8.91 \\ * & * & 35.82 & -11.09 & -23.88 \\ * & * & * & 27.65 & 7.74 \\ * & * & * & * & 43.31 \end{bmatrix}$$

with latent roots

$$(17) \quad \theta_1 = 88.89, \quad \theta_2 = 71.72, \quad \theta_3 = 40.73, \quad \theta_4 = 14.66, \quad \theta_5 = 10.14.$$

Canonical analyses were carried out, and 200 values of nm^2 were

calculated for each sample size using the ROTATE directive in the program GENSTAT [7].

Table 1 shows the means and variances of the sampled values, together with the asymptotic values obtained from (13), (14) and (17). The corresponding approximate asymptotic distribution (equation (15)) was

$$(18) \quad 5.273\chi_{\alpha}^2, \quad \text{with } \alpha=23.680.$$

As a rough indication of the approach of the simulated distributions to (18), Pearson's goodness-of-fit criterion was calculated for each set of 200 values, based on the deciles of (18). By sample size 50, (18) appears to be giving a reasonable approximation to the overall shape of the curve. It would be desirable to investigate the approximation in higher dimensions.

Table 1 Approach of nm^2 to the approximate asymptotic distribution
(10 5-variate populations, 200 trials)

n	Mean (nm^2)	Var (nm^2)	Pearson χ^2	P
50	129.33 (2.86)	1639 (333)	12.8	0.17
100	125.01 (2.51)	1264 (145)	4.0	0.91
200	129.30 (2.65)	1407 (144)	4.2	0.90
∞	124.86	1317		

Figures in brackets denote standard errors.

Table 2 Mean vectors used in Table 1

Popula- tion	Mean Vector				
1	2.1368	4.3781	-0.3241	-2.2281	0.9611
2	-2.6930	-1.3939	-1.7595	-0.1107	3.4107
3	0.1313	3.2676	1.9097	-4.1451	-1.1622
4	3.4693	2.7769	1.4386	-0.9864	-2.1049
5	-2.8560	-1.0431	-3.5809	-0.0487	-0.4427
6	-3.6216	0.4996	2.0489	0.8789	-3.0896
7	0.1676	1.8274	-2.1315	1.5259	3.1300
8	1.8050	-1.7579	1.4265	0.3312	1.4627
9	4.1422	-1.9152	-0.9020	1.5625	0.1427
10	3.3105	1.3766	-1.9964	-0.4910	2.3302

Appendix

(a) *Derivation of (12)*: Write

$$(A1) \quad S_x^{*-1} = I_p + C\mathcal{V}_x, \quad S_x^{*-1/2} = I_p + C\bar{\mathcal{V}}_x (C\mathcal{V}_x = (C\mathcal{V}_{ij}^x), C\bar{\mathcal{V}}_x = (C\bar{\mathcal{V}}_{ij}^x)),$$

with a similar notation for S_y^* . Equations (9) and (A1) are to be substituted in the starred version of (3), retaining terms of order 1, $n^{-1/2}$ and n^{-1} ; thus

$$(A2) \quad n^{-1} \text{trace} [B_{xx}^* S_x^{*-1}] = \sum_{i=1}^p \theta_i + pq/n + \sum_{i=1}^p [U_{ii} + (\theta_i + q/n) C\mathcal{V}_{ii}^x] \\ + \sum_{i=1}^p \sum_{j=1}^p U_{ij} C\mathcal{V}_{ij}^*,$$

and a similar expression is obtained for $n^{-1} \text{trace} [B_{yy}^* S_y^{*-1}]$.

To derive a corresponding result for the final term in (3), let

$$(A3) \quad C^2 = (I_p + C\bar{\mathcal{V}}_x)(A + U_{xy})(I_p + C\mathcal{V}_y)(A + U'_{xy})(I_p + C\bar{\mathcal{V}}_x) \\ = A^2 + A, \quad (A = (A_{ij})),$$

say, where

$$C = \sum_{l=0}^{\infty} C_l, \quad C_l = (C_{ij}^{(l)}) = O(n^{-l/2}),$$

($l=0, 1, 2, \dots$) and each C_l is symmetric. Equating terms of like order,

$$C_0 = A,$$

$$C_1 C_0 + C_0 C_1 = A,$$

$$C_2 C_0 + C_1^2 + C_0 C_2 = O_p, \quad \text{etc.}$$

Hence

$$C_{ij}^{(1)} = (\theta_i + \theta_j)^{-1} A_{ij},$$

$$C_{ij}^{(2)} = -(\theta_i + \theta_j)^{-1} \sum_{l=1}^p [A_{il} A_{jl} / (\theta_i + \theta_l)(\theta_j + \theta_l)],$$

etc., so that

$$(A4) \quad 2 \text{trace } C = 2 \sum_{i=1}^p \theta_i + \sum_{i=1}^p A_{ii} / \theta_i - \sum_{i=1}^p \sum_{j=1}^p A_{ij}^2 / \theta_i (\theta_i + \theta_j)^2 + \dots$$

To terms of order $n^{-1/2}$ and n^{-1} ,

$$(A5) \quad A_{ij} = [\theta_i U_{i,j+p} + \theta_j U_{j,i+p} + (\theta_i^2 + \theta_j^2) C\mathcal{V}_{ij}^x + \theta_i \theta_j C\mathcal{V}_{ij}^y] \\ + \sum_{l=1}^p [U_{i,l+p} U_{j,l+p} + \theta_i U_{i,l+p} C\bar{\mathcal{V}}_{jl}^x + \theta_j U_{j,l+p} C\bar{\mathcal{V}}_{il}^x]$$

$$+ \theta_i U_{j, l+p} \mathcal{V}_{il}^y + \theta_j U_{i, l+p} \mathcal{V}_{jl}^y + \theta_i \{ U_{i, l+p} \bar{\mathcal{V}}_{jl}^x \\ + U_{j, l+p} \bar{\mathcal{V}}_{il}^x + \theta_i \bar{\mathcal{V}}_{jl}^x \mathcal{V}_{il}^y + \theta_j \bar{\mathcal{V}}_{il}^x \mathcal{V}_{jl}^y + \theta_i \bar{\mathcal{V}}_{il}^x \mathcal{V}_{jl}^x \} ,$$

and substituting in (A4), the resulting expression takes its most convenient form if we replace \mathcal{V}_x by $\mathcal{V}_x/2 - \mathcal{V}_x^2/8$:

$$(A6) \quad 2 \text{ trace } C = 2 \sum_{i=1}^p \theta_i + \sum_{i=1}^p [2 U_{i, i+p} + \theta_i (\mathcal{V}_{ii}^x + \mathcal{V}_{ii}^y)] \\ + \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p (\theta_i + \theta_j)^{-1} [(U_{i, j+p} - U_{j, i+p})^2 \\ + 4 U_{i, j+p} (\theta_j \mathcal{V}_{ij}^x + \theta_i \mathcal{V}_{ij}^y) - \theta_i \theta_j (\mathcal{V}_{ij}^x - \mathcal{V}_{ij}^y)^2] .$$

Finally, (12) is obtained from (A2) and (A6). Quantities of order 1 and $n^{-1/2}$ cancel, and to order n^{-1} it has been possible to replace \mathcal{V}_x , \mathcal{V}_y by $-V_x$, $-V_y$ respectively.

(b) *Derivation of mean and variance*: For (13), we require only

$$(A7) \quad E[U_{ij} U_{kl}] = \kappa(ij, kl)/n - q\lambda(ij, kl)/n^2$$

$$(A8) \quad E[V_{ij}^x V_{kl}^x] = \lambda(ij, kl)/\nu ,$$

(Jensen [4]), where

$$(A9) \quad \kappa(ij, kl) = \delta_{ik}\omega_{jl} + \delta_{il}\omega_{jk} + \delta_{jk}\omega_{il} + \delta_{jl}\omega_{ik} , \\ \lambda(ij, kl) = \delta_{ik}\delta_{jl} + \delta_{il}\delta_{jk} ,$$

and δ_{ij} is Kronecker's delta. Note that U , V_x and V_y are independently distributed.

In deriving the variance of (12), the $2pq/n$ term may be ignored, and the remainder consists of five terms, 1, 2, ..., 5 say, defined by the bracketing. Let (i, j) denote the covariance of terms i and j . Then from the independence, these vanish apart from the five variances, (1, 3), and (2, 4), and to evaluate these we require the following formulae in addition to (A7) and (A8):

$$(A10) \quad E[U_{ij} U_{kl} U_{mn}] = n^{-2} \sum \omega_{ik}\lambda(jl, mn) + O(n^{-3}) ,$$

where the summation extends over the twelve possible selections of the subscripts of ω from distinct pairs of U 's; and

$$(A11) \quad E[U_{ij} U_{kl} U_{mn} U_{pq}] = n^{-2} \sum \kappa(ij, kl)\kappa(mn, pq) + O(n^{-3}) ,$$

$$(A12) \quad E[V_{ij}^x V_{kl}^x V_{mn}^x V_{pq}^x] = \nu^{-2} \sum \lambda(ij, kl)\lambda(mn, pq) + O(n^{-3}) ,$$

where the summations extend over the three distinct arrangements of (ij) , (kl) , (mn) and (pq) into pairs. These results may be obtained from

the moment generating function of the non-central Wishart distribution ([6], p. 175). The required covariances are found to be

$$\begin{aligned}
 (1, 1) &= 8pq/n^2, & (2, 2) &= (8/n\nu)(p+1) \sum_{i=1}^p \theta_i \\
 (3, 3) &= 4p(p-1)/n^2 = -(1, 3)/2, \\
 (4, 4) &= (8/n\nu) \sum_{1 \leq i \leq j \leq p} (\theta_i^2 + \theta_j^2)/(\theta_i + \theta_j) = -(2, 4)/2 \\
 (5, 5) &= (8/\nu^2) \sum_{1 \leq i \leq j \leq p} \theta_i^2 \theta_j^2 / (\theta_i + \theta_j)^2.
 \end{aligned}
 \tag{A13}$$

On summing we obtain (14).

Acknowledgments

The author's thanks are due to Mr. J. C. Gower for suggesting this investigation and to Mr. L. G. Veitch for assistance with the programming.

CSIRO DIVISION OF MATHEMATICS AND STATISTICS, ADELAIDE

REFERENCES

- [1] Gower, J. C. (1966). Some distance properties of latent and vector methods used in multivariate analysis, *Biometrika*, **53**, 325-338.
- [2] Gower, J. C. (1971). Statistical methods of comparing different multivariate analyses of the same data, *Mathematics in the Archaeological and Historical Sciences*, Hodson, F. R., Kendall, D. G., and Tautu, P. (Eds.), Edinburgh University Press, 138-149.
- [3] Gower, J. C. (1975). Generalized Procrustes Analysis, *Psychometrika*, **40**, (in the press).
- [4] Jensen, D. R. (1972). The limiting form of the non-central Wishart distribution, *Aust. J. Statist.*, **14**, 10-16.
- [5] Johnson, N. L. and Kotz, S. (1970). *Distributions in Statistics: Continuous Univariate Distributions-2*, Houghton Mifflin, Boston.
- [6] Johnson, N. L. and Kotz, S. (1972). *Distributions in Statistics: Continuous Multivariate Distributions*, Wiley, New York.
- [7] Nelder, J. A. and Members of the Rothamsted Statistics Department (1975). *Genstat Reference Manual*, Inter-University/Research Council Series, Report No. 3, Third Edition, Edinburgh Regional Computing Centre.
- [8] Seal, H. (1974). *Multivariate Statistical Analysis for Biologists*, Methuen, London.