

诚信声明

我声明，所呈交的毕业论文是本人在老师指导下进行的研究工作及取得的研究成果。据我查证，除了文中特别加以标注和致谢的地方外，论文中不包含其他人已经发表或撰写过的研究成果，也不包含为获得其他教育机构的学位或证书而使用过的材料。我承诺，论文中的所有内容均真实、可信。

毕业论文作者签名： 签名日期： 年 月 日

高维协方差矩阵的检验综述及其推广

[摘要]随着科技和计算机的发展，在线收集数据得以实现，使得收集的数据在形式上、结构上和数量上发生了革命性的变化。为了适应数据发展的要求，统计研究也应呈现相应的变化。传统的统计研究方法关注的是 $p < n$ 的情形，当今的研究更多的是关注 p 与 n 相当（高维， $p/n \rightarrow c$ ）和 $p \gg c$ （超高维）的情形。本文回顾综述了高维协方差矩阵的检验方法。并针对单个高维总体协方差矩阵等于单位矩阵的假设检验问题，考虑 $p/n \rightarrow c \in (0,1]$ 的情形，提出了一个检验统计量，给出其在原假设下的极限分布，并将此检验方法推广到单个高维总体协方差矩阵等于对角矩阵和对称正定矩阵的假设检验情形。在数值模拟过程中，该统计量比传统似然比检验统计量表现出更好的效果，在保证检验显著性水平前提下，拥有更高的检验功效。

[关键词]高维数据；协方差矩阵；假设检验

A Review and Extension of High Dimensional Testing for Covariance Matrices

Abstract: With the development of technology and computers, more and more data are collected online which is always heterogeneous, varied structure and usually called as big data. In order to meet the requirements of data development, statistical research should also show corresponding changes. Because the classical statistics focused on the data where the dimension of the variable p is less than the sample size n . But today's data is always high dimensional($p/n \rightarrow c$) or ultrahigh dimensional($p \gg c$). This paper reviews and summarizes the hypothesis test methods for high dimensional covariance matrices. A statistics is proposed for the hypothesis test problem that a single high dimensional population covariance matrix Σ is equal to the unit matrix, considering the case of $p/n \rightarrow c \in (0,1]$, and its limit distribution under the null hypothesis is given. This method is extended to the hypothesis test case where Σ is equal to the diagonal matrix or the symmetric positive definite matrix. In the numerical simulation, the proposed statistics show better effect than the traditional likelihood ratio test statistics, and has a higher test power at the ensured level of significance.

Keywords: High dimensional data; Covariance matrix; Hypothesis test

目录

1 绪论	1
1.1 研究背景和目的	1
1.2 国内外研究现状综述	1
1.2.1 国内外研究现状	1
1.2.2 研究趋势和导向	4
1.3 研究思路与框架	4
1.3.1 研究思路	4
1.3.2 研究框架	4
1.3.3 创新之处	5
2 单个高维总体协方差矩阵的检验	6
2.1 单个高维总体协方差矩阵等于给定单位矩阵的检验	6
2.2 单个高维总体协方差矩阵等于给定对角矩阵的检验	8
2.3 单个高维总体协方差矩阵等于给定非负定矩阵的检验	11
3 数值模拟	13
3.1 极限分布的模拟	13
3.2 检验显著性水平和功效模拟（单个高维总体协方差矩阵等于给定单位矩阵）	14
3.3 检验显著性水平和功效模拟（单个高维总体协方差矩阵等于给定对角矩阵）	15
3.4 检验功效变化的数值模拟（单个高维总体协方差矩阵等于给定对角矩阵）	18
4 结论和展望	20
致谢	21
附录	22
定理 1 证明	22
定理 3 证明	25
(10)式证明	27
数值模拟代码	29
参考文献	37

1 绪论

1.1 研究背景和目的

高维数据的出现,对一些传统的多元统计分析方法和理论(如假设检验)来说是一个挑战,因为传统的假设检验理论是在假定变量个数 p 小于样本量 n 的基础上建立起来的,这就需要我们更新或改写传统的假设检验理论,以面对现代数据 p 与 n 相当和 p 远大于 n 的情形。

本文旨在对现有的高维协方差矩阵的假设检验方法进行回顾总结。并针对单个高维总体协方差矩阵的假设检验问题,在 $p/n \rightarrow c \in (0,1]$ 的特定条件下,提出检验效果比传统似然比检验更好的检验统计量。并使用 R 语言进行数值模拟,比较所提检验统计量和已有几个统计量在检验显著性水平和检验功效上的差异,从而得到不同的情形使用不同的检验统计量的结论,也为其他类型检验的检验统计量的扩展和改进提供参考,使之适用于模型构建前对前提条件的假设检验问题。

1.2 国内外研究现状综述

1.2.1 国内外研究现状

在协方差矩阵的等价的假设检验问题上,分为两种情况 (Chang et al. (2015)^[1]): 第一种考虑单个总体,检验方差矩阵等于一个给定矩阵或具有特定结构的假设,如对角等;第二种考虑两个或两个以上总体,如两总体有相同的协方差矩阵 ($H_0: \Sigma_0 = \Sigma_1$)。而后的研究都是针对非精确检验方法或修正的多元统计方法的前提假设进行的,不一样的前提假设有不同的检验效

率，也就有不同方法对应的优缺点。

对于单个总体协方差矩阵的假设检验问题，传统常规的检验统计量有 Nagao(1973)^[2]提出的 V 统计量 $\frac{1}{p}tr(S - I_p)^2$ 和似然比检验统计量 (LTR , 参考 Anderson(2004)^[3]): $n(trS - \log|S| - p)$, 其中 S 指样本协方差矩阵。而随着 p 增大, 且 $p/n \rightarrow c \in (0, +\infty)$, 在原假设下似然比检验统计量对应的卡方检验不再适用, Lediot & Wolf(2002)^[4]对 V 统计量进行改进并提出统计量 $\frac{1}{p}tr(S - I_p)^2 - \frac{p}{n-1}[\frac{1}{p} \cdot tr(S)]^2$, 在 $p \rightarrow \infty$ 和多元正态分布假设等条件下, Lediot & Wolf 证明中心极限定理如下:

$$n \cdot \left\{ \frac{1}{p}tr(S - I_p)^2 - \frac{p}{n-1} \left[\frac{1}{p} \cdot tr(S) \right]^2 \right\} \xrightarrow{D} N(1, 4).$$

放松正态假定, 在总体分布峰度为3且 $p/n \rightarrow c \in (0, 1)$ 等前提下, Bai et al.(2009)^[5]利用随机矩阵理论提出统计量: $\frac{1}{p}tr(B) - \frac{1}{p}\log|B| - 1$, 其中 $B = \frac{n-1}{n}S$ 。Jiang et al.(2012)^[6]利用 Selberg 不等式将其扩展到 $p < n$ 且 $p/n \rightarrow c \in (0, 1]$ 的情形。Wang et al.(2013)^[7]在 $p/n \rightarrow c \in (0, +\infty)$ 前提下提出检验统计量 $\frac{1}{p}tr(S) - \frac{1}{p}\log|S| - 1$, 并且相比于 Bai et al.(2009)^[5]的检验统计量, Wang et al.(2013)^[7]指出前者更适用于均值向量未知且没有正态分布假定的假设检验问题。在 $p > n$ 的情形, 样本协方差矩阵是奇异的, 由于上述检验统计量包含 $\log|S|$ 项不再适用, Chen et al.(2010)^[8]在 $tr(\Sigma^4) = o(tr^2(\Sigma^2))$ 假设下, 提出统计量 $\frac{1}{p}T_1 - \frac{1}{p}T_2 - 1$, 其中 T_1 和 T_2 分别为 $tr(\Sigma)$ 和 $tr(\Sigma^2)$ 的无偏估计。Ahmad & Rosen(2015)^[9]改进了 Chen et al.(2010)^[8]的方法, 进一步提出统计量 $\frac{1}{p}C_1 - \frac{1}{p}C_2 - 1$, 其中 C_1 和 C_2 分别为 $tr(\Sigma)$ 和 $tr(\Sigma^2)$ 的一致无偏估计。此外 Wu & Li (2015)^[10]在概念上提出一种简单的方法, 使用随机投影将数据投影

到一维随机子空间上,便可以使用传统的统计方法。Jiang & Yang(2013)^[11], Jiang & Qi(2015)^[12]和 Chen & Jiang(2018)^[13]针对传统的似然比检验统计量是在固定 p 且 $n \rightarrow \infty$ 情形提出的,在正态假定和 $p/n \rightarrow c \in (0,1]$ 情形,在多元假设检验问题上提出了似然比检验统计量的改进方法。

对两个或两个以上总体的高维假设检验问题, Li & Chen(2012)^[14]提出,对任意的 $i, j, k, l \in \{1, 2\}$, 满足:

$$tr(\Sigma_i \Sigma_j) \rightarrow \infty, tr[(\Sigma_i \Sigma_j)(\Sigma_k \Sigma_l)] = o(tr(\Sigma_i \Sigma_j)tr(\Sigma_k \Sigma_l))$$

的前提下,可用基于 $\Sigma_1 - \Sigma_2$ 的 Frobenius 范数的 U 统计量检验 $H_0: \Sigma_0 = \Sigma_1$ 。Cai et al.(2013)^[15]在 Σ_1 和 Σ_2 差异微小,且 $\Sigma_1 - \Sigma_2$ 是稀疏的前提下,对 $H_0: \Sigma_0 = \Sigma_1$ 进行检验很有效力,且该方法在检验 $H_0: \Sigma_0 = \Sigma_1$ 的同时,可以逐行检验 Σ_1 和 Σ_2 的差异所在,在基因组学中,当 $H_0: \Sigma_0 = \Sigma_1$ 的零假设被拒绝时,进一步研究它们彼此之间的差异具体所在有重要意义。Hong & Kim(2018)^[16]指出在假设协方差矩阵密集的前提下,可以使用 Cai et al.(2013)^[15]的方法,然而,考虑交互作用后 $\Sigma_1^{-1} - \Sigma_2^{-1}$ 稀疏,而 $\Sigma_1 - \Sigma_2$ 不稀疏,根本原因在于边际相关性和条件相关性导致对 $H_0: \Sigma_1 = \Sigma_2$ 的检验方法并不适用于对 $H_0: \Sigma_1^{-1} = \Sigma_2^{-1}$ 进行检验。Cai(2017)^[17]在对以往协方差矩阵估计方法的回顾中指出, Li & Chen(2012)^[14]所用的基 Frobenius 范数的检验方法,在 $\Sigma_1 - \Sigma_2$ 非稀疏的情况下表现较好,而在稀疏的情形下检验效力较低。Cai & Wu(2018)^[18]从统计和计算角度探究了被高斯白噪声污染的高维稀疏矩阵的基本限制,指出构建好的检验统计量的关键步骤是稀疏矩阵的结构性质。

1.2.2 研究趋势和导向

国内外对高维协方差矩阵的假设检验通常需要正态假定，而对超高维矩阵的假设检验方法则通常需要稀疏条件，除了多种类型的假设检验问题，除了是在各种不同前提下有好的性质，研究的目的也是很重要因素。如基因表达谱数据、单个核苷酸多态性数据、金融衍生数据分析、无线网络通信技术等领域对研究趋势本身有重要导向。例如对于基因表达谱数据，检验两总体是否有相同的协方差矩阵，若不相同，进一步了解它们彼此差异所在有重要作用。

下文将针对特定的假设检验，即单个高维总体协方差矩阵等于单位矩阵的检验问题，并将其推广至单个高维总体协方差矩阵等于给定对称正定矩阵的检验。从这特定的问题中寻找高维众多类型假设检验中检验统计量构建的灵感。

1.3 研究思路与框架

1.3.1 研究思路

为了达到研究目的，本文采用文献研究法和模拟仿真法来进行研究。特别的，针对单个高维总体协方差矩阵的假设检验问题，通过改进似然比检验框架下统计量的构建，提出适用于 $p/n \rightarrow c \in (0,1]$ 情形的检验统计量，采用R语言编程进行仿真模拟，比较提出的检验统计量的检验效果。

1.3.2 研究框架

本文具体的研究框架如下：

第一部分：介绍研究的背景、目的、思路等内容；

第二部分：介绍国内外学者对高维协方差矩阵假设检验方法的研究现状、成果和发展趋势；

第三部分：在似然比检验框架下，对单个高维总体协方差矩阵的假设检验的检验统计量进行改进和推广，并使用 R 语言对其进行数值模拟，比较几个检验统计量的检验显著性和检验功效；

第四部分：结论和展望。

1.3.3 创新之处

本文针对单个高维总体协方差矩阵等于单位矩阵的假设检验问题，考虑 $p/n \rightarrow c \in (0,1]$ 的情形，给出一个新检验统计量及其极限分布，并将其推广到单个高维总体协方差矩阵等于对角矩阵和对称正定矩阵的假设检验情形。在特定假设检验情形下，探讨检验功效函数的特性，并在数值模拟部分进行验证。此外，本文指出和强调在一般高维协方差矩阵的假设检验问题中，为了方便进行假设检验而对样本数据进行线性变换时需要注意的地方，并在数值模拟部分予以验证和凸显。

2 单个高维总体协方差矩阵的检验

2.1 单个高维总体协方差矩阵等于给定单位矩阵的检验

假设 $X = (X_1, X_2, \dots, X_n)$ 是独立同分布于 $N_p(\mu, \Sigma)$ 的实值随机向量，其中 $\mu \in R^p$ 为均值向量， $\Sigma_{p \times p}$ 为协方差矩阵。考虑假设检验问题：

$$H_0: \Sigma = I_p \text{ v.s. } H_1: \Sigma \neq I_p, \quad (1)$$

其中 I_p 为 $p \times p$ 单位矩阵。记：

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \\ A = \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})' = X'(I_{n \times n} - \frac{1}{n} J_{n \times n})X.$$

其中， $I_{n \times n}$ 为 $n \times n$ 单位矩阵， $J_{n \times n} = (1, \dots, 1)'_{n \times 1} \cdot (1, \dots, 1)_{1 \times n}$ 。对于假设检验 (1)，显著性水平为 α 的似然比检验的拒绝域为 $\{\Lambda^* \leq c_\alpha\}$

(Muirhead(1982)^[19], Theorem 8.4.5)，其中：

$$\Lambda^* = \left(\frac{e}{m}\right)^{\frac{mp}{2}} e^{-\frac{1}{2} \text{tr}(A)} |A|^{\frac{m}{2}}, \quad (2)$$

其中 $m = n - 1$ 为自由度（修正作用，使似然比检验具有无偏性）， $\text{tr}(A)$ 和 $|A|$ 分别表示离差阵 A 的迹和行列式。

由于在 $p \geq n$ 时样本离差阵 A 不是满秩的，导致行列式为零，因而 (2) 中的统计量仅适用于 $p < n$ 情形。Muirhead(1982)^[19] 的 Theorem 8.4.9 表明，在原假设 $H_0: \Sigma = I_p$ 成立，固定 p 且 $n \rightarrow \infty$ 时，有：

$$-2\rho \cdot \log \Lambda^* \xrightarrow{D} \chi_f^2. \quad (3)$$

其中，

$$\rho = 1 - \frac{2p^2 + 3p - 1}{6m(p+1)},$$

$$f = \frac{1}{2}p(p+1).$$

显然, 在此种情况下 $\rho = \rho_n \rightarrow 1$ 。

下面我们提出在 p 随着 n 的增大而增大, 具体来说当 $p = p_n \rightarrow \infty, n \rightarrow \infty$ 且 $\lim_{n \rightarrow \infty} \frac{p_n}{n} \rightarrow c \in (0,1]$ 时, 有如下中心极限定理。

定理 1: 假设 $p = p_n$, 对任意 $n \geq 3$, 满足 $n > p + 3$, 且 $\lim_{n \rightarrow \infty} \frac{p}{n} \rightarrow c \in (0,1]$ 。

$\Lambda_m^* = \Lambda^*$ 按(2)中定义所示, 则对于(1)中假设检验问题, 当 $n \rightarrow \infty$, 在原假设 $H_0: \Sigma = I_p$ 成立时, 有:

$$\frac{\log \Lambda_m^* - \mu_{m,0}}{m \cdot \sigma_{m,0}} \xrightarrow{D} N(0,1). \quad (4)$$

其中, \xrightarrow{D} 表示依分布收敛, $N(0,1)$ 表示标准正态分布,

$$\mu_{m,0} = \frac{1}{4} m [(2m - 2p - 1) \cdot r_m^2 - 2p],$$

$$\sigma_{m,0}^2 = \frac{1}{2} \left(r_m^2 - \frac{p}{m} \right),$$

$$r_m = [-\log \left(1 - \frac{p}{m} \right)]^{\frac{1}{2}}.$$

定理 1 的证明见附录。类似的定理可参考 Chen & Jiang(2018):

定理 2 (Chen & Jiang(2018)^[1], Corollary 1): 假设 $n > p_n + 1$ 对任意 $n \geq 3$ 成立, 且 $\lim_{n \rightarrow \infty} p_n \rightarrow \infty$, $\Lambda_n^* = \Lambda^*$ 按(2)中定义所示, 则对于(1)中假设检验问题, 当 $n \rightarrow \infty$, 在原假设 $H_0: \Sigma = I_p$ 成立时, 有:

$$\frac{\log \Lambda_n^* - \mu_{n,0}}{n \cdot \sigma_{n,0}} \xrightarrow{D} N(0,1). \quad (5)$$

其中,

$$\begin{aligned}\mu_{n,0} &= -\frac{1}{4}(n-1) \left[2p + (2n-2p-3) \log \left(1 - \frac{p}{n-1} \right) \right], \\ \sigma_{n,0}^2 &= -\frac{1}{2} \left[\frac{p}{n-1} + \log \left(1 - \frac{p}{n-1} \right) \right].\end{aligned}$$

定理 1 和定理 2 的推导过程相似,在这两个检验统计量的构建过程中, $\mu_{m,0}$ 和 $\mu_{n,0}$, $\sigma_{m,0}^2$ 和 $\sigma_{n,0}^2$ 是等价的,但两者的微小差异在于构建的检验统计量的分母中,回顾定理 1 检验统计量为 $\frac{\log \Lambda_m^* - \mu_{m,0}}{m \cdot \sigma_{m,0}}$, 定理 2 检验统计量为 $\frac{\log \Lambda_n^* - \mu_{n,0}}{n \cdot \sigma_{n,0}}$, 当 $n \rightarrow \infty$ 时检验统计量分母包含的是 m 还是 n 的差异是微乎其微的,但是在 n “较小” 时,检验统计量的前者要大于后者,这就导致定理 2 模拟的显著性水平及其方差和定理 1 有明显差异。

2.2 单个高维总体协方差矩阵等于给定对角矩阵的检验

假设 $X = (X_1, X_2, \dots, X_n)$ 是独立同分布于 $N_p(\mu, \Sigma)$ 的实值随机向量,其中 $\mu \in R^p$ 为均值向量, $\Sigma_{p \times p}$ 为协方差矩阵。考虑假设检验问题:

$$H_0: \Sigma = \Sigma_0 \text{ v.s. } H_1: \Sigma \neq \Sigma_0, \quad (6)$$

其中 Σ_0 为给定 $p \times p$ 对角矩阵, 对角线元素为 $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_p < \infty$ 。

下面我们提出在 p 随着 n 的增大而增大,具体来说当 $p = p_n \rightarrow \infty, n \rightarrow \infty$ 且 $\lim_{n \rightarrow \infty} \frac{p_n}{n} \rightarrow c \in (0,1]$ 时, 有如下中心极限定理。

定理 3: 假设 $p = p_n$, 对任意 $n \geq 3$, 满足 $n > p + 3$, 且 $\lim_{n \rightarrow \infty} \frac{p}{n} \rightarrow c \in (0,1]$ 。 $\Lambda_m^* = \Lambda^*$ 按(2)中定义所示, 则对于(6)中假设检验问题, 当 $n \rightarrow \infty$, 在原假设 $H_0: \Sigma = \Sigma_0$ 成立时, 有:

$$\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m} \xrightarrow{D} N(0,1). \quad (7)$$

其中,

$$\begin{aligned}\mu_m &= \frac{1}{4}m \left[(2m - 2p - 1) \cdot r_m^2 + 2 \sum_{i=1}^p \log \lambda_i - 2 \sum_{i=1}^p \lambda_i \right], \\ \sigma_m^2 &= \frac{1}{2} \left[r_m^2 - \frac{2}{m} \sum_{i=1}^p \left(\lambda_i - \frac{\lambda_i^2}{2} \right) \right], \\ r_m &= \left[-\log \left(1 - \frac{p}{m} \right) \right]^{\frac{1}{2}}.\end{aligned}$$

定理 3 的证明见附录。显然，令 $\lambda_1 = \lambda_2 = \dots = \lambda_p = 1$ 可知，定理 1 是定理 3 的特殊情形。类似的定理可参考 Chen & Jiang(2018)提出的如下定理：

定理 4 (Chen & Jiang(2018)^[1], Theorem 1): 假设 $n > p_n + 1$ 对任意 $n \geq 3$ 成立，且 $\lim_{n \rightarrow \infty} p_n \rightarrow \infty$ ， $\Lambda_n^* = \Lambda^*$ 按(2)中定义所示，如果 Σ_0 为非奇异矩阵，当 $n \rightarrow \infty$ ，则有：

$$\frac{\log \Lambda_n^* - \mu_n}{n \cdot \sigma_n} \xrightarrow{D} N(0, 1). \quad (8)$$

其中，

$$\begin{aligned}\mu_n &= -\frac{1}{4}(n-1)(2n-2p-3) \log \left(1 - \frac{p}{n-1} \right) \\ &\quad + \frac{1}{2}(n-1)[\log |\Sigma_0| - \text{tr}(\Sigma_0)], \\ \sigma_n^2 &= -\frac{1}{2} \left[\frac{p}{n-1} + \log \left(1 - \frac{p}{n-1} \right) \right] + \frac{1}{2n} \text{tr}[(\Sigma_0 - I)^2].\end{aligned}$$

定理 3 和定理 4 的推导过程相似，对比可知 μ_m 和 μ_n 是等价的，但两者的微小差异在于构建的检验统计量的分母中，类似 2.1 节最后一段所述，两者除了有受到分母中取 m 还是 n 的影响，此时还受到 σ_m^2 和 σ_n^2 差异的影响。回顾定理 3 中 σ_m^2 为 $\frac{1}{2} [r_m^2 - \frac{2}{m} \sum_{i=1}^p (\lambda_i - \frac{\lambda_i^2}{2})]$ ，定理 4 中 σ_n^2 整理后为 $\frac{1}{2} [r_m^2 + \frac{1}{n} \sum_{i=1}^p (\lambda_i - 1)^2 - \frac{p}{m}]$ ，当 $n \rightarrow \infty$ 时两者的差异是微小得可以忽略不计的、但是在 n “较小” 时，两者差异会很大，而且受到具体样本协方差矩阵特征值

$\lambda_1, \lambda_2, \dots, \lambda_p$ 的影响, 这就导致定理 4 模拟的显著性水平及其方差和定理 3 有明显差异 (见表 3-2)。

对于定理 3, 检验功效函数如下 (更广泛的定理可参考 Chen & Jiang(2018)^[1]):

$$\beta(\Sigma) = P(\log \Lambda_m^* \leq c_\alpha | \Sigma) \sim \Phi\left(\frac{c_\alpha - \mu_m}{m \cdot \sigma_m}\right) \quad (9)$$

其中 $c_\alpha = \mu_{m,0} + m \cdot \sigma_{m,0} \cdot \Phi^{-1}(\alpha)$, 对应定理 1 中显著性水平为 α 时的拒绝域 $\{\log \Lambda_m^* \leq c_\alpha\}$, $\Phi^{-1}(\cdot)$ 为标准正态分布的密度函数。在 $0 < \lambda_1 = \lambda_2 = \dots = \lambda_p = \lambda < \infty$, 即 $\Sigma = \lambda I_p$ 前提下, 进一步我们推得如下结论:

$$\frac{c_\alpha - \mu_m}{m \cdot \sigma_m} = \sqrt{\frac{r_m^2 - \frac{p}{m}}{r_m^2 - \frac{p}{m} \cdot 2(\lambda - \frac{\lambda^2}{2})}} \cdot \Phi^{-1}(\alpha) + \frac{p \cdot (\lambda - \log \lambda - 1)}{\sqrt{2} \cdot \sqrt{r_m^2 - \frac{p}{m} \cdot 2(\lambda - \frac{\lambda^2}{2})}} \quad (10)$$

(10)式的证明见附录。

从(10)式可看出, 当 $\lambda = 1$ 时, 有 $2\left(\lambda - \frac{\lambda^2}{2}\right) = 1$, $\lambda - \log \lambda - 1 = 0$, 从而 $\frac{c_\alpha - \mu_m}{m \cdot \sigma_m} = \Phi^{-1}(\alpha)$, 这表明 $\Sigma_0 = \lambda I_p = I_p$ 时, 有下式成立:

$$\beta(I_p) \sim \Phi\left(\frac{c_\alpha - \mu_m}{m \cdot \sigma_m}\right) = \Phi(\Phi^{-1}(\alpha)) = \alpha.$$

说明检验功效接近显著性水平 α 。

此外, 当 $p = p_n \rightarrow \infty, n \rightarrow \infty$ 且 $\lim_{n \rightarrow \infty} \frac{p_n}{n} \rightarrow c \in (0, C_1]$ 时 ($C_1 < 1$ 为常数), (10)式中 r_m^2 和 $\frac{p}{m}$ 可视作常数, $\frac{c_\alpha - \mu_m}{m \cdot \sigma_m}$ 主要受 p 和 λ 影响。显然的, 在 λ 为有界量,

我们可以把(10)整理得:

$$\frac{c_\alpha - \mu_m}{m \cdot \sigma_m} = C_1 + p \cdot C_2.$$

其中,

$$C_1 = \sqrt{\frac{r_m^2 - \frac{p}{m}}{r_m^2 - \frac{p}{m} \cdot 2(\lambda - \frac{\lambda^2}{2})}} \cdot \Phi^{-1}(\alpha),$$

$$C_2 = \frac{\lambda - \log \lambda - 1}{\sqrt{2} \cdot \sqrt{r_m^2 - \frac{p}{m} \cdot 2(\lambda - \frac{\lambda^2}{2})}}.$$

视为常数。表明随着 $p = p_n \rightarrow \infty$, $\beta(\Sigma_0) \rightarrow 1$. (见表 3-3)

值得注意的是, 当 $\lambda \in (0, 1) \cup (1, +\infty)$, 有 $\sqrt{\frac{r_m^2 - \frac{p}{m}}{r_m^2 - \frac{p}{m} \cdot 2(\lambda - \frac{\lambda^2}{2})}} < 1$, 所以当 p

有限时, 着眼于第二项 $\frac{\lambda - \log \lambda - 1}{\sqrt{2} \cdot \sqrt{r_m^2 - \frac{p}{m} \cdot 2(\lambda - \frac{\lambda^2}{2})}}$ 而增大 λ 并不一定能使整体 $\frac{C\alpha - \mu_m}{m \cdot \sigma_m}$ 提升,

因为 λ 增大的同时第一项 $\sqrt{\frac{r_m^2 - \frac{p}{m}}{r_m^2 - \frac{p}{m} \cdot 2(\lambda - \frac{\lambda^2}{2})}}$ 减小了。这可能是导致模拟过程中,

限定变量个数 p , 增大 λ 但检验的功效依然很小的原因。

2.3 单个高维总体协方差矩阵等于给定非负定矩阵的检验

假设 $X = (X_1, X_2, \dots, X_n)$ 是独立同分布于 $N_p(\mu, \Sigma)$ 的实值随机向量, 其中 $\mu \in R^p$ 为均值向量, $\Sigma_{p \times p}$ 为协方差矩阵。考虑假设检验问题:

$$H_0: \Sigma = \Sigma_1 \text{ v.s. } H_1: \Sigma \neq \Sigma_1, \quad (11)$$

其中 Σ_1 为给定 $p \times p$ 非负定矩阵。则由奇异值分解 $\Sigma_1 = U \Sigma_0 U'$, 其中 Σ_0 为对角矩阵, U 为正交矩阵, 可将(11)中假设检验问题转化为(6)中针对对角矩阵的假设检验问题。从而有如下方法 1:

方法 1: 对 X 左乘 U^{-1} , 得到 $U^{-1} \cdot X \sim N_p(U^{-1}\mu, U^{-1} \cdot \Sigma \cdot U^{-1T})$, 进而转化为检验 $U^{-1} \cdot X \sim N_p(U^{-1}\mu, U^{-1} \cdot U\Sigma_0U' \cdot U^{-1T}) = N_p(U^{-1}\mu, \Sigma_0)$ 的问题。

另外还可将(11)中的假设检验问题转化为(1)中的假设检验问题, 即如下

方法 2:

方法 2: 对 X 左乘 $\Sigma_1^{-\frac{1}{2}}$, 得到 $\Sigma_1^{-\frac{1}{2}} \cdot X \sim N_p(\Sigma_0^{-\frac{1}{2}}\mu, \Sigma_1^{-\frac{1}{2}} \cdot \Sigma \cdot \Sigma_1^{-\frac{1}{2}T})$, 进而转化为检验 $\Sigma_1^{-\frac{1}{2}} \cdot X \sim N_p(\mu\Sigma_1^{-\frac{1}{2}}, I_p)$ 的问题。

在上述奇异值分解中, 正交矩阵 U 仅起到旋转的作用, 对角矩阵 Σ_0 仅起到拉伸的作用。对于方法 1, 回顾定理 4, $U^{-1} \cdot X \sim N_p(U^{-1}\mu, \Sigma_0)$, 对(6)中对角矩阵的假设检验问题, 我们构建的检验统计量 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 中, $\log \Lambda_m^*$ 由样本离差阵 $A = (U^{-1} \cdot X)'(I_{n \times n} - \frac{1}{n}J_{n \times n})(U^{-1} \cdot X)$ 的特征值决定, μ_m 和 σ_m 由原假设 $H_0: \Sigma = \Sigma_0$ 对应的 Σ_0 的特征值决定。对于方法 2, 回顾定理 1, $\Sigma_1^{-\frac{1}{2}} \cdot X \sim N_p(\mu\Sigma_1^{-\frac{1}{2}}, I_p)$, 对(1)中单位矩阵的假设检验问题, 我们构建的检验统计量 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 中, $\log \Lambda_m^*$ 由样本离差阵 $A = (\Sigma_1^{-\frac{1}{2}} \cdot X)'(I_{n \times n} - \frac{1}{n}J_{n \times n})(\Sigma_1^{-\frac{1}{2}} \cdot X)$ 的特征值决定, μ_m 和 σ_m 由原假设 $H_0: \Sigma = I_p$ 对应的 I_p 的特征值决定。

当 Σ_1 中含有较大的特征值时, 方法 2 中的样本离差阵 A 会接近奇异阵(不满秩), 因为 $(\Sigma_1^{-\frac{1}{2}} \cdot X)$ 接近奇异阵, 这对检验统计量中 $\log \Lambda_m^*$ 的构建产生大的影响, 由于 A 的最小特征值很小很小, 对数函数在区间 $(0, 1]$ 要比 $(1, +\infty)$ 陡峭, 这就导致在 Σ_1 中含有较大的特征值时, 方法 2 检验效果比方法 1 不稳定, 如检验显著性水平波动比方法 1 的要大, 且功效要低(见表 3-2)。

3 数值模拟

3.1 极限分布的模拟

为比较定理 1 和(3)式两个检验统计量在特定 n 和 p 情形下的整体表现，我们选取样本量 $n = 100$ 和变量个数 $p = 5, 30, 60, 90$.生成均值向量为零向量，协方差矩阵为单位矩阵的 $N_p(\mu, I)$ 随机数，模拟次数都为100,000次，得到图 3-1 如下所示（橙色直方图表示100,000次模拟得到的检验统计量的频率分布，蓝色曲线为检验统计量在特定的 n 和 p 取值下的极限分布。上层四幅小图对应(3)式所示的检验统计量 $-2\rho \cdot \log \Lambda^*$ ，下层四幅小图对应定理 1 所示的检验统计量 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ ）：

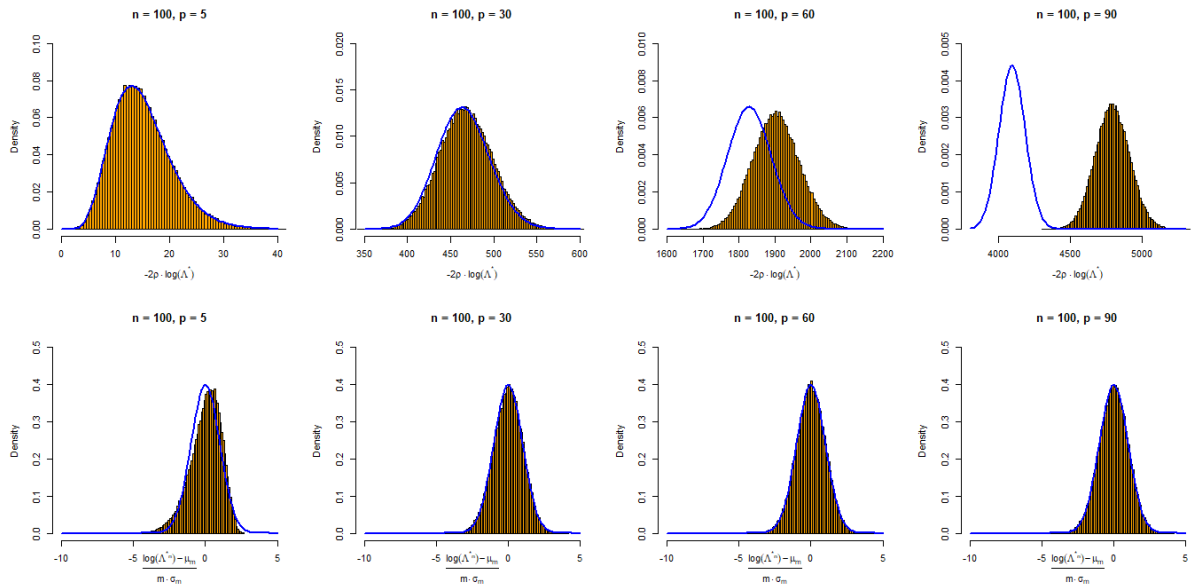


图 3-1 两检验统计量及其极限分布比较

从图 3-1 可知，在 $n = 100, p = 5$ 的情形，传统的似然比检验统计量 $-2\rho \cdot \log \Lambda^*$ 的直方图和其理论上极限分布（ χ^2 分布）契合程度比定理 1 的检验统计量 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 要好，但随着变量个数的增加（从上层四幅小图可知），检验统计量 $-2\rho \cdot \log \Lambda^*$ 对 $\frac{p}{n} \rightarrow c \in (0, 1]$ ，特别是 c 越接近于1的情形，和其对应的

χ^2 分布相差较大。而检验统计量 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 在变量个数越来越接近样本量个数的情形（从下层四幅小图可知），其频率分布直方图和理论上的极限分布（标准正态分布）仍然契合。

因为传统似然比检验统计量是在固定变量个数 p ，让样本量 $n \rightarrow \infty$ 推导得到，所以在 $p = p_n \rightarrow \infty$ ， $n \rightarrow \infty$ 且 $\lim_{n \rightarrow \infty} \frac{p_n}{n} \rightarrow c \in (0,1]$ 情形下表现欠佳。

3.2 检验显著性水平和功效模拟（单个高维总体协方差矩阵等于给定单位矩阵）

为比较定理 1 和(3)式两个检验统计量在特定 n 和 p 情形下的显著性水平和检验功效的差异，我们选取样本量 $n = 100$ 和变量个数 $p = 5, 30, 60, 90$ 。模拟次数都为100,000次，对于检验显著性水平的模拟，生成均值向量为零向量，协方差矩阵为单位矩阵的 $N_p(\mu, I)$ 随机数，检验水平设置在 $\alpha = 0.05$ ；对于检验功效的模拟，生成均值向量为零向量，协方差矩阵如下

$$\Sigma^* = \{\sigma_{ij}\} = \begin{cases} 1, & \text{if } i = j \\ 0.1, & \text{if } 0 < |i - j| \leq 3 \\ 0, & \text{if } |i - j| > 3 \end{cases}$$

所示的 $N_p(\mu, \Sigma^*)$ 随机数。模拟结果如下表 3-1 所示，括号内为相应的标准差估计值。

表 3-1 两检验统计量的检验显著性水平和检验功效对比

	检验显著性水平（在 H_0 下）		检验功效（在 H_1 下）	
	$-2\rho \cdot \log \Lambda^*$	$\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$	$-2\rho \cdot \log \Lambda^*$	$\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$
$n = 100, p = 5$	0.0497 (0.00069)	0.0619 (0.00076)	0.2527 (0.00137)	0.3566 (0.00151)

$n = 100, p = 30$	0.0528 (0.00071)	0.0541 (0.00072)	0.5705 (0.00157)	0.5604 (0.00157)
$n = 100, p = 60$	0.2289 (0.00133)	0.0523 (0.00070)	0.8593 (0.00110)	0.4694 (0.00158)
$n = 100, p = 90$	1.0000 (0.00000)	0.0531 (0.00071)	1.0000 (0.00000)	0.2705 (0.00140)

从表 3-1 可以看出, 在 $n = 100, p = 5$ (变量个数很小) 的情形, 检验统计量 $-2\rho \cdot \log \Lambda^*$ 的效果要比 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 要好, 其显著性水平更接近真实显著性水平 $\alpha = 0.05$, 不过检验效力略低于后者。在 $n = 100, p = 30$ 的情形, 两检验统计量的检验效果相近。随着 p 增加到 60 和 90, 我们看到检验统计量 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 的检验效果要好于 $-2\rho \cdot \log \Lambda^*$ (前者的显著性水平在 p 变大时仍保持接近真实显著性水平), 而检验统计量 $-2\rho \cdot \log \Lambda^*$ 不再适用。值得注意的是, 随着变量个数 p 越来越接近样本量 n , $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 的检验效力明显下降, 这从定理 1 中 σ_m^2 的表达式可以看出, 随着 $n \rightarrow \infty$, $\sigma_m^2 \rightarrow \infty$, 这导致在样本量 n 不够大情形下 (如表中 $n = 100, p = 90$ 情形), 检验的整体上的准确性下降 (检验显著性水平仍接近 0.05, 但检验效力下降明显)。

3.3 检验显著性水平和功效模拟 (单个高维总体协方差矩阵等于给定对角矩阵)

为比较定理 3 和定理 4 (由 Chen & Jiang(2018)^[1]提出) 的检验统计量检验效果的优劣, 和在 2.3 节对于单个高维总体协方差矩阵等于给定非负定矩阵的检验时两个概括性方法的优劣, 我们选取样本量 $n = 100$ 和变量个数 $p = 60$, 对于第 k 个数值实验 (其中编号 $k = 1, \dots, 60$, 模拟次数都为 10,000

次), 对于检验显著性水平的模拟, 生成均值向量为零向量, 协方差矩阵为 $\Sigma_0 = \text{diag}\{1, \dots, 1000\}$ (其中第 k 个数值实验 Σ_0 对角线中 1000 的个数为 k) 的 $N_p(\mu, \Sigma_0)$ 随机数, 检验水平设置在 $\alpha = 0.05$; 对于检验功效的模拟, 生成均值向量为零向量, 协方差矩阵为 $\Sigma^{**} = \text{diag}\{1, \dots, 1200\}$ (其中第 k 个数值实验 Σ_1 对角线中 1200 的个数为 k) $N_p(\mu, \Sigma^{**})$ 随机数。

$T_1 \text{ size}$ 表示定理 3 对应的检验统计量 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 的显著性水平, $T_2 \text{ size}$ 表示定理 4 对应的检验统计量 $\frac{\log \Lambda_n^* - \mu_n}{n \cdot \sigma_n}$ 的显著性水平, $T_3 \text{ size}$ 表示使用 2.3 节方法 2 (将假设检验问题 $H_0: \Sigma = \Sigma_0$ 转换为 $H_0: \Sigma = I_p$) 对应检验统计量的显著性水平, 类似的, $T_1 \text{ power}$ 、 $T_2 \text{ power}$ 和 $T_3 \text{ power}$ 分别表示它们的检验功效。模拟结果如下表 3-2 所示:

表 3-2 三种方法的检验统计量的检验显著性水平和检验功效对比

编号 k	检验显著性水平 (在 H_0 下)			检验功效 (在 H_1 下)		
	$T_1 \text{ size}$	$T_2 \text{ size}$	$T_3 \text{ size}$	$T_1 \text{ power}$	$T_2 \text{ power}$	$T_3 \text{ power}$
1	0.0499	0.0471	0.0494	0.3117	0.3104	0.0545
2	0.0486	0.0481	0.0497	0.4880	0.4969	0.0520
3	0.0510	0.0478	0.0540	0.6416	0.6405	0.0521
4	0.0504	0.0498	0.0530	0.7489	0.7521	0.0549
5	0.0518	0.0425	0.0521	0.8366	0.8365	0.0564
6	0.0472	0.0476	0.0584	0.8961	0.8878	0.0535
7	0.0481	0.0475	0.0539	0.9346	0.9359	0.0604
8	0.0504	0.0455	0.0526	0.9565	0.9571	0.0607
9	0.0526	0.0503	0.0540	0.9725	0.9715	0.0581
10	0.0466	0.0475	0.0539	0.9814	0.9835	0.0608
11	0.0501	0.0468	0.0527	0.9881	0.9900	0.0641
12	0.0479	0.0519	0.0541	0.9941	0.9932	0.0666
13	0.0529	0.0511	0.0518	0.9960	0.9962	0.0616
14	0.0517	0.0464	0.0541	0.9972	0.9976	0.0629
15	0.0492	0.0488	0.0505	0.9986	0.9987	0.0713
16	0.0471	0.0478	0.0533	0.9988	0.9989	0.0689

17	0.0507	0.0468	0.0496	0.9993	0.9992	0.0678
18	0.0490	0.0524	0.0533	0.9998	0.9995	0.0796
19	0.0503	0.0508	0.0492	0.9997	0.9999	0.0753
20	0.0521	0.0464	0.0530	0.9998	0.9998	0.0748
21	0.0509	0.0489	0.0509	1.0000	0.9998	0.0778
22	0.0499	0.0492	0.0552	1.0000	1.0000	0.0859
23	0.0485	0.0507	0.0568	1.0000	1.0000	0.0880
24	0.0515	0.0492	0.0544	1.0000	1.0000	0.0874
25	0.0508	0.0516	0.0530	1.0000	1.0000	0.0895
26	0.0496	0.0492	0.0515	1.0000	1.0000	0.0944
27	0.0495	0.0524	0.0532	1.0000	1.0000	0.1025
28	0.0481	0.0457	0.0559	1.0000	1.0000	0.0989
29	0.0470	0.0489	0.0558	1.0000	1.0000	0.1000
30	0.0453	0.0470	0.0534	1.0000	1.0000	0.1083
31	0.0472	0.0506	0.0503	1.0000	1.0000	0.1154
32	0.0487	0.0501	0.0522	1.0000	1.0000	0.1180
33	0.0469	0.0480	0.0544	1.0000	1.0000	0.1251
34	0.0498	0.0492	0.0515	1.0000	1.0000	0.1250
35	0.0509	0.0466	0.0554	1.0000	1.0000	0.1254
36	0.0514	0.0488	0.0523	1.0000	1.0000	0.1284
37	0.0525	0.0487	0.0515	1.0000	1.0000	0.1345
38	0.0510	0.0484	0.0555	1.0000	1.0000	0.1400
39	0.0496	0.0471	0.0527	1.0000	1.0000	0.1467
40	0.0472	0.0466	0.0550	1.0000	1.0000	0.1506
41	0.0529	0.0500	0.0512	1.0000	1.0000	0.1576
42	0.0493	0.0531	0.0520	1.0000	1.0000	0.1666
43	0.0478	0.0489	0.0538	1.0000	1.0000	0.1578
44	0.0527	0.0486	0.0528	1.0000	1.0000	0.1716
45	0.0513	0.0491	0.0536	1.0000	1.0000	0.1716
46	0.0555	0.0512	0.0469	1.0000	1.0000	0.1812
47	0.0473	0.0476	0.0549	1.0000	1.0000	0.1871
48	0.0485	0.0507	0.0551	1.0000	1.0000	0.1980
49	0.0497	0.0471	0.0513	1.0000	1.0000	0.2005
50	0.0506	0.0487	0.0515	1.0000	1.0000	0.2071
51	0.0487	0.0473	0.0496	1.0000	1.0000	0.2066
52	0.0459	0.0460	0.0540	1.0000	1.0000	0.2158
53	0.0480	0.0485	0.0552	1.0000	1.0000	0.2229
54	0.0512	0.0554	0.0556	1.0000	1.0000	0.2278
55	0.0504	0.0519	0.0534	1.0000	1.0000	0.2365
56	0.0497	0.0499	0.0523	1.0000	1.0000	0.2355
57	0.0509	0.0426	0.0527	1.0000	1.0000	0.2460
58	0.0509	0.0484	0.0517	1.0000	1.0000	0.2502

59	0.0514	0.0455	0.0544	1.0000	1.0000	0.2663
60	0.0530	0.0466	0.0539	1.0000	1.0000	0.2698

为看出三种方法对应检验显著性水平和理论显著性水平 $\alpha = 0.05$ 的差异，可对 T_1 size、 T_2 size 和 T_3 size 分别对应的 60 个值取均值，分别得到 0.04983 (0.0020)、0.0486 (0.0024) 和 0.0530 (0.0021)。因而 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ (标为 T_1 ，定理 3) 的显著性水平是更接近理论显著性水平的，而且也更稳定。在保证显著性水平情况下看三种方法的检验功效，从表 3-2 可以看出，当样本协方差矩阵含有较高的特征值时，方法 2 的检验功效很低，而且随着所含较高特征值的个数的增加，检验功效的波动也很大，而 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ (标为 T_1 ，定理 3) 和 $\frac{\log \Lambda_n^* - \mu_n}{n \cdot \sigma_n}$ (标为 T_2 ，定理 4) 的检验功效都很高，且随着所含较高特征值个数的增加而迅速增加。

3.4 检验功效变化的数值模拟(单个高维总体协方差矩阵等于给定对角矩阵)

回顾 2.2 节 (10) 式，为验证在 $p = p_n \rightarrow \infty, n \rightarrow \infty, \lim_{n \rightarrow \infty} \frac{p_n}{n} \rightarrow c \in (0, C_1]$ 且样本协方差矩阵的特征值不变时 (r_m^2 和 $\frac{p}{m}$ 可视作常数)， $\beta(\Sigma_0) \rightarrow 1$ 主要是由 $p = p_n \rightarrow \infty$ 而拉动的，我们进行如下模拟。

选取样本量 $m = 100, 150, 200, 250, 300, 350$ 和对应变量个数 $p = 60, 90, 120, 150, 180, 210$ (保证 $\frac{p}{m}$ 不变)，模拟次数都为 10,000 次，对于检验显著性水平的模拟，生成均值向量为零向量，协方差矩阵为对角线元素全为 2 的对角矩阵 Σ_0 的 $N_p(\mu, \Sigma_0)$ 随机数 (固定 (10) 式中的 λ 保持不变)，检验水平设置在 $\alpha = 0.05$ ；对于检验功效的模拟，生成均值向量为零向量，协方差矩阵

为对角线值全为2.1的对角矩阵 Σ'' 的 $N_p(\mu, \Sigma'')$ 随机数。计算定理3所示检验统计量 $\frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m}$ 对应的检验显著性水平和检验功效，结果如表3-3所示：

表 3-3 协方差矩阵等于对角矩阵检验的检验功效和变量个数的关系对比

	检验显著性水平（在 H_0 下）	检验功效（在 H_1 下）
$m = 100, p = 60$	0.0506 (0.00219)	0.6066 (0.00489)
$m = 150, p = 90$	0.0509 (0.00220)	0.9115 (0.00284)
$m = 200, p = 120$	0.0484 (0.00215)	0.9924 (0.00087)
$m = 250, p = 150$	0.0529 (0.00224)	0.9997 (0.00017)
$m = 300, p = 180$	0.0494 (0.00217)	1.0000 (0.00000)
$m = 350, p = 210$	0.0491 (0.00216)	1.0000 (0.00000)

从表3-3可以看出，当我们固定协方差矩阵的特征值，保持变量个数 p 和与样本量相关的 m 的比不变，检验统计量的模拟显著性水平和真实显著性水平 $\alpha = 0.05$ 和接近，但检验功效会随着 p 的增长而被迅速拉高。这也表明，在变量个数和样本量成比例增加的过程中，若能控制协方差矩阵的特征值相近，或通过变换使之相近，则在保证控制协方差矩阵的假设检验显著性水平前提下，还能有较高的检验功效。

4 结论和展望

本文对现有的高维总体协方差矩阵假设检验的方法进行综述与回顾，并针对单个高维总体协方差矩阵等于单位矩阵、对角矩阵和对称正定矩阵的假设检验问题，分别提出检验方法。结果表明，该检验方法比传统的似然比检验有更好的检验显著性水平和检验功效。并且对于协方差矩阵等于特定对称正定矩阵的检验问题，在样本协方差矩阵含有较大特征值时，提出了更稳健的检验流程。我们有如下几点结论与展望：

- (a) 对多种假设检验类型，似然比检验框架下推得的检验统计量的组成成分在超高维的情形或有启发，如文中检验统计量含有样本协方差矩阵特征值的和、特征值的平方和，这在超高维便可用 $tr(S)$ 等来浓缩检验对象的信息。
- (b) 对于两个或两个以上总体协方差矩阵的假设检验问题，推广要比本文的推广方法复杂，因为涉及 Zonal 多项式，在复杂原假设下（如检验两个总体协方差矩阵是否等于特定对角矩阵、对称正定矩阵等），检验统计量的化简和计算机运算需要更多工具。
- (c) 本文假定样本服从多元正态分布假定，并且有 $\lim_{n \rightarrow \infty} \frac{p}{n} \rightarrow c \in (0,1]$ 的限制，进一步研究可考虑放松正态假定，或改为更易验证的接近正态假定的条件；或改进检验统计量使之适用于 $p > n$ 的情形。
- (d) 基于本文的研究，可进一步探索该检验统计量的统计性质，如是否具有无偏性。也可使用判决分析的方法比较假设检验，而不使用功效函数，使检验统计量具有损失函数最优性。

致谢

陈寅恪《海宁王静安先生纪念碑铭》曰：“士之读书治学，盖将以脱心志于俗谛之桎梏，真理因得以发扬。”读书先与治学，心志重于真理，暨南园求学四年，文亮喟然叹曰：“博我以文，约我以礼。暨南园所遇循循然善诱者唯二，一约堂，一孤狼。”自 2015 年以来，恩师一直关心本人的学业与生活，并寄予充分的信任与期许，愧无以报，为终身之憾。就读暨南园期间，有幸结识明诚堂诸兄，明诚明之性，诚明诚之教，认真生活，深得切磋之乐。旧友振兴、杨强、丘梵诸兄，时通音问，天涯若比邻。知我者，二三子。

纵观暨南园课堂内外至略窥专业堂奥，及此论文开题、中期考核、答辩诸环节，柳向东、陈光慧、伍业锋、杨广仁、姜云卢、伍海军诸位老师费心指导，深受启益。在此谨致谢忱。

需要说明的是，拙稿论述架构的确立，乃受孤狼启发。数月来，本人随明诚堂诸兄几度蹭约堂的课，随组员几度蹭孤狼的饭，得识前辈学人中正温粹、淡泊自守的气象，窃奉以为步趋之仪型。

家父常忧吾业国学、经济统计无以谋生，而后竟得自立，实属侥幸耳。

附录

定理 1 证明

对于两数列 $\{a_n; n \geq 1\}$ 和 $\{b_n; n \geq 1\}$, $a_n = O(b_n)$ 表示当 $n \rightarrow \infty$ 时, 有 $\limsup_{n \rightarrow \infty} \left| \frac{a_n}{b_n} \right| < \infty$ 。 $a_n = o(b_n)$ 表示当 $n \rightarrow \infty$ 时, 有 $\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = 0$ 。对于两函数 $f(x)$ 和 $g(x)$, 在 $x \rightarrow x_0 \in [-\infty, \infty]$ 时, $f(x) = O(g(x))$ 和 $f(x) = o(g(x))$ 也是类似定义。

作如下定义 (参考 Muirhead(1982)^[19], Theorem 2.1.12):

$$\Gamma_p(z) := \pi^{\frac{p(p-1)}{4}} \prod_{i=1}^p \Gamma(z - \frac{1}{2}(i-1)) \quad (12)$$

其中 z 为满足 $\text{Re}(z) > \frac{1}{2}(p-1)$ 的复数。

引理 1 (Jiang & Yang(2013)^[12], Lemma 5.4): $\Gamma_p(z)$ 由(12)式定义, 记 $n > p = p_n$, $r_n = (-\log(1 - \frac{p}{n}))^{\frac{1}{2}}$, 若 $n \rightarrow \infty$ 时满足 $\frac{p}{n} \rightarrow c \in (0,1]$, $s = s_n = O(1/r_n)$, $t = t_n = O(1/r_n)$ 。则当 $n \rightarrow \infty$ 时, 有下式成立:

$$\log \frac{\Gamma_p(\frac{n}{2}+t)}{\Gamma_p(\frac{n}{2}+s)} = p(t-s)(\log n - 1 - \log 2) + r_n^2 \left[(t^2 - s^2) - \left(p - n + \frac{1}{2} \right) (t-s) \right] + o(1) \quad (13)$$

引理 2 (Muirhead(1982)^[19], Theorem 8.4.7 和 Corollary 8.4.8 P384): (2)式定义的修正的似然比检验统计量 Λ^* 的 t 阶矩, 由下式给出:

$$E(\Lambda^{*t}) = \left(\frac{2e}{m} \right)^{\frac{mpt}{2}} \cdot \frac{(\det \Sigma)^{-\frac{mt}{2}}}{\det(I+t\Sigma)^{\frac{m(1+t)}{2}}} \cdot \frac{\Gamma_p(\frac{m}{2}(1+t))}{\Gamma_p(\frac{m}{2})} \quad (14)$$

从而, 在原假设 $H_0: \Sigma = I_p$ 成立时, Λ^* 的 t 阶矩满足:

$$E(\Lambda^{*t}) = \left(\frac{2e}{m}\right)^{\frac{mpt}{2}} \cdot (1+t)^{-\frac{mp(1+t)}{2}} \cdot \frac{\Gamma_p(\frac{m}{2}(1+t))}{\Gamma_p(\frac{m}{2})} \quad (15)$$

定理 1 证明:

由 $\lim_{n \rightarrow \infty} \frac{p}{n} \rightarrow c \in (0,1]$, 可知 $\lim_{m \rightarrow \infty} \frac{p}{m} \rightarrow c \in (0,1]$ 。对于任意 $x < 1$, 有 $\log(1-x) < -x$, 从而 $\sigma_{m,0}^2 > 0$ 对于任意 $m \geq 2$ 成立 ($m \geq 2$ 是为了使 r_m 有意义)。进而, 有:

$$\lim_{m \rightarrow \infty} \sigma_{m,0}^2 = \begin{cases} -\frac{1}{2}(\log(1-c) + c), & \text{if } c < 1; \\ +\infty, & \text{if } c = 1. \end{cases} \quad (16)$$

表明 $\lim_{m \rightarrow \infty} \sigma_m^2$ 这极限总为正的, 因而有:

$$\delta_0 := \inf\{\sigma_{m,0}; m \geq 2\} > 0.$$

固定 h , 其中 $|h| < \delta_0$, 则对于任意 $m \geq 2$, 满足 $h > -\delta_0 \geq -\sigma_{m,0}$, 进而有:

$$\frac{p}{m} - 1 \leq \frac{m-1}{m} - 1 = -\frac{1}{m} < \frac{h}{m \cdot \sigma_{m,0}} \quad (17)$$

设 $t := t_m = \frac{h}{m \cdot \sigma_{m,0}}$, ($m \geq 2$), 则

$$t > \frac{p}{m} - 1, (m \geq 2) \quad (18)$$

由(16)知, $\{t_m; m \geq 2\}$ 有界, 从而根据(18)和引理 2, 可得:

$$E e^{t \cdot \log \Lambda_m^*} = E \Lambda_m^{*t} = \left(\frac{2e}{m}\right)^{\frac{mpt}{2}} \cdot (1+t)^{-\frac{mp(1+t)}{2}} \cdot \frac{\Gamma_p(\frac{m}{2}(1+t))}{\Gamma_p(\frac{m}{2})} \quad (19)$$

其中记号 Λ_m^* 表示 Λ^* 依赖于 m .

为证定理 1, 只需证明当 $m \rightarrow \infty$, 对于任意满足 $|h| < \delta_0$ 的 h , 有下式成立:

$$E \exp\left\{\frac{\log \Lambda_m^* - \mu_{m,0}}{m \cdot \sigma_{m,0}} h\right\} \rightarrow e^{\frac{h^2}{2}} \quad (20)$$

由(16)和 t 的定义, 可知 $\frac{mt}{2} = \frac{h}{2\sigma_m} = O(\frac{1}{r_{m-1}})$, 从而在引理 1 中取 $s = 1, t = \frac{mt}{2}$ 可以得到, 当 $m \rightarrow \infty$:

$$\begin{aligned}
 & \log \frac{\Gamma_p(\frac{m}{2}(1+t))}{\Gamma_p(\frac{m}{2})} \\
 &= \frac{mpt}{2}(\log m - 1 - \log 2) + r_m^2 \left[\frac{m^2 t^2}{4} - \left(p - m + \frac{1}{2} \right) \frac{mt}{2} \right] + o(1) \\
 &= \frac{mpt}{2}(\log m - 1 - \log 2) + \frac{m^2 r_m^2}{4} t^2 + \frac{1}{4}(2m - 2p - 1)mr_m^2 t + o(1)
 \end{aligned} \tag{21}$$

由 $\log(1+s) = s - \frac{s^2}{2} + o(s^2)$, 可知当 $m \rightarrow \infty$, $\lim_{m \rightarrow \infty} t = 0$, 有:

$$\begin{aligned}
 & (1+t) \log(1+t) \\
 &= (1+t) \left(t - \frac{t^2}{2} + o(t^2) \right) \\
 &= t - \frac{t^2}{2} + o(t^2) + t^2 - \frac{t^3}{2} + to(t^2) \\
 &= t + \frac{t^2}{2} + o(t^2)
 \end{aligned}$$

从而有:

$$\begin{aligned}
 & \log(1+t) - \frac{mp(1+t)}{2} = -\frac{1}{2}mp(1+t) \log(1+t) \\
 &= -\frac{1}{2}mp \left(t + \frac{t^2}{2} + o(t^2) \right) \\
 &= -\frac{1}{2}mp \left(t + \frac{t^2}{2} \right) + o(1)
 \end{aligned} \tag{22}$$

由(20), (21)和(22), 当 $m \rightarrow \infty$, 对于任意满足 $|h| < \delta_0$ 的 h , 可得:

$$\begin{aligned}
& \log E(\Lambda_m^{*t}) \\
&= \frac{mpt}{2} \log \frac{2e}{m} - \frac{1}{2}mp \left(t + \frac{t^2}{2} \right) + \frac{mpt}{2} (\log m - 1 - \log 2) + \frac{m^2 r_m^2}{4} t^2 \\
&\quad + \frac{1}{4} (2m - 2p - 1) m r_m^2 t + o(1) \\
&= \frac{1}{4} (m^2 r_m^2 - mp) t^2 + \frac{1}{4} [(2m - 2p - 1) m r_m^2 - 2mp] t + o(1) \\
&= \frac{1}{2} m^2 \sigma_{m,0}^2 t^2 + \mu_{m,0} t + o(1) \\
&= \frac{h^2}{2} + \mu_{m,0} t + o(1)
\end{aligned}$$

从而得到：

$$E \exp \left\{ \frac{\log \Lambda_m^* - \mu_{m,0}}{m \cdot \sigma_{m,0}} h \right\} = \log E(\Lambda_m^{*t}) - \mu_{m,0} t \rightarrow \frac{h^2}{2}$$

当 $m \rightarrow \infty$ ，对任意满足 $|h| < \delta_0$ 的 h 成立。定理 1 得证。

定理 3 证明

定理 3 证明：

仿照定理 1 的证明，由于 $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_p < \infty$ 为有界量，我们可类似定义 δ_0 和 t ，我们只需对(19)和(22)式稍作修改。

为证定理 3，只需证明当 $m \rightarrow \infty$ ，对于任意满足 $|h| < \delta_0$ 的 h ，有下式成立：

$$E \exp \left\{ \frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m} h \right\} \rightarrow e^{\frac{h^2}{2}} \quad (23)$$

由引理 2，在原假设 $H_0: \Sigma = \Sigma_0$ 成立时 (Σ_0 为实对角矩阵，对角元素 $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_p < \infty$)，则 Λ^* 的 t 阶矩满足：

$$E(\Lambda^{*t}) = \left(\frac{2e}{m}\right)^{\frac{mpt}{2}} \cdot \frac{(\prod_{i=1}^p \lambda_i)^{\frac{mt}{2}}}{[\prod_{i=1}^p (\lambda_i t + 1)]^{\frac{m(t+1)}{2}}} \cdot \frac{\Gamma_p(\frac{m}{2}(1+t))}{\Gamma_p(\frac{m}{2})} \quad (24)$$

由 $\log(1+s) = s - \frac{s^2}{2} + o(s^2)$, 可知当 $m \rightarrow \infty$, $\lim_{m \rightarrow \infty} \lambda_i t = 0$, 有:

$$\begin{aligned} & (1+t) \log(\lambda_i t + 1) \\ &= (1+t) \left(\lambda_i t - \frac{\lambda_i^2 t^2}{2} + o(\lambda_i^2 t^2) \right) \\ &= \lambda_i t - \frac{\lambda_i^2 t^2}{2} + o(\lambda_i^2 t^2) + \lambda_i t^2 - \frac{\lambda_i^2 t^3}{2} + t o(\lambda_i^2 t^2) \\ &= \lambda_i t + (\lambda_i - \frac{\lambda_i^2}{2}) t^2 + o(1) \end{aligned}$$

结合(24)式, 进而有:

$$\begin{aligned} & \log \frac{(\prod_{i=1}^p \lambda_i)^{\frac{mt}{2}}}{[\prod_{i=1}^p (\lambda_i t + 1)]^{\frac{m(t+1)}{2}}} \\ &= \frac{mt}{2} \cdot \left(\sum_{i=1}^p \log \lambda_i \right) - \frac{m(t+1)}{2} \cdot \left[\sum_{i=1}^p \log(\lambda_i t + 1) \right] \\ &= \frac{m}{2} \cdot \left(\sum_{i=1}^p \log \lambda_i \right) \cdot t - \frac{m}{2} \cdot \left(\sum_{i=1}^p \lambda_i \right) \cdot t - \frac{m}{2} \cdot \left(\sum_{i=1}^p \left(\lambda_i - \frac{\lambda_i^2}{2} \right) \right) \cdot t^2 + o(1) \end{aligned} \quad (25)$$

由(24), (21)和(25), 当 $m \rightarrow \infty$, 对于任意满足 $|h| < \delta_0$ 的 h , 可得:

$$\begin{aligned}
& \log E(\Lambda_m^{*t}) \\
&= \frac{mpt}{2} \log \frac{2e}{m} + \frac{m}{2} \cdot \left(\sum_{i=1}^p \log \lambda_i \right) \cdot t - \frac{m}{2} \cdot \left(\sum_{i=1}^p \lambda_i \right) \cdot t - \frac{m}{2} \cdot \sum_{i=1}^p \left(\lambda_i - \frac{\lambda_i^2}{2} \right) \\
&\quad \cdot t^2 + \frac{mpt}{2} (\log m - 1 - \log 2) + \frac{m^2 r_m^2}{4} t^2 \\
&\quad + \frac{1}{4} (2m - 2p - 1) m r_m^2 t + o(1) \\
&= \frac{1}{4} \left(m^2 r_m^2 - 2m \cdot \sum_{i=1}^p \left(\lambda_i - \frac{\lambda_i^2}{2} \right) \right) t^2 \\
&\quad + \frac{1}{4} \left[(2m - 2p - 1) m r_m^2 + 2m \cdot \sum_{i=1}^p \log \lambda_i - 2m \cdot \sum_{i=1}^p \lambda_i \right] t \\
&\quad + o(1) \\
&= \frac{1}{2} m^2 \sigma_m^2 t^2 + \mu_m t + o(1) \\
&= \frac{h^2}{2} + \mu_m t + o(1)
\end{aligned}$$

从而得到：

$$E \exp \left\{ \frac{\log \Lambda_m^* - \mu_m}{m \cdot \sigma_m} h \right\} = \log E(\Lambda_m^{*t}) - \mu_m t \rightarrow \frac{h^2}{2}$$

当 $m \rightarrow \infty$ ，对任意满足 $|h| < \delta_0$ 的 h 成立。定理 3 得证。

(10) 式证明

回顾定理 1 和定理 3， $\mu_{m,0}$ 、 $\sigma_{m,0}^2$ 、 μ_m 和 σ_m^2 分别定义如下：

$$\begin{aligned}
\mu_{m,0} &= \frac{1}{4} m [(2m - 2p - 1) \cdot r_m^2 - 2p], \\
\sigma_{m,0}^2 &= \frac{1}{2} \left(r_m^2 - \frac{p}{m} \right),
\end{aligned}$$

$$\mu_m = \frac{1}{4}m \left[(2m - 2p - 1) \cdot r_m^2 + 2 \sum_{i=1}^p \log \lambda_i - 2 \sum_{i=1}^p \lambda_i \right],$$

$$\sigma_m^2 = \frac{1}{2} [r_m^2 - \frac{2}{m} \sum_{i=1}^p (\lambda_i - \frac{\lambda_i^2}{2})],$$

从而有：

$$\begin{aligned} \frac{c_\alpha - \mu_m}{m \cdot \sigma_m} &= \frac{\mu_{m,0} + m \cdot \sigma_{m,0} \cdot \Phi^{-1}(\alpha) - \mu_m}{m \cdot \sigma_m} \\ &= \frac{\sigma_{m,0}}{\sigma_m} \cdot \Phi^{-1}(\alpha) + \frac{\sum_{i=1}^p \lambda_i - \sum_{i=1}^p \log \lambda_i - p}{2\sigma_m} \end{aligned}$$

进而，在在 $0 < \lambda_1 = \lambda_2 = \dots = \lambda_p = \lambda < \infty$ ，即 $\Sigma = \lambda I_p$ 前提下， λ 代入上式便可推得(10)式。

数值模拟代码

power = FALSE 对应 size, 输入 sigma0; power = TRUE 对应 power, 输入 sigma0 和 sigma1。针对对角矩阵的 size 和 power 计算。

```
jwen_test <- function(power = FALSE, n, mu, sigma0, sigma1, cycles) {
  require(MASS)
  get_jwen_test <- function(data, n, sigma0) {
    m <- n - 1
    p <- ncol(data)
    r_m <- sqrt(-log(1 - p/m))
    A <- t(data) %*% (diag(1, n) - matrix(rep(1, n*n), nrow = n)/n) %*% data
    lambda_A <- eigen(A, symmetric = TRUE, only.values = TRUE)$values
    lambda <- diag(sigma0)
    mu_m <- 0.25 * ((2*m-2*p-1)*m*r_m^2 + 2*m*sum(log(lambda)) - 2*m*sum(lambda))
    sigma_m <- sqrt(0.5 * (r_m^2 - 2*sum(lambda)/m + sum(lambda^2)/m))
    return((m*p*(1 - log(m))/2 - 0.5*sum(lambda_A) + m*sum(log(lambda_A))/2 - mu_m) / (m*sigma_m))
  }
  test_collect <- c()
  if (power == FALSE) {
    for (i in 1:cycles) {
      data <- mvrnorm(n = n, mu = mu, Sigma = sigma0)
      test <- get_jwen_test(data, n, sigma0)
      test_collect <- c(test_collect, test)
    }
    alpha <- mean(test_collect < qnorm(0.025, 0, 1) | test_collect > qnorm(0.975, 0, 1))
    return(list(test = test_collect, alpha = alpha))
  }
  if (power == TRUE) {
    for (i in 1:cycles) {
      data <- mvrnorm(n = n, mu = mu, Sigma = sigma1)
      test <- get_jwen_test(data, n, sigma0)
      test_collect <- c(test_collect, test)
    }
    alpha <- mean(test_collect < qnorm(0.025, 0, 1) | test_collect > qnorm(0.975, 0, 1))
    return(list(test = test_collect, power = alpha))
  }
}

# 示例
# jwen_test(power = FALSE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), cycles = 100)
# jwen_test(power = TRUE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), sigma1 = diag(100, 60), cycles = 100)
```

power = FALSE 对应 size, 输入 sigma0; power = TRUE 对应 power, 输入 sigma0 和 sigma1。针对任意对称正定矩阵转换为单位阵假设检验的 size 和 power 计算。

```
translate_test <- function(power = FALSE, n, mu, sigma0, sigma1, cycles) {
  require(MASS)
  matrix0.5 <- function(Matrix) {
    Lambda <- solve(eigen(Matrix)$vectors) %*% Matrix %*% eigen(Matrix)$vectors
    result <- eigen(Matrix)$vectors %*% diag(diag(Lambda)^(-0.5)) %*% solve(eigen(Matrix)$vectors)
    return(result)
  }
  get_translate_test <- function(data, n, sigma0) {
    m <- n - 1
    p <- ncol(data)
    r_m <- sqrt(-log(1 - p/m))
    data1 <- data %*% matrix0.5(sigma0)
    A <- t(data1) %*% (diag(1, n) - matrix(rep(1, n*n), nrow = n)/n) %*% data1
    lambda_A <- eigen(A, symmetric = TRUE, only.values = TRUE)$values
    r_m <- sqrt(-log(1 - p/m))
    mu_m <- 0.25 * ((2*m-2*p-1)*m*r_m^2 - 2*m*p)
    sigma_m <- sqrt(0.5 * (r_m^2 - p/m))
    return((m*p*(1 - log(m))/2 - 0.5*sum(lambda_A) + m*sum(log(lambda_A))/2 - mu_m) / (m*sigma_m))
  }
  test_collect <- c()
  if (power == FALSE) {
    for (i in 1:cycles) {
      data <- mvrnorm(n = n, mu = mu, Sigma = sigma0)
      test <- get_translate_test(data, n, sigma0)
      test_collect <- c(test_collect, test)
    }
    alpha <- mean(test_collect < qnorm(0.025, 0, 1) | test_collect > qnorm(0.975, 0, 1))
    return(list(test = test_collect, alpha = alpha))
  }
  if (power == TRUE) {
    for (i in 1:cycles) {
      data <- mvrnorm(n = n, mu = mu, Sigma = sigma1)
      test <- get_translate_test(data, n, sigma0)
      test_collect <- c(test_collect, test)
    }
    alpha <- mean(test_collect < qnorm(0.025, 0, 1) | test_collect > qnorm(0.975, 0, 1))
    return(list(test = test_collect, power = alpha))
  }
}
# 示例
# translate_test(power = FALSE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), cycles = 100)
```

```

# translate_test(power = TRUE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), sigma1 = diag(100, 60),
cycles = 100)

# power = FALSE 对应 size，输入 sigma0； power = TRUE 对应 power，输入 sigma0 和 sigma1。针对
和我类似的 Chen（2018）针对任意对称正定协方差阵的 size 和 power 计算。
chen_test <- function(power = FALSE, n, mu, sigma0, sigma1, cycles) {
  require(MASS)
  get_chen_test <- function(data, n, sigma0) {
    m <- n - 1
    p <- ncol(data)
    A <- t(data) %*% (diag(1, n) - matrix(rep(1, n*n), nrow = n)/n) %*% data
    lambda_A <- eigen(A, symmetric = TRUE, only.values = TRUE)$values
    lambda <- eigen(sigma0, symmetric = TRUE, only.values = TRUE)$values
    mu_n <- -0.25*m*(2*n - 2*p - 3)*log(1 - p/m) + 0.5*m*(sum(log(lambda)) - sum(lambda))
    sigma_n <- sqrt(-0.5*(p/m + log(1-p/m)) + sum(diag((sigma0-diag(1, p))%*%(sigma0-diag(1,
p)))/(2*n))
    return((m*p*(1-log(m))/2 - 0.5*sum(lambda_A) + m*sum(log(lambda_A))/2 - mu_n)/(n*sigma_n))
  }
  test_collect <- c()
  if (power == FALSE) {
    for (i in 1:cycles) {
      data <- mvrnorm(n = n, mu = mu, Sigma = sigma0)
      test <- get_chen_test(data, n, sigma0)
      test_collect <- c(test_collect, test)
    }
    alpha <- mean(test_collect < qnorm(0.025, 0, 1) | test_collect > qnorm(0.975, 0, 1))
    return(list(test = test_collect, alpha = alpha))
  }
  if (power == TRUE) {
    for (i in 1:cycles) {
      data <- mvrnorm(n = n, mu = mu, Sigma = sigma1)
      test <- get_chen_test(data, n, sigma0)
      test_collect <- c(test_collect, test)
    }
    alpha <- mean(test_collect < qnorm(0.025, 0, 1) | test_collect > qnorm(0.975, 0, 1))
    return(list(test = test_collect, power = alpha))
  }
}
# 示例
# chen_test(power = FALSE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), cycles = 100)
# chen_test(power = TRUE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), sigma1 = diag(100, 60), cycles
= 100)

```

power = FALSE 对应 size, 输入 sigma0; power = TRUE 对应 power, 输入 sigma0 和 sigma1。针对传统似然比卡方检验的 size 和 power 计算。

```
chi_test <- function(power = FALSE, n, mu, sigma0, sigma1, cycles) {
  require(MASS)
  p <- length(mu)
  get_chi_test <- function(data, n, p, sigma0) {
    m <- n - 1
    A <- t(data) %*% (diag(1, n) - matrix(rep(1, n*n), nrow = n)/n) %*% data
    lambda_A <- eigen(A, symmetric = TRUE, only.values = TRUE)$values
    rho <- 1 - (2*p^2 + 3*p - 1)/(6*m*(p + 1))
    return(-2 * rho * (m*p*(1 - log(m))/2 - 0.5*sum(lambda_A) + m*sum(log(lambda_A))/2))
  }
  test_collect <- c()
  if (power == FALSE) {
    for (i in 1:cycles) {
      data <- mvrnorm(n = n, mu = mu, Sigma = sigma0)
      test <- get_chi_test(data, n, p, sigma0)
      test_collect <- c(test_collect, test)
    }
    alpha <- mean(test_collect < qchisq(0.025, df = 0.5*p*(p+1)) | test_collect > qchisq(0.975, df =
0.5*p*(p+1)))
    return(list(test = test_collect, alpha = alpha))
  }
  if (power == TRUE) {
    for (i in 1:cycles) {
      data <- mvrnorm(n = n, mu = mu, Sigma = sigma1)
      test <- get_chi_test(data, n, p, sigma0)
      test_collect <- c(test_collect, test)
    }
    alpha <- mean(test_collect < qchisq(0.025, df = 0.5*p*(p+1)) | test_collect > qchisq(0.975, df =
0.5*p*(p+1)))
    return(list(test = test_collect, power = alpha))
  }
}
# 示例
# chi_test(power = FALSE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), cycles = 100)
# chi_test(power = TRUE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), sigma1 = diag(100, 60), cycles =
100)
```

power = FALSE 对应 size, 输入 sigma0; power = TRUE 对应 power, 输入 sigma0 和 sigma1。针对对称正定矩阵正交变换后进行对角矩阵检验的 size 和 power 计算。

```
jwen_test2 <- function(power = FALSE, n, mu, sigma0, sigma1, cycles) {
  require(MASS)
```

```

get_jwen_test2 <- function(data, n, sigma0) {
  m <- n - 1
  p <- ncol(data)
  r_m <- sqrt(-log(1 - p/m))
  eigen_sigma0 <- eigen(sigma0)
  data1 <- data %*% eigen_sigma0$vectors
  A <- t(data1) %*% (diag(1, n) - matrix(rep(1, n*n), nrow = n)/n) %*% data1
  lambda_A <- eigen(A, symmetric = TRUE, only.values = TRUE)$values
  lambda <- eigen_sigma0$values
  mu_m <- 0.25 * ((2*m-2*p-1)*m*r_m^2 + 2*m*sum(log(lambda)) - 2*m*sum(lambda))
  sigma_m <- sqrt(0.5 * (r_m^2 - 2*sum(lambda)/m + sum(lambda^2)/m))
  return((m*p*(1 - log(m))/2 - 0.5*sum(lambda_A) + m*sum(log(lambda_A))/2 - mu_m) / (m*sigma_m))
}

test_collect <- c()
if (power == FALSE) {
  for (i in 1:cycles) {
    data <- mvrnorm(n = n, mu = mu, Sigma = sigma0)
    test <- get_jwen_test2(data, n, sigma0)
    test_collect <- c(test_collect, test)
  }
  alpha <- mean(test_collect < qnorm(0.025, 0, 1) | test_collect > qnorm(0.975, 0, 1))
  return(list(test = test_collect, alpha = alpha))
}

if (power == TRUE) {
  for (i in 1:cycles) {
    data <- mvrnorm(n = n, mu = mu, Sigma = sigma1)
    test <- get_jwen_test2(data, n, sigma0)
    test_collect <- c(test_collect, test)
  }
  alpha <- mean(test_collect < qnorm(0.025, 0, 1) | test_collect > qnorm(0.975, 0, 1))
  return(list(test = test_collect, power = alpha))
}
}

# 示例
# jwen_test2(power = FALSE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), cycles = 100)
# jwen_test2(power = TRUE, n = 100, mu = rep(0, 60), sigma0 = diag(1, 60), sigma1 = diag(100, 60), cycles
= 100)

diag_p_q <- function(lambda, p, q) {
  return(diag(lambda, p) + diag(q, (p+1), p)[2:(p+1), 1:p] + diag(q, (p+2), (p+1))[3:(p+2), 1:p] + diag(q,
(p+3), (p+2))[4:(p+3), 1:p] + diag(q, p, (p+1))[1:p, 2:(p+1)] + diag(q, (p+1), (p+2))[1:p, 3:(p+2)] + diag(q,
(p+2), (p+3))[1:p, 4:(p+3)])
}

```

图 3-1

```

chi_size <- chi_power <- jwen_size <- jwen_power <- list()
parameters <- c(5, 30, 60, 90)
for (i in 1:length(parameters)) {
  p <- parameters[i]
  chi_size[[i]] <- chi_test(power = FALSE, n = 100, mu = rep(0, p), sigma0 = diag(1, p), cycles = 100000)
  jwen_size[[i]] <- jwen_test(power = FALSE, n = 100, mu = rep(0, p), sigma0 = diag(1, p), cycles = 100000)
  chi_power[[i]] <- chi_test(power = TRUE, n = 100, mu = rep(0, p), sigma0 = diag(1, p), sigma1 =
diag_p_q(1, p, 0.1), cycles = 100000)
  jwen_power[[i]] <- jwen_test(power = TRUE, n = 100, mu = rep(0, p), sigma0 = diag(1, p), sigma1 =
diag_p_q(1, p, 0.1), cycles = 100000)
}
windows()
par(mfcol = c(2, 4))
hist(chi_size[[1]]$test, freq = FALSE, col = "orange", xlim=c(0,40),
     main = paste0("n = 100, p = 5"), ylim = c(0, 0.1), breaks = 100,
     xlab = expression(paste("-2", rho %.% log(Lambda^paste("*")))))
curve(dchisq(x, df = 0.5*5*(5+1)), add = TRUE, col = "blue", lwd = 2)
hist(jwen_size[[1]]$test, freq = FALSE, col = "orange", xlim=c(-10,5),
     main = paste0("n = 100, p = 5"), ylim = c(0, 0.5), breaks = 100,
     xlab = expression(frac(log(Lambda^paste("*"))[m]) - mu[m], m %.% sigma[m])))
curve(dnorm(x, mean = 0, sd = 1), add = TRUE, col = "blue", lwd = 2)
hist(chi_size[[2]]$test, freq = FALSE, col = "orange", xlim=c(350, 600),
     main = paste0("n = 100, p = 30"), ylim = c(0, 0.02), breaks = 100,
     xlab = expression(paste("-2", rho %.% log(Lambda^paste("*")))))
curve(dchisq(x, df = 0.5*30*(30+1)), add = TRUE, col = "blue", lwd = 2)
hist(jwen_size[[2]]$test, freq = FALSE, col = "orange", xlim=c(-10,5),
     main = paste0("n = 100, p = 30"), ylim = c(0, 0.5), breaks = 100,
     xlab = expression(frac(log(Lambda^paste("*"))[m]) - mu[m], m %.% sigma[m])))
curve(dnorm(x, mean = 0, sd = 1), add = TRUE, col = "blue", lwd = 2)
hist(chi_size[[3]]$test, freq = FALSE, col = "orange", xlim=c(1600, 2200),
     main = paste0("n = 100, p = 60"), ylim = c(0, 0.01), breaks = 100,
     xlab = expression(paste("-2", rho %.% log(Lambda^paste("*")))))
curve(dchisq(x, df = 0.5*60*(60+1)), add = TRUE, col = "blue", lwd = 2)
hist(jwen_size[[3]]$test, freq = FALSE, col = "orange", xlim=c(-10,5),
     main = paste0("n = 100, p = 60"), ylim = c(0, 0.5), breaks = 100,
     xlab = expression(frac(log(Lambda^paste("*"))[m]) - mu[m], m %.% sigma[m])))
curve(dnorm(x, mean = 0, sd = 1), add = TRUE, col = "blue", lwd = 2)
hist(chi_size[[4]]$test, freq = FALSE, col = "orange", xlim=c(3800, 5300),
     main = paste0("n = 100, p = 90"), ylim = c(0, 0.005), breaks = 100,
     xlab = expression(paste("-2", rho %.% log(Lambda^paste("*")))))
curve(dchisq(x, df = 0.5*90*(90+1)), add = TRUE, col = "blue", lwd = 2)
hist(jwen_size[[4]]$test, freq = FALSE, col = "orange", xlim=c(-10,5),

```

```

main = paste0("n = 100, p = 90"), ylim = c(0, 0.5), breaks = 100,
xlab = expression(frac(log(Lambda^paste(" ")[m]) - mu[m], m %>% sigma[m])))
curve(dnorm(x, mean = 0, sd = 1), add = TRUE, col = "blue", lwd = 2)

```

表 3-1

```

(table1 <- cbind(chi_size = sapply(chi_size, function(x) x$alpha),
                jwen_size = sapply(jwen_size, function(x) x$alpha),
                chi_power = sapply(chi_power, function(x) x$power),
                jwen_power = sapply(jwen_power, function(x) x$power)))
(table1_sd <- sqrt(table1*(1-table1)/100000))

```

表 3-2

```

jwen_size <- jwen_power <- chen_size <- chen_power <- tran_size <- tran_power <- list()
for (i in 1:60) {
  n <- 100; p <- 60; mu <- rep(0, p); cycles <- 10000;
  sigma0 <- diag(c(rep(1, p-i), rep(1000, i)))
  sigma1 <- diag(c(rep(1, p-i), rep(1200, i)))
  jwen_size[[i]] <- jwen_test(power = FALSE, n = n, mu = mu, sigma0 = sigma0, cycles = cycles)
  chen_size[[i]] <- chen_test(power = FALSE, n = n, mu = mu, sigma0 = sigma0, cycles = cycles)
  tran_size[[i]] <- translate_test(power = FALSE, n = n, mu = mu, sigma0 = sigma0, cycles = cycles)
  jwen_power[[i]] <- jwen_test(power = TRUE, n = n, mu = mu, sigma0 = sigma0, sigma1 = sigma1, cycles
= cycles)
  chen_power[[i]] <- chen_test(power = TRUE, n = n, mu = mu, sigma0 = sigma0, sigma1 = sigma1, cycles
= cycles)
  tran_power[[i]] <- translate_test(power = TRUE, n = n, mu = mu, sigma0 = sigma0, sigma1 = sigma1,
cycles = cycles)
  print(i)
}

```

```

(table2 <- cbind(jwen_size = sapply(jwen_size, function(x) x$alpha),
                chen_size = sapply(chen_size, function(x) x$alpha),
                tran_size = sapply(tran_size, function(x) x$alpha),
                jwen_power = sapply(jwen_power, function(x) x$power),
                chen_power = sapply(chen_power, function(x) x$power),
                tran_power = sapply(tran_power, function(x) x$power)))
(table2_sd <- sqrt(table2*(1-table2)/100))

```

表 3-3

```

jwen_size <- jwen_power <- list()
n <- 100*c(1, 1.5, 2, 2.5, 3, 3.5, 4)
p <- 60*c(1, 1.5, 2, 2.5, 3, 3.5, 4)
for (i in 1:length(n)) {
  jwen_size[[i]] <- jwen_test(power = FALSE, n = n[i], mu = rep(0, p[i]), sigma0 = diag(2, p[i]), cycles =

```

```
10000)
  jwen_power[[i]] <- jwen_test(power = TRUE, n = n[i], mu = rep(0, p[i]), sigma0 = diag(2, p[i]), sigma1
= diag(2.1, p[i]), cycles = 10000)
  print(i)
}
(table3 <- cbind(jwen_size = sapply(jwen_size, function(x) x$alpha),
                jwen_power = sapply(jwen_power, function(x) x$power)))
(table3_sd <- sqrt(table3*(1-table3)/10000))
```


参考文献

- [1] Chang J, Guo B, Yao Q. High dimensional stochastic regression with latent factors, endogeneity and nonlinearity[J]. *Journal of Econometrics*, 2015, 189(2): 297-312.
- [2] Nagao H. On some test criteria for covariance matrix[J]. *The Annals of Statistics*, 1973: 700-709.
- [3] Anderson-Cook C M. An introduction to multivariate statistical analysis[J]. *Journal of the American Statistical Association*, 2004, 99(467): 907-909.
- [4] Ledoit O, Wolf M. Some hypothesis tests for the covariance matrix when the dimension is large compared to the sample size[J]. *The Annals of Statistics*, 2002, 30(4): 1081-1102.
- [5] Bai Z, Jiang D, Yao J F, et al. Corrections to LRT on large-dimensional covariance matrix by RMT[J]. *The Annals of Statistics*, 2009, 37(6B): 3822-3840.
- [6] Jiang D, Jiang T, Yang F. Likelihood ratio tests for covariance matrices of high-dimensional normal distributions[J]. *Journal of Statistical Planning and Inference*, 2012, 142(8): 2241-2256.
- [7] Wang C, Yang J, Miao B, et al. Identity tests for high dimensional data using RMT[J]. *Journal of Multivariate Analysis*, 2013, 118: 128-137.
- [8] Chen S X, Zhang L X, Zhong P S. Tests for high-dimensional covariance matrices[J]. *Journal of the American Statistical Association*, 2010, 105(490): 810-819.
- [9] Ahmad M R, Rosen D. Tests for high-dimensional covariance matrices using the theory of U-statistics[J]. *Journal of Statistical Computation and Simulation*, 2015, 85(13): 2619-2631.
- [10] Wu T L, Li P. Tests for high-dimensional covariance matrices using random matrix projection[J]. *arXiv preprint arXiv:1511.01611*, 2015.
- [11] Jiang T, Yang F. Central limit theorems for classical likelihood ratio tests for high-dimensional normal distributions[J]. *The Annals of Statistics*, 2013, 41(4): 2029-2074.
- [12] Jiang T, Qi Y. Likelihood Ratio Tests for High-Dimensional Normal Distributions[J]. *Scandinavian Journal of Statistics*, 2015, 42(4): 988-1009.
- [13] Chen H, Jiang T. A study of two high-dimensional likelihood ratio tests under alternative hypotheses[J]. *Random Matrices: Theory and Applications*, 2018, 7(01): 1750016.
- [14] Li J, Chen S X. Two sample tests for high-dimensional covariance matrices[J]. *The Annals of Statistics*, 2012, 40(2): 908-940.
- [15] Cai T, Liu W, Xia Y. Two-sample covariance matrix testing and support recovery in high-dimensional and sparse settings[J]. *Journal of the American Statistical Association*, 2013, 108(501): 265-277.
- [16] Hong Y, Kim C. Recent developments in high dimensional covariance

- estimation and its related issues, a review[J]. Journal of the Korean Statistical Society, 2018.
- [17] Cai T T. Global testing and large-scale multiple testing for high-dimensional covariance structures[J]. Annual Review of Statistics and Its Application, 2017, 4: 423-446.
- [18] Cai T T, Wu Y. Statistical and computational limits for sparse matrix detection[J]. arXiv preprint arXiv:1801.00518, 2018.
- [19] Muirhead R J. Aspects of multivariate statistical theory Wiley[J]. New York, 1982.