

Profile analysis in high-dimensional data

Nobumichi Shutoh^{*,1}, Sho Takahashi²

¹ Department of Health Sciences, Oita University of Nursing and Health Sciences
2944-9, Megusuno, Oita-shi, Oita 870-1201, JAPAN

² Department of Mathematical Information Science, Graduate School of Science,
Tokyo University of Science
1-3, Kagurazaka, Shinjuku-ku, Tokyo 162-8601, JAPAN

Abstract

This paper provides test statistics for the three hypotheses in profile analysis under high-dimensional data. The existing methods in profile analysis have suffered from the curse of high-dimensionality of the datasets, i.e., the singularity of the sample covariance matrix. We propose the new test statistics without the inverse of the sample covariance matrix via Cauchy-Schwarz inequality. The performed simulation evaluates the proposed procedures.

AMS 2010 Mathematics Subject Classification. 62H15, 62E20.

Key Words and Phrases: profile analysis, high-dimensional data, asymptotic distribution.

1 Introduction

In the practical applications, we may have the statistical inference for the existence of the interaction between k treatments and p responses or may be interested in the main effects of k treatments and p responses, respectively. Then we may apply two-way analysis of variance if the both of k treatments and p responses are mutually independent.

In this paper, let us consider the case that p responses are not independent. For $k = 2$, we connect the lines between the plots $(1, \mu_{i1}), (2, \mu_{i2}), \dots, (p, \mu_{ip})$, which is called mean profile, where $\mu_{i1}, \dots, \mu_{ip}$ are the p mean components from the i -th group for $i = 1, 2$.

No interaction can be described as Figure 1(a), i.e., $H_1 : \mu_{11} - \mu_{21} = \mu_{12} - \mu_{22} = \dots = \mu_{1p} - \mu_{2p} \equiv \gamma$, which is called “the parallelism hypothesis”. If there is no interaction between

^{*}Corresponding author. E-mail: shutoh@oita-nhs.ac.jp Telephone: +81-97-586-4473 Fax: +81-97-586-4382

the treatments and the responses, then we may also estimate the level difference γ . The fact that $H_2|H_1 : \gamma = 0$, which is called “the level hypothesis”, implies that no treatments’ main effect exists, as described in Figure 1(b). On the other hand, mean profile with flatness implies no responses’ main effect described as Figure 1(c) with the formulation: $H_3|H_1 : \mu_{11} = \mu_{12} = \dots = \mu_{1p}$, $\mu_{21} = \mu_{22} = \dots = \mu_{2p}$, which is called “the flatness hypothesis”.

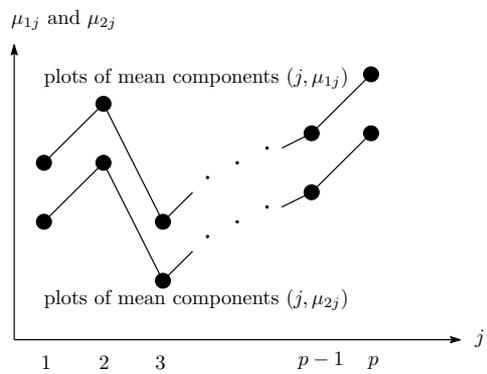
If the p -dimensional sample vector \mathbf{x}_{ij} ($i = 1, 2$, $j = 1, \dots, N_i$) can be observed from two groups Π_i under multivariate normality with mean vector $\boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{ip})'$ and the covariance matrix Σ , then the statistical inference for the interaction or main effects can be considered. Using the estimators of mean vector and covariance matrix:

$$\bar{\mathbf{x}}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} \mathbf{x}_{ij}, \quad S = \frac{1}{n} \sum_{i=1}^2 \sum_{j=1}^{N_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)',$$

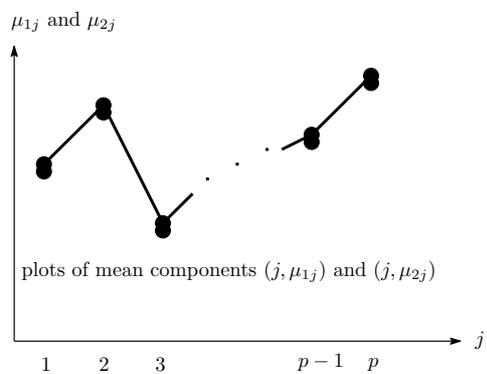
where $n = N_1 + N_2 - 2$, as mentioned in Greenhouse and Geisser (1959), we can have the exact tests that follows F distribution under null hypotheses. For general k groups, Srivastava (1987) derived the likelihood ratio tests for the three hypotheses. In growth curve model, Fujikoshi (2009) also derived a test for the parallelism hypothesis H_1 via the likelihood ratio for k groups using a canonical form. Yokoyama and Fujikoshi (1993) discussed the profile analysis with random effects model. In recent years, some authors may also be interested in the effects of non-normality in profile analysis, see e.g., Okamoto et al. (2006) and Maruyama (2007).

The above-mentioned classical methods suffer from the curse of high-dimensionality. They cannot be applicable to the data set, such as microarrays data, for $p > n$ owing to the singularity of S . Thereby, in the last decade, several authors have the motivation to derive another statistical procedures without S^{-1} and $|S|$. For example, the tests for covariance structure in high-dimensional data could be derived in Ledoit and Wolf (2002) and Srivastava (2005). Fujikoshi et al. (2010) have the invaluable reviews of the research based on the asymptotic theory in high-dimensional data.

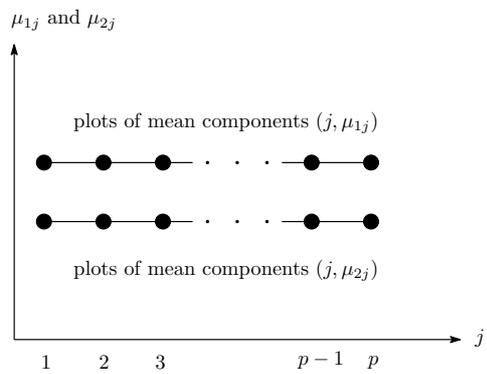
In this paper, we derive new test statistics in profile analysis without S^{-1} and $|S|$. The testing procedures for H_1 and $H_3|H_1$ are derived under the following high-dimensional asymptotic



(a) The parallelism hypothesis H_1



(b) The level hypothesis $H_2|H_1$



(c) The flatness hypothesis $H_3|H_1$

Figure 1: The mean profiles under three hypotheses

frameworks:

$$A1 : \text{tr}[(\Sigma A)^i] = O(p), \quad i = 1, 2, 3, 4,$$

$$A2 : n \rightarrow \infty, \quad p \rightarrow \infty, \quad p/n \rightarrow \xi \in (0, \infty),$$

where $A = I_p - (1/p)\mathbf{1}_p\mathbf{1}_p'$ and $\mathbf{1}_p$ denotes p -dimensional vector with the components of 1's. On the other hand, the testing procedure for $H_2|H_1$ is derived under the asymptotic frameworks A2 and

$$A3 : (\mathbf{1}_p'\Sigma\mathbf{1}_p)^i = O(p), \quad i = 1, 2.$$

This paper is organized as follows. Section 2 presents the test statistics for H_1 , $H_2|H_1$ and $H_3|H_1$, respectively. Section 3 proves the lemmas and corollaries. Section 4 has the simulation studies and presents the asymptotic behavior of the proposed results. Finally, Section 5 concludes the paper.

2 The proposed test statistics

In this section, we propose the test statistics for the three hypotheses: “the parallelism hypothesis”, “the level hypothesis” and “the flatness hypothesis”. The first and third test statistics are derived via Cauchy-Schwarz inequality

$$\|\mathbf{a}\|^2 \cdot \|\mathbf{b}\|^2 - (\mathbf{a}, \mathbf{b})^2 \geq 0, \quad (2.1)$$

where \mathbf{a} and \mathbf{b} are p -dimensional vectors, $(\mathbf{a}, \mathbf{b}) = \mathbf{a}'\mathbf{b}$ and $\|\mathbf{a}\|^2 = (\mathbf{a}, \mathbf{a})$.

2.1 Test for the parallelism hypothesis H_1

The parallelism hypothesis can be also formulated by $H_1 : \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2 = \gamma\mathbf{1}_p$, where γ is unknown parameter for the level difference. Applying $\mathbf{a} = \mathbf{1}_p$, $\mathbf{b} = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$ to (2.1), we can obtain a new formulation for the parallelism hypothesis:

$$H_1 : (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)'A(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = 0 \quad \text{vs.} \quad A_1 : (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)'A(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) > 0 \quad (2.2)$$

because the equality of (2.1) holds if and only if $\mathbf{b} = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$ is proportional to $\mathbf{a} = \mathbf{1}_p$. For the naive estimator of the left-side of H_1 stated in (2.2): $T_1 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' A(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$, we have

$$E(T_1) = \left(\frac{1}{N_1} + \frac{1}{N_2} \right) \text{tr}[\Sigma A], \quad \text{Var}(T_1) = 2 \left(\frac{1}{N_1} + \frac{1}{N_2} \right)^2 \text{tr}[(\Sigma A)^2]$$

and

$$U_1 = \frac{1}{\sqrt{p}} \left\{ \left(\frac{1}{N_1} + \frac{1}{N_2} \right)^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' A(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - \text{tr}[\Sigma A] \right\} \xrightarrow{d} N(0, 2c_2)$$

under the asymptotic frameworks A1 and A2, where “ \xrightarrow{d} ” denotes the convergence in distribution, as shown in Lemma 3.1. Therefore, by standardizing T_1 , we have $T_1^* = (\sqrt{p}/\sigma_1)U_1 \xrightarrow{d} N(0, 1)$ under the asymptotic frameworks A1 and A2, where $\sigma_1 = \sqrt{2pc_2}$. However, T_1^* depends on the unknown constant c_i 's for $i = 1, 2$. Instead we propose the following test statistic \widehat{T}_1^* obtained by using unbiased and consistent estimator of c_i obtained in Lemma 3.3.

Theorem 2.1. *If the parallelism hypothesis H_1 holds, then*

$$\widehat{T}_1^* = \frac{1}{\widehat{\sigma}_1} \left\{ \left(\frac{1}{N_1} + \frac{1}{N_2} \right)^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' A(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) - p\widehat{c}_1 \right\} \xrightarrow{d} N(0, 1)$$

under the high-dimensional asymptotic frameworks A1 and A2, where $\widehat{\sigma}_1 = \sqrt{2p\widehat{c}_2}$, $\widehat{c}_1 = \text{tr}(SA)/p$ and

$$\widehat{c}_2 = \frac{n^2}{(n+2)(n-1)} \left\{ \text{tr}[(SA)^2] - \frac{1}{n} \{ \text{tr}(SA) \}^2 \right\}.$$

Proof. Using Lemma 3.3, under the high-dimensional asymptotic frameworks A1 and A2, it also holds that $1/\sqrt{2\widehat{c}_2} \xrightarrow{p} 1/\sqrt{2c_2}$. On the other hand, it follows that $U_1 \xrightarrow{d} N(0, 2c_2)$ and $V_1 = -\sqrt{p}(\widehat{c}_1 - c_1) \xrightarrow{p} 0$ under the asymptotic frameworks A1 and A2. Finally, applying Slutsky's theorem completes the proof. \square

By the testing procedure based on \widehat{T}_1^* , if $\widehat{T}_1^* > z_\alpha$, then the parallelism hypothesis H_1 is rejected, where z_α is the upper 100 α % point of $N(0, 1)$.

2.2 Test for the level hypothesis $H_2|H_1$

The level hypothesis can be also formulated by $H_2|H_1 : \gamma = 0$. We have an unbiased estimator of γ : $T_2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{1}_p / p$ and its variance

$$\text{Var}(T_2) = \frac{1}{p} \left(\frac{1}{N_1} + \frac{1}{N_2} \right) \mathbf{1}'_p \Sigma \mathbf{1}_p.$$

It follows from standardizing T_2 that

$$T_2^* = \left\{ \frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{1}_p}{p} - \gamma \right\} / \sqrt{\frac{1}{p} \left(\frac{1}{N_1} + \frac{1}{N_2} \right) d_1} \sim N(0, 1), \quad (2.3)$$

where $d_1 = \mathbf{1}'_p \Sigma \mathbf{1}'_p / p$. Combining (2.3) and Corollary 3.4, we obtain the following theorem and corollary.

Theorem 2.2. *If the level hypothesis $H_2|H_1$ holds, then*

$$\widehat{T}_2^* = \left(\frac{p}{N_1} + \frac{p}{N_2} \right)^{-\frac{1}{2}} \frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{1}_p}{\sqrt{\widehat{d}_1}} \xrightarrow{d} N(0, 1)$$

under the high-dimensional asymptotic frameworks A2 and A3, where $\widehat{d}_1 = \mathbf{1}'_p S \mathbf{1}_p / p$.

By the testing procedure based on \widehat{T}_2^* , if $|\widehat{T}_2^*| > z_{\alpha/2}$, then the level hypothesis $H_2|H_1$ is rejected. If $H_2|H_1$ is rejected, we may be interested in the interval estimation for the level difference γ .

Corollary 2.3. *Under asymptotic frameworks A2 and A3, an approximate interval estimation for the level difference γ with the significant level $1 - \alpha$ can be obtained by*

$$\left[\frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{1}_p}{p} - z_{\frac{\alpha}{2}} \sqrt{\frac{\widehat{d}_1 (N_1 + N_2)}{p N_1 N_2}}, \frac{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{1}_p}{p} + z_{\frac{\alpha}{2}} \sqrt{\frac{\widehat{d}_1 (N_1 + N_2)}{p N_1 N_2}} \right].$$

2.3 Test for the flatness hypothesis $H_3|H_1$

The flatness hypothesis can be also formulated by $H_3|H_1 : \boldsymbol{\mu}_{12} = \delta \mathbf{1}_p$, where $\boldsymbol{\mu}_{12} = \{N_1 / (N_1 + N_2)\} \boldsymbol{\mu}_1 + \{N_2 / (N_1 + N_2)\} \boldsymbol{\mu}_2$ and δ is unknown constant. In a similar manner to Subsection 2.1, we can obtain a new formulation for the flatness hypothesis:

$$H_3|H_1 : \boldsymbol{\mu}'_{12} A \boldsymbol{\mu}_{12} = 0 \quad \text{vs.} \quad A_3|H_1 : \boldsymbol{\mu}'_{12} A \boldsymbol{\mu}_{12} > 0.$$

Corollary 3.2 shows that

$$U_3 = \frac{1}{\sqrt{p}} \left\{ \left(\frac{1}{N_1 + N_2} \right)^{-1} \bar{\mathbf{x}}'_{12} A \bar{\mathbf{x}}_{12} - \text{tr}[\Sigma A] \right\} \xrightarrow{d} N(0, 2c_2)$$

under asymptotic frameworks A1 and A2, where $\bar{\mathbf{x}}_{12} = \{N_1/(N_1 + N_2)\}\bar{\mathbf{x}}_1 + \{N_2/(N_1 + N_2)\}\bar{\mathbf{x}}_2$.

Then, we also propose the following test statistic.

Theorem 2.4. *If the flatness hypothesis $H_3|H_1$ holds, then*

$$\widehat{T}_3^* = \frac{1}{\widehat{\sigma}_3} \left\{ \left(\frac{1}{N_1 + N_2} \right)^{-1} \bar{\mathbf{x}}'_{12} A \bar{\mathbf{x}}_{12} - p\widehat{c}_1 \right\} \xrightarrow{d} N(0, 1)$$

under the high-dimensional asymptotic frameworks A1 and A2, where $\widehat{\sigma}_3 = \sqrt{2p\widehat{c}_2}$.

Proof. The proof is completed by a similar manner to Theorem 2.1. □

By the testing procedure based on \widehat{T}_3^* , if $\widehat{T}_3^* > z_\alpha$, then the flatness hypothesis $H_3|H_1$ is rejected.

3 Lemmas

We list the key results which complete the proofs of the main theorems.

Lemma 3.1. *If the parallelism hypothesis H_1 holds, then $U_1 \xrightarrow{d} N(0, 2c_2)$ under the high-dimensional asymptotic frameworks A1 and A2.*

Proof. U_1 equals $\{\sum_{j=1}^r \lambda_j Y_j - \text{tr}[\Sigma A]\}/\sqrt{p}$, where r is the rank of $A\Sigma A'$, λ_j is the characteristic roots of $A\Sigma A'$ and Y_j is independently and identically distributed as a chi-square distribution with 1 degree of freedom for $j = 1, \dots, r$. It should be noted that $r = p - 1$ or $p - 2$ because $\text{rank}[A\Sigma] + \text{rank}[A'] - p \leq r \leq \min\{\text{rank}[A\Sigma], \text{rank}[A']\}$. Further we have the following characteristic function:

$$\varphi(t) = \text{E}[\exp(itU_1)] = \prod_{j=1}^r \left(1 - 2i \frac{\lambda_j}{\sqrt{p}} t \right)^{-\frac{1}{2}} \exp \left\{ -\frac{it}{\sqrt{p}} \text{tr}[\Sigma A] \right\},$$

where $i = \sqrt{-1}$. Therefore, we have

$$\ln \varphi(t) = -\frac{1}{2} \sum_{j=1}^r \ln \left(1 - 2i \frac{\lambda_j}{\sqrt{p}} t \right) - \frac{it}{\sqrt{p}} \text{tr}[\Sigma A] = \frac{(it)^2}{2} \cdot 2c_2 + O \left(\frac{|t|^3}{\sqrt{p}} c_3 \right),$$

by applying Taylor expansion to $\ln\{1 - 2it(\lambda_j/\sqrt{p})\}$. □

In a similar manner to Lemma 3.1, we have the following corollary.

Corollary 3.2. *If the flatness hypothesis $H_3|H_1$ holds, then $U_3 \xrightarrow{d} N(0, 2c_2)$ under the high-dimensional asymptotic frameworks A1 and A2.*

The next lemma and corollary prove the unbiasedness and consistency of \widehat{c}_i ($i = 1, 2$) and \widehat{d}_1 by using the moments derived by Himeno (2007) and Hyodo et al. (2012).

Lemma 3.3. *For $i = 1, 2$, \widehat{c}_i is an unbiased and consistent estimator of c_i under the high-dimensional asymptotic frameworks A1 and A2.*

Proof. It holds that $E(\widehat{c}_i) = c_i$, $\text{Var}(\widehat{c}_1) = (2c_2)/(np) = O(n^{-2})$, and

$$\text{Var}(\widehat{c}_2) = -\frac{8c_2^2}{n^4} + \frac{4c_2^2}{n^3} + \frac{4c_2^2}{n^2} + \frac{48c_4}{n^5p} - \frac{48c_4}{n^4p} - \frac{28c_4}{n^3p} + \frac{20c_4}{n^2p} + \frac{8c_4}{np} = O(n^{-2}),$$

for $i = 1, 2$. Under the high-dimensional frameworks A1 and A2, applying \widehat{c}_i 's for $i = 1, 2$ to Chebyshev's inequality completes the proof. \square

Corollary 3.4. *\widehat{d}_1 is an unbiased and consistent estimator of d_1 under the high-dimensional asymptotic frameworks A2 and A3.*

4 Simulation studies

Under the null hypotheses, we conduct Monte Carlo simulation with 10,000 replications in order to investigate the attained significance level of \widehat{T}_i^* ($i = 1, 2, 3$) (i.e., type I error), respectively. Then the selected parameters are as follows: $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \mathbf{0}$, $\alpha = 0.05$, $\Sigma = (1 - \rho)I_p + \rho\mathbf{1}_p\mathbf{1}_p'$ for $\rho = 0, 0.2, 0.5$. The dimensionality and the sample sizes are set as $\xi = 4$; $(p, N_1, N_2) = (100, 14, 13)$, $(200, 26, 26)$, $(400, 51, 51)$ and $\xi = 2$; $(p, N_1, N_2) = (100, 26, 26)$, $(200, 51, 51)$, $(400, 101, 101)$. The results are listed in Tables 1 and 2.

5 Conclusion

We have the testing procedures for three hypotheses for profile analysis in high-dimensional data via Cauchy-Schwarz inequality. In profile analysis, our paper provides the testing procedures

which are applicable to high-dimensional data. We also perform Monte Carlo simulation in order to observe the type I error of the proposed results under some selected parameters. The proposed procedures have the type I error that is close to the desired level of significance when both of the dimensionality and the total sample size are larger.

However, the type I errors in the both of the tests for the parallelism hypothesis and the flatness hypothesis increase when the off-diagonal elements of Σ are larger. It can be also observed that the type I error also increases in all the three procedures when the ratio of the dimensionality to the total sample size denoted by ξ is larger. Thereby, in some cases, the proposed testing procedures should be improved.

References

- Dempster, A. P., 1958. A high dimensional two samples significance test. *Annals of Mathematical Statistics* 29, 995–1010.
- Dempster, A. P., 1960. A significance test for separation of two highly multivariate small samples. *Biometrics*, 16, 41–50.
- Fujikoshi, Y., 2009. Statistical inference for parallelism hypothesis in growth curve model. *SUT Journal of Mathematics* 45, 137–148.
- Fujikoshi, Y., Ulyanov, V. V. and Shimizu, R., 2010. *Multivariate Statistics High-Dimensional and Large-Sample Approximations*. John Wiley & Sons, Inc., Hoboken, New Jersey.
- Greenhouse, S. W., Geisser, S., 1959. On methods in the analysis of profile data. *Psychometrika* 24, 95–112.
- Himeno, T., 2007. Asymptotic expansions of the null distributions for the Dempster trace criterion. *Hiroshima Mathematical Journal*, 37, 431–454.
- Hyodo, M., Takahashi, S., Nishiyama, T., 2012. Multiple comparisons among mean vectors when the dimension is larger than the total sample size. Technical Report 12–01, Hiroshima Statistical Research Group, Hiroshima University.
- Ledoit, O., Wolf, M., 2002. Some hypothesis tests for the covariance matrix when the dimension is large compared to the sample size.
- Okamoto, N., Miura, N., Seo, T., 2006. On the distributions of some test statistics for profile analysis in elliptical populations. *American Journal of Mathematical and Management Sciences* 26, 1–31.
- Maruyama, Y., 2007. Asymptotic expansions of the null distributions of some test statistics for profile analysis in general distributions. *Journal of Statistical Planning and Inference* 137, 506–526.

Srivastava, M. S., 1987. Profile analysis of several groups. *Communications in Statistics – Theory and Methods* 16, 909–926.

Srivastava, M. S., 2005. Some tests concerning the covariance matrix in high dimensional data. *Journal of Japan Statistical Society* 35, 251–272.

Yokoyama, T., Fujikoshi, Y., 1993. A parallel profile model with random-effects covariance structure. *Journal of Japan Statistical Society* 23, 83–89.

