



A two sample test in high dimensional data

Muni S. Srivastava^a, Shota Katayama^{b,*}, Yutaka Kano^b

^a Department of Statistics, University of Toronto, 100 St. George Street, Toronto, Ontario M5S 3G3, Canada

^b Graduate School of Engineering Science, Osaka University, 1-3 Machikaneyama, Toyonaka, Osaka 560-8531, Japan

ARTICLE INFO

Article history:

Received 6 August 2011

Available online 27 August 2012

AMS subject classifications:

primary 62H15

secondary 62E20

Keywords:

High-dimensional data

Hypothesis testing

Behrens–Fisher problem

Asymptotic theory

ABSTRACT

In this paper we propose a test for testing the equality of the mean vectors of two groups with unequal covariance matrices based on N_1 and N_2 independently distributed p -dimensional observation vectors. It will be assumed that N_1 observation vectors from the first group are normally distributed with mean vector μ_1 and covariance matrix Σ_1 . Similarly, the N_2 observation vectors from the second group are normally distributed with mean vector μ_2 and covariance matrix Σ_2 . We propose a test for testing the hypothesis that $\mu_1 = \mu_2$. This test is invariant under the group of $p \times p$ nonsingular diagonal matrices. The asymptotic distribution is obtained as $(N_1, N_2, p) \rightarrow \infty$ and $N_1/(N_1 + N_2) \rightarrow k \in (0, 1)$ but N_1/p and N_2/p may go to zero or infinity. It is compared with a recently proposed non-invariant test. It is shown that the proposed test performs the best.

© 2012 Elsevier Inc. All rights reserved.

1. Introduction

Let \mathbf{x}_{ij} be independently distributed as the multivariate normal distribution with the mean vector μ_i and the positive definite covariance matrix Σ_i for $i = 1, 2$ and $j = 1, 2, \dots, N_i$. For notational convenience, we shall denote it as $N_p(\mu_i, \Sigma_i)$, $i = 1, 2$, where p denotes the dimension of the random vectors \mathbf{x}_{ij} . In this article, we consider the problem of testing the hypothesis

$$H : \mu_1 = \mu_2 \quad (1.1)$$

against the alternative

$$A : \mu_1 \neq \mu_2, \quad (1.2)$$

when the covariance matrices Σ_1 and Σ_2 of the two groups may be unequal. This problem has recently been considered by Chen and Qin [2] who proposed a test which we denote by T_{cq} . The test T_{cq} will be described in Section 2 from which it will be clear that it is a rather complicated test and requires considerable terms in programming and computing. Also, it is shown that the T_{cq} test is almost identical to a test that can be obtained by generalizing the Bai and Saranadasa [1] test when $\Sigma_1 \neq \Sigma_2$. In addition, the test T_{cq} , although invariant under the group of orthogonal transformations, is not invariant under the units of measurements. That is, if we consider $\mathbf{D}\mathbf{x}_{ij}$ instead of \mathbf{x}_{ij} , where \mathbf{D} is a nonsingular $p \times p$ diagonal matrix, the test T_{cq} changes, which is an undesirable feature. It may be noted that when N_i is less than p , no fully affine invariant test exists. Thus, in this article, we propose a test that is invariant under the transformation of the observation vector \mathbf{x}_{ij} by nonsingular

* Corresponding author.

E-mail addresses: srivasta@utstat.utoronto.ca (M.S. Srivastava), sfujimoto@sigmath.es.osaka-u.ac.jp (S. Katayama), kano@sigmath.es.osaka-u.ac.jp (Y. Kano).

$p \times p$ diagonal matrices. It will be shown that this new test, denoted by T , performs better than T_{cq} . To describe this new test T , we introduce some notations with $n_i = N_i - 1$, $i = 1, 2$:

$$\bar{\mathbf{x}}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} \mathbf{x}_{ij} \quad \text{and} \quad \mathbf{S}_i = \frac{1}{n_i} \sum_{j=1}^{N_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)'. \tag{1.3}$$

In high dimensional data, since N_i may be less than p , the sample covariance matrices \mathbf{S}_i may be singular. However, the diagonal matrices consisting of only the diagonal elements of $\mathbf{S}_i = (s_{ijk})$, $i = 1, 2$, namely,

$$\hat{\mathbf{D}}_i = \text{diag}(s_{i11}, \dots, s_{ipp}), \quad i = 1, 2, \tag{1.4}$$

are non-singular matrices. Let

$$\hat{\mathbf{D}} = \frac{\hat{\mathbf{D}}_1}{N_1} + \frac{\hat{\mathbf{D}}_2}{N_2} = (\hat{d}_{ij}). \tag{1.5}$$

Then

$$\mathbf{R} = \hat{\mathbf{D}}^{-1/2} \left(\frac{\mathbf{S}_1}{N_1} + \frac{\mathbf{S}_2}{N_2} \right) \hat{\mathbf{D}}^{-1/2} = (r_{ij}) \tag{1.6}$$

is the sample correlation matrix, while \mathbf{S}_i may not converge to Σ_i in probability since N_i may be less than p , $\hat{\mathbf{D}}_i$ converges in probability to \mathbf{D}_i , where

$$\mathbf{D}_i = \text{diag}(\sigma_{i11}, \dots, \sigma_{ipp}), \quad \Sigma_i = (\sigma_{ijk}), \quad i = 1, 2, \tag{1.7}$$

if $\max_{1 \leq k \leq p} \sigma_{ikk} < \infty$ uniformly in p . Let

$$\mathbf{D} = \frac{\mathbf{D}_1}{N_1} + \frac{\mathbf{D}_2}{N_2} = (d_{ij}). \tag{1.8}$$

Then, $\hat{\mathbf{D}} \rightarrow \mathbf{D}$ in probability. Similar to the sample correlation matrix \mathbf{R} , we define the population correlation matrix \mathcal{R} by

$$\mathcal{R} = \mathbf{D}^{-1/2} \left(\frac{\Sigma_1}{N_1} + \frac{\Sigma_2}{N_2} \right) \mathbf{D}^{-1/2} = (\rho_{ij}). \tag{1.9}$$

We note that under the null hypothesis H in (1.1),

$$\begin{aligned} E[(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{D}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)] &= \text{tr} \mathbf{D}^{-1} \left(\frac{\Sigma_1}{N_1} + \frac{\Sigma_2}{N_2} \right) \\ &= \text{tr} \mathcal{R} = p. \end{aligned}$$

Also, under the null hypothesis H in (1.1),

$$\begin{aligned} \text{Var}[(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{D}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)] &= \text{Var}(\bar{\mathbf{x}}_1' \mathbf{D}^{-1} \bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2' \mathbf{D}^{-1} \bar{\mathbf{x}}_2 - 2\bar{\mathbf{x}}_1' \mathbf{D}^{-1} \bar{\mathbf{x}}_2) \\ &= \text{Var}(\bar{\mathbf{x}}_1' \mathbf{D}^{-1} \bar{\mathbf{x}}_1) + \text{Var}(\bar{\mathbf{x}}_2' \mathbf{D}^{-1} \bar{\mathbf{x}}_2) + 4\text{Var}(\bar{\mathbf{x}}_1' \mathbf{D}^{-1} \bar{\mathbf{x}}_2) \\ &= \frac{2\text{tr}(\mathbf{D}^{-1} \Sigma_1)^2}{N_1^2} + \frac{2\text{tr}(\mathbf{D}^{-1} \Sigma_2)^2}{N_2^2} + \frac{4\text{tr} \mathbf{D}^{-1} \Sigma_1 \mathbf{D}^{-1} \Sigma_2}{N_1 N_2} \\ &= 2\text{tr} \left[\left(\frac{\mathbf{D}^{-1/2} \Sigma_1 \mathbf{D}^{-1/2}}{N_1} \right) + \left(\frac{\mathbf{D}^{-1/2} \Sigma_2 \mathbf{D}^{-1/2}}{N_2} \right) \right]^2 \\ &= 2\text{tr} \mathcal{R}^2. \end{aligned}$$

Following Corollary 2.6 of [5], we have for $i = 1, 2$ and $j = 1, \dots, p$ that $E(s_{ij}^{-1}) = \sigma_{ij}^{-1} + O(N_i^{-1})$. Hence, $s_{ij}^{-1} = \sigma_{ij}^{-1} + O_p(N_i^{-1})$. Thus,

$$\hat{\mathbf{D}}^{-1} = \mathbf{D}^{-1} [1 + O_p(N_m^{-1})], \quad N_m = \min(N_1, N_2),$$

which implies

$$\begin{aligned} \hat{q}_n &= \frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \hat{\mathbf{D}}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - p}{\sqrt{p}} \\ &= \frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{D}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - p [1 + O_p(N_m^{-1})]}{\sqrt{p}} [1 + O_p(N_m^{-1})] \\ &= \tilde{q}_n + O_p \left(\frac{\sqrt{p}}{N_m} \right), \end{aligned} \tag{1.10}$$

Download English Version:

<https://daneshyari.com/en/article/1145937>

Download Persian Version:

<https://daneshyari.com/article/1145937>

[Daneshyari.com](https://daneshyari.com)