

**Statistica Sinica Preprint No: SS-2023-0117**

|                                 |                                                                                                                    |
|---------------------------------|--------------------------------------------------------------------------------------------------------------------|
| <b>Title</b>                    | Combining P-values Using Heavy Tailed Distributions and Their Asymptotic Results with Applications to Genomic Data |
| <b>Manuscript ID</b>            | SS-2023-0117                                                                                                       |
| <b>URL</b>                      | <a href="http://www.stat.sinica.edu.tw/statistica/">http://www.stat.sinica.edu.tw/statistica/</a>                  |
| <b>DOI</b>                      | 10.5705/ss.202023.0117                                                                                             |
| <b>Complete List of Authors</b> | Junsik Kim and Junyong Park                                                                                        |
| <b>Corresponding Authors</b>    | Junyong Park                                                                                                       |
| <b>E-mails</b>                  | Junyongpark@snu.ac.kr                                                                                              |

## Combining $p$ -values using heavy tailed distributions and their asymptotic results with applications to genomic data

Junsik Kim\* and Junyong Park†

\**Department of Statistics, Duksung Women's University, Seoul, Korea*

†*Department of Statistics, Seoul National University, Seoul, Korea*

*Abstract:* Combining individual  $p$ -values to handle large scale inferences or to aggregate results of different studies is one of major interest in meta-analysis which has been traditionally based on independent  $p$ -values. In contrast to combining methods that are constructed when  $p$ -values are independent, recently proposed combinations of  $p$ -values transformed into heavy-tailed distribution are known to be robust to the dependence structure of  $p$ -values. In this paper, we investigate theoretical properties of combining  $p$ -value methods for different heaviness of transformation under a wider class of correlation structures compared to existing studies from the view point of controlling Type I error and obtaining powers. We also investigate relationships between harmonic mean type combination methods and combining methods that use transformation of  $p$ -values into stable distribution including Cauchy and Lévy combination methods. We provide extensive numerical studies supporting theoretical results. We also apply these  $p$ -value combining methods to real example of Crohn's disease data and present some idea on how to validate these methods.

*Key words and phrases:* Combining  $p$ -values; Meta-Analysis, GWAS; Cauchy transformation; Lévy transformation

## 1. Introduction

Meta-analysis has been used in various fields as a statistical technique to draw more reliable conclusions by combining the results of different studies or experiments. Most of research has focused on how to aggregate  $p$ -values generated from the assumption that all experiments or studies are independent. Among the methods for combining such independent  $p$ -values, well known methods including Fisher's method in Fisher (1934) and Stouffer's method in Stouffer et al. (1949) are based on aggregating some transformed  $p$ -values leading to chi-square distribution and normal distribution, respectively. In particular, the distribution of the sum of independent  $p$ -values or other forms of statistics needs to be derived relatively easily for practical use. For example, if the  $p$ -value is transformed into an infinitely divisible random variable, the distribution of the sum of those transformed  $p$ -values is the same type of distribution as the transformed distribution as in the case of Fisher's method and Stouffer's method.

On the other hand, there are limited studies on combining dependent  $p$ -values. When  $p$ -values are dependent, Kost and McDermott (2002) and Hartung (1999) aggregated the dependent  $p$ -values by estimating the correlation coeffi-

---

cients. Recently, Liu and Xie (2020) showed that combining dependent  $p$ -values transformed into the Cauchy distribution is robust to the overall dependency when a significance level is fairly small. Numerical studies in Wilson (2021) demonstrated the robustness of Lévy combination test for dependent  $p$ -values. Fang et al. (2023) showed that combinations of dependent  $p$ -values have the same tail behavior as that from independent  $p$ -values, however their result is restricted to the case when the number of  $p$ -values is fixed. Liang and Rho (2022) used stable combination tests under long range or short range dependence which is somewhat restricted dependence in meta analysis. Wilson (2019) proposed to use a harmonic mean of  $p$ -values to combine dependent  $p$ -values.

In this paper, we consider a broader class of dependence structure induced from Gaussian copula used in Liu and Xie (2020) and investigate the asymptotic performance of existing test statistics from the view point of controlling a given Type I error and obtaining power. Liu and Xie (2020) assumed that eigenvalues of correlation matrix are bounded which cannot include highly correlated  $p$ -values. The dependence structure in Liang and Rho (2022) is also somewhat restricted in practice. We focus on figuring out the scopes of dependence structure in which existing methods control a given Type I error and obtain powers.

Among these two criteria in hypothesis testing, the first requirement is controlling Type I error while obtaining power is pursued thereafter. In order to

obtain the robustness to dependence of  $p$ -values, it is advantageous to have test statistics dominated by a small number of small  $p$ -values. However this property is undesirable to obtain testing power under non-null hypothesis especially when there are many  $p$ -values generated from non-null hypothesis. Test statistics depending on a relatively small number of  $p$ -values may lose testing power since transformation of fairly small values of  $p$ -values tend to be extremely large, so they dominate transformation of some other  $p$ -values which are also generated from non-null hypotheses.

With this motivation, for dependent  $p$ -values, we investigate the properties of methods of combining  $p$ -values in controlling Type I error and obtaining power for different tail behaviors from different transformations. In particular, it is highlighted that all test statistics considered in this paper are shown to have a trade-off relationship between size and power. In other words, if the size is controlled stably at strongly dependent  $p$ -values, the power will be low, whereas it becomes difficult to control the size of a test statistic with an advantage in power at  $p$ -values with strong dependence. We demonstrate the trade-off theoretically and numerically via presenting asymptotic results, simulations and real data examples.

Test statistics considered in this paper are classified by the heaviness of the tails of the transformed  $p$ -values. Depending on the heaviness of the tail, we

---

consider heavier tailed distributions such as the Cauchy distribution and the Lévy distribution which have been studied in Liu and Xie (2020), Chen et al. (2023) and Wilson (2021), in addition to well known methods such as Stouffer's and Fisher's methods. We provide various results on how different heaviness of tails affects the performance of test statistics for dependent  $p$ -values.

We also investigate asymptotic relations between  $p$ -values combining methods for heavy-tailed distribution and for a type of harmonic mean of  $p$ -values. Combining methods using the Cauchy and Lévy distribution and harmonic mean are designed to rely on a small subset of transformed  $p$ -values which dominates all other transformed  $p$ -values due to the heavy tail of Cauchy or Lévy distributions and inverse of  $p$ -values in harmonic mean. This property results in diminishing the effect of dependence of  $p$ -values since a few  $p$ -values dominating all other  $p$ -values avoid the accumulation of dependence. Chen et al. (2023) showed that Cauchy combination method and harmonic mean method are asymptotically equivalent when the number of  $p$ -values is fixed and the smallest of  $p$ -value goes to zero. On the other hand, Fisher's method and Stouffer's method are not explained by a few number of  $p$ -values, since these two methods use chi-square and normal distributions which do not have heavy tails. As an extreme case, the Tippett's method (Tippett, 1931) which uses only the smallest  $p$ -value is also considered as the heaviest transformation.

---

This paper is organized as follows; Section 2 includes the introduction of test statistics for different transformations. Section 3 presents theoretical results from the view point of asymptotic size and power of test statistics. In Section 4, relations between  $p$ -values combining methods for heavy-tailed distribution are investigated. Numerical studies and real data examples are provided in Section 5 and Section 6, respectively. Section 7 presents concluding remarks.

## 2. Test Statistics for combining $p$ -values

We consider the following global hypothesis testing ;

$$H_0 : H_{0i} \text{ are all true for } 1 \leq i \leq d \quad \text{vs.} \quad H_1 : \text{at least one } H_{0i} \text{ is false} \quad (2.1)$$

where  $H_{0i} : \mu_i = 0$  for two-sided test and  $H_{0i} : \mu_i \leq (\geq)0$  for one-sided test including right-tailed and left-tailed tests. For each hypothesis  $H_{0i}$ , we obtain  $p$ -values,  $p_1, \dots, p_d$  by marginally testing each hypothesis. In meta analysis, it may not be known about the marginal data as well as the joint distribution of those test statistics, but Liu and Xie (2020) assumed that the original data or the test statistic  $\mathbf{X} = (X_1, \dots, X_d)^T$  follows the multivariate normal distribution,  $\mathbf{X} \sim N_d(\mu, \Sigma)$ , where  $\mu = (\mu_1, \dots, \mu_d)^T$  and  $\Sigma$  is unknown.  $p$ -values are represented by  $p_i = 1 - \Phi(X_i)$  or  $\Phi(X_i)$  depending on right-tailed or left-tailed test and  $p_i = 2\{1 - \Phi(|X_i|)\}$  for two-sided test.  $\Sigma$  characterizes dependence

---

structure of  $p$ -values, however, since it is unknown, it is not feasible to identify such dependence.

Various methods of combining  $p$ -values in meta-analysis are constructed through an appropriate transformation of  $p$ -values such as  $T = \sum_{i=1}^d h(p_i)$  where  $h(\cdot)$  is a decreasing function. Typical examples are Stouffer's method,  $T_{\text{Stouffer}} = \frac{1}{\sqrt{d}} \sum_{i=1}^d \Phi^{-1}(p_i)$ , and Fisher's method,  $T_{\text{Fisher}} = -2 \sum_{i=1}^d \log(p_i)$ . When  $p$ -values are independent,  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  have the standard normal distribution and chi-square distribution with  $2d$  degrees of freedom under  $H_0$ , respectively. We reject  $H_0$  if  $T_{\text{Stouffer}} > z_\alpha$  and if  $T_{\text{Fisher}} > \chi_\alpha^2(2d)$ , where  $z_\alpha$  and  $\chi_\alpha^2(2d)$  are the upper  $\alpha$  quantiles of a standard normal distribution and chi-square distribution with degrees of freedom  $2d$ , respectively. However, when  $p$ -values are dependent, distributions of  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  are not guaranteed to be the standard normal and chi-square distributions, respectively, so such decision rules are not valid any more. If there exists relatively strong dependency among  $p$ -values,  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  with the critical values  $z_\alpha$  and  $\chi_\alpha^2(2d)$  are not able to control Type I error.

Recently, methods for combining dependent  $p$ -values using heavy-tailed distributions have been studied intensively since combination of dependent heavy tailed distributions may have similar behavior to the combination of those independent distributions. For example, the Cauchy transformation,  $T_{\text{Cauchy}} =$

---

$\frac{1}{d} \sum_{i=1}^d \tan(\pi(1/2 - p_i))$ , in Liu and Xie (2020), follows Cauchy distribution under  $H_0$  and independent  $p$ -values. Liu and Xie (2020) showed that, under dependent  $p$ -values with some conditions, the tail behavior of  $T_{Cauchy}$  is asymptotically equivalent to the standard Cauchy distribution. Then  $H_0$  is rejected when  $T_{Cauchy} > c_\alpha$  where  $c_\alpha$  is the upper  $\alpha$  quantile of Cauchy distribution. With a similar motivation, we can also consider a heavier tailed distribution than Cauchy distribution such as Lévy distribution (Wilson, 2021). We define  $T_{Lévy} = \frac{1}{d^2} \sum_{i=1}^d \left\{ \Phi^{-1} \left( \frac{1}{2}(1 + p_i) \right) \right\}^{-2}$  and reject  $H_0$  if  $T_{Lévy} > l_\alpha$  where  $l_\alpha$  is the upper  $\alpha$  quantile of Lévy distribution.

Note that tail behaviors of  $T_{Stouffer}$ ,  $T_{Fisher}$ ,  $T_{Cauchy}$  and  $T_{Lévy}$  are different in that the thickness of the tail part has an order of  $T_{Stouffer} < T_{Fisher} < T_{Cauchy} < T_{Lévy}$ . Lastly, we also consider the Tippett's method which uses the minimum of  $p$ -values,  $T_{\min P} = \max_i \Phi^{-1}(1 - p_i)$ , or  $T'_{\min P} = \min_i p_i$ , equivalently. As a critical value which is based on independence, we reject  $H_0$  if  $T_{\min P} \geq \Phi^{-1}((1 - \alpha)^{1/d})$  or equivalently  $T'_{\min P} = \min_i p_i \leq 1 - (1 - \alpha)^{1/d}$ , for a given significant level  $\alpha$ .

Throughout this paper, we evaluate different properties of  $T_{Stouffer}$ ,  $T_{Fisher}$ ,  $T_{Cauchy}$ ,  $T_{Lévy}$  and  $T_{\min P}$  by their different distributions with different tail behaviors. Each of test statistics uses critical values derived under independent  $p$ -values such as  $z_\alpha$ ,  $\chi_\alpha^2(2d)$ ,  $c_\alpha$ ,  $l_\alpha$  and  $\Phi^{-1}((1 - \alpha)^{1/d})$ , respectively. These critical

---

values may depart from the true critical values of  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$ ,  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  under dependent  $p$ -values. In practice, since we do not know the dependent structure induced from unknown  $\Sigma$ , we have no choice but to use them derived under independent  $p$ -values. We highlight that, as the tail of the transformed distribution  $h(\cdot)$  is heavier, a smaller number of  $p$ -values dominate the variation of the test statistics which avoids the accumulation of dependence of  $p$ -values. We also address the issue of testing power of these test statistics which have reverse order of performances of robustness in controlling Type I error.

In the following sections, we provide theoretical studies on the properties of  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$ ,  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  and present numerical studies including simulations and real data examples. All proofs of theoretical studies are presented in the supplementary materials.

### 3. Type I error and testing power

In this section, we present theoretical results on controlling a given Type I error and obtaining power of  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$ ,  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$ . Throughout this section, we see that the heavier the distribution of the transformed  $p$ -values, the more advantageous it is to control Type I error. However, we also see that such heavy tailed distributions can be detrimental to the achievement of testing power which is  $1 - \text{Type II error}$ . As an overall performance, we investigate the

---

conditions for each of test statistics under which we have the following properties: Let  $\alpha$  be a given level of Type I error. Then, as  $d \rightarrow \infty$  and  $\alpha \rightarrow 0$ ,

$$\text{Type I error} \leq \alpha \rightarrow 0 \text{ and Type II error} \rightarrow 0. \quad (3.1)$$

We compare conditions satisfying (3.1) for different types of testing procedures and observe there exists trade-off relationship to obtain these two criteria. To implement dependent  $p$ -values, we use the Gaussian copula as in Liu and Xie (2020) rather than considering arbitrary dependence in Vovk and Wang (2020) and Chen et al. (2023). As seen in Vovk and Wang (2020), a modification of some test statistics for arbitrary dependent  $p$ -values become impractical due to their fairly low powers.

We use the Gaussian copula to model dependent  $p$ -values which are not that restrictive and we derive conditions for realistic powers of some existing test statistics. More specifically, the raw data or statistics  $X_i$ 's for  $1 \leq i \leq d$  are normally distributed with a correlation matrix  $\Sigma = (\rho_{ij})_{1 \leq i, j \leq d}$  which is an element in a collection of correlation matrices,  $\mathcal{F}_{d, \rho} = \{\Sigma : 0 \leq \rho_{ij} \leq \rho < 1 \text{ for } 1 \leq i \neq j \leq d\}$ .

Although our dependence is implemented via the Gaussian copula,  $\mathcal{F}_{d, \rho}$  covers various dependent structures including a correlation matrix in Liu and Xie (2020). In particular, as  $d$  increases, the correlation matrices in Liu and Xie (2020) are limited in that they represent weak dependence among all  $p$ -values

---

due to the condition  $\lambda_{\max}(\Sigma) \leq C$  for some constant  $C$  where  $\lambda_{\max}(\Sigma)$  is the maximum eigenvalue of  $\Sigma$ . For example, they exclude a case of equally correlated data. Specifically, since it is known that  $1 + \frac{1}{d} \sum_{i \neq j} \rho_{ij} \leq \lambda_{\max}(\Sigma)$ , under a high dimensional setting or a case where the dimension increases, an assumption of bounded eigenvalues requires each correlation coefficient becomes very small or diminishes to zero order of reciprocal of dimension. In practice, especially for applications to genomic data, there are many cases that the assumption of bounded eigenvalues is not realistic due to dependence among SNPs in genes in high dimensionality. On the other hand,  $\mathcal{F}_{d,\rho}$  covers highly correlated data including equally correlated case since  $\mathcal{F}_{d,\rho}$  allows  $\lambda_{\max}(\Sigma)$  to diverge as  $d$  increases.  $\mathcal{F}_{d,\rho}$  includes correlation matrices since it allows a case where  $\lambda_{\max}(\Sigma)$  diverges as  $d$  increases. Similarly, Liang and Rho (2022) used the weakly dependent structures such as serially correlated cases also included in  $\mathcal{F}_{d,\rho}$ . And in the real data, since true correlation coefficients are unknown, the condition in this paper can include the worst case scenario. From this point, we deal with a broader class of correlation matrices,  $\mathcal{F}_{d,\rho}$ .

We introduce some notations used in this paper. Denote  $f(x) \sim g(x)$  as  $x \rightarrow c$  if  $\lim_{x \rightarrow c} f(x)/g(x) = 1$ . The statement  $a \lesssim b$  means that  $a \leq \gamma \cdot b$ , where  $\gamma > 0$  is a fixed constant independent of dimension  $d$ . Let  $\phi(x)$ ,  $\Phi(x)$  and  $\bar{\Phi}(x)$  be the density function, the cumulative distribution function and the

### 3.1 Type I error control

survival function of standard normal distribution, respectively. For  $f(x) > 0$  and  $g(x) > 0$ , we define  $f(x) \ll g(x)$  if  $\limsup_{x \rightarrow \infty} f(x)/g(x) = C$  for a constant  $C \in (1, \infty)$ . The notation  $[x]$  means the integer part of a positive value  $x$ .

#### 3.1 Type I error control

We present theorems showing that each of test statistics,  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$ ,  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  controls Type I error under different scenarios such as  $\rho_{ij}$  or  $\alpha$  is fixed while other parameters depends on  $d$ . We first present a theorem implying that for an exchangeable case, both  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  fail to control a given Type I error when the correlation coefficient is fixed. This shows the weakness of  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  which do not use the heavy tailed distribution. On the other hand, Type I errors of  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  are approximately a given level  $\alpha$  if the correlation coefficient diminishes to zero under some conditions. In the following theorem, we define  $P_{H_0}$  as the probability measure under  $H_0$ .

**Theorem 1.** 1. Suppose that  $0 < \rho = \text{Corr}(p_i, p_j)$  for  $1 \leq i \neq j \leq d$  is fixed.

Then for any  $\alpha \in (0, 1)$ , under  $H_0$  in (2.1), we have  $\lim_{d \rightarrow \infty} P_{H_0}(T_{\text{Stouffer}} \geq z_\alpha) = 1/2$ . However, if  $\rho = o(1/d)$ , then  $\lim_{d \rightarrow \infty} P_{H_0}(T_{\text{Stouffer}} \geq z_\alpha) = \alpha$ .

2. Suppose that  $0 < \xi = \text{Corr}(-2 \log p_i, -2 \log p_j)$  for  $1 \leq i \neq j \leq d$  is

fixed. Then for any  $\alpha \in (0, 1)$ , under  $H_0$  in (2.1), we have  $\lim_{d \rightarrow \infty} P_{H_0}(T_{\text{Fisher}} \geq \chi_\alpha^2(2d)) =$

$1/e$ . However, if  $\xi = o(1/d)$  and  $d\sqrt{\xi} \rightarrow c \in [0, +\infty]$ , then  $\lim_{d \rightarrow \infty} P_{H_0}(T_{\text{Fisher}} \geq \chi_\alpha^2(2d)) =$

$\alpha$ .

We note that  $T_{\text{Fisher}}$  can be expressed as a special combining method of generalized averaging of  $p$ -values. Let  $p_1, \dots, p_d$  be any  $p$ -values following the uniform distribution marginally and  $M_r = d^{-1/r}(p_1^r + \dots + p_d^r)^{1/r}$  be a combining function. From Vovk and Wang (2020), the combining function is precise if for each  $\epsilon \in (0, 1)$ ,  $\sup P(M_r \leq \epsilon) = \epsilon$ , where the supremum encompasses arbitrary dependence structure between  $p_1, \dots, p_d$ . Vovk and Wang (2020) showed that, if a constant  $a_d$  depending on  $d$  and  $a_d \rightarrow e$  as  $d \rightarrow \infty$ , the supremum of adjusted statistic  $a_d M_0$  for  $M_0 := \left(\prod_{i=1, \dots, d} p_i\right)^{1/d}$  achieves precise probability, that is,  $\sup P(a_d M_0 \leq \epsilon) = \epsilon$ . Therefore,  $T_{\text{Fisher}}$  also can be considered to be precise. However, in the perspective of hypothesis testing procedure for Fisher's method, the adjusted method of Vovk and Wang (2020) becomes too conservative. Indeed, from the precise probability, for a given significance level  $\alpha \in (0, 1)$ ,  $\alpha = \sup P(a_d M_0 \leq \alpha) = \sup P(T_{\text{Fisher}} \geq 2d \log(a_d/\alpha))$ . However, for large  $d$ , since  $\log(a_d/\alpha) > 1$ , the rejection region,  $[2d \log(a_d/\alpha), +\infty)$  is narrower than that of Fisher's method. Hence, the adjusted statistic becomes conservative which is a cost for validity under arbitrary dependence structure. Numerical studies in Chen et al. (2023) showed that under arbitrary dependence structure, powers of methods based on the generalized averaging are lower than power of Bonferroni method so that the adjusted statistics are hard to be used

### 3.1 Type I error control

in practice. Therefore, in this paper, we adapt the quantile constructed under independence assumption that is called VI (Valid for Independent) method for a practical balance between power and size. In the following theorem, we define  $P^\Sigma(\cdot)$  as the probability under a given correlation matrix  $\Sigma \in \mathcal{F}_{d,\rho}$ .

**Theorem 2.** 1. For any  $\Sigma \in \mathcal{F}_{d,\rho}$  and any  $\alpha \in (0, 1)$ , under  $H_0$  in (2.1), we have

$$\inf_{\Sigma \in \mathcal{F}_{d,\rho}} P_{H_0}^\Sigma(T_{\text{Stouffer}} \geq z_\alpha) = P_{H_0}^I(T_{\text{Stouffer}} \geq z_\alpha) = \alpha. \quad (3.2)$$

Furthermore, for any given  $\alpha \in (0, 1)$  and if  $\sum_{i \neq j} \rho_{ij} = o(d)$ , under  $H_0$  in (2.1), we have

$$\lim_{d \rightarrow \infty} P_{H_0}^\Sigma(T_{\text{Stouffer}} \geq z_\alpha) = \alpha. \quad (3.3)$$

2. For any  $\Sigma \in \mathcal{F}_{d,\rho}$  and any  $\alpha \in (0, 1)$ , under  $H_0$  in (2.1), we have

$$\inf_{\Sigma \in \mathcal{F}_{d,\rho}} P_{H_0}^\Sigma(T_{\text{Fisher}} \geq \chi_\alpha^2(2d)) = P_{H_0}^I(T_{\text{Fisher}} \geq \chi_\alpha^2(2d)) = \alpha. \quad (3.4)$$

Furthermore, let  $\xi_{ij} = \text{Corr}(-2 \log p_i, -2 \log p_j)$  for  $i \neq j = 1, \dots, d$ . if  $\sum_{i \neq j} \xi_{ij} = o(1/d)$ , then

$$\lim_{d \rightarrow \infty} P_{H_0}^\Sigma(T_{\text{Fisher}} \geq \chi_\alpha^2(2d)) = \alpha. \quad (3.5)$$

3. For any  $\Sigma \in \mathcal{F}_{d,\rho}$  and any  $\alpha \in (0, 1)$ , under  $H_0$  in (2.1), we have

$$\sup_{\Sigma \in \mathcal{F}_{d,\rho}} P_{H_0}^\Sigma(T_{\text{minP}} \geq x_\alpha) = P_{H_0}^I(T_{\text{minP}} \geq x_\alpha) = \alpha \quad (3.6)$$

### 3.1 Type I error control

where  $x_\alpha = \Phi^{-1}((1 - \alpha)^{1/d})$ . Furthermore, for any  $\Sigma$  with  $\rho_{ij} \geq \epsilon > 0$ , under  $H_0$  in (2.1), we have

$$\lim_{d \rightarrow \infty} P_{H_0}^\Sigma(T_{\min P} \geq x_\alpha) = 0. \quad (3.7)$$

Note that (3.2), (3.4) and (3.6) are non-asymptotic results while (3.3), (3.5) and (3.7) are asymptotic results for varying correlation coefficients depending on  $d$ . Theorem 1 and 2 show the weakness of  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$  and  $T_{\min P}$  in the sense that  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  may fail in controlling a given Type I error and  $T_{\min P}$  is too conservative when there exist some dependence among  $p$ -values. Liu and Xie (2020) and Wilson (2019) presented numerical studies showing that Stouffer's test and Fisher's test fail in controlling Type I error when  $p$ -values are correlated. However, the results of Theorem 1 and 2 have the novelty in that they provide theoretical reasons of the failures of Stouffer's and Fisher's tests in controlling Type I error. These results are related to how many components in  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$  and  $T_{\min P}$  contribute to explain the whole variation of those three tests, respectively. In other words,  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  are aggregating all transformed  $p$ -values while  $T_{\min P}$  relies on only one term, maximum of  $X_{i.s}$ . We leave the following remarks as implications of Theorem 2.

**Remark 1.** 1. For any given  $\Sigma$ ,  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  control Type I error non-asymptotically only when all  $p$ -values are independent, i.e.,  $\Sigma = I$ . Otherwise,

### 3.1 Type I error control

$T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  cannot control Type I error for dependent  $p$ -values. However, when the number of  $p$ -values increases,  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  may control Type I error asymptotically when the effect of overall correlations is not large enough such as  $\sum_{i \neq j} \rho_{ij} = o(d)$  and  $\sum_{i \neq j} \xi_{ij} = o(1/d)$ , respectively.

2. For exchangeable  $p$ -values under  $H_0$ , if, for all  $1 \leq i \neq j \leq d$ ,  $\text{Corr}(p_i, p_j) = \rho$  is positive and fixed,  $P_{H_0}(T_{\text{Stouffer}} > z_\alpha) \rightarrow 1/2$  under the null, as  $d \rightarrow \infty$ , whereas the Type I error of  $T_{\text{Fisher}}$  converges to  $e^{-1}$ , if  $0 < \xi = \text{Corr}(-2 \log p_i, -2 \log p_j)$  is fixed. Since  $1/e \approx 0.37$ , Type I error of  $T_{\text{Stouffer}}$  is larger than that of  $T_{\text{Fisher}}$  so that  $T_{\text{Stouffer}}$  obtains more seriously inflated Type I error than  $T_{\text{Fisher}}$ . These results can be confirmed by Figure 1 in Section 5.

3.  $T_{\text{minP}} = \max_i X_i$  or  $\min_i p_i$  can control Type I error for any  $\Sigma \in \mathcal{F}_{d,\rho}$  and exact  $\alpha$  is obtained only when all  $p$ -values are independent. In particular, when there exist serious correlations such that  $\rho_{ij} \geq \epsilon > 0$  for some  $\epsilon$ ,  $T_{\text{minP}}$  is fairly conservative in the sense that Type I error converges to 0 as  $d$  increases.

Compared to the weakness of  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$  and  $T_{\text{minP}}$  which are either too liberal or conservative for dependent  $p$ -values, the motivation of  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  is developing test statistics which are intermediate positions between  $T_{\text{Stouffer}}$  (or  $T_{\text{Fisher}}$ ) and  $T_{\text{minP}}$  in the sense that some small subset of transformed  $p$ -values can explain most of variations.

Before we present the asymptotic properties of  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$ , we pro-

### 3.1 Type I error control

vide the following lemma. Roughly speaking, Lemma 1 shows that the lower bound of  $\max_{1 \leq i \leq d} X_i$  is decreasing in  $\rho$  which leads the corresponding  $p$ -value  $\min_{1 \leq i \leq d} p_i$  to be larger since  $\min_{1 \leq i \leq d} p_i$  is derived under the critical value with  $\rho = 0$ . In fact, this result is also related to the third result in Theorem 2 since Type I error of  $\max_{1 \leq i \leq d} X_i$  converges to 0 when  $\rho_{ij} \geq \epsilon > 0$  for some  $\epsilon$ .

**Lemma 1.** *Let  $\mathbf{X} \sim N_d(\mathbf{0}, \Sigma)$  for  $\Sigma \in \mathcal{F}_{d,\rho}$ , then  $\max_{1 \leq i \leq d} X_i \geq \sqrt{1 - \rho} \sqrt{2 \log d} + o_p(1)$ .*

Now we present asymptotic properties of  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  under varying quantities such as the nominal level of Type I error  $\alpha$  and the dimension  $d$ . Our asymptotic results of  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  in the following theorem are different from Liu and Xie (2020) and Liang and Rho (2022) in that they dealt with the case of weakly dependent cases such as covariance matrix with bounded eigenvalues and serially correlated cases while we discuss the cases of  $\Sigma \in \mathcal{F}_{d,\rho}$  which also includes exchangeable matrix as an example of strong dependence. Furthermore, the following theorem shows an answer to the robustness of  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  for  $\Sigma \in \mathcal{F}_{d,\rho}$  as  $d$  diverges while Fang et al. (2023) considered the case of given  $\Sigma$  with fixed  $d$ .

**Theorem 3.** *Suppose  $X \sim N_d(0, \Sigma)$  where  $\Sigma \in \mathcal{F}_{d,\rho}$ . Let  $c_\alpha$  and  $l_\alpha$  be upper  $\alpha$  quantiles of standard Cauchy and Lévy distributions, respectively.*

### 3.2 Comparison of $T_{\text{Cauchy}}$ and $T_{\text{Lévy}}$ in controlling Type I error

1. For  $0 < a_C < \frac{1-\rho^2}{(\sqrt{3}\rho+\sqrt{2})^2}$  and  $d = [(c_\alpha)^{a_C}]$ , we have,

$$\lim_{\alpha \rightarrow 0} \sup_{\Sigma \in \mathcal{F}_{d,\rho}} \frac{P_{H_0}^\Sigma(T_{\text{Cauchy}} \geq c_\alpha)}{P_{H_0}^I(T_{\text{Cauchy}} \geq c_\alpha)} = \lim_{\alpha \rightarrow 0} \sup_{\Sigma \in \mathcal{F}_{d,\rho}} \frac{P_{H_0}^\Sigma(T_{\text{Cauchy}} \geq c_\alpha)}{\alpha} = 1.$$

2. For  $0 < a_L < \frac{1-\rho^2}{5\rho^2+4\sqrt{3}\rho+2}$  and  $d = [(\ell_\alpha)^{a_L}]$ , we have,

$$\lim_{\alpha \rightarrow 0} \sup_{\Sigma \in \mathcal{F}_{d,\rho}} \frac{P_{H_0}^\Sigma(T_{\text{Lévy}} \geq \ell_\alpha)}{P_{H_0}^I(T_{\text{Lévy}} \geq \ell_\alpha)} = \lim_{\alpha \rightarrow 0} \sup_{\Sigma \in \mathcal{F}_{d,\rho}} \frac{P_{H_0}^\Sigma(T_{\text{Lévy}} \geq \ell_\alpha)}{\alpha} = 1.$$

### 3.2 Comparison of $T_{\text{Cauchy}}$ and $T_{\text{Lévy}}$ in controlling Type I error

Theorem 3 shows that  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  are robust to dependence under given dimension  $d_C = [(c_{\alpha_C})^{a_C}]$  and  $d_L = [(\ell_{\alpha_L})^{a_L}]$  for some constraints on  $a_C$  and  $a_L$ . Such constraints on  $d$  can be considered as the scope of the robustness to the accumulation of dependence among  $p$ -values which is due to increasing  $d$ . From the relation between  $d$  and  $\alpha$ , the robustness in controlling Type I error under the dependence consists of conditions on  $d$  and  $\alpha$ . If  $d \rightarrow \infty$ , to handle the accumulation of dependence, the robustness of  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  in controlling Type I error is attained when  $\alpha \rightarrow 0$ . From this motivation, we compare  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  from the following view points : for Type I error and dimension  $(\alpha_C, d_C)$  and  $(\alpha_L, d_L)$  corresponding to  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$ , (i) under  $\alpha_C = \alpha_L \rightarrow 0$ , we compare divergence rates of  $d_C$  and  $d_L$  which are required to control Type I error and (ii) under  $d_C = d_L \rightarrow \infty$ , we compare the rates of  $\alpha_C$  and  $\alpha_L$  converging to 0 which are required to control Type I error. Based on these two criteria, we

### 3.2 Comparison of $T_{\text{Cauchy}}$ and $T_{\text{Lévy}}$ in controlling Type I error

now show that  $T_{\text{Lévy}}$  is more robust than  $T_{\text{Cauchy}}$  in controlling Type I error for  $\Sigma \in \mathcal{F}_{d,\rho}$ . We first present the following lemma which is used in comparison of  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$ .

**Lemma 2.** For  $\alpha \in (0, 1)$ , we have  $2\pi c_\alpha^2 \sim \ell_\alpha$  as  $\alpha \rightarrow 0$ .

In Theorem 3, the dimension  $d$  is expressed based on the  $a_C$  and  $a_L$  such that the polynomial orders of the upper  $\alpha$  quantiles are  $d_C = [(c_\alpha)^{a_C}]$  and  $d_L = [(l_\alpha)^{a_L}]$ . In the proof of Lemma 2, approximations of tail probabilities of Cauchy and Lévy distributions are obtained by  $\alpha \sim 1/\pi c_\alpha$  and  $\alpha \sim \sqrt{2/\pi \ell_\alpha}$ . Then, three parameters  $(d, \alpha, a)$  for  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  have the following relationship :

$$\alpha_L \sim \sqrt{\frac{2}{\pi}} \left( \frac{1}{d_L} \right)^{\frac{1}{2a_L}}, \quad \alpha_C \sim \frac{1}{\pi} \left( \frac{1}{d_C} \right)^{\frac{1}{a_C}} \quad (3.8)$$

where  $(d_C, \alpha_C, a_C)$  and  $(d_L, \alpha_L, a_L)$  are three parameters corresponding to  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$ .

Theorem 3 shows that  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  control Type I error  $\alpha$  asymptotically under  $0 < a_C < \frac{1-\rho^2}{(\sqrt{3}\rho+\sqrt{2})^2}$  and  $0 < a_L < \frac{1-\rho^2}{5\rho^2+4\sqrt{3}\rho+2}$ , respectively. It is obvious that  $\{(2\alpha_L, d_L)\}$  satisfying (3.8) includes  $\{(\alpha_C, d_C)\}$  since  $0 < 2a_L < \frac{2(1-\rho^2)}{5\rho^2+4\sqrt{3}\rho+2}$  includes  $0 < a_C < \frac{1-\rho^2}{(\sqrt{3}\rho+\sqrt{2})^2}$  due to  $\frac{1-\rho^2}{(\sqrt{3}\rho+\sqrt{2})^2} < \frac{2(1-\rho^2)}{5\rho^2+4\sqrt{3}\rho+2}$ . Figure S3 in the supplementary materials shows the upper bounds of  $2a_L$  and  $a_C$  for the values of  $0 \leq \rho < 1$  and it is shown that the upper bound of  $2a_L$  is larger than that of  $a_C$ .

### 3.2 Comparison of $T_{\text{Cauchy}}$ and $T_{\text{Lévy}}$ in controlling Type I error

Based on these, we define ranges of  $a_C$  and  $2a_L$  which are

$$\begin{aligned} \mathcal{I}_C &= \left\{ a_C : 0 < a_C < \frac{1 - \rho^2}{(\sqrt{3}\rho + \sqrt{2})^2} \right\}, \\ \mathcal{I}_L &= \left\{ 2a_L : 0 < 2a_L < \frac{2(1 - \rho^2)}{5\rho^2 + 4\sqrt{3}\rho + 2} \right\} \end{aligned}$$

and regions of  $(\alpha, d)$  for  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$

$$\begin{aligned} \mathcal{G}_C &= \left\{ (\alpha_C, d_C) : \alpha_C \sim \frac{1}{\pi} \left( \frac{1}{d_C} \right)^{\frac{1}{a_C}}, a_C \in \mathcal{I}_C \right\}, \\ \mathcal{G}_L &= \left\{ (\alpha_L, d_L) : \alpha_L \sim \sqrt{\frac{2}{\pi}} \left( \frac{1}{d_L} \right)^{\frac{1}{2a_L}}, a_L \in \mathcal{I}_L \right\} \end{aligned}$$

where  $\mathcal{G}_C$  is the collection of sequences  $(\alpha, d)$  satisfying  $\alpha \sim c \cdot \left(\frac{1}{d}\right)^{\frac{1}{a}}$  for  $a \in \mathcal{I}_C$

where  $c$  is a constant. and  $\mathcal{G}_L$  is similarly interpreted.

We note that, to control Type I error, Liu and Xie (2020) presents the condition of dimension  $d$  in  $T_{\text{Cauchy}}$  corresponding to  $0 < a_C < 1/2$  which is wider than  $0 < a_C < \frac{1-\rho^2}{(\sqrt{3}\rho+\sqrt{2})^2}$  in Theorem 3. Indeed,  $0 < a_C < \frac{1-\rho^2}{(\sqrt{3}\rho+\sqrt{2})^2}$  equals to  $0 < a_C < 1/2$  when  $\rho = 0$ . Hence, the conditions in Theorem 3 can be considered as a cost to encompass arbitrary structures of correlations matrix.

$\mathcal{G}_C$  and  $\mathcal{G}_L$  provide the scopes of  $\alpha$  and  $d$  under which  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  control Type I error for dependent  $p$ -values. We have the following corollary derived from Theorem 3 which states some advantage of  $T_{\text{Lévy}}$  over  $T_{\text{Cauchy}}$  in controlling Type I errors.

**Corollary 1.** *Under the conditions in Theorem 3, we have the following results*

### 3.3 Power Studies

for any  $0 < \rho < 1$ : (i)  $\mathcal{G}_C \subset \mathcal{G}_L$ . (ii) For  $a_L \in \mathcal{I}_L - \mathcal{I}_C$ , if  $\alpha_L \sim \alpha_C$ , we have  $d_C/d_L \rightarrow 0$  for any  $a_C \in \mathcal{I}_C$  and the corresponding  $d_C$ . If  $d_L \sim d_C$ , we have  $\alpha_C/\alpha_L \rightarrow 0$  for any  $a_C \in \mathcal{I}_C$  and the corresponding  $\alpha_C$ .

The first result in Corollary 1 shows that  $T_{\text{Lévy}}$  controls Type I error  $\alpha$  in a wider class of  $(\alpha, d)$  than  $T_{\text{Cauchy}}$ . The second result shows the behavior of  $\alpha$  and  $d$  of  $T_{\text{Lévy}}$  in  $\mathcal{G}_L$  but outside of  $\mathcal{G}_C$ . From these results,  $T_{\text{Lévy}}$  can control Type I error asymptotically in additional region of  $(\alpha, d)$  such as larger values of dimension  $d$  and  $\alpha$  compared to those of  $T_{\text{Cauchy}}$ .

### 3.3 Power Studies

To find sufficient conditions under which combination test statistics have asymptotically power converging to 1, we consider Type II error of each test which is  $1 - \text{power}$ . If Type II error of a test goes to zero asymptotically, it can be concluded that the test is asymptotically powerful. Denote  $\mathcal{N}$  and  $\mathcal{S}$  as index sets of null and signal part under non-null hypothesis, respectively, i.e., for  $i \in \mathcal{N}$ ,  $X_i \stackrel{d}{=} Z_i$  and for  $i \in \mathcal{S}$ ,  $X_i = Z_i + \mu_i$  where  $Z_i \sim N(0, 1)$  and  $\mu_i \equiv \mu = \sqrt{2\tau \log d}$  for  $\tau > 0$ . Assume that  $|\mathcal{S}| = d^\beta$ ,  $0 < \beta < 1$ . Donoho and Jin (2008) called the case of  $0 < \beta < 1/2$  a strongly sparse case and the case of  $1/2 < \beta < 1$  a moderately sparse case.

The following theorem states that neither of  $T_{\text{Cauchy}}$  nor  $T_{\text{Lévy}}$  dominates the

other in terms of controlling Type I error and obtaining asymptotic power 1.

**Theorem 4.** Assume that  $|\mathcal{S}| = d^\beta$  for  $0 < \beta < 1$ . (i) If  $\frac{1}{\beta+\tau-1} < a_C$ , Type II error of  $T_{Cauchy}$  converges to 0, as  $d \rightarrow \infty$ . (ii) If  $\frac{1}{\beta+2\tau-2} < a_L$ , Type II error of  $T_{Lévy}$  converges to 0, as  $d \rightarrow \infty$ .

From Theorem 3 and 4, we have sufficient conditions to achieve two criteria in (3.1). The following corollary presents conditions for  $T_{Cauchy}$  and  $T_{Lévy}$  to achieve (3.1) as  $\alpha \rightarrow 0$  and  $d \rightarrow \infty$ .

**Corollary 2.** Under the settings of Theorem 3 and 4, (i)  $T_{Cauchy}$  has Type I error + Type II error  $\rightarrow 0$  under  $d^{1-\beta-\tau} \ll \alpha \ll d^{-\frac{(\sqrt{3}\rho+\sqrt{2})^2}{1-\rho^2}}$  if  $\beta + \tau > (\sqrt{2}\rho + \sqrt{3})^2/(1-\rho^2)$ . (ii)  $T_{Lévy}$  has Type I error + Type II error  $\rightarrow 0$  under  $d^{1-\beta/2-\tau} \ll \alpha \ll d^{-\frac{5\rho^2+4\sqrt{3}\rho+2}{2(1-\rho^2)}}$  if  $\beta + 2\tau > (\sqrt{3}\rho + 2)^2/(1-\rho^2)$ .

Corollary 2 presents the sufficient asymptotic conditions of  $\alpha$  to attain the robustness and obtain the power for  $T_{Cauchy}$  and  $T_{Lévy}$ . For a given  $\rho$ , the upper bounds of  $\alpha$  implies that  $\alpha$  of  $T_{Cauchy}$  converges to 0 faster than  $T_{Lévy}$  to attain the robustness. On the other hand, for given moderately large values of  $\beta$  and  $\tau$ , the lower bounds of  $\alpha$  implies that  $T_{Lévy}$  requires  $d$  to increase faster than  $T_{Cauchy}$  to obtain the power. Therefore, the relationship between  $\alpha$  and  $d$  with  $\beta$ ,  $\tau$  and  $\rho$  shows the trade-off between robustness and power for  $T_{Cauchy}$  and  $T_{Lévy}$ .

---

#### 4. Asymptotic equivalence between heavy-tailed stable distribution and harmonic mean in combining $p$ -values

In this section, we discuss that hypothesis testing procedures based on  $p$ -value combining methods using heavy-tailed stable distribution are asymptotically equivalent to those using a type of harmonic mean. Wilson (2019) showed that the harmonic mean of  $p$ -values can be considered as a global  $p$ -value of hypothesis testing since it approximately tends to the Landau distribution which is a specific form of heavy-tailed stable distribution. Fang et al. (2023) and Chen et al. (2023) showed that Cauchy combination is asymptotically equivalent to the harmonic mean. Also, Wilson (2021) discussed the equivalence of Lévy combination and squared type of harmonic mean.

We investigate more general cases than those in Fang et al. (2023) and Wilson (2021). Fang et al. (2023) used the Pareto distribution which has different tail behaviors to represent heavy-tailed distribution for different parameters such as  $P(X > t) = C(\eta)/t^\eta$  for some constant  $C(\eta)$ . In particular, the Pareto distributions with  $\eta = 1$  and  $\eta = 1/2$  have the same asymptotic tail behavior as Cauchy and Lévy distributions, respectively except some constants. Indeed, the Cauchy and Lévy distributions are special cases of a family of stable distribution, referred to as the Lévy alpha-stable distribution of which the tail probability is given by  $P(X > t) \sim C^*(\eta, \psi, \gamma)/t^\eta$ , as  $t \rightarrow \infty$  where  $C^*$  is a constant de-

---

pending on parameters  $\eta$ ,  $\psi$  and  $\gamma$ . Here,  $\eta$  and  $\psi$  are parameters reflecting the heaviness of tail and the skewness, respectively, and  $\gamma$  is a scale parameter. With these parameters, the Pareto and the Lévy alpha-stable distributions reflect the properties of the heaviness of the tail probability.

Although Fang et al. (2023) and Wilson (2021) analyzed tail behaviors of heavy tailed distributions and their properties, in this section, we investigate the connection between combining method based on transforming into the Lévy alpha-stable distribution and the harmonic mean type method in detail. Concretely, we highlight the result of the equivalence of combining method of heavy tailed transformation and generalized harmonic mean method when dimension  $d$  goes to infinity while Fang et al. (2023), Wilson (2021) and Chen et al. (2023) used the assumption of fixed  $d$ .

First, we define the generalized harmonic mean of  $p$ -values by  $p_{global}^{gHMP(\eta)}$ : for  $p_1, \dots, p_d$  and  $\eta > 0$ ,  $p_{global}^{gHMP(\eta)} = d / \left( \sum_{j=1}^d p_j^{-1/\eta} \right)^\eta$ . The generalized harmonic mean includes typical cases, for example, when  $\eta = 1$ ,  $p_{global}^{gHMP(1)}$  is a harmonic mean of  $p$ -values. Note that the generalized harmonic mean is also used in Vovk and Wang (2020), Chen et al. (2023) and Wilson (2021) with the notation  $M_{r,K}$ .

Using the theory of regularly varying functions as in Wilson (2019), we can

derive the following result: for  $\eta > 0$  and weights  $w_i > 0, i = 1, \dots, d$ , as  $\epsilon \rightarrow 0$

$$P \left( \left( \sum_{j=1}^d w_j p_j^{-1/\eta} \right)^\eta \geq \frac{1}{\epsilon} \right) \sim \left( \sum_{j=1}^d w_j \right) P \left( p_i^{-1/\eta} \geq \epsilon^{-1/\eta} \right) = \epsilon,$$

where  $p_i$ 's are allowed to be arbitrary dependent and  $d$  is fixed. That is,  $p_{global}^{gHMP(\eta)}$  can be considered as a global  $p$ -value approximately.

In the perspective of testing procedure, we define a test statistic  $T_{gHMP(\eta)}$  based on the generalized harmonic mean of  $p$ -values as a global  $p$ -value. Let  $\alpha > 0$  and  $\eta > 0$  be given. Define  $T_{gHMP(\eta)} = I \left( p_{global}^{gHMP(\eta)} < \alpha \right)$ , where  $I(\cdot)$  is an indicator function. If  $T_{gHMP(\eta)} = 1$ , the null hypothesis is rejected and if  $T_{gHMP(\eta)} = 0$ , the null hypothesis cannot be rejected.

The following Theorem 5 and 6 present the robustness of Type I error and asymptotic power for  $T_{gHMP(\eta)}$ , respectively. In the below theorem, let  $E_{H_0}^\Sigma$  be the expectation under  $H_0$  with a given correlation matrix  $\Sigma \in \mathcal{F}_{d,\rho}$ .

**Theorem 5.** Let  $X \sim N_d(0, \Sigma)$  where  $\Sigma \in \mathcal{F}_{d,\rho}$ . Let  $d = \lceil \alpha^{-a_G} \rceil$ . Then for  $0 <$

$a_G < \frac{1-\rho^2}{(\eta+2)\rho^2+2\eta+2\rho\sqrt{(\eta+1)(2\eta+1)}}$ , under  $H_0$ , we have  $\lim_{d \rightarrow \infty} \sup_{\Sigma \in \mathcal{F}_{d,\rho}} \frac{E_{H_0}^\Sigma T_{gHMP(\eta)}}{\alpha} = 1$ , as  $\alpha \rightarrow 0$ .

**Theorem 6.** Assume that  $|\mathcal{S}| = d^\beta$  for  $0 < \beta < 1$ . If  $\frac{1}{\eta\beta+\tau-1} < a_G$ , Type II error of  $T_{gHMP(\eta)}$  converges to 0 as  $d \rightarrow \infty$ .

From Theorem 5 and 6, the following corollary presents the condition of  $\alpha$  and  $d$  under which both Type I and Type II errors of  $T_{gHMP(\eta)}$  converge to 0.

**Corollary 3.** Under settings of Theorem 5 and 6,  $T_{gHMP(\eta)}$  has Type I error + Type II error  $\rightarrow 0$  under  $d^{1-\eta\beta-\tau} \ll \alpha \ll d^{-\frac{(\eta+2)\rho^2+2\eta+2\rho\sqrt{(\eta+1)(2\eta+1)}}{1-\rho^2}}$ , if  $\tau + \eta\beta > \frac{(\eta+1)\rho^2+2\eta+1+2\rho\sqrt{(\eta+1)(2\eta+1)}}{1-\rho^2}$ .

Similar to  $p_{global}^{gHMP(\eta)}$ , we also generalize the transformation of  $p$ -values to heavy tailed stable distribution. For this, let  $h_{\eta,\psi}(p_i)$  be a transformation function of  $p$ -value to a random variable of heavy-tailed stable distribution with  $0 < \eta \leq 1$  and  $-1 \leq \psi \leq 1$ .  $h_{\eta,\psi}(p_i)$  also can be considered as a form of inverse cumulative distribution function of heavy-tailed stable distribution. If  $\eta = 1$  and  $\psi = 0$ , then  $h_{1,0}$  is a Cauchy transformation and if  $\eta = 1/2$  and  $\psi = 1$ , then  $h_{1/2,1}$  is a Lévy transformation. We denote a global  $p$ -value of transformation of stable distribution by  $g_{global}^{S(\eta,\psi)} = P_{H_0} \left( T_{\eta,\psi} > \sum_{i=1}^d d^{-1/\eta} h_{\eta,\psi}(p_i) \right)$ , where  $T_{\eta,\psi}$  is a combination test statistic of heavy-tailed transformation. It follows from Nolan (2020) that, for observed  $p$ -values,  $p_1, \dots, p_d$ ,

$$p_{global}^{S(\eta,\psi)} = P_{H_0} \left( T_{\eta,\psi} > \sum_{i=1}^d d^{-1/\eta} h_{\eta,\psi}(p_i) \right) \sim d \cdot c_{\eta,\psi} \cdot \left( \sum_{i=1}^d h_{\eta,\psi}(p_i) \right)^{-\eta},$$

where  $c_{\eta,\psi}$  denotes a generic constant depending on  $\eta$  and  $\psi$ .

We now present the following Theorem 7 which investigates the asymptotic equivalence between  $p_{global}^{gHMP(\eta)}$  and  $p_{global}^{S(\eta,\psi)}$  under some conditions when  $d \rightarrow \infty$ .

**Theorem 7.** For  $\Sigma \in \mathcal{F}_{d,\rho}$ ,  $0 < \eta \leq 1$  and  $-1 \leq \psi \leq 1$ , if one of the following conditions: (i)  $\max_j \sum_{i=1}^d \rho_{ij}^2 \leq C_0$  for some constant  $C_0 > 0$ , (ii) for  $0 < \eta <$

---

1,  $0 < \rho < 1 - \eta$  and for  $\eta = 1$ ,  $0 \leq \rho_{ij} \leq \rho = o\left(\frac{\log \log d}{\log d}\right)$  is satisfied, then we have  $\frac{p_{global}^{S(\eta, \psi)}}{p_{global}^{gHMP(\eta)}} \xrightarrow{p} 1$  as  $d \rightarrow \infty$ .

We note that Fang et al. (2023) proved the equivalence between transformation functions of Cauchy combination and harmonic mean for one  $p$ -value. On the other hand, Vovk and Wang (2020) showed the asymptotic equivalence of Cauchy combination and harmonic mean when the minimum of  $p$ -values goes to zero for fixed dimension. However, Theorem 7 presents the equivalence of combining method of heavy tailed transformation and generalized harmonic mean method when  $d \rightarrow \infty$  which is more practical assumption than that of Fang et al. (2023) and Vovk and Wang (2020).

**Remark 2.** 1. We request some conditions on  $\Sigma$  in Theorem 7 such as (i)  $\max_j \sum_{i=1}^d \rho_{ij}^2 \leq C_0$  for some constant  $C_0 > 0$  or (ii)  $\rho = o\left(\frac{\log \log d}{\log d}\right)$ . Neither (i) nor (ii) implies the other. The first one allows sparse  $\rho_{ij}^2$ , such that a small number of  $\rho_{ij}^2$  can have values not diminishing to 0 and all others are fairly negligible. This does not satisfy the condition (ii) since there exist some  $\rho_{ij}$  which does not diminish to zero. On the other hand, (ii) can include the case of dense such that most of  $\rho_{ij}$  diminishing to zero, but decreasing order is faster than  $\frac{\log \log d}{\log d}$ . This does not satisfy (i) since  $\sum_{j=1}^d \rho_{ij}^2 = o\left(\frac{d \log \log d}{\log d}\right)$  may diverge.

2. The asymptotic equivalence between  $p_{global}^{S(\eta, \psi)}$  and  $p_{global}^{gHMP(\eta)}$  shown in Theorem 7 implies that the condition of robustness of two methods are also equiva-

---

lent. In Theorem 5 and 6, when  $\eta = 1/2$ , the upper and lower bounds of  $a_G$  can be expressed by 2 times of the upper and lower bounds of  $a_L$  in Theorem 3 and 4, respectively. On the other hand, when  $\eta = 1$ , the upper and lower bounds of  $a_G$  equal to the upper and lower bounds of  $a_C$  in Theorem 3 and 4, respectively.

## 5. Numerical Studies

In this section, we evaluate Type I error and power for various dependency structures to compare  $p$ -values combining methods according to heaviness of transformation functions and simulate the relationships between  $p$ -values combining methods. To investigate properties of  $p$ -values combining methods, we construct  $p$ -values from the multivariate normal distribution with mean  $\boldsymbol{\mu}$  and correlation matrix  $\Sigma$  in  $\mathbf{X} \sim N_d(\boldsymbol{\mu}, \Sigma)$  and then derive one-sided right-tailed  $p$ -values. The correlation matrix  $\Sigma = (\rho_{ij})_{i,j=1,\dots,d}$  is defined by as follows:

1. **Exchangeable:** For  $0 < \rho < 1$ ,  $\rho_{ij} = \rho$  if  $i \neq j$ , and 1 if  $i = j$ .
2. **Polynomial decay:** For  $r > 0$ ,  $\rho_{ij} = 1/(1 + |i - j|^r)$  for  $i \neq j$  and 1 if  $i = j$ .
3. **Exponential decay, AR(1):** For  $0 < \rho < 1$ ,  $\rho_{ij} = \rho^{|i-j|}$ .

For each dependency structure,  $\rho$  and  $r$  decide the strength of the dependency.

Thus we vary  $\rho$  and  $r$  to evaluate effects of the dependencies.

## 5.1 Asymptotic Equivalence between Heavy Tailed combination test and Generalized Harmonic mean method

### 5.1 Asymptotic Equivalence between Heavy Tailed combination test and Generalized Harmonic mean method

In this subsection, we present simulation studies for three types of correlation matrices which demonstrate the similarity of  $p_{\text{global}}^{S(\eta,\psi)}$  and  $p_{\text{global}}^{gHMP(\eta)}$  for finite  $d$ . We consider  $(\eta, \psi) = (1, 0)$  and  $(1/2, 1)$  in  $p_{\text{global}}^{S(\eta,\psi)}$  which are corresponding to the Cauchy method and Lévy method, respectively.

Figures S4 and S5 in the supplementary materials provide plots of  $p_{\text{global}}^{S(1,0)}$  *vs.*  $p_{\text{global}}^{gHMP(1)}$  and  $p_{\text{global}}^{S(1/2,1)}$  *vs.*  $p_{\text{global}}^{gHMP(1/2)}$ . As mentioned, we observe that the Cauchy method and the harmonic mean method are quite close to each other for the cases of polynomial decay with large coefficient (Bottom panel) and exponentially decay (Middle panel) since the scatter plots are almost located around the diagonal line. However, under exchangeable dependency structures indicating strong dependence, the top panels in Figure S4 show that these two methods tend to be different. All plots in Figure S5 show that the Lévy method  $p_{\text{global}}^{S(1/2,1)}$  looks fairly close to  $p_{\text{global}}^{gHMP(1/2)}$  for all three correlation matrices. From Figure S4 and S5, we see that  $p_{\text{global}}^{S(1/2,1)}$  and  $p_{\text{global}}^{gHMP(1/2)}$  are close to each other compared to  $p_{\text{global}}^{S(1,0)}$  and  $p_{\text{global}}^{gHMP(1)}$ . This is due to the fact that  $p_{\text{global}}^{S(1/2,1)}$  and  $p_{\text{global}}^{gHMP(1/2)}$  are dominated by a smaller number of  $p$ -values which are affected less severely by correlation structures compared to  $p_{\text{global}}^{S(1,0)}$  and  $p_{\text{global}}^{gHMP(1)}$ .

## 5.2 Type I error

We present Figure 1 showing Type I error of each method at significant level 0.05 under exchangeable cases where  $\rho$  varies from 0 to 0.9 and  $d$  is fixed at 200. As expected from Theorem 1,  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  fail in controlling a given Type I error 0.5 except  $\rho = 0$ , the case of independent  $p$ -values. We also see that as  $\rho$  increases, Type I errors of  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  converge to  $\frac{1}{2}$  and  $\frac{1}{e}$ , respectively, shown in Theorem 1. In Figure 1, patterns of Type I errors of  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  are similar, since two methods are influenced by a small number of dominating  $p$ -values. Additionally, as  $\rho$  increases, both  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  are getting more conservative. On the other hand,  $T_{\text{Cauchy}}$  fails to control Type I error at the mid-level dependency structures while Type I error is well controlled for independent or extremely correlated  $p$ -values. Since  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  tend to have inflated Type I errors, we focus on the rest of tests such as  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$ .

Figure 2 shows Type I errors of  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  under exchangeable, polynomially and exponentially decaying dependency structures at significance levels 0.05, respectively. As seen in Figure 1, Figure 2 shows almost similar patterns of Type I errors under three types of correlation matrices. In addition, to evaluate effects of dimension and significance level, we present Type I errors of each method with  $d = 2000$  and 3000 in Figure S7 and at significance level 0.01 in Figure S6 in the supplementary materials showing also similar patterns.

## 5.2 Type I error

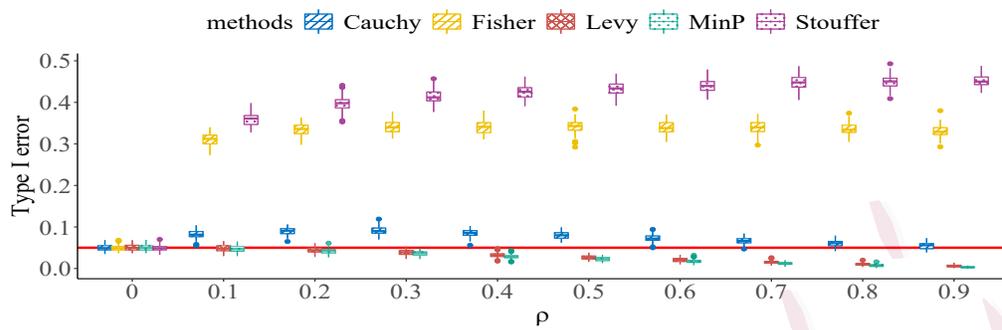


Figure 1: Box plots of Type I errors of  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$ ,  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  under exchangeable case. Horizontal line indicates the significance level 0.05.

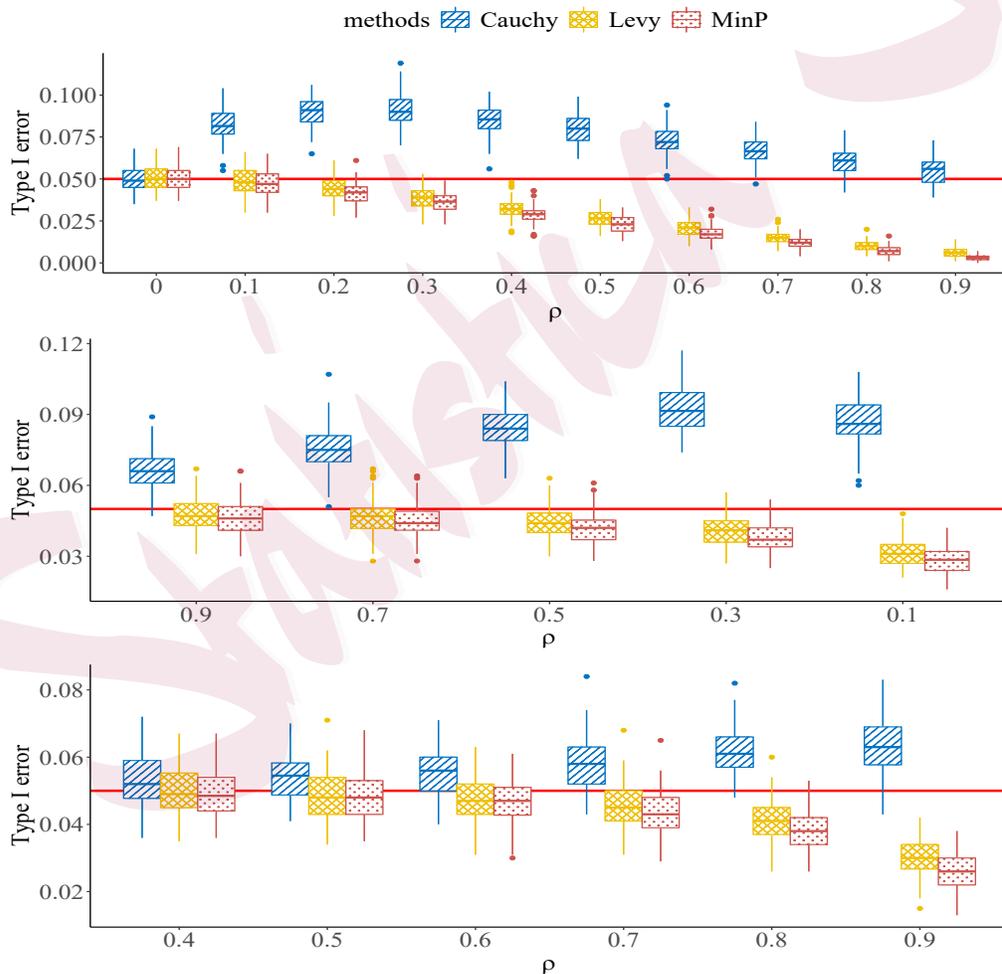


Figure 2: Box plots of Type I error of  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  at 0.05 with exchangeable, polynomially and exponentially decaying dependency structures.

### 5.3 Power

In order to compare the finite-sample performances of  $p_{\text{global}}^{gHMP(\eta)}$  and  $p_{\text{global}}^{S(\eta)}$ , Figure S8 in the supplementary materials represents Type I error of  $T_{\text{Cauchy}}$ ,  $T_{\text{gHMP}(1)}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{gHMP}(1/2)}$  under the three types of correlation. As mentioned in Section 5.1 and as shown in Figure S4 and S5,  $p_{\text{global}}^{S(1/2,1)}$  and  $p_{\text{global}}^{gHMP(1/2)}$  are close to each other compared to  $p_{\text{global}}^{S(1,0)}$  and  $p_{\text{global}}^{gHMP(1)}$  in the finite samples.

### 5.3 Power

We compare powers of  $T_{\text{Stouffer}}$ ,  $T_{\text{Fisher}}$ ,  $T_{\text{Cauchy}}$ ,  $T_{\text{gHMP}(1)}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  under the exchangeable case in Figure 3. We also provide Figure S9 and S10 in the supplementary materials for polynomially and exponentially decaying dependency structures at significance level  $\alpha = 0.05$ . Simulation settings are similar to Liu and Xie (2020) under the model  $\mathbf{X} \sim N_d(\mu, \Sigma)$ , where  $\mu = (\mu_i)_{1 \leq i \leq d}$  and  $\Sigma$  is a correlation matrix. All signals for non-null are defined to have the same strength and to account for effects of sparsity of signals such that  $\mu_i := \sqrt{2 \log(d)}/s^{1/3}$  for  $i \in \mathcal{S}$ , where  $\mathcal{S}$  denotes the index set of signals and  $s$  is the cardinality of  $\mathcal{S}$ . The range of  $d$  is 20 to 500.

In Figure 3, powers of  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  are unreliable since Type I errors are inflated under  $\rho > 0$  shown in Figure 1. On the other hand, powers of  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  are reliable regardless of  $\rho > 0$ . Powers of  $T_{\text{Cauchy}}$  depend on  $\rho$  since its Type I errors are inflated except the cases where correlations are close to 0 or

1. Figure S9 and S10 show similar patterns to Figure 3.

Compared to  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$ ,  $T_{\text{minP}}$  has the lowest power which is related to the conservativeness represented in numerical studies of Type I errors. There are additional patterns in the powers of each method. For example, in exchangeable case in Figure 3,  $T_{\text{Cauchy}}$  has more power than  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$ . In addition, for more weak dependency structures such that correlations are polynomially and exponentially decaying,  $T_{\text{Lévy}}$  achieves comparable power compared to  $T_{\text{Cauchy}}$  (Figure S9 and S10). As signals become dense such that the cardinality of  $\mathcal{S}$  increases, the power of  $T_{\text{Lévy}}$  is getting lower than that of  $T_{\text{Cauchy}}$ . This is from the fact that  $T_{\text{Lévy}}$  is dominated by a smaller number of components than  $T_{\text{Cauchy}}$ . This indicates the trade-off relationship between Type I error and power discussed in section 3 that  $T_{\text{Lévy}}$  has advantage in controlling Type I error under dependent  $p$ -values at the cost of losing powers under dense signals.

We also compare size-adjusted powers of each method in finite samples under exchangeable, polynomially and exponentially decaying cases in Figure S11, S12 and S13, respectively, in the supplementary materials.

## 6. Real Data Example

For a real-data analysis, we apply the combination methods to the Crohn's disease GWAS (Duerr et al., 2006) which is also used in Liu and Xie (2020). For

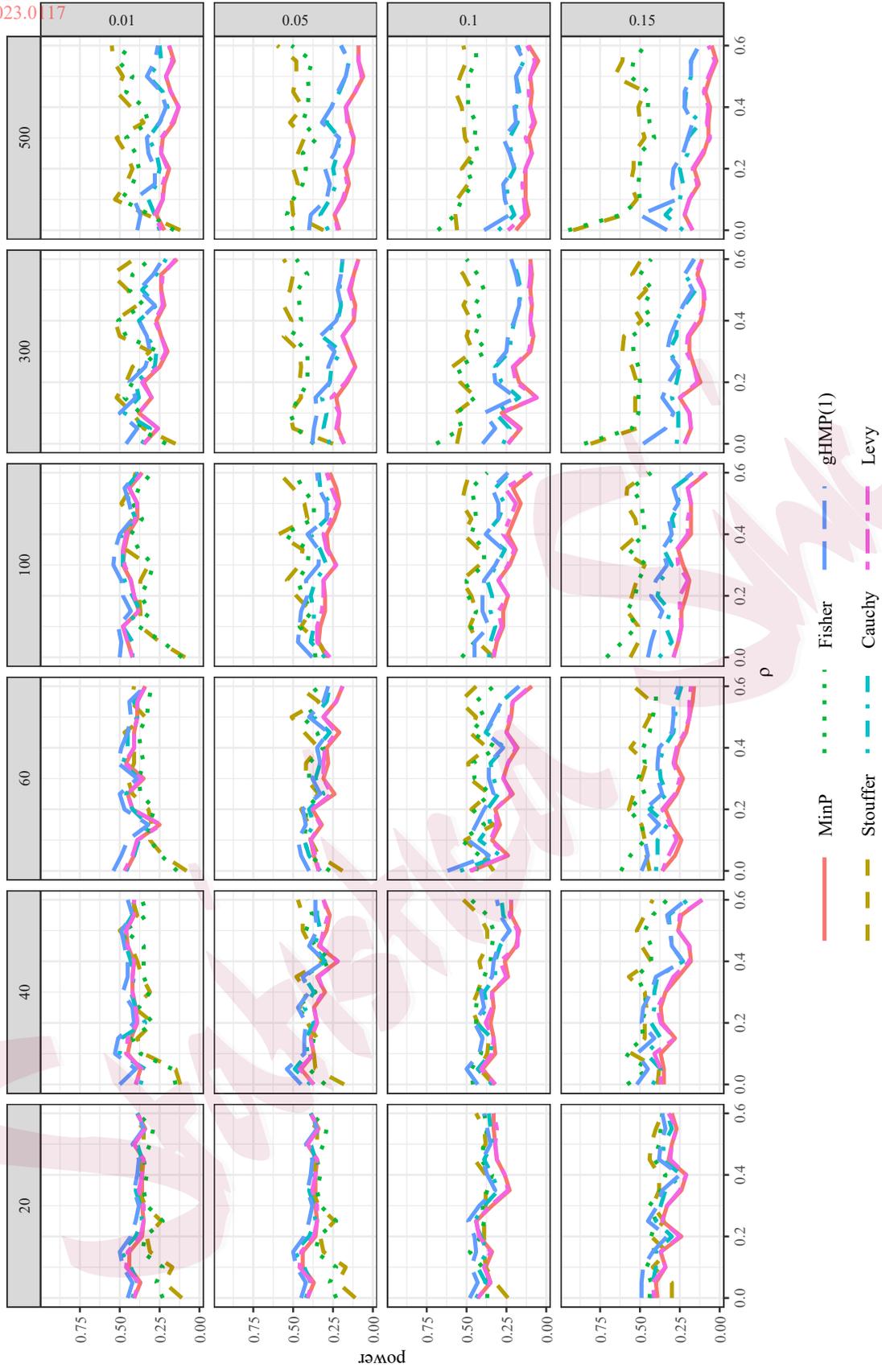


Figure 3: Power comparison for exchangeable case. The columns from left to right correspond to the dimension  $d = 20, 40, 60, 100, 300, 500$ . The rows from top to bottom correspond to the signal proportion 1%, 5%, 10%, 15%.

---

this data, the goal is identifying genes that are associated with the disease.

The dataset is based on 968 cases and 995 controls using the Illumina HumanHap300 Genotyping BeadChip. The cases were selected to have Crohn's disease, and the controls were matched to the cases based on sex and year of birth. Subjects were drawn from two cohorts: 1. persons with non-Jewish, European ancestry (561 cases and 563 controls), 2. persons with Jewish ancestry (407 cases and 432 controls). To analyze the association between genes and the disease, we grouped SNPs to genes via Genome Browser in a Box. Each SNP can be contained in multiple genes. As a result, all 242,535 SNPs are grouped into 19,769 genes according to the Genome Reference Consortium Human Build 38. The number of SNPs in each gene ranges from 1 to 676. Among all genes, only 4,969 genes having more than 10 SNPs are used.  $p$ -values are constructed by using a Cochran-Mantel Haenszel chi-square test separately for each SNP. The dataset is downloaded from the database of Genotypes and Phenotypes.

Figure 4 shows histograms of  $p$ -values of SNPs in two genes, "NLGN1" and "CDH4". We present these two genes since they show some strong deviation from the uniform distribution which may be expected for highly correlated  $p$ -values. There are 118 and 123 SNPs in each gene, respectively. The histogram of  $p$ -values in each gene deviate from the uniform distribution since there are more null  $p$ -values, which is  $p$ -values from the null distribution than those of

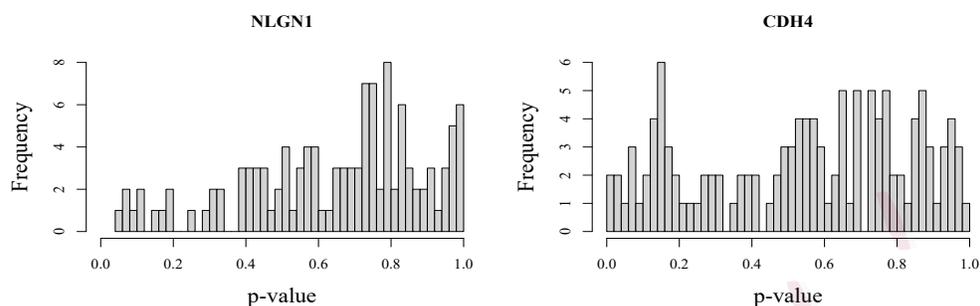


Figure 4: Histogram of  $p$ -values of SNPs for the “NLGN1” and “CDH4” genes ideal uniform distribution. There may be possible cases such that either the  $p$ -values are highly correlated or there are many weak signals which are hard to be detected. However, it is more reasonable to assume that  $p$ -values in a specific gene are highly correlated based on biological properties of GWAS. From this point of view, the objective of the study is to test if there are any significant signals among the highly correlated SNPs in each gene and to compare combination methods. Since the dataset consists of only  $p$ -values of individual SNPs, a dependency structure between  $p$ -values is unknown and can not be estimated.

Table 1 and Table 2 in the supplementary materials represent  $p$ -values generated from  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$ . Table 1 contains 15 genes rejected by  $T_{\text{Cauchy}}$  and  $T_{\text{Lévy}}$  simultaneously. Table 2 contains 15 genes rejected by  $T_{\text{Cauchy}}$  but cannot be rejected by  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$ . Table 1 and 2 show that  $p$ -values from  $T_{\text{Cauchy}}$ ,  $T_{\text{Lévy}}$  and  $T_{\text{minP}}$  tend to be ordered for each gene which matches to our theoretical and numerical results.

---

On the other hand,  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  produce unstable  $p$ -values since their  $p$ -values have large variability compared to the other three methods. Although it is difficult to validate this for real data examples, we conjecture that  $T_{\text{Stouffer}}$  and  $T_{\text{Fisher}}$  are not robust to dependent  $p$ -values and all the other three methods have robustness to dependent  $p$ -values in controlling Type I error. We also discuss Type I error of each method using simulated null  $p$ -values which is discussed in detail in the supplementary materials. In addition, we present Quantile-Quantile plots (Figure S14) of real dataset to check Type I errors of each methods.

## 7. Concluding Remarks

In this paper, we analyze properties of different  $p$ -value combining methods according to their heaviness of transformation functions of individual  $p$ -values, when  $p$ -values are correlated and the dependence structure is unknown. We investigated Type I error and power of several methods under a wide class of dependence structures both theoretically and numerically. We also provide the theoretical study on the asymptotic equivalence between two types of recently proposed  $p$ -value combining methods, stable combination and harmonic mean.

Throughout this paper, we provide intensive results on characteristic of different  $p$ -value transformations from the view point of the thickness of the tail which affects the robustness to dependence structure. Our contributions are more

extensive compared to existing studies in the following sense : Firstly, we introduce a class of dependence structure based on Gaussian copula with increasing dimension while Chen et al. (2023) used arbitrary dependence under fixed dimension and Liu and Xie (2020) and Liang and Rho (2022) used a correlation matrix with bounded eigenvalues and serially correlated  $p$ -values. Secondly, we present theoretical studies on controlling Type I error for different dependent structures and provide insight on the reason for such failures for different transformations. Thirdly, we provide theoretical results showing the trade-off relationship in controlling Type I error and obtaining powers depending on different heaviness of the tail. Lastly, we investigate relationships between harmonic mean type methods and combining methods with stable distribution.

### **Supplementary Materials**

All proofs of theoretical studies and additional numerical results are contained in the supplementary materials.

### **Acknowledgements**

Research of J. Park was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. 2020R1A2C1A01100526).

## References

- Chen, Y., P. Liu, K. S. Tan, and R. Wang (2023). Trade-off between validity and efficiency of merging  $p$ -values under arbitrary dependence. *Statistica Sinica* 33, 851–872.
- Donoho, D. and J. Jin (2008). Higher criticism thresholding: Optimal feature selection when useful features are rare and weak. *Proceedings of the National Academy of Sciences* 105(39), 14790–14795.
- Duerr, R. H., K. D. Taylor, and S. R. Brant (2006). A genome-wide association study identifies *il23r* as an inflammatory bowel disease gene. *Science* 314(5804), 1461–1463.
- Fang, Y., G. C. Tseng, and C. Chang (2023). Heavy-tailed distribution for combining dependent  $p$ -values with asymptotic robustness. *Statistica Sinica* 33, 1115–1142.
- Fisher, R. A. (1934). *Statistical methods for research workers*. Edinburgh: Oliver & Boyd.
- Hartung, J. (1999). A note on combining dependent tests of significance. *Biometrical Journal* 45(7), 849–855.
- Kost, J. T. and M. P. McDermott (2002). Combining dependent  $p$ -values. *Statis-*

## REFERENCES

---

- tics and Probability Letters* 60(2), 183–190.
- Liang, X. and Y. Rho (2022). Stable combination tests. *Statistica Sinica* 32, 641–644.
- Liu, Y. and J. Xie (2020). Cauchy combination test: A powerful test with analytic  $p$ -value calculation under arbitrary dependency structures. *Journal of the American Statistical Association* 115(529), 393–402.
- Nolan, J. P. (2020). *Univariate Stable Distributions; Models for Heavy Tailed Data*. Switzerland: Springer.
- Stouffer, S. A., E. A. Suchman, L. C. Edvinney, S. A. Star, and R. M. Williams (1949). *The American Soldier: Adjustment During Army Life*. New Jersey: Princeton University Press.
- Tippett, L. H. C. (1931). *The Methods of Statistics*. London: Williams and Norgate.
- Vovk, V. and R. Wang (2020). Combining  $p$ -values via averaging. *Biometrika* 107(4), 791–808.
- Wilson, D. J. (2019). The harmonic mean  $p$ -value for combining dependent tests. *Proceedings of the National Academy of Sciences* 116(4), 1195–1200.
- Wilson, D. J. (2021). The lévy combination test. *arXiv:2105.01501v1*.

## REFERENCES

---

Department of Statistics, Duksung Women's University, Seoul, Korea

E-mail: (junsik@duksung.ac.kr)

Department of Statistics, Seoul National University, Seoul, Korea

E-mail: (junyongpark@snu.ac.kr)

Statistica Sinica