Minimax Optimality of Sign Test for Paired Heterogeneous Data

Martin J. Zhang Stanford University jinye@stanford.edu Meisam Razaviyayn University of Southern California razaviya@usc.edu **David Tse** Stanford University dntse@stanford.edu

Abstract

Comparing two groups under different conditions is ubiquitous in the biomedical sciences. In many cases, samples from the two groups can be naturally paired; for example a pair of samples may come from the same individual under the two conditions. However samples across different individuals may be highly heterogeneous. Traditional methods often ignore such heterogeneity by assuming the samples are identically distributed. In this work, we study the problem of comparing paired heterogeneous data by modeling the data as Gaussian distributed with different parameters across the samples. We show that in the minimax setting where we want to maximize the worst-case power, the sign test, which only uses the signs of the differences between the paired sample, is optimal in the one-sided case and near optimal in the two-sided case. The superiority of the sign test over other popular tests for paired heterogeneous data is demonstrated using both synthetic data and a real-world RNA-Seq dataset.

1 INTRODUCTION

A common form of scientific experimentation is the comparison of two groups. Suppose we collected 2n samples $\{X_i^A\}_{i=1}^n$, $\{X_i^B\}_{i=1}^n$ under two conditions A and B. The conditions may be sick v.s. healthy, pre- v.s. post- treatment, etc. In the traditional homogeneous setting, samples within each group are assumed to be independently and identically distributed (i.i.d.), i.e.

$$X_i^A \overset{\text{i.i.d.}}{\sim} \mathbb{P}_A, \ X_i^B \overset{\text{i.i.d.}}{\sim} \mathbb{P}_B, \ \forall \ i=1,\cdots,n,$$

where the distributions \mathbb{P}_A and \mathbb{P}_B typically come from some common distribution families like Gaussian or Poisson. The goal is to infer whether there is a difference between the mean of the two groups, i.e. if $\mathbb{E}[X_i^A] \neq \mathbb{E}[X_i^B]$. One of the most commonly used test in this case is the two-sample t-test, which assumes the distributions \mathbb{P}_A and \mathbb{P}_B are Gaussian and is based on the t-statistic [5]. However, in many real-world applications, the data are *paired* and *heterogeneous*. The paired structure means that for each i, the samples X_i^A , X_i^B are similar due to some shared properties. The heterogeneity means that the data within the same group, $\{X_i^A\}$ or $\{X_i^B\}$, may be non-identically distributed. As a result, the paired differences $\{X_i^B - X_i^A\}$ may also be non-identically distributed.

Such paired heterogeneous data may occur in many scenarios. For example, in pre- v.s. post-treatment studies [16], samples were taken before and after the treatment from the same individuals. Samples from the same person are similar and thus can be paired, while samples from different individuals may be very different due to individual-level heterogeneity. In another study, samples were obtained from the same person over a long period time to study viral infection disease (VID) [3]. Samples under different conditions (sick/healthy) can be paired if they are close to each other in time. As the person may change a lot over time, within-pair samples are more similar than within-group samples that are far from one another in time¹. See the following example.

¹One may argue that a time-series analysis is more appropriate [1]. However, when we are not interested in the time-series pattern, the differential expression analysis by two-group comparison is still a valid method and has been used in various studies. Second, many RNA-Seq datasets including VID are noisy and have very few

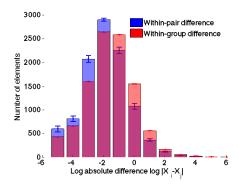


Figure 1: Visualization of within-pair difference and within-group different for VID.

Example 1. (Heterogeneity in VID data) We present a visualization on VID to illustrate the presence of the heterogeneity in data. In this dataset, a sample $\mathbf{X}_i \in \mathbb{R}^{23,231}$ is a measurement of the gene expression level of 23, 231 genes, and samples are taken from one person over 1124 days under two conditions, healthy and sick (see Fig. 3). We match samples close to each other in time, one from each group (healthy/infected), as pairs. We plot the histogram of the pairwise difference of within-pair samples and of within-group samples (see the details in Supp. Sec. 1). As can be seen in Figure 1, the within-pair difference is systematically smaller than the within-group difference, indicating that the within-pair samples are more similar than within-group samples.

A popular treatment to the paired heterogeneous data is to allow each data point to have a different distribution but assume that the paired differences $X_i^B - X_i^A$ are i.i.d. If the paired differences are Gaussian, then the paired t-test is most powerful [19]. The Gaussian assumption is quite reasonable and used in many applications, e.g. the RNA-Seq data [11]. However, given the heterogeneity across data pairs, it is hard to justify that the paired differences are actually identically distributed. For example, in the pre- and post- treatment studies, different individuals may have different responses to the treatment, and hence the paired differences may not be i.i.d.

In this work we keep the Gaussian assumption but allow the paired differences to be non-identically distributed. To characterize the *systematic* difference between the two group means, we assume the probability of increase/decrease from one group to the other is the same across all pairs. Specifically, we assume:

• Independently Gaussian (possibly non-identical):

$$X_i^A \sim \mathcal{N}(\nu_i^A, (\sigma_i^A)^2), \ X_i^B \sim \mathcal{N}(\nu_i^B, (\sigma_i^B)^2), \ \forall i.$$

• The "tendency of shift" is the same across all data pairs: $\mathbb{P}(X_i^B \geq X_i^A) = \theta, \forall i$.

The tendency assumption is *weaker* than the previous i.i.d. assumption and allows the paired differences to be non-identically distributed. In the Gaussian case, it implies a natural scaling for the paired differences $X_i^B - X_i^A$; their means being proportional to their standard deviations. Even if the tendency assumption violated, the proposed test in this manuscript will still maintain a good power and be minimax optimal under certain conditions. See Remark 7 for details.

Then the problem of interest is to test if $\theta=0.5$. We seek robust tests that consistently produce high power under different levels of heterogeneity. A natural approach is to consider the minimax setting, where we fix the level of shift θ and maximize the worst-case power over all values of nuisance parameters.

Contributions. The main contribution of the paper is to identify the optimal test for our minimax setting, which turns out to be the sign test [12]. The sign test uses the number of times of the paired differences being positive as the summary statistics, i.e. $W = \sum_i \mathbb{I}_{\{X_i^B - X_i^A > 0\}}$. Our result shows that the sign test is maximin in the one-sided case where we want to test $\theta = 0.5$ against $\theta > 0.5$. In

samples. In those cases, the two-group comparison can produce more stable results compared to the time-series analysis.

the two-sided case where the alternative hypothesis becomes $\theta \neq 0.5$, we show that the worst-case power of any test can be upper bounded by that of the sign test plus a negligible additive term, implying the sign test is near optimal. In addition, we verify our theoretical analysis using both synthetic data and a real-world RNA-Seq dataset.

Let us explain our contributions within the context of prior art. Let the paired difference be $Y_i = X_i^B - X_i^A$. Prior to this work, it is known that if we restrict our attention to Y_i 's only, which is natural given the paired structure, then the sign test is maximin over the class of all distributions where the differences Y_i 's are independent and the "tendency of shift" $\mathbb{P}(Y_i)$ is the same across all pairs [12]. In this case, the worst-case distribution is any distribution pair \mathbb{P}_{H_0} , \mathbb{P}_{H_1} satisfying $\forall i, y_+ \geq 0, y_- < 0$,

$$\frac{\mathbb{P}_{H_1}(Y_i = y_+ | Y_i \ge 0)}{\mathbb{P}_{H_0}(Y_i = y_+ | Y_i \ge 0)} = \frac{\mathbb{P}_{H_1}(Y_i = y_- | Y_i < 0)}{\mathbb{P}_{H_0}(Y_i = y_- | Y_i < 0)} = 1.$$
(1)

The result is a direct consequence of plugging the above worst-case distribution in Theorem 8.1.1 in [12]. The above argument considers a class of distributions so general that it may yield an overly pessimistic result. Indeed, the worst-case distribution \mathbb{P}_{H_1} is very artificial in that it is continuous everywhere else but at y=0. Hence, it is natural to restrict ourselves to a smaller and more natural class.

According to empirical studies, the RNA-Seq data can be modeled as Gaussian random variables after variance-stabilization transformation [11]. In this work, therefore, we restrict ourselves to the family of normal distributions. As a result, the paired differences Y_i 's now have normal distributions. Ideally, this distribution information should be properly utilized. The question is that whether it leads to a test more powerful than distribution-free tests. Our result states that even this extra information does not help us to go beyond the sign test, indicating the importance of the sign information for robust testing. We also note that our result is not a straight forward extension of existing results. In fact, it is not even clear that whether the sign test remains optimal after restricting to the Gaussian class because the Gaussian assumption excludes the worst-case distribution (1) in the general class.

In terms of the novelty of the proof, the proof techniques here are completely different from the older proof. Theorem 8.1.1 [12] cannot be applied here because it is extremely difficult to find the worst-case distribution in the Gaussian case (See the discussion after Theorem 3). In fact, our conjecture is that the worst-case distribution does not even exist. We used a different strategy in our proof: we first show that the sign test is maximin among the family of "simple tests". Then we show that "simple tests" can approximate the Borel measurable tests arbitrarily well. This is inspired by the widely used techniques in measure theory. However, we made two changes here: first, the notion of "simple tests" is tailored to fit the location-scale invariance property, different from simple functions in measure theory; second, the approximation is in terms of the testing performance (size and power), rather than some function norms. From a theoretical point of view, our result fills in a missing piece in the minimax analysis (of the adaptivity of sign test to the Gaussian family) in the classical statistical literature.

Related works. This work has a very classical flavor. Some related topics are testing within-group heterogeneity [6], robust tests for paired data [9], and rank-based tests [13]. In the literature, paired t-test and the Wilcoxon signed-rank test are compared most often to the sign test [19]. The paired t-test assumes the paired differences $X_i^B - X_i^A$ are i.i.d. Gaussian and uses the t-statistic (2). Let S_i be the sign of $X_i^B - X_i^A$ and R_i be the rank of $|X_i^B - X_i^A|$ among all pairs. The Wilcoxon test uses the sign-rank statistic (3).

Paired t-test :
$$T = \sqrt{n} \frac{\text{mean}(X_i^B - X_i^A)}{\text{std}(X_i^B - X_i^A)},$$
 (2)

Paired t-test :
$$T = \sqrt{n} \frac{\text{mean}(X_i^B - X_i^A)}{\text{std}(X_i^B - X_i^A)},$$
 (2)
Wilcoxon test : $U = \sum_{i=1}^n S_i R_i$.

Most of the existing results for the above tests are based on the i.i.d. (homogeneous) scenario. For example, when the differences $X_i^B - X_i^A$ are i.i.d. Gaussian, the paired t-test is known to be most powerful and the relative efficiency of the sign test v.s. the paired t-test is $2/\pi$ (Table 14.1,[19]), and $3/\pi$ for Wilcoxon test (an extension of the former). Results are rare on the heterogeneous case, especially under the minimax setting.

The motivating application for the present work is RNA-Seq experiments. In RNA-Seq experiments, the gene expression level of people under different conditions are measured and the task is to

identify genes differentially expressed under the two conditions. In related works, within-group heterogeneity is usually modeled by assuming some prior distribution on the expression level, e.g. gamma distribution [4, 10]. The paired structure is modeled by the design matrix in the generalized linear model [11, 14], or by assigning same expression level parameters to samples in the same pair [4, 10]. All above methods assume some complex models for the data, e.g. Bayesian hierarchical model in [4] or some mean-variance function shared across genes in [17, 14]. This leads to the lack of thorough theoretical understanding and, consequently, difficulty in establishing theoretical guarantees.

"Those methods treat the estimated parameters as if they were known parameters, without allowing for the uncertainty of estimation, and this leads to statistical tests that are overly liberal in some situations" [18].

In fact, no theoretical result is yet available for the RNA-Seq data analysis [11]. On the contrary, the sign test is theoretically justified in this manuscript. As shown in the experiments, it can be easily applied to the RNA-Seq data after simple normalization. Moreover, despite its simplicity, it yields reasonable results compared to other much more complex methods.

The rest of the paper is organized as follows. After the problem formulation in Sec. 2, we prove the optimality of the sign test in Sec. 3, followed by a theoretical comparison of the sign test with the paired t-test in SubSec. 3.3. Finally, we present numerical experiments on both synthetic data and real-world data in Sec. 4. We postpone the proofs to the supplementary materials.

PROBLEM FORMULATION

Consider n paired data points $\{X_i^A, X_i^B\}_{i=1}^n$, where (X_i^A, X_i^B) denotes the i-th sample pair in groups A and B. Our goal is to detect whether there is a systematic difference between samples in two groups. We assume that 1. the samples are independently and normally distributed; 2. the "tendency of shift" is the same across all sample pairs. Let [n] denote the set $\{1, 2, \dots, n\}$. Then mathematically the above assumptions can be written as

$$\begin{split} X_i^A &\sim \mathcal{N}(\nu_i^A, (\sigma_i^A)^2), \ X_i^B \sim \mathcal{N}(\nu_i^B, (\sigma_i^B)^2), \ \forall i \in [n], \\ \mathbb{P}(X_i^B \geq X_i^A) &= \mathbb{P}(X_i^B \geq X_i^A) \triangleq \theta, \ \forall i, j \in [n], \end{split}$$

The range of parameters are

$$\{\nu_i^A, \nu_i^B \in \mathbb{R}, \ \sigma_i^A, \sigma_i^B \in \mathbb{R}_{\geq 0}, \ \text{s.t.} \ \mathbb{P}(X_i^B \geq X_i^A) = \theta\},$$

where $\mathbb{R}_{>0}$ denotes the set of non-negative real numbers. Clearly, $\theta=0.5$ means that there is no systematic shift from group A to group B, and $\theta \neq 0.5$ indicates that such systematic difference exists. Hence, our null hypothesis is $\theta = 0.5$. For the alternative hypothesis, if we have some prior knowledge on the shifting direction, we can test a one-sided alternative $\theta > 0.5$. Otherwise, a two-sided alternative, $\theta \neq 0.5$, is appropriate.

Let us represent our statistical model in a more tractable way. Let $\Phi(\cdot)$ be the cumulative density

Let us represent our statistical model in a more tractable way. Let
$$\Phi(\cdot)$$
 be the cuminative density function of the standard normal distribution. Since $\forall i, \mathbb{P}(X_i^B \geq X_i^A) = \Phi(-\frac{\nu_i^B - \nu_i^A}{\sqrt{(\sigma_i^B)^2 + (\sigma_i^A)^2}}) = \theta$, we obtain $\forall i, j, \frac{\nu_i^B - \nu_i^A}{\sqrt{(\sigma_i^B)^2 + (\sigma_i^A)^2}} = \frac{\nu_j^B - \nu_j^A}{\sqrt{(\sigma_j^B)^2 + (\sigma_j^A)^2}}$. Let $\delta \triangleq \frac{\nu_i^B - \nu_i^A}{\sqrt{(\sigma_i^B)^2 + (\sigma_i^A)^2}}$. Clearly, $\theta = \Phi(-\delta)$, and by defining $\mu_i \triangleq \sqrt{(\sigma_i^A)^2 + (\sigma_i^B)^2}$, $\rho_i \triangleq \frac{(\sigma_i^A)^2}{(\sigma_i^A)^2 + (\sigma_i^B)^2}$, and $\nu_i \triangleq \nu_i^A$, the above model becomes

$$X_i^A \sim \mathcal{N}(\nu_i, \rho_i \mu_i^2),$$

$$X_i^B \sim \mathcal{N}(\nu_i + \delta \mu_i, (1 - \rho_i) \mu_i^2), \ i \in [n],$$
(4)

where $\nu_i \in \mathbb{R}$, $\mu_i \in \mathbb{R}_{>0}$, $\rho_i \in [0, 1]$. Here $\mathbb{R}_{>0}$ is the set of positive real numbers and the tendency assumption prevents μ_i to be 0. It is not hard to see that in this equivalent representation, we test the null hypothesis $\mathcal{H}_0: \delta = 0$ against the alternative hypothesis $\mathcal{H}_1: \delta > 0$ (one-sided) or $\mathcal{H}_1: \delta \neq 0$ (two-sided). To quantify the size and the power, we let δ have some fixed unknown magnitude. For the two-sided case, it may be either positive or negative with the sign s_{δ} . Since we have no knowledge about the nuisance parameters $\{\nu_i, \mu_i, \rho_i, s_\delta\}_{i=1}^n$, a natural formulation is to look for the maximin test that maximizes the worst-case power over all possible values of the nuisance parameters.

Let the data vectors $\mathbf{X}^A \triangleq \{X_i^A\}_{i=1}^n$ and $\mathbf{X}^B \triangleq \{X_i^B\}_{i=1}^n$ be data points obtained by model (4). We use $\phi(\mathbf{X}^A, \mathbf{X}^B) : \mathbb{R}^n \times \mathbb{R}^n \mapsto [0,1]$ to denote a test that rejects the null hypothesis with probability $\phi(\mathbf{X}^A, \mathbf{X}^B)$ when the data are $\mathbf{X}^A, \mathbf{X}^B$. Given the nuisance parameters $\gamma \triangleq \{\nu_i, \mu_i, \rho_i, s_\delta\}_{i=1}^n$, the size is given by $\mathbb{E}_{\mathbb{P}_0(\gamma)}[\phi(\mathbf{X}^A, \mathbf{X}^B)]$, where the expectation is taken with respect to the null distribution $\mathbb{P}_0(\gamma)$ with the given nuisance parameters γ . Similarly, the power of the test is given by $\mathbb{E}_{\mathbb{P}_1(\gamma)}[\phi(\mathbf{X}^A, \mathbf{X}^B)]$. We call a test $\phi^*(\cdot, \cdot)$, a level- α maximin test if for any other test ϕ ,

$$\inf_{\gamma} \mathbb{E}_{\mathbb{P}_{1}(\gamma)}[\phi^{*}(\mathbf{X}^{A}, \mathbf{X}^{B})] \geq \inf_{\gamma} \mathbb{E}_{\mathbb{P}_{1}(\gamma)}[\phi(\mathbf{X}^{A}, \mathbf{X}^{B})],$$

$$\sup_{\gamma} \mathbb{E}_{\mathbb{P}_{0}(\gamma)}[\phi^{*}(\mathbf{X}^{A}, \mathbf{X}^{B})] \leq \alpha.$$

In other words, ϕ^* has the best worst-case power among all tests with size smaller than α over all values of the nuisance parameters γ . Equivalently, the problem can be stated as

$$\phi^* \in \arg \max_{\phi} \inf_{\gamma} \mathbb{E}_{\mathbb{P}_1(\gamma)}[\phi(\mathbf{X}^A, \mathbf{X}^B)],$$
s.t.
$$\sup_{\gamma} \mathbb{E}_{\mathbb{P}_0(\gamma)}[\phi(\mathbf{X}^A, \mathbf{X}^B)] \leq \alpha.$$
 (5)

For the sake of notational simplicity, we abbreviate the expressions as follows. Given two vectors \mathbf{a} and \mathbf{b} , let $\mathbf{a} \circ \mathbf{b}$ be the vector of the same dimension that contains the element-wise product of \mathbf{a} and \mathbf{b} . Also by the inequalities $\mathbf{a} > \mathbf{b}$ we refer to element-wise comparison, i.e $a_i > b_i$, $\forall i$. For two tests ϕ and ψ , we use $\phi = \psi$ to denote that the two tests have ϵ -similar performance,

$$|\inf_{\gamma} \mathbb{E}_{\mathbb{P}_{1}(\gamma)}[\phi(\mathbf{X}^{A}, \mathbf{X}^{B})] - \inf_{\gamma} \mathbb{E}_{\mathbb{P}_{1}(\gamma)}[\psi(\mathbf{X}^{A}, \mathbf{X}^{B})]| \leq \epsilon$$

$$|\sup_{\gamma} \mathbb{E}_{\mathbb{P}_{0}(\gamma)}[\phi(\mathbf{X}^{A}, \mathbf{X}^{B})] - \sup_{\gamma} \mathbb{E}_{\mathbb{P}_{0}(\gamma)}[\psi(\mathbf{X}^{A}, \mathbf{X}^{B})]| \leq \epsilon.$$
(6)

As a natural extension, $\phi \doteq \psi$ means that the two tests have the same performance, i.e., $\phi \stackrel{.}{=} \psi$.

Similarly, $\phi \leq \psi$ means ϕ has no better performance as ψ ; in other words,

$$\begin{split} &\inf_{\gamma} \mathbb{E}_{\mathbb{P}_{1}(\gamma)}[\phi(\mathbf{X}^{A}, \mathbf{X}^{B})] \leq \inf_{\gamma} \mathbb{E}_{\mathbb{P}_{1}(\gamma)}[\psi(\mathbf{X}^{A}, \mathbf{X}^{B})] \\ &\sup_{\gamma} \mathbb{E}_{\mathbb{P}_{0}(\gamma)}[\phi(\mathbf{X}^{A}, \mathbf{X}^{B})] \geq \sup_{\gamma} \mathbb{E}_{\mathbb{P}_{0}(\gamma)}[\psi(\mathbf{X}^{A}, \mathbf{X}^{B})]. \end{split}$$

Finally, we will use \mathbb{E}_0 and \mathbb{E}_1 as shorthand representations for $\mathbb{E}_{\mathbb{P}_0(\gamma)}$ and $\mathbb{E}_{\mathbb{P}_1(\gamma)}$, respectively.

3 OPTIMALITY OF SIGN TEST

Let us start by showing the location-scale invariance of the minimax problem (5). Generally speaking, the worst-case performance of a test $\phi(\cdot,\cdot)$ does not change under any shifting or scaling of the input: **Fact 2.** Let $\phi(\mathbf{X}^A, \mathbf{X}^B) : \mathbb{R}^n \times \mathbb{R}^n \mapsto [0,1]$ be an arbitrary test. For any $\mathbf{a} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^n$, let $\psi(\mathbf{X}^A, \mathbf{X}^B) = \phi(\mathbf{a} + \mathbf{b} \circ \mathbf{X}^A, \mathbf{a} + \mathbf{b} \circ \mathbf{X}^B)$. Then, $\phi \doteq \psi$.

See Supp. Sec 2 for the proof. Fact 2 leads to two conjectures. First, any data shifting does not affect the performance, suggesting that the absolute offset may be redundant and we should only focus on the relative difference $\mathbf{Y} \triangleq \mathbf{X}^B - \mathbf{X}^A$. Second, any data scaling does not affect the performance either. This suggests that the magnitude of the data, $|\mathbf{Y}|$, may also be redundant. Then, intuitively what remains, namely the signs of the difference $S_i \triangleq \operatorname{sgn}(Y_i)$, is the actual informative part for the minimax problem.

For the most powerful test using only the sign information, a sufficient statistic is the number of positive signs $W = \sum_{i=1}^n \mathbb{I}_{\{Y_i>0\}}$. Clearly, under the null distribution, $W \sim Bin(n,0.5)$. Then the sign test $\phi^S(\mathbf{X}^A,\mathbf{X}^B)$ can be written as

one-sided:
$$\mathbb{I}_{\{W>c_1\}} + p_1 \mathbb{I}_{\{W=c_1\}},$$
 (7)

two-sided:
$$\mathbb{I}_{\{|W-n/2|>c_2\}} + p_2 \mathbb{I}_{\{|W-n/2|=c_2\}},$$
 (8)

where c_1, c_2, p_1, p_2 are some constants calculated according to $W \sim Bin(n, 0.5)$ and the level α . Note that $p_1, p_2 \in [0, 1]$. In addition, the distribution of W and the power of the sign test does not depend on the values of the nuisance parameters. We next prove the optimality of the sign test in the one-sided case and near optimality in the two-sided case.

3.1 One-sided Case

In the one-sided case, let us assume $\theta \geq 0.5$ ($\delta \geq 0$). Then the nuisance parameters are just $\gamma \triangleq \{\nu_i, \mu_i, \rho_i\}_{i=1}^n$. Our main result below confirms that the sign test is indeed maximin in the sense of (5).

Theorem 3. Let $\mathcal{B}_{n\times n}$ be the class of Borel measurable functions that maps $\mathbb{R}^n \times \mathbb{R}^n$ to [0,1]. Then the one-sided sign test, as given in (7), is maximin among all tests $\phi(\mathbf{X}^A, \mathbf{X}^B) \in \mathcal{B}_{n\times n}$.

Recall from Section 2 that our statistical model is

$$X_i^A \sim \mathcal{N}(\nu_i, \rho_i \mu_i^2), \quad X_i^B \sim \mathcal{N}(\nu_i + \delta \mu_i, (1 - \rho_i) \mu_i^2).$$

Paired tests literatures suggest us to look at only the difference $Y_i = X_i^B - X_i^A \sim \mathcal{N}(\delta\mu_i, \mu_i^2)$ to get rid of the nuisance parameters ν_i 's and ρ_i 's, which is not surprising. Here let us take the case n=1 as an example to give some high-level intuitions why we can further reduce the sufficient statistics from Y_i 's to S_i 's.

According to Theorem 8.1.1 and the Neyman-Pearson lemma in (author?) [12], to show the sign test is maximin, it suffices to find a prior on μ_1 where S_1 is a sufficient statistic. However, a careful inspection of the proof of Lemma 5, especially on (25), reveals that there is no such single prior on μ_1 for which S_1 is a sufficient statistic. However, there is in fact a sequence of priors on μ_1 such that fixing the observation Y_1 , S_1 is asymptotically a sufficient statistic. For $k=1,2,\cdots$, consider the sequence of priors $g_k(\mu_1)=c_k/\mu_1$ for $\mu_1\in(1/k,k)$ and some normalizing constant c_k .

Let $f_0(\cdot;\mu)$, $f_1(\cdot;\mu)$ be the densities of $\mathcal{N}(0,\mu^2)$, $\mathcal{N}(\delta\mu,\mu^2)$ respectively and let f_0 , f_1 be that of $\mathcal{N}(0,1)$, $\mathcal{N}(\delta,1)$. A direct calculation (by change of variable $\mu'=\frac{Y_1}{\mu}$) shows that as $k\to\infty$, the likelihood ratio

$$\frac{f(Y_1|H_1)}{f(Y_1|H_0)} = \frac{\int_{\frac{1}{k}}^k f_1(Y_1;\mu) \frac{1}{\mu} d\mu}{\int_{\frac{1}{k}}^k f_0(Y_1;\mu) \frac{1}{\mu} d\mu} = \frac{\int_{\frac{1}{k}}^k f_1(\frac{Y_1}{\mu}) \frac{1}{\mu^2} d\mu}{\int_{\frac{1}{k}}^k f_0(\frac{Y_1}{\mu}) \frac{1}{\mu^2} d\mu} \stackrel{\mu' = \frac{Y_1}{\mu}}{=} \frac{\int_{\frac{1}{k}}^{k} f_1(\frac{Y_1}{\mu}) \frac{1}{\mu^2} d\mu}{\int_{\frac{1}{k}}^{k} f_1(\mu') d\mu'} \stackrel{k \to \infty}{\longrightarrow} 2 \left[\theta \mathbb{I}_{\{S_1 = +\}} + (1 - \theta) \mathbb{I}_{\{S_1 = -\}}\right].$$

Then asymptotically, the likelihood ratio of the observation $f(Y_1|H_1)/f(Y_1|H_0)$ depends only on S_1 , implying that S_1 is indeed asymptotically a sufficient statistic.

We note the above argument serves *only* as a high-level intuition and is by no means rigorous. More specifically, here we only show that for each Y_1 , S_1 is asymptotically a sufficient statistic for the sequence of priors g_1, g_2, \cdots , which is essentially a point-wise convergence result. However, a uniform convergence is needed to actually prove the final result. The rigorous proof is as follows.

Proof. (Proof sketch of Theorem 3) The idea is to show any test in $\mathcal{B}_{n\times n}$ performs no better than the sign test. The proof is composed of three main steps. We state the lemmas being used right after the corresponding steps, and relegate their proofs to the supplementary material. We recall that inequalities on vectors are element-wise, e.g. $\mathbf{a} \geq \mathbf{b}$ means $a_i \geq b_i$ for all i.

Step 1: Define the set of Borel measurable tests:

$$\mathcal{B}_n = \{ f : \mathbb{R}^n \mapsto [0,1], \ f \ is \ Borel \ measurable \}.$$

Lemma 4 implies that it suffices to show that the sign test is maximin among all tests $\phi(\mathbf{Y}) \in \mathcal{B}_n$.

Lemma 4. For any test $\phi(\mathbf{X}^A, \mathbf{X}^B) \in \mathcal{B}_{n \times n}$, there exists a Borel measurable test $\psi(\mathbf{Y}) \in \mathcal{B}_n$, such that $\phi \leq \psi$. Moreover, if ϕ is symmetric, then ψ is also symmetric².

Step 2: We show the sign test is maximin over the set of "simple test" $\mathcal{S} \subset \mathcal{B}_n$, which is defined as follows. Let $\mathcal{O} = \{-,+\}^n$ be the set of 2^n orthants in \mathbb{R}^n and let $\mathbf{o} = (o_1, \cdots, o_n) \in \mathcal{O}$. Consider any $\omega > 0$. For any $b \in \mathbb{Z}$, let the 1-D intervals be $I_b^+ = ((1+\omega)^b, (1+\omega)^{b+1}]$ and

²The symmetry is used for proving Theorem 8.

 $I_b^- = [-(1+\omega)^{b+1}, -(1+\omega)^b)$. Define the n-D box, specified by the orthant index \mathbf{o} and the interval index $\mathbf{b} = (b_1, \cdots, b_n)$, as $I_\mathbf{b}^\mathbf{o} = I_{b_1}^{o_1} \times \cdots \times I_{b_n}^{o_n}$. Then define the set of simple tests, denoted by $\mathcal{S}(\omega)$, to be the test that are piece-wise constant on the boxes $I_\mathbf{b}^\mathbf{o}$'s:

$$S(\omega) = \{ \phi : \phi = \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \phi_{\mathbf{b}}^{\mathbf{o}} \mathbb{I}_{I_{\mathbf{b}}^{\mathbf{o}}} + \phi_{0} \mathbb{I}_{\{0\}},$$

$$for \ some \ 0 \le \phi_{\mathbf{b}}^{\mathbf{o}} \le 1, 0 \le \phi_{0} \le 1 \},$$

$$(9)$$

and let $S = \bigcup_{\omega > 0} S(\omega)$. By Lemma 5, the sign test is maximin among all tests in S.

Lemma 5. The one-sided sign test (7) is maximin among all α -level tests in S.

Step 3: We show that \mathcal{S} approximates \mathcal{B}_n arbitrarily well in terms of testing performance as defined in (6), and hence establish the optimally of the one-sided sign test in \mathcal{B}_n . Specifically, by Lemma 6, $\forall \ \phi \in \mathcal{B}_n, \ \epsilon > 0, \ \exists \ \psi \in \mathcal{S}, \ \text{s.t.} \ \phi \stackrel{.}{=} \psi$. Letting $\epsilon \downarrow 0$ we have that ϕ^S is maximin among all tests in \mathcal{B}_n , concluding the proof.

Lemma 6. Let \mathcal{B}_n be the set of Borel measurable functions $f : \mathbb{R}^n \mapsto [0,1]$. For any $\phi(\mathbf{Y}) \in \mathcal{B}_n$ and any $\epsilon > 0$, there exists a measurable function $\psi \in \mathcal{S}$ such that $\phi = \psi$.

Remark 7. Under the alternative distribution, the testing statistics W will follow a binomial distribution $Bin(n, \theta)$. If the tendency assumption is violated, i.e. each pair has a different θ_i , it will instead follow a Poisson binomial distribution with parameter $(\theta_1, \dots, \theta_n)$, which has a tail property similar to that of the binomial distribution. Hence the sign test will still maintain a good power.

Moreover, when θ_i 's are different, one can consider a minimax setting over θ_i 's by defining the nuisance parameters to be $\{\mu_i, \rho_i, \nu_i\} \cup \{\theta_i : \theta_i \geq \theta_0\}$ for some $\theta_0 > 0.5$. It is not hard to see that the one-sided sign test is still maximin in this case. Specifically, the power will increase with the increase of any θ_i , and hence the worst-case is when $\theta_i = \theta_0 \ \forall i$, which reduces to the setting where θ_i 's are same.

3.2 Two-sided Case

Now we extend our result to the two-sided case, where we want to test $\theta=0$ v.s. $\theta\neq0$. Recall that in this case, we can no longer assume $\delta\geq0$ for the distribution in (4). So we modify the formulation in Section 2 by letting $s_{\delta}\in\{-1,1\}$ to be the sign of δ , and letting the nuisance parameter be $\gamma=\{\nu_i,\mu_i,\rho_i,s_{\delta}\}_{i=1}^n$. We fix the magnitude $|\delta|$ and consider the maximin problem (5). Without loss of generality assume $\alpha<0.5$. Let $\tilde{\phi}^S$ be the $\frac{\alpha}{2}$ -level one-sided sign test. The α -level two-sided sign test ϕ^S can be written as

$$\phi^S = \tilde{\phi}^S(\mathbf{Y}) + \tilde{\phi}^S(-\mathbf{Y}) \tag{10}$$

The following theorem shows in the two-sided case, the sign test is near optimal. See Supp. Subsec. 3.1 for the proof.

Theorem 8. (Two-sided case) Let $\mathcal{B}_{n \times n}$ be the class of Borel measurable functions that maps $\mathbb{R}^n \times \mathbb{R}^n$ to [0,1], and let ϕ^S be the two-sided sign test as defined in (10). For any α -level test $\phi \in \mathcal{B}_{n \times n}$, the worst-case power satisfies

$$\inf_{\gamma} \mathbb{E}_1[\phi] \le \inf_{\gamma} \mathbb{E}_1[\phi^S] + \frac{\alpha}{2} \exp(-\frac{n\delta^2}{2}).$$

If $\alpha = 0.05$ and $\delta = \frac{3}{\sqrt{n}}$, the additive term is 2.7e-4, almost negligible.

Remark 9. The proofs in both cases mainly use two properties of the Gaussian distribution. The first is the location-scale invariance, i.e., if we scale or shift the data points, the distribution still lies within the family of interest. This is used in the proof of Fact 2, Lemma 4, and Lemma 5. The second is the sub-Gaussian tail property. This is used in the proof of Lemma 6 and Theorem 8. Specifically, if we have a heavy-tailed distribution but the tail probability still vanishes, Lemma 6 will still hold, while the exponential term on the RHS of Theorem 8 will become a term with a slower vanishing speed that depends on the tail property of the distribution under consideration.

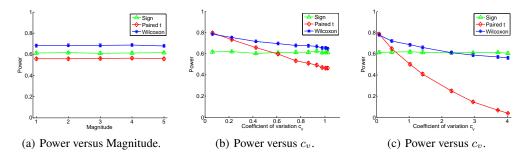


Figure 2: Effect of various parameters on the statistical power of the tests.

To generalize to result, in order for the sign test to be maximin in some family of distribution, the family needs to have location-scale invariance and sub-Gaussian tail property. Our conjecture is that these two are also sufficient for the minimaxity of the sign test.

3.3 Comparison with Paired T-test

Besides the minimaxity of the sign test, it is interesting to identify (realistic) conditions on the nuisance parameters such that the sign test outperforms other popular tests. We demonstrate this by comparing the asymptotic power of the sign test with that of the paired t-test, whose test statistic is given in (2). Since both the sign test and the paired t-test only use $\mathbf{Y} = \mathbf{X}^B - \mathbf{X}^A$ to compute the test statistics, we can only consider \mathbf{Y} as the input, which is generated by μ_i 's according to $Y_i \sim \mathcal{N}(\delta\mu_i, \mu_i^2)$, $\forall i$. Let $m_1 = \frac{1}{n} \sum_i \mu_i$ and $m_2 = \frac{1}{n} \sum_i (\mu_i - m_1)^2$ be the mean and the variance of μ_i 's. In the homogeneous case, μ_i 's are the same and thus $m_2 = 0$. In the presence of within-group heterogeneity, however, μ_i 's are different and m_2 may be large. Hence, intuitively, it is reasonable to look at the *coefficient of variation* [8], $c_v = m_2/m_1^2$, as a measure for within group heterogeneity. In fact, as shown below, it is the determining factor for the testing performance.

Theorem 10. Let $n \to \infty$ and scale δ with n such that $\delta \sqrt{n}$ remains constant³. Also assume that by increasing n the values of m_1 and m_2 remain constant. Then, the asymptotic power of the α -level two-sided sign test and the α -level two-sided paired t-test are given by (11), respectively. Moreover, the two-sided sign test has a larger asymptotic power if $c_v \ge \pi/2 - 1$.

Power of sign test:
$$Q\left(z_{\frac{\alpha}{2}} - \sqrt{\frac{2}{\pi}}\sqrt{n}\delta\right)$$
, (11)

Power of paired t-test:
$$Q\left(z_{\frac{\alpha}{2}} - \frac{\sqrt{n\delta}}{\sqrt{1+c_v}}\right)$$
, (12)

where $Q(\cdot)$ is the tail function of the standard normal distribution and $z_{\frac{\alpha}{2}}$ is the $(1-\frac{\alpha}{2})$ -th quantile of the standard normal distribution.

Remark. As shown in Theorem 10, c_v is the key quantity that determines the performance of the paired t-test as compared to the signed test. The condition $c_v \ge \pi/2 - 1$ is quite general, implying under a variety of the nuisance parameters, the sign test can outperform the paired t-test.

4 NUMERICAL EXPERIMENTS

In this section, we provide numerical evidence on the theoretical optimality and the practicality of our results. First, we evaluate our results on the synthetic data. Then, the viral infection disease dataset [3] is used to further evaluate the practicality of our theoretical findings.

³Such asymptotic scaling makes the power converge to some constant between 0 and 1, and thus making the power of the tests comparable [19].

	Sign	Wil	Pair T	DESeq2	Voom
Sign	225	156	140	66	193
Wil		267	223	95	260
Pair T			292	97	282
DESeq2				170	163
Voom					628

Table 1: Number of common discoveries across various methods: the sign test (Sign), Wilcoxon signed-rank test (Wil), paired t-test (Pair T), paired-mode DESeq2 (DESeq 2), paired-mode Voom (Voom).

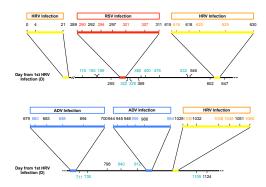


Figure 3: The description of the VID dataset [3].

4.1 Synthetic Data

Here, we compare the performance of sign test with two other popular tests for paired data: the paired t-test and the Wilcoxon signed rank test [19]. All three tests calculate the test statistic using only Y. Hence we only consider different values of $\{\mu_i\}_{i=1}^n$ and generate the samples Y according to $Y_i \sim \mathcal{N}(\delta\mu_i, \mu_i^2)$, for i=1,...,n. In all experiments, we fix the sample size $n=20, \delta=3/\sqrt{n}$, and the size of the tests $\alpha=0.05$. We repeat experiments under each parameter setting 10,000 times, and plot the corresponding 3 std confidence intervals.

Recall that, according to Theorem 10, the coefficient of variation $c_v = m_2/m_1^2$ determines the power of the paired t-test. Our first experiment examines if c_v can quantify the within-group heterogeneity level reasonably for finite values of n. In this experiment, we generated μ using the two-group model [7], where the μ_i 's are 50/50 with values $1\times$ and $10\times$ of some given magnitudes. We plot the corresponding powers for different values of the given magnitudes while the corresponding c_v 's are kept fixed to the value 0.7. As shown in Fig. 2 (a), for all 3 tests, the powers are the same for experiments under different magnitudes, implying that same value of c_v always results in the same power regardless of the values of other parameters. Hence, it is reasonable to assume that c_v can well quantify the level of heterogeneity.

In the other two experiments, we use c_v to represent heterogeneity level and plot the power of different test versus different values of c_v . In Fig. 2 (b), the value of μ is also generated by the two-group model as before. As can be seen in this figure, the sign test outperforms the paired t-test when c_v exceeds 0.58. This phenomenon is consistent with the condition in Theorem 10 where the threshold is computed as $c_v \geq \frac{\pi-2}{2} \approx 0.57$. In Fig. 2 (c), μ is generated according to the multi-group model with 5 groups of different values. As can be seen in the figure, the sign test has a better power than the Wilcoxon test when c_v exceeds 2.3. In addition, as c_v increases, the power of both the paired t-test and the Wilcoxon signed test decreases, but the later decreases much slower than the former. This is in line with our intuition since the Wilcoxon signed-rank test statistic uses the sign information and is more robust to the within-group heterogeneity than the paired t-test statistic.

4.2 The Viral Infection Dataset

In VID [3], one subject went through 6 viral infection periods in an overall time period of 1124 days. During this period, 57 RNA-Seq blood samples were collected under two conditions, healthy and

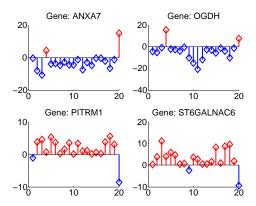


Figure 4: Count difference of the top four genes discovered only by the sign test. X-axis: sample index; Y-axis: gene expression level.

sick; see Fig. 3 for more details. The task is to find differentially-expressed genes under the two conditions. We manually pair the samples under the two conditions that are close to each other in time; and altogether acquire 20 data pairs. As shown in Fig. 1 in the beginning of this manuscript, within-pair samples are more similar to each other than within-group samples. We compare the performance of the two-sided sign test (8), the Wilcoxon test (3), the paired t-test (2). We also report the performance of two popular differential expression analysis packages paired-mode DESeq2 [14] and paired-mode Voom [11]. According to [15], they have the most promising performance among differential expression analysis tools. Prior to testing, genes with the total number of counts less than 50 or having some counts less than or equal to 1 are removed since we do not have enough observations of them. For the sign test and the paired t-test, we used the size factor normalization method as in DESeq2. For all methods, after the p-value calculation, the BH procedure [2] is used to control the false discovery rate (FDR) at the level of 0.1.

We present our results in Table 1, where the ij-th entry of this table is the number of genes discovered by both methods i and j. There are 26 genes discovered *only* by the sign test. We plot the the paired differences $\mathbf{Y} = \mathbf{X}^B - \mathbf{X}^A$ for 4 of these 26 genes with smallest p-values in Fig. 4, where the signals for differential expression are very strong. For genes ANXA7, PITRM1, ST6GLLNAC6, sample 20 has the opposite direction and a magnitude much larger than others. This prevents other methods that use the magnitude from discovering these genes. But the sign test is robust to the heterogeneity in the magnitude and discovers them.

In Table 1, the paired t-test has more discoveries than the sign test and the Wilcoxon test. A possible explanation is that the within-group heterogeneity level is not high enough due to the fact that the samples were all drawn from the same person. Voom makes more discoveries because it makes a strong assumption that a mean-variance function is shared across the genes (which also could lead to false discoveries).

References

- [1] Tarmo Äijö, Vincent Butty, Zhi Chen, Verna Salo, Subhash Tripathi, Christopher B Burge, Riitta Lahesmaa, and Harri Lähdesmäki. Methods for time series analysis of rna-seq data with application to human th17 cell differentiation. *Bioinformatics*, 30(12):i113–i120, 2014.
- [2] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the royal statistical society. Series B (Methodological)*, pages 289–300, 1995.
- [3] Rui Chen, George I Mias, Jennifer Li-Pook-Than, Lihua Jiang, Hugo YK Lam, Rong Chen, Elana Miriami, Konrad J Karczewski, Manoj Hariharan, Frederick E Dewey, et al. Personal omics profiling reveals dynamic molecular and medical phenotypes. *Cell*, 148(6):1293–1307, 2012.
- [4] Lisa M Chung, John P Ferguson, Wei Zheng, Feng Qian, Vincent Bruno, Ruth R Montgomery, and Hongyu Zhao. Differential expression analysis for paired rna-seq data. *BMC bioinformatics*, 14(1):110, 2013.
- [5] NAC Cressie and HJ Whitford. How to use the two sample t-test. *Biometrical Journal*, 28(2):131–148, 1986.
- [6] AC Davison. Treatment effect heterogeneity in paired data. Biometrika, pages 463-474, 1992.
- [7] Bradley Efron. Microarrays, empirical bayes and the two-groups model. *Statistical science*, pages 1–22, 2008.
- [8] Brian S Everitt. The Cambridge dictionary of statistics. Cambridge University Press, 2006.
- [9] Patricia M Grambsch. Simple robust tests for scale differences in paired data. *Biometrika*, pages 359–372, 1994.
- [10] Thomas J Hardcastle and Krystyna A Kelly. Empirical bayesian analysis of paired high-throughput sequencing data with a beta-binomial distribution. *BMC bioinformatics*, 14(1):135, 2013.
- [11] Charity W Law, Yunshun Chen, Wei Shi, and Gordon K Smyth. Voom: precision weights unlock linear model analysis tools for rna-seq read counts. *Genome biology*, 15(2):R29, 2014.
- [12] Erich L Lehmann and Joseph P Romano. *Testing statistical hypotheses*. Springer Science & Business Media, 2006.
- [13] Erich Leo Lehmann and Howard JM D'Abrera. *Nonparametrics: statistical methods based on ranks*. Springer New York, 2006.
- [14] Michael I Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome biology*, 15(12):550, 2014.
- [15] Harold J Pimentel, Nicolas Bray, Suzette Puente, Páll Melsted, and Lior Pachter. Differential analysis of rna-seq incorporating quantification uncertainty. bioRxiv, page 058164, 2016.
- [16] Feng Qian, Lisa Chung, Wei Zheng, Vincent Bruno, Roger P Alexander, Zhong Wang, Xiaomei Wang, Sebastian Kurscheid, Hongyu Zhao, Erol Fikrig, et al. Identification of genes critical for resistance to infection by west nile virus using rna-seq analysis. *Viruses*, 5(7):1664–1681, 2013.
- [17] Mark D Robinson, Davis J McCarthy, and Gordon K Smyth. edger: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140, 2010.
- [18] Charlotte Soneson and Mauro Delorenzi. A comparison of methods for differential expression analysis of rna-seq data. *BMC bioinformatics*, 14(1):91, 2013.
- [19] Aad W Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.

Supplemental Materials

1 The details of the visualization

For any two samples $\mathbf{X}_i, \mathbf{X}_j$, we can calculate the histogram of the k values $\log |\mathbf{X}_i - \mathbf{X}_j| \in \mathbb{R}^k$. Concentration in small values indicates $\mathbf{X}_i, \mathbf{X}_j$ are similar. The histogram is averaged over all sample pairs for the within-pair difference and over all size-2 subsets of samples from the same group for the within-group difference.

2 Proof of Fact 2

Proof. (Proof of Fact 2) It suffices to show that

$$\sup_{\gamma} \mathbb{E}_0[\phi] = \sup_{\gamma} \mathbb{E}_0[\psi], \quad and \quad \inf_{\gamma} \mathbb{E}_1[\phi] = \inf_{\gamma} \mathbb{E}_1[\psi]. \tag{13}$$

Now we prove the first equation in (13). Let $\{Z_i^A, Z_i^B\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathcal{N}(0,1)$. Then by (4) we can write

$$\mathbf{X}^A = \boldsymbol{\nu} + \boldsymbol{\mu} \circ \sqrt{\boldsymbol{\rho}} \circ \mathbf{Z}^A, \ \mathbf{X}^B = \boldsymbol{\nu} + \delta \boldsymbol{\mu} + \boldsymbol{\mu} \circ \sqrt{1 - \boldsymbol{\rho}} \circ \mathbf{Z}^B.$$

Similarly, we can write the transformed data as

$$\mathbf{a} + \mathbf{b} \circ \mathbf{X}^A = \nu' + \mu' \circ \sqrt{\rho} \circ \mathbf{Z}^A$$
$$\mathbf{a} + \mathbf{b} \circ \mathbf{X}^B = \nu' + \delta \mu' + \mu' \circ \sqrt{1 - \rho} \circ \mathbf{Z}^B,$$

for $\nu' \triangleq \mathbf{a} + \mathbf{b} \circ \nu$ and $\mu' \triangleq \mathbf{b} \circ \mu$.

The key idea for proving the first equation in (13) is that $\{\nu, \mu\}$ and $\{\nu', \mu'\}$ actually consist of the same parameter space $\mathbb{R}^n \times \mathbb{R}^n_{>0}$. To be more exact, notice that under the null hypothesis $\delta = 0$. We have

$$\sup_{\boldsymbol{\nu},\boldsymbol{\mu},\boldsymbol{\rho}} \mathbb{E}_0[\psi(\mathbf{X}^A,\mathbf{X}^B)]$$

$$= \sup_{\boldsymbol{\nu},\boldsymbol{\mu},\boldsymbol{\rho}} \mathbb{E}[\phi(\boldsymbol{\nu}' + \boldsymbol{\mu}' \circ \sqrt{\boldsymbol{\rho}} \circ \mathbf{Z}^A, \boldsymbol{\nu}' + \boldsymbol{\mu}' \circ \sqrt{1-\boldsymbol{\rho}} \circ \mathbf{Z}^B)]$$

$$= \sup_{\boldsymbol{\nu}',\boldsymbol{\mu}',\boldsymbol{\rho}} \mathbb{E}[\phi(\boldsymbol{\nu}' + \boldsymbol{\mu}' \circ \sqrt{\boldsymbol{\rho}} \circ \mathbf{Z}^A, \boldsymbol{\nu}' + \boldsymbol{\mu}' \circ \sqrt{1-\boldsymbol{\rho}} \circ \mathbf{Z}^B)]$$

$$= \sup_{\boldsymbol{\nu},\boldsymbol{\mu},\boldsymbol{\rho}} \mathbb{E}[\phi(\boldsymbol{\nu} + \boldsymbol{\mu} \circ \sqrt{\boldsymbol{\rho}} \circ \mathbf{Z}^A, \boldsymbol{\nu} + \boldsymbol{\mu} \circ \sqrt{1-\boldsymbol{\rho}} \circ \mathbf{Z}^B)]$$

$$= \sup_{\boldsymbol{\nu},\boldsymbol{\mu},\boldsymbol{\rho}} \mathbb{E}[\phi(\mathbf{X}^A,\mathbf{X}^B)],$$

where in the second equation we use the fact that $\{\nu, \mu\}$ and $\{\nu', \mu'\}$ consist of the same parameter space. The second equation of (13) can be shown similarly. Then $\phi \doteq \psi$.

3 Proof of Theorems

3.1 Proof of Theorem 8

Proof. (Proof of Theorem 8) First, for this problem, we should restrict ourselves to symmetrical tests, i.e. any ϕ such that $\phi(\mathbf{X}^A, \mathbf{X}^B) = \phi(-\mathbf{X}^A, -\mathbf{X}^B)$. This is because for any $(\mathbf{X}^A, \mathbf{X}^B)$, the distribution $(-\mathbf{X}^A, -\mathbf{X}^B)$ is also valid for the maximin problem. Second, according to Lemma 4, it suffices to consider the symmetrical tests using only \mathbf{Y} .

Now consider any α -level symmetrical test that uses only \mathbf{Y} , namely $\phi(\mathbf{Y})$. For any set of nuisance parameters γ , let $f_{0,\gamma}(\cdot)$, $f_{1,\gamma}(\cdot)$ be the density function under the null and the alternative hypothesis respectively. Define

$$\mathcal{Y}^+ = \{ \mathbf{y} : f_{1,\gamma}(\mathbf{y}) > f_{1,\gamma}(-\mathbf{y}) \},$$

$$\mathcal{Y}^- = \{ \mathbf{y} : f_{1,\gamma}(\mathbf{y}) < f_{1,\gamma}(-\mathbf{y}) \}.$$

Then by explicitly writing out the density function, it is not hard to see that for any $\mathbf{y} \in \mathcal{Y}^+$, we have $-\mathbf{y} \in \mathcal{Y}^-$. As the null distribution is symmetrical around 0, we have

$$\sup_{\gamma} \mathbb{E}_{0}[\phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y}\in\mathcal{Y}^{+}\}}] = \sup_{\gamma} \mathbb{E}_{0}[\phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y}\in\mathcal{Y}^{-}\}}] \le \frac{\alpha}{2}$$
(14)

Next consider the power of ϕ . As $\phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y}\in\mathcal{Y}^+\}}$ is a $\frac{\alpha}{2}$ -level test, by the optimality of the one-sided sign test, $\forall \epsilon > 0, \exists \gamma^*, \text{ s.t.,}$

$$\mathbb{E}_{1,\gamma^*}[\phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y}\in\mathcal{Y}^+\}}] - \epsilon < \inf_{\gamma} \mathbb{E}_1[\phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y}\in\mathcal{Y}^+\}}] \le \inf_{\gamma} \mathbb{E}_1[\tilde{\phi}^S(\mathbf{Y})]. \tag{15}$$

Also we have

$$\mathbb{E}_{1,\gamma^*}[\phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y}\in\mathcal{Y}^-\}}] = \int_{\mathcal{Y}^-} \phi(\mathbf{y})f_{1,\gamma^*}(\mathbf{y})d\mathbf{y} = \int_{\mathcal{Y}^-} \phi(\mathbf{y})f_{0,\gamma^*}(\mathbf{y})\frac{f_{1,\gamma^*}(\mathbf{y})}{f_{0,\gamma^*}(\mathbf{y})}d\mathbf{y}
\leq \sup_{\mathbf{y}\in\mathcal{Y}^-} \left(\frac{f_{1,\gamma^*}(\mathbf{y})}{f_{0,\gamma^*}(\mathbf{y})}\right) \sup_{\gamma} \mathbb{E}_0[\phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y}\in\mathcal{Y}^-\}}] \leq \frac{\alpha}{2} \exp(-\frac{\delta^2 n}{2}).$$
(16)

In the last inequality of (16), the first term is due to the fact that if $\mathbf{y} \in \mathcal{Y}^-$, we have $\frac{f_{1,\gamma^*}(\mathbf{y})}{f_{1,\gamma^*}(-\mathbf{y})} \leq 1$, giving $\delta \sum_{i=1}^i \frac{y_i}{\mu_i} \leq 0$, which further gives $\forall \mathbf{y} \in \mathcal{Y}^-$,

$$\frac{f_{1,\gamma^*}(\mathbf{y})}{f_{0,\gamma^*}(\mathbf{y})} = \exp(-\frac{n\delta^2}{2} + \delta \sum_{i=1}^i \frac{y_i}{\mu_i}) \le \exp(-\frac{n\delta^2}{2}).$$

The second term is due to (14). Finally, combining (15) and (16), we reach that the power of ϕ

$$\inf_{\gamma} \mathbb{E}_1[\phi(\mathbf{Y})] \leq \mathbb{E}_{1,\gamma^*}[\phi(\mathbf{Y})] = \mathbb{E}_{1,\gamma^*}[\phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y} \in \mathcal{Y}^+\}} + \phi(\mathbf{Y})\mathbb{I}_{\{\mathbf{Y} \in \mathcal{Y}^-\}}]$$

$$\leq \inf_{\gamma} \mathbb{E}_1[\tilde{\phi}^S(\mathbf{Y})] + \frac{\alpha}{2} \exp(-\frac{n\delta^2}{2}) + \epsilon \leq \inf_{\gamma} \mathbb{E}_1[\phi^S(\mathbf{Y})] + \frac{\alpha}{2} \exp(-\frac{n\delta^2}{2}) + \epsilon,$$

where we recall that ϕ^S is the two-sided sign test as defined in (10). Finally, let $\epsilon \to 0$ to complete the proof.

3.2 Proof of Theorem 10

Proof. (Proof of Theorem 10) Let us start by defining the following notations: we will use $\stackrel{d}{\to}$ and $\stackrel{p}{\to}$ to denote convergence in distribution and convergence in probability, respectively. Next, we first compute the power of the two-sided sign test. Let us consider the sufficient statistics $W = \sum_i \mathbb{I}_{\{Y_i>0\}}$ for the sign test. Clearly, under the null distribution, $\frac{W-0.5n}{\sqrt{0.25n}} \stackrel{d}{\to} \mathcal{N}(0,1)$ due to the standard central limit theorem. Thus the two-sided sign test asymptotically rejects when $|\frac{W-0.5n}{\sqrt{0.25n}}| \geq z_{\frac{\alpha}{2}}$. Furthermore, under the alternative hypothesis, by Taylor's expansion, we have

$$\theta = \mathbb{P}(Y_i \ge 0) = Q(-\delta) = 0.5 + \frac{1}{\sqrt{2\pi}}\delta + O(\delta^2).$$

Without loss of generality assume that $\delta > 0$. Then the power of the test can be computed as

$$\mathbb{P}(|\frac{W-0.5n}{\sqrt{0.25n}}| \ge z_{\frac{\alpha}{2}}) \approx \mathbb{P}(\frac{W-0.5n}{\sqrt{0.25n}} \ge z_{\frac{\alpha}{2}}),$$

where by using the approximately equality we neglect the lower tail, a small quantity that decreases exponentially with α . This power can be further computed as

$$Q\left(z_{\frac{\alpha}{2}} - \sqrt{\frac{2}{\pi}}\sqrt{n}\delta\right). \tag{17}$$

On the other hand, for the paired t-test, we can asymptotically write the numerator and the denominator of the T-statistics as

$$\sqrt{n}\bar{Y} \stackrel{d}{\to} \mathcal{N}\left(\sqrt{n}\delta m_1, m_1^2 + m_2\right)$$
$$\frac{1}{n-1} \sum_i (Y_i - \bar{Y})^2 \stackrel{p}{\to} m_1^2 + m_2,$$

where we recall that $\sqrt{n}\delta$ is some given constant by appropriate scaling of δ . Therefore, under the null distribution, we have $T \stackrel{d}{\to} \mathcal{N}(0,1)$. Consequently, the paired t-test rejects when $|T| \geq z_{\frac{\alpha}{2}}$. Again by assuming $\delta > 0$ under the alternative hypothesis, the power of the test can be written as

$$\mathbb{P}(|T| \ge z_{\frac{\alpha}{2}}) \approx \mathbb{P}(T \ge z_{\frac{\alpha}{2}}) \to Q\left(z_{\frac{\alpha}{2}} - \frac{\sqrt{n\delta}}{\sqrt{1 + c_v}}\right). \tag{18}$$

Combining (17) and (18) will complete the proof.

4 Proofs of Lemmas

4.1 Proof of Lemma 4

Proof. Since $\mathbf{X}^B = \mathbf{X}^A + \mathbf{Y}$, ϕ can be equivalently represented as $\tilde{\phi}(\mathbf{Y}, \mathbf{X}^A) = \phi(\mathbf{X}^A, \mathbf{Y} + \mathbf{X}^A)$. Let $\psi(\mathbf{Y}) = \tilde{\phi}(\mathbf{Y}, 0)$. If ϕ is symmetric, then

$$\psi(\mathbf{Y}) = \tilde{\phi}(\mathbf{Y}, 0) = \phi(0, \mathbf{X}^B - \mathbf{X}^A)$$

$$= \phi(0, \mathbf{X}^A - \mathbf{X}^B) = \tilde{\phi}(-\mathbf{Y}, 0) = \psi(-\mathbf{Y}),$$
(19)

giving that ψ is also symmetric.

We next show that $\tilde{\phi} \leq \psi$. First notice that for each $i, (X_i^A, Y_i)$ follows a joint Gaussian distribution:

$$\left[\begin{array}{c} X_i^A \\ Y_i \end{array}\right] \sim \mathcal{N} \left(\begin{array}{c} \nu_i \\ \delta \mu_i \end{array}\right], \quad \left[\begin{array}{cc} \rho_i \mu_i^2 & -\rho_i \mu_i^2 \\ -\rho_i \mu_i^2 & \mu_i^2 \end{array}\right] \right),$$

and the samples are independent across all indices $i=1,\ldots,n$. Then $X_i^A|Y_i=y_i\sim\mathcal{N}(\nu_i-\rho_iy_i+\rho_i\delta\mu_i,\rho_i(1-\rho_i)\mu_i^2)$ and $\forall~i\neq j,X_i^A\perp X_j^A|\mathbf{Y}$, where we note that $\left[X_i^A|Y_i=y_i,\nu_i=\rho_i=0\right]=0$. For the worst-case size,

$$\sup_{\boldsymbol{\mu},\boldsymbol{\nu},\boldsymbol{\rho}} \mathbb{E}_{0}[\boldsymbol{\phi}] = \sup_{\boldsymbol{\mu},\boldsymbol{\nu},\boldsymbol{\rho}} \mathbb{E}_{0,\mathbf{Y}}[\mathbb{E}_{0,\mathbf{X}^{A}|\mathbf{Y}}[\tilde{\boldsymbol{\phi}}(\mathbf{Y},\mathbf{X}^{A})|\mathbf{Y}]]$$

$$\geq \sup_{\boldsymbol{\mu},\boldsymbol{\nu}=\boldsymbol{\rho}=\mathbf{0}} \mathbb{E}_{0,\mathbf{Y}}[\mathbb{E}_{0,\mathbf{X}^{A}|\mathbf{Y}}[\tilde{\boldsymbol{\phi}}(\mathbf{Y},\mathbf{X}^{A})|\mathbf{Y}]]$$

$$= \sup_{\boldsymbol{\mu},\boldsymbol{\nu}=\boldsymbol{\rho}=\mathbf{0}} \mathbb{E}_{0,\mathbf{Y}}[\mathbb{E}_{0,\mathbf{X}^{A}|\mathbf{Y}}[\tilde{\boldsymbol{\phi}}(\mathbf{Y},0)|\mathbf{Y}]]$$

$$= \sup_{\boldsymbol{\mu},\boldsymbol{\nu}=\boldsymbol{\rho}=\mathbf{0}} \mathbb{E}_{0,\mathbf{Y}}[\mathbb{E}_{0,\mathbf{X}^{A}|\mathbf{Y}}[\boldsymbol{\psi}(\mathbf{Y})|\mathbf{Y}]] = \sup_{\boldsymbol{\mu},\boldsymbol{\nu},\boldsymbol{\rho}} \mathbb{E}_{0}[\boldsymbol{\psi}],$$

$$(20)$$

where the last equality is because $\psi(\mathbf{Y})$ does not depend on ν, ρ . Similarly, the power of the two tests can be related by

$$\inf_{\boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{\rho}, s_{\delta}} \mathbb{E}_{1}[\phi] = \inf_{\boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{\rho}, s_{\delta}} \mathbb{E}_{1, \mathbf{Y}}[\mathbb{E}_{1, \mathbf{X}^{A} | \mathbf{Y}}[\tilde{\phi}(\mathbf{Y}, \mathbf{X}^{A}) | \mathbf{Y}]]$$

$$\leq \inf_{\boldsymbol{\mu}, \boldsymbol{\nu} = \boldsymbol{\rho} = \mathbf{0}, s_{\delta}} \mathbb{E}_{1, \mathbf{Y}}[\mathbb{E}_{1, \mathbf{X}^{A} | \mathbf{Y}}[\tilde{\phi}(\mathbf{Y}, \mathbf{X}^{A}) | \mathbf{Y}]]$$

$$= \inf_{\boldsymbol{\mu}, \boldsymbol{\nu} = \boldsymbol{\rho} = \mathbf{0}, s_{\delta}} \mathbb{E}_{1, \mathbf{Y}}[\mathbb{E}_{1, \mathbf{X}^{A} | \mathbf{Y}}[\tilde{\phi}(\mathbf{Y}, 0) | \mathbf{Y}]]$$

$$= \inf_{\boldsymbol{\mu}, \boldsymbol{\nu} = \boldsymbol{\rho} = \mathbf{0}, s_{\delta}} \mathbb{E}_{1, \mathbf{Y}}[\mathbb{E}_{1, \mathbf{X}^{A} | \mathbf{Y}}[\psi(\mathbf{Y}) | \mathbf{Y}]] = \inf_{\boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{\rho}, s_{\delta}} \mathbb{E}_{1}[\psi],$$
(21)

where $s_{\delta} = 1$ for the one-sided case and $s_{\delta} \in \{1, -1\}$ for the two-sided case. Combining (20) and (21) we conclude $\phi \leq \psi$.

4.2 Proof of Lemma 5

Proof. Let $Z_1,\cdots,Z_n \overset{i.i.d.}{\sim} \mathcal{N}(0,1)$, and let P_0,P_1 denote the n-D product probability measure of $\mathbf{Z}=(Z_1,\cdots,Z_n)$ and $\mathbf{Z}+\delta=(Z_1+\delta,\cdots,Z_n+\delta)$ respectively. In addition, denote the probability of orthant \mathbf{o} under measure P_0 and P_1 by $P_0^{\mathbf{o}}$ and $P_1^{\mathbf{o}}$ respectively. Notice that $\sum_{\mathbf{o}\in\mathcal{O}}P_0^{\mathbf{o}}=\sum_{\mathbf{o}\in\mathcal{O}}P_1^{\mathbf{o}}=1$, and for any $\mathbf{o}\in\mathcal{O},P_0^{\mathbf{o}}=2^{-n}$.

Depending only on the sign of the data, the sign test has constant value on each orthant. In fact, it is not hard to see that the level- α one-sided sign test $\phi^S(\mathbf{Y})$ (7) maximizes its power by assigning 1 to orthants with larger values of P_1^{o} until the corresponding size reaches α . We next show that such procedure has a power that upper bounds the power of any level- α test in S, which proves the lemma.

We first define some useful quantities. Now consider any $\phi \in \mathcal{S}(\omega)$. Define the discretized space of μ to be $D = \{ \mu = (\mu_1, \dots, \mu_n) : \forall i, \mu_i = (1 + \omega)^{d_i}, d_i \in \mathbb{Z} \}$. Consider any $\mu \in D$ with $\mathbf{d} = (d_1, \dots, d_n)$. For any n-D box $I_{\mathbf{b}}^{\mathbf{o}}$, the element-wise multiplication by μ maps it to $I_{\mathbf{b}+\mathbf{d}}^{\mathbf{o}}$. Define $f(\mathbf{z})$ to be the density function of \mathbf{Z} . Let $\alpha(\mu) = \mathbb{E}_0[\phi(\mathbf{Y})]$ and $\beta(\mu) = \mathbb{E}_1[\phi(\mathbf{Y})]$ be the size and power of ϕ when the nuisance parameters have value μ . We have

$$\alpha(\boldsymbol{\mu}) = \mathbb{E}[\phi(\boldsymbol{\mu} \circ \mathbf{Z})] = \int_{\mathbb{R}^n} \phi(\boldsymbol{\mu} \circ \mathbf{z}) f(\mathbf{z}) d\mathbf{z} = \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \phi_{\mathbf{d} + \mathbf{b}}^{\mathbf{o}} P_0(I_{\mathbf{b}}^{\mathbf{o}}),$$

$$\beta(\boldsymbol{\mu}) = \mathbb{E}[\phi(\boldsymbol{\mu} \circ (\mathbf{Z} + \delta))] = \int_{\mathbb{R}^n} \phi(\boldsymbol{\mu} \circ \mathbf{z}) f(\mathbf{z} - \delta) d\mathbf{z} = \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \phi_{\mathbf{d} + \mathbf{b}}^{\mathbf{o}} P_1(I_{\mathbf{b}}^{\mathbf{o}}).$$

Moreover, we can write the size and power corresponding to orthant o as

$$\alpha^{\mathbf{o}}(\boldsymbol{\mu}) = \sum_{-\infty < \mathbf{b} < \infty} \phi_{\mathbf{d}+\mathbf{b}}^{\mathbf{o}} P_0(I_{\mathbf{b}}^{\mathbf{o}}), \qquad \beta^{\mathbf{o}}(\boldsymbol{\mu}) = \sum_{-\infty < \mathbf{b} < \infty} \phi_{\mathbf{d}+\mathbf{b}}^{\mathbf{o}} P_1(I_{\mathbf{b}}^{\mathbf{o}}), \tag{22}$$

where we note that $\alpha(\mu) = \sum_{\mathbf{o} \in \mathcal{O}} \alpha^{\mathbf{o}}(\mu)$ and $\beta(\mu) = \sum_{\mathbf{o} \in \mathcal{O}} \beta^{\mathbf{o}}(\mu)$. Next, for each orthant \mathbf{o} and any positive integer m, define the m-th approximation of the probability measure $P_0^{\mathbf{o}}$, $P_1^{\mathbf{o}}$, the size $\alpha^{\mathbf{o}}(\boldsymbol{\mu})$ and the power $\beta^{\mathbf{o}}(\boldsymbol{\mu})$ to be

$$P_{0,m}^{\mathbf{o}} = \sum_{-m \le \mathbf{b} \le m} P_0(I_{\mathbf{b}}^{\mathbf{o}}), \qquad P_{1,m}^{\mathbf{o}} = \sum_{-m \le \mathbf{b} \le m} P_1(I_{\mathbf{b}}^{\mathbf{o}})$$
(23)

$$P_{0,m}^{\mathbf{o}} = \sum_{-m \le \mathbf{b} \le m} P_0(I_{\mathbf{b}}^{\mathbf{o}}), \qquad P_{1,m}^{\mathbf{o}} = \sum_{-m \le \mathbf{b} \le m} P_1(I_{\mathbf{b}}^{\mathbf{o}}) \qquad (23)$$

$$\alpha_m^{\mathbf{o}}(\boldsymbol{\mu}) = \sum_{-m \le \mathbf{b} \le m} \phi_{\mathbf{d}+\mathbf{b}}^{\mathbf{o}} P_0(I_{\mathbf{b}}^{\mathbf{o}}), \qquad \beta_m^{\mathbf{o}}(\boldsymbol{\mu}) = \sum_{-m \le \mathbf{b} \le m} \phi_{\mathbf{d}+\mathbf{b}}^{\mathbf{o}} P_1(I_{\mathbf{b}}^{\mathbf{o}}), \qquad (24)$$

where the summation is over all indices $\mathbf{b} \in \mathbb{Z}^n$ with elements all between -m and m. It is not hard to see that $P_{0,m}^{\mathbf{o}} \uparrow P_0^{\mathbf{o}}$, $P_{1,m}^{\mathbf{o}} \uparrow P_1^{\mathbf{o}}$, $\alpha_m^{\mathbf{o}}(\boldsymbol{\mu}) \uparrow \alpha^{\mathbf{o}}(\boldsymbol{\mu})$, and $\beta_m^{\mathbf{o}}(\boldsymbol{\mu}) \uparrow \beta^{\mathbf{o}}(\boldsymbol{\mu})$.

We next show the key step of the proof: for every m,

$$\inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_{1,m}^{\mathbf{o}}} \le \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_{0,m}^{\mathbf{o}}}.$$
 (25)

Recall that $\mu = (\mu_1, \dots, \mu_n)$, for $\mu_i = (1 + \omega)^{d_i}$. We can write,

$$\begin{split} \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_{m}^{\mathbf{o}}(\boldsymbol{\mu})}{P_{0,m}^{\mathbf{o}}} - \inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_{m}^{\mathbf{o}}(\boldsymbol{\mu})}{P_{1,m}^{\mathbf{o}}} &= \lim_{l \to \infty} \left[\sup_{-l \le \mathbf{d} \le l} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_{m}^{\mathbf{o}}(\boldsymbol{\mu})}{P_{0,m}^{\mathbf{o}}} - \inf_{-l \le \mathbf{d} \le l} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_{m}^{\mathbf{o}}(\boldsymbol{\mu})}{P_{1,m}^{\mathbf{o}}} \right] \\ &= \lim_{l \to \infty} \left[\sup_{-l \le \mathbf{d} \le l} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\sum_{-m \le \mathbf{b} \le m} \phi_{\mathbf{d} + \mathbf{b}}^{\mathbf{o}} P_{0}(I_{\mathbf{b}}^{\mathbf{o}})}{P_{0,m}^{\mathbf{o}}} - \inf_{-l \le \mathbf{d} \le l} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\sum_{-m \le \mathbf{b} \le m} \phi_{\mathbf{d} + \mathbf{b}}^{\mathbf{o}} P_{1}(I_{\mathbf{b}}^{\mathbf{o}})}{P_{1,m}^{\mathbf{o}}} \right] \\ &\ge \lim_{l \to \infty} \frac{1}{(2l+1)^{n}} \left[\sum_{-l \le \mathbf{d} \le l} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\sum_{-m \le \mathbf{b} \le m} \phi_{\mathbf{d} + \mathbf{b}}^{\mathbf{o}} P_{0}(I_{\mathbf{b}}^{\mathbf{o}})}{P_{0,m}^{\mathbf{o}}} - \sum_{-l \le \mathbf{d} \le l} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\sum_{-m \le \mathbf{b} \le m} \phi_{\mathbf{d} + \mathbf{b}}^{\mathbf{o}} P_{1}(I_{\mathbf{b}}^{\mathbf{o}})}{P_{1,m}^{\mathbf{o}}} \right] \\ &\ge \lim_{l \to \infty} \mathcal{O}(\frac{m}{l}) = 0, \end{split}$$

which gives (25). The first equality is because $D = \lim_{l \to \infty} {\{ \mu : \mu_i = (1 + \omega)^{d_i}, -l \leq \mathbf{d} \leq l \}}$. The second equality is because of (24). The first inequality is obtained by replacing sup and inf by averaging. The second inequality is a little tricky. Inside the square brackets, both the first and the second big term can be rearranged by $\phi_{\mathbf{j}}^{\mathbf{o}}$ to have the form $\sum_{-l-m\leq\mathbf{j}\leq l+m}c_{\mathbf{j}}\phi_{\mathbf{j}}^{\mathbf{o}}$. A careful inspection reveals that $c_{\mathbf{j}} = 1$ for any $-l + m \leq \mathbf{j} \leq l - m$. Indeed, e.g. for the first big term, as long as $-l + m \leq \mathbf{j} \leq l - m$, $\phi_{\mathbf{j}}^{\mathbf{o}}$ will be multiplied by $P_0(I_{\mathbf{b}}^{\mathbf{o}})$ once for each $-m \leq \mathbf{b} \leq m$, and the sum of those coefficients, $\sum_{-m < \mathbf{b} < m} P_0(I_{\mathbf{b}}^{\mathbf{o}})$, exactly equals the denominator $P_{0,m}^{\mathbf{o}}$ according to (23). Thus,

for all $-l+m \leq \mathbf{j} \leq l-m$, $\phi^{\mathbf{o}}_{\mathbf{j}}$ will have the same coefficients for both the first big term and the second big term, resulting altogether $(2l-2m+1)^n$ terms canceling each other. As a result, there remain $(2l+2m+1)^n-(2l-2m+1)^n=O(ml^{n-1})$ terms, whose corresponding summation can be upper bounded by $O(\frac{m}{l})$ because of the multiplicative factor $\frac{1}{(2l+1)^n}$ outside the square brackets. Finally, $O(\frac{m}{l})$ vanishes as $l\to\infty$.

Next we prove the limiting case of (25):

$$\inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} \le \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}$$
 (26)

by showing

$$\lim_{m \to \infty} \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_{0,m}^{\mathbf{o}}} = \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}$$
(27)

$$\lim_{m \to \infty} \inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_{1,m}^{\mathbf{o}}} = \inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}}.$$
 (28)

Proving (27). Because $P_{0,m}^{\mathbf{o}} \uparrow P_0^{\mathbf{o}}$ not depending on μ , $P_0^{\mathbf{o}}$ bounded away from 0, and $\alpha_m^{\mathbf{o}}(\mu) \leq 1$, we have that for any $\epsilon_0 > 0$, there exists m_0 such that $\forall m > m_0, \mu \in D$,

$$\sum_{\mathbf{o}\in\mathcal{O}}\frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}\leq \sum_{\mathbf{o}\in\mathcal{O}}\frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_{0,m}^{\mathbf{o}}}<\sum_{\mathbf{o}\in\mathcal{O}}\frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}+\epsilon_0,$$

giving

$$\lim_{m \to \infty} \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_{0,m}^{\mathbf{o}}} = \lim_{m \to \infty} \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}.$$
 (29)

Next, as $\alpha_m^{\mathbf{o}}(\boldsymbol{\mu}) \uparrow \alpha^{\mathbf{o}}(\boldsymbol{\mu})$, we have $\sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} \uparrow \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}$, giving

$$\lim_{m \to \infty} \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}_m(\boldsymbol{\mu})}{P^{\mathbf{o}}_0} \leq \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P^{\mathbf{o}}_0}.$$

On the other hand, for any $\epsilon_1 > 0$, let

$$E = \{ \boldsymbol{\mu} : \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} > \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} - \epsilon_1 \}$$

$$E_m = \{ \boldsymbol{\mu} : \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} > \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} - \epsilon_1 \}.$$

Notice that $\sum_{\mathbf{o}\in\mathcal{O}}\frac{\alpha_{m}^{\mathbf{o}}(\boldsymbol{\mu})}{P_{0}^{\mathbf{o}}}\uparrow\sum_{\mathbf{o}\in\mathcal{O}}\frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_{0}^{\mathbf{o}}}$, we have $E_{m}\subset E_{m+1}$, \forall n, and $E=\bigcup_{m}E_{m}$. Therefore, for any $\epsilon_{1}>0$, \exists m_{1} such that for all $m>m_{1}$, $E_{m}\neq\emptyset$, which implies

$$\lim_{m \to \infty} \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} \ge \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}.$$

Therefore,

$$\lim_{m \to \infty} \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} = \sup_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}.$$
 (30)

Combining (29) and (30) to have (27).

Proving (28). Due to the same reason of (29) we have

$$\lim_{m \to \infty} \inf_{\mu \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\mu)}{P_{1,m}^{\mathbf{o}}} = \lim_{m \to \infty} \inf_{\mu \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\mu)}{P_1^{\mathbf{o}}}.$$
 (31)

Since $\beta_m(\mu) \uparrow \beta(\mu)$, we have $\sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\mu)}{P_1^{\mathbf{o}}} \uparrow \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\mu)}{P_1^{\mathbf{o}}}$, giving

$$\lim_{m \to \infty} \inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}} \le \inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}}.$$
 (32)

Define $c_m = \sum_{\mathbf{o} \in \mathcal{O}} \sum_{\mathbf{b} < -m \text{ or } \mathbf{b} > m} \frac{P_1(I_{\mathbf{o}}^{\mathbf{o}})}{P_1^{\mathbf{o}}}$. Then $c_m \downarrow 0$ and $\forall \mu \in D$, $\sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\mu)}{P_1^{\mathbf{o}}} + c_m \downarrow \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\mu)}{P_1^{\mathbf{o}}}$. Furthermore, for any $\epsilon_1 > 0$, let

$$E = \{ \boldsymbol{\mu} : \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}} < \lim_{m \to \infty} \inf_{\boldsymbol{\mu} \in D} \left(\sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}} + c_m \right) + \epsilon_1 \}$$

$$E_m = \{ \mu : \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}} + c_m < \lim_{m \to \infty} \inf_{\boldsymbol{\mu} \in D} \left(\sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}} + c_m \right) + \epsilon_1 \}.$$

Since $\sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_m^{\mathbf{o}}(\mu)}{P_1^{\mathbf{o}}} + c_m \downarrow \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\mu)}{P_1^{\mathbf{o}}}$, we have $E_m \subset E_{m+1}$, $\forall m$, and $E = \bigcup_m E_m$. For any $\epsilon_1 > 0$, since for any $m, E_m \neq \emptyset$, we have $E \neq \emptyset$ and hence

$$\inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_{\mathbf{i}}^{\mathbf{o}}} \le \lim_{m \to \infty} \inf_{\boldsymbol{\mu} \in D} \left(\sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_{m}^{\mathbf{o}}(\boldsymbol{\mu})}{P_{\mathbf{i}}^{\mathbf{o}}} + c_{m} \right) = \lim_{m \to \infty} \inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_{m}^{\mathbf{o}}(\boldsymbol{\mu})}{P_{\mathbf{i}}^{\mathbf{o}}}.$$
 (33)

Combining (32) and (33) we obtain

$$\inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_{\mathbf{o}}^{\mathbf{o}}} = \lim_{m \to \infty} \inf_{\boldsymbol{\mu} \in D} \sum_{\mathbf{o} \in \mathcal{O}} \frac{\beta_{m}^{\mathbf{o}}(\boldsymbol{\mu})}{P_{\mathbf{o}}^{\mathbf{o}}}.$$
 (34)

Combining (31) and (34) to have (28). Then finally we proved (26).

Finally we show the power of the sign test upper bounds that of any test in $S(\omega)$ using (26). (26) further gives

$$\inf_{\boldsymbol{\mu}>0} \sum_{\mathbf{o}\in\mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}} \le \inf_{\boldsymbol{\mu}\in D} \sum_{\mathbf{o}\in\mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}} \le \sup_{\boldsymbol{\mu}\in D} \sum_{\mathbf{o}\in\mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}} \le \sup_{\boldsymbol{\mu}>0} \sum_{\mathbf{o}\in\mathcal{O}} \frac{\alpha^{\mathbf{o}}(\boldsymbol{\mu})}{P_0^{\mathbf{o}}}$$
(35)

Recall that for every orthant o, $P_0^o = 2^{-n}$. Multiplying both sides of (35) by 2^{-n} we have

$$\inf_{\boldsymbol{\mu}>0} 2^{-n} \sum_{\mathbf{o}\in\mathcal{O}} \frac{\beta^{\mathbf{o}}(\boldsymbol{\mu})}{P_1^{\mathbf{o}}} \le \sup_{\boldsymbol{\mu}>0} \sum_{\mathbf{o}\in\mathcal{O}} \alpha^{\mathbf{o}}(\boldsymbol{\mu}) = \sup_{\boldsymbol{\mu}>0} \alpha(\boldsymbol{\mu}) \le \alpha, \tag{36}$$

where α is the size of the test. For any $\epsilon>0$, there exists an μ' such that $2^{-n}\sum_{\mathbf{o}\in\mathcal{O}}\frac{\beta^{\mathbf{o}}(\mu')}{P_1^{\alpha}}\leq \alpha+\epsilon$. For this specific μ' , the power for each orthant is weighted by $1/P_1^{\alpha}$ in the upper bound. To maximize the overall power $\beta(\mu')=\sum_{\mathbf{o}\in\mathcal{O}}\beta^{\mathbf{o}}(\mu')$, we start from the orthant with the largest P_1^{α} , maximizing its power $\beta^{\mathbf{o}}(\mu')$ by letting the test ϕ to be 1 for that orthant. Then we go to the second largest and keep doing it until the inequality becomes equal. This is indeed the sign test, giving that $\beta(\mu')$ is no larger than the power of the sign test. Together with the fact that $\beta(\mu')$ serves as an upper bound for the worse-case power of ϕ , it implies that ϕ^S is maximin among all tests in $S(\omega)$. Furthermore, for any $\phi \in S$, $\exists \omega$ such that $\phi \in S(\omega)$. By noting that ϕ^S is a maximin α -level test in any $S(\omega)$ we complete the proof.

4.3 Proof of Lemma 6

Proof. Consider any $\phi(\mathbf{Y}) \in \mathcal{B}_n$. For any $\omega > 0$, define $\psi_{(\omega)} = \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \psi_{\mathbf{b}}^{\mathbf{o}} \mathbb{I}_{I_{\mathbf{b}}^{\mathbf{o}}}$, where $\psi_{\mathbf{b}}^{\mathbf{o}} = \frac{1}{|I_{\mathbf{b}}^{\mathbf{o}}|} \int_{I_{\mathbf{b}}^{\mathbf{o}}} \phi(\mathbf{y})$ and $|I_{\mathbf{b}}^{\mathbf{o}}|$ is the volume of the box $I_{\mathbf{b}}^{\mathbf{o}}$. We next show that there exists a small enough ω such that $\phi \doteq \psi_{(\omega)}$.

Notice that by fixing $\omega \in \mathbb{R}$, any real positive vector $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)$ can be written as $\mu_j = (1+\omega)^{d_j}(1+\omega_j')$, for $\mathbf{d} = (d_1, \dots, d_n)$ and $\boldsymbol{\omega}' = (\omega_1', \dots, \omega_n')$, where $0 \leq \boldsymbol{\omega}' < \omega$ and we recall that the inequality of the vector $\boldsymbol{\omega}'$ is element-wise. Similar to Step 2 in the proof sketch

of Theorem 3, define $\tilde{I}_b^+(\omega_j') = \frac{1}{1+\omega_j'}I_b^+ = (\frac{(1+\omega)^b}{1+\omega'}, \frac{(1+\omega)^{b+1}}{1+\omega_j'}]$ and $\tilde{I}_b^-(\omega_j') = \frac{1}{1+\omega_j'}I_b^-(\omega_j') = [-\frac{(1+\omega)^{b+1}}{1+\omega_j'}, -\frac{(1+\omega)^b}{1+\omega_j'}]$. Then we can define the rescaled box $\tilde{I}_b^{\mathbf{o}} = \tilde{I}_{b_1}^{o_1}(\omega_1') \times \cdots \times \tilde{I}_{b_n}^{o_n}(\omega_n')$. Note that the element-wise multiplication by $\boldsymbol{\mu}$ maps the rescaled box $\tilde{I}_b^{\mathbf{o}}$ to the original box $I_{\mathbf{b}+\mathbf{d}}^{\mathbf{o}}$. First, for any index \mathbf{b} , we have

$$\int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[\phi(\boldsymbol{\mu} \circ \mathbf{z}) - \psi(\boldsymbol{\mu} \circ \mathbf{z}) \right] d\mathbf{z} = \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \phi(\boldsymbol{\mu} \circ \mathbf{z}) d\mathbf{z} - |\tilde{I}_{\mathbf{b}}^{\mathbf{o}}| \psi_{\mathbf{b}+\mathbf{d}}^{\mathbf{o}} \right]
= \frac{1}{\prod_{j} \mu_{j}} \int_{I_{\mathbf{b}+\mathbf{d}}^{\mathbf{o}}} \phi(\mathbf{z}) d\mathbf{z} - |\tilde{I}_{\mathbf{b}}^{\mathbf{o}}| \psi_{\mathbf{b}+\mathbf{d}}^{\mathbf{o}} = \frac{1}{\prod_{j} \mu_{j}} |I_{\mathbf{b}+\mathbf{d}}^{\mathbf{o}}| \psi_{\mathbf{b}+\mathbf{d}}^{\mathbf{o}} - |\tilde{I}_{\mathbf{b}}^{\mathbf{o}}| \psi_{\mathbf{b}+\mathbf{d}}^{\mathbf{o}} = 0.$$

Let us define $\underline{f}_{\mathbf{b}}^{\mathbf{o}}(\boldsymbol{\omega}') = \inf_{\mathbf{z} \in \tilde{I}_{\mathbf{b}}^{\mathbf{o}}} f(\mathbf{z})$. For any $\boldsymbol{\mu}$ specified by \mathbf{d} and $\boldsymbol{\omega}'$, the difference of the sizes of the two tests ϕ and ψ can be upper bounded as

$$\begin{aligned} &|\mathbb{E}_{0}[\phi(\mathbf{Y}) - \psi(\mathbf{Y})]| = \left| \int_{\mathbb{R}^{n}} \left[\phi(\boldsymbol{\mu} \circ \mathbf{z}) - \psi(\boldsymbol{\mu} \circ \mathbf{z}) \right] f(\mathbf{z}) d\mathbf{z} \right| \\ &= \left| \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[\phi(\boldsymbol{\mu} \circ \mathbf{z}) - \psi(\boldsymbol{\mu} \circ \mathbf{z}) \right] f(\mathbf{z}) d\mathbf{z} \right| \\ &= \left| \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[\phi(\boldsymbol{\mu} \circ \mathbf{z}) - \psi(\boldsymbol{\mu} \circ \mathbf{z}) \right] \left[f(\mathbf{z}) - \underline{f}_{\mathbf{b}}^{\mathbf{o}}(\boldsymbol{\omega}') \right] d\mathbf{z} \right| \\ &\leq \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z}) - \underline{f}_{\mathbf{b}}^{\mathbf{o}}(\boldsymbol{\omega}') \right] d\mathbf{z} \\ &\leq \sum_{\mathbf{o} \in \mathcal{O}} \left[\sum_{-m \le \mathbf{b} \le m} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z}) - \underline{f}_{\mathbf{b}}^{\mathbf{o}}(\boldsymbol{\omega}') \right] d\mathbf{z} + \sum_{\mathbf{b} : \exists |b_{j}| > m} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z}) - \underline{f}_{\mathbf{b}}^{\mathbf{o}}(\boldsymbol{\omega}') \right] d\mathbf{z} \right| \end{aligned}$$

As $\omega \to 0$ and $m \to \infty$, the above expression goes uniformly to zero for all $\forall \omega' < \omega$. Therefore $\exists \omega_0, m_0$, such that $\forall \omega < \omega_0, m > m_0$, we have

$$\sum_{\mathbf{o} \in \mathcal{O}} \sum_{-m \le \mathbf{b} \le m} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z}) - \underline{f}_{\mathbf{b}}^{\mathbf{o}}(\boldsymbol{\omega}') \right] d\mathbf{z} \le \frac{\epsilon}{2}, \quad \forall \boldsymbol{\omega}' < \omega,$$

$$\sum_{\mathbf{o} \in \mathcal{O}} \sum_{\mathbf{b} : \exists |b_{i}| > m} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z}) - \underline{f}_{\mathbf{b}}^{\mathbf{o}}(\boldsymbol{\omega}') \right] d\mathbf{z} \le \frac{\epsilon}{2}.$$

Thus $\forall \mu > 0$, $|\mathbb{E}_0[\phi(\mathbf{Y}) - \psi(\mathbf{Y})]| \leq \epsilon$, which implies that the sizes of the two tests are within ϵ -distance of each other.

Similarly, we can bound the difference between the power of the two tests. Let $\underline{f}_{\underline{\mathbf{b}}}^{\mathbf{o}}(\omega') = \inf_{\mathbf{z} \in \tilde{I}_{\mathbf{b}}^{\mathbf{o}}} f(\mathbf{z} - \delta)$. For any μ specified by \mathbf{d} and ω' , the difference of the power of ϕ and $\overline{\psi}$ can be written as

$$\begin{aligned} &|\mathbb{E}_{1}[\phi(\mathbf{Y}) - \psi(\mathbf{Y})]| = \left| \int_{\mathbb{R}^{n}} \left[\phi(\boldsymbol{\mu} \circ \mathbf{z}) - \psi(\boldsymbol{\mu} \circ \mathbf{z}) \right] f(\mathbf{z} - \delta) d\mathbf{z} \right| \\ &= \left| \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[\phi(\boldsymbol{\mu} \circ \mathbf{z}) - \psi(\boldsymbol{\mu} \circ \mathbf{z}) \right] f(\mathbf{z} - \delta) d\mathbf{z} \right| \\ &= \left| \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[\phi(\boldsymbol{\mu} \circ \mathbf{z}) - \psi(\boldsymbol{\mu} \circ \mathbf{z}) \right] \left[f(\mathbf{z} - \delta) - \underline{\underline{f}_{\mathbf{b}}^{\mathbf{o}}}(\boldsymbol{\omega}') \right] d\mathbf{z} \right| \\ &\leq \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-\infty < \mathbf{b} < \infty} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z} - \delta) - \underline{\underline{f}_{\mathbf{b}}^{\mathbf{o}}}(\boldsymbol{\omega}') \right] d\mathbf{z} \\ &\leq \sum_{\mathbf{o} \in \mathcal{O}} \left[\sum_{-m \le \mathbf{b} \le m} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z} - \delta) - \underline{\underline{f}_{\mathbf{b}}^{\mathbf{o}}}(\boldsymbol{\omega}') \right] d\mathbf{z} + \sum_{\mathbf{b} : \exists |b_{j}| > m} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z} - \delta) - \underline{\underline{f}_{\mathbf{b}}^{\mathbf{o}}}(\boldsymbol{\omega}') \right] d\mathbf{z} \right] \end{aligned}$$

Similar to the argument in the previous case, as $\omega \to 0$ and $m \to \infty$, the above expression goes uniformly to zero for all $\omega' < \omega$. Hence there exists ω_1, m_1 , such that $\forall \omega < \omega_1, m > m_1$,

$$\begin{split} & \sum_{\mathbf{o} \in \mathcal{O}} \sum_{-m \leq \mathbf{b} \leq m} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z} - \delta) - \underline{\underline{f}_{\mathbf{b}}^{\mathbf{o}}}(\boldsymbol{\omega}') \right] d\mathbf{z} < \frac{\epsilon}{2}, \quad \forall \boldsymbol{\omega}' < \omega, \\ & \sum_{\mathbf{o} \in \mathcal{O}} \sum_{\mathbf{b} : \exists |b_{j}| > m} \int_{\tilde{I}_{\mathbf{b}}^{\mathbf{o}}} \left[f(\mathbf{z} - \delta) - \underline{\underline{f}_{\mathbf{b}}^{\mathbf{o}}}(\boldsymbol{\omega}') \right] d\mathbf{z} < \frac{\epsilon}{2}. \end{split}$$

Thus $\forall \mu > 0$, $|\mathbb{E}_1[\phi(\mathbf{Y}) - \psi(\mathbf{Y})]| \le \epsilon$ which implies the distance between the power of the two tests is bounded by ϵ . Therefore, for any $\epsilon > 0$, by selecting $\omega < \min(\omega_0, \omega_1)$ and $m > \max(m_0, m_1)$, we have a test $\psi_{(\omega)} \in \mathcal{S}$, such that $\forall \ \mu > 0$, $|\mathbb{E}_0[\phi(\mathbf{Y}) - \psi(\mathbf{Y})]| \le \epsilon$ and $|\mathbb{E}_1[\phi(\mathbf{Y}) - \psi(\mathbf{Y})]| \le \epsilon$. As a result, $\phi \doteq \psi_{(\omega)}$.