

Wiesław W a g n e r (Poznań)

O TESTOWANIU RÓWNOŚCI K MACIERZY KOWARIANCJI

1. Wstęp

W doświadczeniach rolniczych i przyrodniczych obserwuje się zwykle więcej niż jedną cechę na każdej jednostce doświadczalnej. Stąd wyniki z takich doświadczeń mają przeważnie charakter wielowymiarowy. Ponieważ cechy te są najczęściej skorelowane, więc analizowanie każdej z obserwowanych cech oddzielnie może okazać się w wielu przypadkach mało efektywne. Jeśli badacza interesuje kompleksowe ujęcie obserwowanych cech, wówczas naturalnym staje się przejście od jednozmiennnej analizy wariancji (ang. ANOVA) do wielozmiennnej analizy wariancji (ang. MANOVA), badającej wpływ różnych źródeł zmienności na ogólną zmienność wielu cech. Istotnym problemem staje się wykorzystanie metod wielozmiennnej analizy, które umożliwiają uzyskanie informacji dotyczących p skorelowanych zmiennych w jednej lub więcej populacjach, z których pobrano losowe próby. Przy wielu metodach przyjmuje się założenie, że te zmienne mają p-wymiarowe rozkłady normalne o tych samych macierzach kowariancji. Jednak taka jednorodność nie zawsze zachodzi i w związku z tym konieczne jest sprawdzenie hipotez dotyczących równości macierzy kowariancji.

2. Jednokierunkowa klasyfikacja obserwacji

Wyobraźmy sobie, że porównujemy k-populacji pod względem p różnych cech. Z każdej z tych populacji pobrano po jednej próbie losowej. Oznaczmy liczebność próby z g-tej populacji przez N'_g ($g =$

= 1, 2, ..., k), a sumę wszystkich liczebności, to znaczy liczbę wszystkich jednostek doświadczalnych przez N' :

$$N' = N'_1 + N'_2 \dots + N'_k.$$

Na każdej jednostce doświadczalnej obserwowanych jest p cech. Ma my wówczas N' obserwacji p -wymiarowych sklasyfikowanych według jed- nego kryterium - przynależności do jednej z k populacji. Oznaczmy j -tą obserwację ($j = 1, \dots, N'_g$) i -tej cechy ($i = 1, \dots, p$) w pró- bie pobranej w g -tej populacji ($g = 1, \dots, k$) przez X_{gji} . Przy tych oznaczeniach macierz obserwacji dla g -tej próby można zapi- sać w postaci

$$(1) \quad \underline{X}_g = \begin{bmatrix} x_{g11} & x_{g12} & \dots & x_{g1p} \\ x_{g21} & x_{g22} & \dots & x_{g2p} \\ \dots & \dots & \dots & \dots \\ x_{gN'_g1} & x_{gN'_g2} & \dots & x_{gN'_gp} \end{bmatrix},$$

gdzie N'_g wierszy odpowiada kolejnym obserwacjom, a p kolumn od- powiada kolejnym cechom. Macierz

$$(2) \quad \underline{X} = (\underline{X}'_1, \underline{X}'_2, \dots, \underline{X}'_k)',$$

typu $N' \times p$, złożona z podmacierzy określonych w (1), jest wówczas macierzą wszystkich p -wymiarowych obserwacji. Zakładamy, że wier- sze macierzy \underline{X}_g są wzajemnie niezależnymi p -wymiarowymi zmiennymi losowymi o jednakowym rozkładzie normalnym $N(\underline{\mu}_g, \underline{\Sigma}_g)$, gdzie wek- tor wartości średnich

$$(3) \quad \underline{\mu}_g = [\mu_{g1}, \mu_{g2}, \dots, \mu_{gp}]',$$

a macierz kowariancji

$$(4) \quad \underline{\Sigma}_g = \begin{bmatrix} \sigma_{g11} & \sigma_{g12} & \dots & \sigma_{g1p} \\ \sigma_{g21} & \sigma_{g22} & \dots & \sigma_{g2p} \\ \dots & \dots & \dots & \dots \\ \sigma_{gp1} & \sigma_{gp2} & \dots & \sigma_{gpp} \end{bmatrix} \quad (g = 1, 2, \dots, k).$$

Przyjmując takie założenie dopuszczamy istnienie różnych wektorów wartości średnich i różnych macierzy kowariancji dla różnych po- pulacji.

Przy wyborze odpowiedniej metody do analizy obserwacji wielo- wymiarowych sklasyfikowanych według jednego kryterium - według po-

populacji, interesuje nas to, czy dopuszczalne jest przyjęcie upraszczającego założenia, że wszystkie macierze $\underline{\Sigma}_1, \underline{\Sigma}_2, \dots, \underline{\Sigma}_k$ są jednakowe.

3. Rozkład statystyki M

Hipotezę zerową o jednorodności między populacjami w sensie równości macierzy kowariancji możemy zapisać w postaci:

$$H_0 : \underline{\Sigma}_1 = \underline{\Sigma}_2 = \dots = \underline{\Sigma}_k .$$

Przy sprawdzaniu tej hipotezy stosuje się statystykę M Bartletta postaci

$$(5) \quad M = N \ln |\underline{S}| - \sum_{g=1}^k N_g \ln |\underline{S}_g|$$

gdzie $N = N_1 + N_2 + \dots + N_k$, $N_g = N'_g - 1$, macierz \underline{S}_g typu $p \times p$ jest nieobciążonym estymatorem macierzy kowariancji dla g -tej populacji o N_g stopniach swobody, a $N\underline{S} = N_1\underline{S}_1 + N_2\underline{S}_2 + \dots + N_k\underline{S}_k$.

Anderson [1] wyprowadził kryterium M modyfikując test oparty na ilorazie wiarygodności.

Kulback [5] natomiast pokazał, że kryterium M jest statystyką minimalnej informacji dyskryminacyjnej dla hipotezy zerowej i wykazuje, że przy hipotezie alternatywnej o niejednorodności macierzy kowariancji przyjmuje ona wartość minimalną.

W oparciu o wyniki Boxa [2] Korin [3], [4] zbudował tablice 5-procentowych punktów istotności dla rozkładu statystyki M uwzględniając różne wartości parametrów k, p oraz N_0 , gdzie $N_0 = N_1 = N_2 = \dots = N_k$. Podane przez Korina tablice [4] obejmują jednak zbyt mały zakres zmieniającej się wielkości N_0 . Często zachodzi potrzeba odczytania punktu istotności dla większej liczby obserwacji N_0 . Wówczas z konieczności musimy szukać przybliżonych punktów istotności dla M. Jeśli spełniony jest warunek $N_1 = N_2 = \dots = N_k = N_0$ to możemy je znaleźć używając jednego z następujących przybliżeń Boxa [2] dla rozkładu M:

(a) Przybliżenie rozkładem χ^2 . Tutaj korzystamy z tego, że statystyka $(1 - A_1) \cdot M$ ma w przybliżeniu rozkład chi-kwadrat o f_1 stopniach swobody, tzn.

$$(6) \quad M \sim \chi^2_{f_1} / (1 - A_1),$$

gdzie liczba stopni swobody

$$(7) \quad f_1 = \frac{1}{2} p(p+1)(k-1),$$

natomiast

$$(8) \quad A_1 = \frac{2p^2 + 3p - 1}{6(p+1)kN_0}.$$

(b) Przybliżenie rozkładem F. Tutaj wykorzystuje się to, że statystyka M/b ma w przybliżeniu rozkład F o f_1 i f_2 stopniach swobody, tzn.

$$(9) \quad M \sim b F_{f_1; f_2},$$

gdzie f_1 i A_1 są określone w (7) i (8), zaś

$$(10) \quad f_2 = \text{entier} \frac{f_1 + 2}{A_2 - A_1},$$

przy czym

$$(11) \quad A_2 = \frac{(p-1)(p+2)(k^2+k+1)}{6k^2N_0^2},$$

natomiast

$$(12) \quad b = \frac{f_1}{1 - A_1 - f_1/f_2}$$

Obliczenie wartości statystyki M poprzedzone jest znalezieniem macierzy

$$(13) \quad \underline{S} = N^{-1}(N_1 \underline{S}_1 + N_2 \underline{S}_2 + \dots + N_k \underline{S}_k),$$

gdzie macierz \underline{S}_g jest macierzą sum kwadratów i iloczynów odchyłeń dla próby z j-tej populacji ($g = 1, \dots, k$), to znaczy

$$(14) \quad \underline{S}_g = \begin{bmatrix} s_{g11} & s_{g12} & \dots & s_{g1p} \\ s_{g21} & s_{g22} & \dots & s_{g2p} \\ \dots & \dots & \dots & \dots \\ s_{gp1} & s_{gp2} & \dots & s_{gpp} \end{bmatrix}$$

gdzie

$$s_{gii'} = \frac{1}{N_g} \sum_{j=1}^{N_g'} x_{gji} x_{gji'} - \frac{1}{N_g^2} T_{gi} T_{gi'},$$

przy czym

$$T_{gi} = \sum_{j=1}^{N_g'} x_{gji} \quad \text{jest sumą wszystkich obserwacji dotyczących } i\text{-tej}$$

cechy ($i, i' = 1, \dots, p$) dla próby z g-tej populacji.

Wartość funkcji testowej M porównuje się z wartością krytycz-

ną odczytaną z tablic Korina [4] dla poziomu istotności 5%. W przypadku większej liczby obserwacji można posłużyć się przybliżeniem rozkładem χ^2 lub F (patrz wzory (6), (9)). W przypadku, gdy wartość M przekroczy wartość krytyczną przy ustalonym poziomie istotności, wówczas odrzucamy hipotezę zerową o równości macierzy kowariancyjnych.

4. Przykład

Jako przykład niech posłuży analiza danych pochodzących z badań hodowlanych nad słonecznikiem oleistym, przeprowadzonych w SHB Borowo. Na dane te składają się obserwacje dotyczące 5 cech u 7 rodów słonecznika, z których wylosowano po 28 roślin. Korzystając ze wzorów (13) i (14) wyliczono macierze \underline{S} i \underline{S}_g ($g = 1, \dots, 7$), a następnie uzyskano wartość statystyki M równą 26,2056.

Ponieważ tablice [4] nie zawierają wartości krytycznych dla poziomu istotności 5% i parametrów $N_0 = 27$, $p = 5$ i $k = 7$, dlatego też zostały one wyznaczone według wzorów (6) i (9). Wyniki zestawiono w tablicy 1.

T a b l i c a 1. Wartości krytyczne na poziomie 5%
dla parametrów $k = 7$, $p = 5$, $N_0 = 27$

Przybliżenie	χ^2 wg wzoru (6)	122,352
Przybliżenie	F wg wzoru (9)	124,593

Wartość krytyczną dla rozkładów χ^2 i F przy poziomie istotności 5% dla odpowiedniej liczby stopni swobody odczytano z tablic [6]. Ponieważ wartość M jest znacznie mniejsza od wartości krytycznych podanych w tablicy 1, więc nie mamy podstaw do odrzucenia hipotezy zerowej o jednorodności macierzy kowariancyjnych.

Obliczenia do tego przykładu wykonano według programu ułożonego w języku ODRA ALGOL na emc ODRA 1204.

Literatura cytowana

- [1] Anderson, T. W., An Introduction to Multivariate Statistical Analysis, New York 1958.
- [2] Box, G. E. P., A general distribution theory for a class of likelihood criteria, *Biometrika*, 36 (1949), str. 317-346.
- [3] Korin, B. P., On the distribution of a statistic used for testing a covariance matrix, *Biometrika*, 55 (1968), str. 171-179.
- [4] Korin, B. P., On testing the equality of k covariance matrices, *Biometrika*, 56 (1969), str. 216-219.
- [5] Kulback, S., Information Theory and Statistics, New York 1959.
- [6] Pearson, E. S. and Hartley, H. O., *Biometrika Tables for Statisticians*, vol. 1, Cambridge 1966.