CONSISTENCY OF INVARIANCE-BASED RANDOMIZATION TESTS

BY EDGAR DOBRIBAN^a

Department of Statistics and Data Science, The Wharton School, University of Pennsylvania, adobriban@wharton.upenn.edu

Invariance-based randomization tests—such as permutation tests, rotation tests, or sign changes—are an important and widely used class of statistical methods. They allow drawing inferences under weak assumptions on the data distribution. Most work focuses on their type I error control properties, while their consistency properties are much less understood.

We develop a general framework to study the consistency of invariance-based randomization tests, assuming the data is drawn from a signal-plus-noise model. We allow the transforms (e.g., permutations or rotations) to be general compact topological groups, such as rotation groups, acting by linear group representations. We study test statistics with a generalized subadditivity property.

We apply our framework to a number of fundamental and highly important problems in statistics, including sparse vector detection, testing for low-rank matrices in noise, sparse detection in linear regression, and two-sample testing. Comparing with minimax lower bounds we develop, we find perhaps surprisingly that in some cases, randomization tests detect signals at the minimax optimal rate.

1. Introduction. Invariance-based randomization tests—such as permutation tests—are an important, fundamental, and widely used class of statistical methods (see Figure 1 for an illustration). They allow making inferences in general settings, with few assumptions on the data distribution. Most methodological and theoretical work focuses on their validity, studying their type I error (false positive rate) control. There is also work on their robustness properties, but less is known about their power and consistency properties.

Our work develops a general theoretical framework to understand the consistency properties of invariance-based randomization tests. We assume that the data follows a "signal-plus-noise" model, being the sum of a deterministic signal and a random noise component. We allow the randomization distributions to be Haar measures over general compact topological groups, such as rotation groups. We go beyond most prior work, which often focuses on discrete groups (mainly permutation groups), and does not fully develop the technically challenging case of compact groups. Moreover, we allow the action of these groups on the data to be via arbitrary compact linear group representations.

We apply our theoretical framework to a number of fundamental and highly important problems in statistics, including sparse vector detection, low-rank matrix detection, sparse detection in linear regression, and two-sample testing. Perhaps surprisingly, combining with minimax lower bounds that we develop, we find that invariance-based randomization tests for appropriate test statistics are minimax rate optimal in a number of cases. We consider this surprising, because the critical values of randomization tests are determined using the same universal principle. These critical values rely on very little information about the problem, namely a set of symmetries of the noise.

In more detail, our contributions are as follows:

Received January 2022; revised May 2022.

MSC2020 subject classifications. Primary 62G10; secondary 62G09, 62H15.

Key words and phrases. Randomization test, permutation test, nonparametric group invariance, sparse detection.

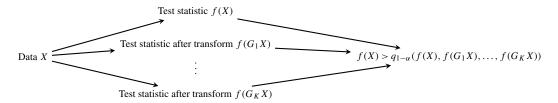


FIG. 1. Flowchart of the invariance-based randomization test: Given the observed data X, we compute the test statistic f(X), as well as the test statistics f(X), $f(G_1X)$,..., $f(G_KX)$ applied to the data X transformed by some transforms G_1, \ldots, G_K , for instance, permutations, sign changes, or rotations. Under the null hypothesis, we assume that the distribution X is invariant under a group \mathcal{G} (e.g., permutations, rotations, or sign changes), and that G_1, \ldots, G_K are sampled i.i.d. from the uniform distribution over this group. We reject the null hypothesis if f(X) is larger than the $1-\alpha$ th quantile of $f(G_1X), \ldots, f(G_KX)$. This includes familiar permutation tests for exchangeable or i.i.d. data, rotation tests assuming spherically distributed data, and sign changes assuming noise symmetric about zero.

- 1. Framework for studying consistency. We develop a framework for studying the consistency of invariance-based randomization tests using the language of group representation theory. In our framework, we have a compact topological group (e.g., permutations, rotations) that acts linearly on the data space. We assume that under the null hypothesis, the distribution of the data is invariant under the action of the group. In the standard randomization test, we sample several group elements chosen at random from the Haar measure on the group, and apply them to the data. This test rejects the null hypothesis when a chosen test statistic is larger than an appropriate quantile of the values of the test statistic applied to the randomly transformed data.
- 2. Consistency results. We develop consistency results for the invariance-based randomization test, assuming the data has been sampled from signal-plus-noise models. We consider sequences of signal-plus-noise models where the signal equals zero under the null hypothesis. We study broad classes of test statistics satisfying the weak requirement of so-called ψ -subadditivity. This includes, for instance, suprema of linear functionals of the data, norms and seminorms, concave nondecreasing functions in one dimension, and convex functions of bounded growth. Further, this class is closed under conic combinations, taking maxima, and compositions with one-dimensional nondecreasing subadditive functions.

We develop a general consistency result, showing that if the sequence of alternatives is such that the value of the test statistic is large enough, then the randomization test rejects with probability tending to unity. We compare this to the corresponding result for the deterministic test based on the same the statistic. The consistency threshold for the randomization test is inflated slightly by a signal-noise interference effect. By randomly transforming the signal, we create additional noise, inflating the effective noise level in the randomized statistic compared to its distribution under the null. However, we later show that in many examples this inflated noise level can be controlled. As part of our consistency theory, we extend to the setting with nuisance parameters, which allows us to handle problems such as two-sample testing.

- 3. New proof techniques. Our proofs are based on novel approaches. For the proofs of the general consistency result, we proceed by a series of reductions, first reducing from the quantile of the randomization distribution to its maximum, then from considering several random transformations to only one transform, and then reducing from a dependent transformed signal and noise to independent ones, via a "deterministic separating sequence" argument.
- 4. *Examples*. We illustrate our results in several important examples. We show that our results provide consistency conditions for invariance-based randomization tests in a number of problems, including sparse vector detection, low-rank matrix detection, sparse detection in linear regression, and two-sample testing.

For sparse vector detection, we consider two settings: where the noise vectors for the different observations are independent and sign symmetric (but not necessarily identically distributed), and where they are rotationally symmetric (spherical). For both cases, we obtain general consistency results, and some matching lower bounds. Specifically for the sign symmetric case where the entries of the noise are independent and identically distributed according to a subexponential distribution, our upper bound for the signflip randomization test matches a lower bound that we obtain. For spherical noise, we obtain general upper bounds as well as specific examples for multivariate t distributions. We also provide similar results for two sample testing.

For low rank matrix detection, we consider the case where each of the columns of the noise matrix has an independent spherical distribution. We obtain a general consistency guarantee for the randomization test based on the operator norm test statistic, using rotation transforms. We show that this result is rate optimal for the special case of normal noise.

For sparse vector detection in linear regression, we study detection based on the ℓ_{∞} norm of the least-squares estimator. We assume that the noise entries of each observation are independent and sign-symmetric. We provide a consistency result for the associated signflip based randomization test, in terms of geometric quantities determined by the feature matrix; namely the suprema of two associated Bernoulli processes.

As a general conclusion, we think it is perhaps surprising that invariance-based randomization tests can sometimes detect signals at the same rate as the optimal tests that assume knowledge about the exact noise distribution. We support our claims with numerical experiments. These experiments can be reproduced with the code provided at https://www.github.com/dobriban/randomization_test, also provided in the Supplementary Material (Dobriban (2022a)).

Note on terminology. We follow the terminology of "randomization tests" from Ch. 15.2 of the standard textbook by Lehmann and Romano (2005): "the term randomization test will refer to tests obtained by recomputing a test statistic over transformations (not necessarily permutations) of the data." This does not consider tests based on randomization of treatments; see, for example, Hemerik and Goeman (2020), Onghena (2018) for discussion. In particular, Hemerik and Goeman (2020) suggest using "randomization tests" only when the treatments are randomized, and suggest using "group invariance tests" for the type of tests we consider. For consistency with the standard textbook by Lehmann and Romano (2005), we will simply use the terminology "invariance-based randomization tests" or "randomization tests." Another well-known example of randomization occurs with discretely distributed tests, to ensure exact type I error control; our work is unrelated to this issue.

Some notations. For a positive integer $m \ge 1$, the m-dimensional all-ones vector is denoted as $1_m = (1, 1, \dots, 1)^{\top}$. We denote $[m] := \{1, 2, \dots, m\}$, and for $j \in [m]$, the jth standard basis vector by $e_j = (0, \dots, 1, \dots, 0)$, where only the jth entry equals unity, and all other entries equal zero. The variance of a random variable X is denoted as Var X or Var[X]. For two random vectors X, Y, we denote by $X =_d Y$ that they have the same distribution. For an index $m = 1, 2, \dots$, and two sequences $(a_m)_{m \ge 1}$, $(b_m)_{m \ge 1}$, $a_m \le b_m$ (and $a_m = O(b_m)$) means that $a_m \le Cb_m$ for some $C \ge 1$ independent of m, but possibly dependent on other problem parameters as specified case by case. We write $a_m \ge b_m$ (or $a_m = \Omega(b_m)$) when $b_m \le a_m$, and $a_m \sim b_m$ (or $a_m = \Theta(b_m)$) when $a_m \le b_m \le a_m$. For a vector $v \in \mathbb{R}^m$, and $p \in (0, \infty)$, $\|v\|_p$ denotes the ℓ_p norm. Unless otherwise specified, $\|v\|$ denotes the Euclidean or ℓ_2 norm, $\|v\| = \|v\|_2$. For a matrix s, the norm $\|s\|_{2,\infty}$ is the maximum of the column ℓ_2 norms of s. For two subsets s, s of a vector space, s, s, s is the maximum of the column s denotes the Minkowski sum. For a s is a vector s, let s is the maximum of the column s denotes the Minkowski sum. For a s is a vector s, let s denotes the Minkowski sum. For a s is a vector s, let s denotes the s denotes the Minkowski sum. For a s is a function s is an odd function if s is a function s is an odd function if s is a function s is an odd function if s is a function s in the position variable is uniform over

the set $\{\pm 1\}$. For a probability distribution Q and a random variable $X \sim Q$, we may write probability statements involving X in several equivalent ways, for instance for the probability that X belongs to a measurable set A, we may write: $P(X \in A)$, $P_X(A)$, $P_Q(X \in A)$, $P_{X \sim Q}(X \in A)$, $Q(X \in A)$, or Q(A). Further, if Q belongs to a collection of probability measures H (e.g., a null or an alternative hypothesis), then we may also write $P_H(A)$ to denote Q(A) for an arbitrary $Q \in H$.

1.1. Related works. There is a large body of important related work. Here we can only review the most closely related ones due to space limitations; see Section 1 of the Supplementary Material (Dobriban (2022b)) for additional related works. The idea of constructing a statistical test based on randomly chosen permutations of datapoints in a dataset dates back at least to Eden and Yates (1933), Fisher (1935); see Berry, Johnston and Mielke (2014), David (2008) for historical details. General references on permutation tests include Anderson and Robinson (2001), Ernst (2004), Good (2006), Hemerik and Goeman (2018a), Kennedy (1995), Pesarin (2001a), Pesarin and Salmaso (2010a), Pesarin and Salmaso (2012). These tests have many applications, for instance, in genomics (Tusher, Tibshirani and Chu (2001)) and neuroscience (Winkler et al. (2014)). For more general discussions of invariance in statistics, see Eaton (1989), Giri (1996), Wijsman (1990); for a general probabilistic reference, see also Kallenberg (2006).

Two-sample permutation tests date back at least to Pitman (1937), and have recently been studied in more general multivariate contexts (Kim, Balakrishnan and Wasserman (2020)). This problem brings special considerations such as issues with using balanced permutations (Southworth, Kim and Owen (2009)).

For the theoretical aspects of invariance-based randomization tests, Lehmann and Stein (1949) develop results for testing a null of equality in distribution $H_{m0}: X_m =_d g_m X_m$ where a transform $g_m \in \mathcal{G}_m$ is chosen from a group \mathcal{G}_m acting on the data. They show that all admissible tests have constant rejection probability equal to the level over each orbit, that is, are similar tests. They use this to show that the most powerful tests against simple alternatives with density f_m reject when $f_m(X_m)$ is greater than the appropriate quantile of $\{f_m(g_m X_m), g_m \in \mathcal{G}_m\}$. They use the Hunt-Stein theorem to derive uniformly most powerful (or most stringent) invariant tests from maximin tests for testing against certain composite alternatives. These are related to our results, but we focus on consistency against special structured signal-plus-noise alternatives instead of maximizing power in a finite sample.

The seminal work by Hoeffding (1952) considers general group transforms, including signflips, for testing symmetry of distributions, but focuses on permutation groups for most part. The main results center on power and consistency of tests. For consistency, Theorem 2.1 in Hoeffding (1952) states that for a test statistic f_m such that $f_m(x) \geq 0$ and $\mathbb{E}_{G_m \sim Q_m} f_m(G_m x) \leq c$ where $G_m \sim Q_m$ denotes that the group element is distributed according to the probability measure Q_m on G_m , we have $q_{1-\alpha,m}(x) \leq c/\alpha$, $\alpha \in (0,1)$, where $q_{1-\alpha,m}$ is the $1-\alpha$ th quantile of the distribution of $f_m(G_m x)$ when $G_m \sim Q_m$. Then, if $f_m \to \infty$ under a sequence of alternatives, the test that rejects when $f_m > q_{1-\alpha,m}(x)$ is consistent, that is, has power tending to unity. These conditions are distinct from ours. Specifically, his conditions require the test statistic to be pointwise bounded (for each datapoint x, they require that $\mathbb{E}_{G_m \sim Q_m} f_m(G_m x) \leq c$), whereas we make assumptions that need to hold with high probability, over the randomness in the data and the random transform. Thus, the two types of conditions are different. Our conditions seem to be more direct. His Theorem 3.1 requires establishing limiting distributions of test statistics, which may be challenging or impossible in certain cases.

For asymptotic power of certain special invariance-based randomization tests, one can obtain results based on contiguity, see, for example, Example 15.2.4 in Lehmann and Romano (2005). In this line of work Albers, Bickel and van Zwet (1976), Bickel and van Zwet

(1978) have studied asymptotic expansions of the power of distribution-free one-sample and two-sample tests, using Edgeworth expansions, showing that the deficiency (in the sense of Hodges and Lehmann) of appropriate linear rank tests and permutation tests compared to optimal nonparametric tests can often be quite small, and even tend to zero in certain cases. However, we are interested in problems where the contiguity of the alternatives may be unknown, or hard to establish.

For permutation tests, Dwass (1957) shows that it is valid to randomly sample permutations—as opposed to using all permutations—to construct the randomization test. Hemerik and Goeman (2018b) provide a general type I error control result for random group transformations under exact invariance, and apply it to false discovery proportion control. See also Hemerik, Solari and Goeman (2019). Hemerik and Goeman (2018a) extend this in various forms, including to sampling transforms without replacement, and giving rigorously justified formulas for *p*-values.

There are a number of works studying the power properties of invariance-based randomization tests. We have already discussed the fundamental work by Hoeffding (1952). Pesarin and Salmaso (2010) develop finite-sample consistency results for certain combination-based permutation tests for multivariate data, when the sample size is fixed and the dimension tends to infinity. They focus on one-sided two-sample tests, and discuss Hotelling's *T*-test as an example. Pesarin and Salmaso (2013) characterize weak consistency of permutation tests for one-dimensional two-sample problems. They study stochastic dominance alternatives assuming the population mean is finite and without assuming existence of population variance.

Pesarin (2015) develops some further theoretical aspects of permutation tests. This includes consistency properties (Property 9), for two-sample tests under some nonparametric assumptions, and alternatives specified by an increased mean of the test statistic. These have different assumptions than the results in our paper, focusing on two-sample problems (while we have general invariance), and nonparametric models (while we focus on parametric ones).

One one of the most closely related papers is that of Kim, Balakrishnan and Wasserman (2020b). They study permutation tests, which are a special case of group invariance tests. The examples in our work mostly concern signflip-based and rotation tests. Two-sample testing is studied in both works, but under different assumptions: we study testing the equality of means in a location model, whereas they study testing the equality of two distributions, such as multinomials and distributions with Hölder densities. Thus, our results are not directly comparable. For instance, our minimax optimality for two-sample testing involves location families with i.i.d. subexponential noise, whereas their examples are multinomial distributions and Hölder densities.

Recent work (posted publicly after our paper) by Koning and Hemerik (2022) develops a framework for improving the power of randomization tests by choosing appropriate subgroups of transformations. Further, they include results showing that the t-test is equivalent to a randomization test where the randomization is performed under the entire orthogonal group. In this light, properties of the t-test also rely on such minimal information; and this may shed some light on the useful properties of randomization tests that we study.

In context. To further put our work in context, we can make the following comparisons:

- The vast majority of works only provide theoretical results for the behavior of randomization tests under the null hypothesis (Eden and Yates (1933), Fisher (1935), Hemerik and Goeman (2018a), Pitman (1939)).
- The seminal work of Hoeffding (1952), already discussed above, provides consistency conditions that are distinct from ours; and does not discuss any of the specific problems that we study. The power analysis based on contiguity (Lehmann and Romano (2005)) does not apply to many of the problems we study. As detailed above, the consistency and

minimax optimality analyses from Kim, Balakrishnan and Wasserman (2020b), Pesarin (2015), Pesarin and Salmaso (2010), Pesarin and Salmaso (2013) all concern different setups from and/or special cases of our results.

Scientific context. For an even broader scientific context, we emphasize that randomization tests are ubiquitous in modern science. Their proper use is crucial for reproducible results; and failure to use them correctly can result in irreproducible results, false scientific discoveries, and ultimately a waste of resources. Here are some examples:

• In neuroscience, the analysis of fMRI data requires testing hypotheses about the activation of regions in the brain. It has been observed that inferences based on models such as Gaussian fields with parametric covariance functions can have massively inflated false-positive rates (Eklund, Nichols and Knutsson (2016)). To mitigate this problem, it has been proposed to use randomization methods such as permutation methods (for two-sample problems) or random sign flips (for one-sample problems) to set critical values. Further randomization methods have been proposed for other problems such as general linear models (Winkler et al. (2014)), or brain network comparison (Simpson et al. (2013)).

The ultimate goal is to report reliable discoveries, which involves analyzing data not from the null distribution, but rather from an alternative distribution that contains signals. Our work can shed light on when randomization tests can succeed in such an analysis of data containing signals.

• In genetics and genomics, hypothesis testing is routinely performed to identify associations between observed phenotypes and genotypes, or between genotypes, etc. Randomization tests, and in particular permutation tests, are widely used to set critical values, in methods such as transmission disequilibrium tests, etc, and are broadly available in popular software such as PLINK; see, for instance, Churchill and Doerge (1994), Epstein et al. (2012), Purcell et al. (2007). Randomization tests are also used for more sophisticated tasks such as gene set enrichment analysis (Barry, Nobel and Wright (2005), Efron and Tibshirani (2007), Subramanian et al. (2005)).

2. General framework.

2.1. Setup. We consider a sequence of statistical models, indexed by an index parameter $m \to \infty$. We observe data X_m from a real vector space V_m , for instance, a vector or a matrix belonging to Euclidean space \mathbb{R}^{p_m} . We assume that we know a group \mathcal{G}_m of the symmetries of the distribution of the data. See Section 3 of the Supplementary Material (Dobriban (2022b)) for a discussion of how such symmetries can arise in practice. A group \mathcal{G}_m has a multiplication operation "·" that satisfies the axioms of associativity, identity, and invertibility. For instance, we could have that the entries of X_m are exchangeable (corresponding to the permutation group), symmetric about zero (corresponding to the group of addition modulo two) or that the density of X_m is spherical (corresponding to the rotation group).

In addition, to transform the data, we have a group representation $\rho_m: \mathcal{G}_m \to \operatorname{GL}_m(V_m)$, acting linearly on $X_m \in V_m$ via $g_m X_m := \rho_m(g_m) \cdot X_m$. The group representation "represents" the elements of the group \mathcal{G}_m as invertible linear operators $V_m \mapsto V_m$ belonging to the general linear group $\operatorname{GL}_m(V_m)$ of such operators. The group representation ρ_m preserves the group multiplication operation, that is, $\rho_m(g_m g_m') = \rho_m(g_m) \rho_m(g_m')$ for all $g_m, g_m' \in \mathcal{G}_m$, and $\rho_m(e_{\mathcal{G}_m}) = I_{V_m}$, where $e_{\mathcal{G}_m}$ is the identity element of the group, and I_{V_m} is the identity operator on V_m . For general references on representation theory, see Eaton (1989), Fulton and Harris (2013), Hall (2015), James and Liebeck (2001), Knapp (2013), Serre (1977), etc. For group representations in statistics, see Diaconis (1988). We will use basic concepts from this area throughout the paper.

Null hypothesis of invariance, and randomization test. We want to use the symmetries of the noise distribution to detect the presence of nonsymmetric signals. Under the null hypothesis, we assume that the distribution of the data is invariant under the action of each group element $g_m \in \mathcal{G}_m$: $X_m =_d g_m X_m$.\frac{1}{2} We study the following invariance-based randomization test (sometimes also called a group invariance test), which at various levels of generality has been considered dating back to Eden and Yates (1933), Fisher (1935), Hoeffding (1952), Lehmann and Stein (1949), Pitman (1937). We sample G_{m1}, \ldots, G_{mK} i.i.d. from \mathcal{G}_m (in a way specified below), and reject the null if for a fixed test statistic $f_m : V_m \mapsto \mathbb{R}_m$, the following event holds:

(1)
$$\mathcal{E}_m = \{ f_m(X_m) > q_{1-\alpha}(f_m(X_m), f_m(G_{m1}X_m), \dots, f_m(G_{mK}X_m)) \},$$

for the $1-\alpha$ th quantile $q_{1-\alpha}$ of the numbers $f_m(X_m), f_m(G_{m1}X_m), \ldots, f_m(G_{mK}X_m)$ and some $\alpha \in (0,1]$. Specifically, let $G_{m0} = I_{V_m}$ be the identity operator on V_m , and $f_{(1)} \leq f_{(2)} \leq \cdots \leq f_{(K+1)}$ be the order statistics of $f_m(G_{mi}X_m), i \in \{0,1,\ldots,K\}$. Let $k = \lceil (1-\alpha)(K+1) \rceil$. Rejecting the null if $f_m(X_m) > f_{(k)}$ is guaranteed to have level at most α , see, for example, Theorem 2 in Hemerik and Goeman (2018a) for an especially clear and rigorous statement.

Noise invariance and robustness. The advantage of randomization tests compared to a rejection region of the form $f_m(X_m) > \tilde{c}_m$ for a fixed \tilde{c}_m is that it does not require the manual specification of the critical value \tilde{c}_m . The critical value needs to account for the set of distributions included the null hypothesis, which may be a very large nonparametric family. In this case, it might be challenging to set the critical value to ensure type I error control. Randomization tests avoid this problem by relying on the symmetries of the noise distributions. To wit, randomization tests are valid under *any* null hypothesis for which the distribution of the noise is invariant under the group. This effectively amounts to the test only depending on the collection of *orbits*, which form a maximal invariant, see Sections 3 and 4 in Eaton (1989) for examples.

Thus our model can be semiparametric, where the nuisance parameter—the distribution of the noise—is infinite-dimensional and is only restricted by an invariance condition. We may add further assumptions on the noise, to enable cleanly stated consistency results, in which case the model may become parametric.

For instance, for the rotation group $O(p_m)$, we get *spherical* distributions, which have a density $p_m(X_m) = \pi_m(\|X_m\|_2)$ with respect to a σ -finite dominating measure on \mathbb{R}^{p_m} only depending on the Euclidean norm of the data X_m (Fang, Kotz and Ng (2018), Fang and Zhang (1990), Gupta and Varga (2012)). This is a nonparametric class that includes in particular distributions such as the multivariate t, multivariate Cauchy, scale mixtures of spherical normals, etc. In particular, it includes heavy tailed distributions, for which tests based on the normal assumption can have an inflated type I error. As another example, consider a stationary field $X_{m,J} = (X_i)_{i \in J}$, for some index set J. Suppose \mathcal{G}_m acts on J, and induces an action on $X_{m,J}$ via its regular representation, that is, $(g_m X_{m,J})_i = X_{g_m^{-1}i}$. For instance, we can have a discrete-time stationary time series where $J = \mathbb{Z}$, and $\mathcal{G}_m = (\mathbb{Z}, +)$. In this example, any translation of the time series keeps the distribution invariant.

While sometimes it is possible to construct test statistics whose distribution does not depend on a broad set of null hypotheses (see, e.g., Section 4.3 "Null robustness" in Eaton (1989)), this may not be possible when the null hypothesis has a great number of nuisance

¹This is called the "Randomization hypothesis," Definition 15.2.1 in Lehmann and Romano (2005). Using the terminology of Bickel et al. (1993), Section 6.3, it is an example of a "nonparametric group model."

parameters. For example, this holds for null hypotheses where each noise entry is independent with a probability density only assumed to be symmetric around zero, in which case sign-flip based methods are applicable; see, for example, Example 15.2.1 of Lehmann and Romano (2005), and also Hemerik, Goeman and Finos (2020), Hong, Sheng and Dobriban (2020).

Haar measure. In the definition of the randomization test, G_{m1}, \ldots, G_{mK} are chosen i.i.d. from the uniform (Haar) measure on \mathcal{G}_m , which is assumed to exist. We refer to Section 2 in Folland (2016) for details; see also Eaton (1989), Fulton and Harris (2013), Wijsman (1990). Thus, \mathcal{G}_m is assumed to be a compact Hausdorff topological group with the Borel sigma-algebra generated by the open sets. For brevity, we will sometimes refer to such groups as compact groups. The Haar probability measure Q_m on \mathcal{G}_m is the unique probability measure such that $Q_m(G_m \in A) = Q_m(G_m \in g'_m A)$ for all $g'_m \in \mathcal{G}_m$ and for all Borel sets A. See for example, Theorems 2.10 & 2.20 in Folland (2016). Thus, in particular, we have the equality in distribution $G_m =_d G_m g'_m$ for $G_m \sim Q_m$, and any fixed $g'_m \in \mathcal{G}_m$.

Choice of K. We remark that, as is well known, choosing K larger, and k as above, can generally lead to a more precise control of the type I error. Indeed, for a given K, the smallest type I error control guaranteed by the randomization test is 1/(K+1), and there are only K possible values of $k \in [K]$ to control the type I error more generally. Thus, for a larger K, we expect that we can control the type I error more accurately. Indeed, we observe this in our experiments. Intuitively, we generally also expect larger K to lead to higher power; and we also observe this in experiments. However, we are not aware of a general theoretical result to this extent.

Alternative hypothesis: Signal-plus-noise model. To study the consistency of the test, we will consider a sequence of alternative hypotheses in the signal-plus-noise model with a deterministic signal s_m and a random noise N_m ,

$$X_m = s_m + N_m.$$

The null hypothesis is specified by $H_{m0}: s_m = 0_{p_m}$, in which case $X_m = N_m$. The alternative hypothesis H_{m1} is specified by a set $\Theta_{m1} \subset V_m$ of signals $s_m \in \Theta_{m1}$. We call $\Theta_m = \{0\} \cup \Theta_{m1}$ the parameter space. The alternative hypothesis is decisively *not* invariant under \mathcal{G}_m . In fact, one can view the test statistic as detecting deviations from invariance.

We view the signal-plus-noise model as quite broad, and we will study a variety of examples as special cases. The breadth of the model arises from two aspects: First, one can choose the signal parameter space Θ_m to be quite general, for instance, a linear subspace, a union of linear subspaces, a convex cone, etc. Second, one can model the family to which the distribution of the noise N_m belongs; and our theory will rely on the symmetries of these distributions. Further, from classical asymptotic statistics, we know that asymptotically any sufficiently regular parametric model is well approximated by a normal observation model, which can be viewed as a signal-plus-noise model like ours if the noise distribution does not depend on the signal.

However, the scope of this model is limited in a few ways. It assumes a specific "structural model" for the data, and it is essentially a submodel of a multidimensional location family. For instance, it requires the distribution of the noise to be functionally independent on the unknown paramater s_m . In some cases, this may be approximately achieved via appropriate variance-stabilizing transforms. In our analysis, this is currently needed to be able to formulate consistency conditions based on only one global distribution of the noise. If the noise distribution can vary in parameter space, we expect that the behavior of randomization tests could be more complex. We discuss this and further limitations of our work in Section 4.

2.2. General consistency. Our basic idea to establish consistency of randomization tests is to find conditions under which the test statistic under the alternative is much larger than the randomized test statistic, that is, (informally) $f_m(s_m + N_m) \gg f_m(G_m[s_m + N_m])$. We wish to do this by introducing only broadly applicable assumptions. The first key step is to find a lower bound on $f_m(s_m + N_m)$. To achieve this, we make assumptions of f_m .

For a given constant $\psi > 0$, we consider ψ -subadditive test statistics, that is, functions $f_m : V_m \mapsto \mathbb{R}$ such that for all $a, b \in V_m$,

$$\psi \cdot f(a+b) \le f(a) + f(b)$$
.

Note that typically $\psi \leq 1$. In the current argument, we will use that $f_m(s_m + N_m) \geq \psi f_m(s_m) - f_m(-N_m)$. This allows us to lower bound the value $f_m(s_m + N_m)$ of the test statistic by a main term $\psi f_m(s_m)$ depending only on the signal, and an error term $-f_m(-N_m)$ depending only on the noise (which we will also control). We will use a similar argument to upper bound the randomized test statistic $f_m(G_m[s_m + N_m])$. These conditions are enough to guarantee the consistency of tests of the form $f_m(X_m) > \tilde{c}_m$ for appropriately chosen "oracle" critical values \tilde{c}_m (which are not practically implementable in general); and we will compare the resulting conditions later in this section.

Examples of subadditive functions include:

1. Given any set $W_m \subset V_m$, the suprema of linear functionals

$$f_m(x) = \sup_{w_m \in W_m} w_m^\top x,$$

assumed to be finite-valued functions, are 1-subadditive. These are the *sublinear functionals* on V_m , see, for example, Section 5.4, Ch. 7, and specifically Exercise 7.103 in Narici and Beckenstein (2010). In particular, affine functions $f(x) = w^{\top}x + c$ are 1-subadditive for any $w \in V_m$ and any $c \ge 0$.

- 2. For instance, for any norm $\|\cdot\|$ on V_m (with the dependence on m suppressed), we can take $f_m(x) = \|x\|$ by choosing $W_m = \{w_m : \|w_m\|_* \le 1\}$, the unit ball in the dual norm $\|\cdot\|_*$ of $\|\cdot\|_*$.
- 3. When $V_m = \mathbb{R}$ is one-dimensional, for any concave nondecreasing function $c: [0, \infty) \to \mathbb{R}$ such that $c(0) \ge 0$, $f: \mathbb{R} \mapsto \mathbb{R}$ given by f(x) = c(|x|) is 1-subadditive. Examples include $f(x) = |x|^q$ for $q \in (0, 1]$. See Section 4 of the Supplementary Material (Dobriban (2022b)) for the argument.
- 4. Convex functions of bounded growth: If $f: \mathbb{R}^p \to \mathbb{R}$ is convex and satisfies $\psi f(2x) \le 2f(x)$, then f is ψ -subadditive. Indeed, $f(a) + f(b) \ge 2f([a+b]/2) \ge \psi f(a+b)$ by convexity and bounded growth. For instance, $f(x) = ||x||_q^q$, for $q \ge 1$ satisfies $f(2x) = 2^q f(x)$, thus it is 2^{1-q} -subadditive.

Nonexamples include functions of very fast growth, for instance, $f : \mathbb{R} \to \mathbb{R}$, $f(x) = \exp(x)$. However, for the purposes of hypothesis testing, only the acceptance and rejection regions are relevant; and thus even for test statistics that are not subadditive, one may—on a case-by-case basis—find subadditive test statistics with the same acceptance and rejection regions; where our theory can be applied. For instance, instead of the exponential map above, one may consider the identity map as a monotone transform. Further, subadditive maps have a number of closure properties, being closed under:

- 1. Conic combinations: If $f_j: V_j \mapsto [0, \infty)$, $j \in [J]$ are ψ_j -subadditive, then for any $\tau_j \geq 0$, $j \in [J]$, $\sum_{j \in [J]} \tau_j f_j$ is $\min_{j \in [J]} \psi_j$ -subadditive.
- 2. Maxima: If $f_j: V_j \mapsto [0, \infty)$, $j \in [J]$ are ψ_j -subadditive, then $\max_{j \in [J]} f_j$ is $J^{-1} \min_{j \in [J]} \psi_j$ -subadditive.

3. Compositions with 1-D functions: if $f_1:[0,\infty)\mapsto\mathbb{R}$ is nondecreasing and ψ_1 -subadditive; and $f_2:\mathbb{R}^p\mapsto[0,\infty)$ is 1-subadditive, then $f_1\circ f_2:\mathbb{R}^p\mapsto\mathbb{R}$ is ψ_1 -subadditive. Indeed,

$$f_1 \circ f_2(x+y) \le f_1[f_2(x) + f_2(y)] \le \psi_1^{-1}[f_1 \circ f_2(x) + f_1 \circ f_2(y)].$$

Our first theorem is a general consistency result for randomization tests with ψ -subadditive test statistics.

THEOREM 2.1 (Consistency of randomization test). Consider a sequence of models indexed by $m \ge 1$, $m \in \mathbb{N}$, such that the data $X_m \in V_m$ follow a p_m -dimensional signal-plusnoise model $X_m = s_m + N_m$, where $s_m \in \Theta_m$ is deterministic and N_m is a random noise vector. Test the sequence of null hypotheses $H_{m0}: s_m = 0$ against a sequence of alternative hypotheses H_{m1} with signal vectors $s_m \in \Theta_{m1}$ for a fixed level $\alpha \in (0, 1]$. Reject the null hypothesis using the randomization test (1). Let f_m be ψ -subadditive. Assume the following:

- 1. Noise invariance. The distribution of the noise is invariant under \mathcal{G}_m : $N_m =_d g_m N_m$ for all $g_m \in \mathcal{G}_m$.
- 2. Signal strength. There is a sequence $(t_m)_{m\geq 1}$, and for any sequence $(s_m)_{m\geq 1}$ such that for all $m\geq 1$, $s_m\in \Theta_{m1}$, there is another sequence $(\tilde{t}_m)_{m\geq 1}$, that may depend on s_m , $\tilde{t}_m=\tilde{t}_m(s_m)$, such that for all large enough integers m,

(2)
$$f_m(s_m) > \psi^{-2} \tilde{t}_m(s_m) + \psi^{-1} (\psi^{-1} + 1) t_m.$$

Further, as $m \to \infty$:

- (a) Noise level. We have $P(f_m(N_m) \le t_m) \to 1$ and $P(f_m(-N_m) \le t_m) \to 1$.
- (b) Bound on randomized statistic. The test statistics evaluated on the randomized signal fall below $\tilde{t}_m(s_m)$, that is, for any sequence $(s_m)_{m\geq 1}$ such that for all $m\geq 1$, $s_m\in\Theta_{m1}$,

$$P_{G_m \sim Q_m} (f_m(G_m s_m) \leq \tilde{t}_m(s_m)) \to 1.$$

Under condition 1, the randomization test has level at most α . Under conditions 1 & 2, the randomization test is consistent, that is, for the event \mathcal{E}_m from (1), for any sequence $(s_m)_{m\geq 1}$ such that $s_m \in \Theta_{m1}$ for all $m \geq 1$, $\lim_{m \to \infty} P_{G_{m1}, \dots, G_{mK} \sim Q_m, N_m}(\mathcal{E}_m) = 1$.

Some comments on the assumptions are in order:

- 1. The noise invariance condition is required to ensure the exact type I error control, as discussed above.
- 2. Our analysis relies on comparing the size of the test statistic on the data and the randomized data. The subadditivity assumption allows us to reduce this to comparing the size of the test statistics on the signal, the noise, and the randomized signal. The remaining conditions are meant to capture high-probability deterministic bounds on the statistic over the randomness in the remaining stochastic quantities: the noise and random group elements.
- 3. The sequence t_m controls the size of the statistic f_m evaluated on the noise N_m . The sequence $\tilde{t}_m(s_m)$ controls the size of the statistic evaluated on the randomized signal $G_m s_m$.

See Section 4 of the Supplementary Material (Dobriban (2022b)) for the proof, which is novel. For the consistency result, we proceed by a series of reductions, first reducing from the quantile test to a max-based test, then from considering several random transformations to only one transform, and then reducing from a dependent transformed signal and noise to independent ones. In particular, our proof relies on the existence of a "separating sequence," which deterministically separates the values of the test statistic $f_m(X_m)$ from the randomized versions $f_m(G_mX_m)$. For completeness, we record this result below.

PROPOSITION 2.2 (Separating sequence). Under the conditions of Theorem 2.1, assume only Condition 1, without ψ -subadditivity and without Condition 2. Suppose in addition that there is a separating sequence $(t'_m)_{m\geq 1}$ such that $f_m(X_m) > t'_m$ and $f_m(G_mX_m) \leq t'_m$ with probability tending to unity as $m \to \infty$. Then, the randomization test is consistent.

See the end of Section 4 of the Supplementary Material (Dobriban (2022b)) for the proof, which is already contained in the proof of Theorem 2.1. While our analysis relies on the existence of a deterministic separating sequence, it turns out that this is not necessary for the consistency of randomization tests, and we provide a discussion at the end of Section 4 of the Supplementary Material (Dobriban (2022b)).

Conventions. To lighten notation, we will often omit the dependence of $\tilde{t}_m(s_m)$ on s_m , writing simply \tilde{t}_m . Further, when it is clear from context what the sequence of tests is, we will simply say that the "test is consistent," as opposed to saying that the "sequence of tests is consistent."

Consistency of deterministic test. As mentioned, ψ -subadditivity is enough to guarantee the consistency of deterministic tests of the form $f_m(X_m) > \tilde{c}_m$ for appropriately chosen deterministic critical values \tilde{c}_m . We state this result below and compare it as a "baseline" result with the conditions for the consistency of randomization tests.

PROPOSITION 2.3 (Consistency of deterministic test). *In the setting of Theorem* 2.1, *suppose that condition* 2(*a*) *holds, along with the following condition*:

1. Signal strength. There is a sequence $(t_m)_{m\geq 1}$ such that for all large enough integers $m\geq 1$,

$$f_m(s_m) > 2\psi^{-1}t_m.$$

Then, for any sequence $(\tilde{c}_m)_{m\geq 1}$ such that $\tilde{c}_m \leq t_m$ for all $m\geq 1$, the sequence of deterministic tests that rejects when $f_m(X_m) > \tilde{c}_m$ is consistent, that is,

$$\lim_{m \to \infty} P_{H_{m1}} \left(f_m(X_m) > \tilde{c}_m \right) = 1.$$

See Section 4 of the Supplementary Material (Dobriban (2022b)) for the proof. To ensure type I error control at level α , the sequence of critical values $(\tilde{c}_m)_{m\geq 1}$ needs to be chosen such that $\sup_{P\in H_{m0}} P(f_m(X_m) > \tilde{c}_m) \leq \alpha$. As we discussed, this can be difficult when the class of null hypotheses is large and has many nuisance parameters. Thus, the test with deterministic critical values may not be practically implementable. However we can still consider it as an idealized "baseline," to understand the conditions on the signal strength that our approach provides to ensure consistency. Comparing the conditions for data signal strength, (2) and (3), and recalling that typically $\psi \leq 1$, we see that the requirement for the randomization test is stronger. The factor in front the noise level t_m is larger, and in addition the randomization test also has the additional term $\psi^{-2}\tilde{t}_m$ controlling the size of the randomized signal.

Thus, our requirements for the randomization test are more stringent. However, as explained above, the deterministic test requires a method to set the critical value, which may be very hard or impossible in practice in certain problems where the null hypothesis is very large. Specifically, every distribution in the null hypothesis leads to a constraint for the critical value; and thus, in general, makes setting the critical value more challenging.

Nuisance parameters. We next develop a generalization of our consistency results allowing nuisance parameters. This allows handling problems such as two-sample testing where the global mean is a nuisance. Let $X_m = \nu_m + s_m + N_m$, where ν_m is a nuisance parameter, $s_m \in \Theta_m$ is the signal. Suppose ν_m belongs to a known linear space U_m , $\nu_m \in U_m$. We can reduce this to the previous setting by projecting into the orthogonal complement of U_m . Let $P_m = P_{U_m^{\perp}}$ be the orthogonal projection operator into the orthogonal complement of U_m . Then $P_m \nu_m = 0$, so by projecting with P_m , we have

$$P_m X_m = P_m s_m + P_m N_m.$$

Let $\tilde{X}_m = P_m X_m$ be the new observation, $\tilde{S}_m = P_m s_m$ be the new signal, and $\tilde{N}_m = P_m N_m$ be the new noise. Then, this reduces to the standard signal-parameter model, with the signal parameter space $\tilde{\Theta}_m = P_m \Theta_m = \{P_m s_m : s_m \in \Theta_m\}$, and a new induced noise distribution.

- 2.3. Review of tools to obtain concrete results. To analyze concrete examples, we will rely on a few technical tools, reviewed in the following sections.
- 2.3.1. Rate optimality. In this section, we review some basic results on minimax rate optimality for hypothesis testing that we will use, focusing on Ingster's (or the chi-squared) method (Ingster (1987), Ingster and Suslina (2012)). This result allows randomized tests ϕ : $V_m \mapsto [0, 1]$, where $\phi(x)$ is the probability of rejecting the null for data x. Denote the set of all level $\alpha \in (0, 1)$ tests by

$$\Phi_m(\alpha) = \Big\{ \phi : V_m \mapsto [0,1] : \sup_{P \in H_{m0}} \mathbb{E}_P[\phi] \le \alpha \Big\}.$$

Define the minimax type II error as

$$R_m = \inf_{\phi \in \Phi_m(\alpha)} \sup_{P \in H_{m1}} \mathbb{E}_P[1 - \phi].$$

Suppose that $P_{m0} \in H_{m0}$ and $P_{m1}, \ldots, P_{mM_m} \in H_{m1}$. Define the average likelihood ratio between P_{m0} and P_{m1}, \ldots, P_{mM_m} as

$$L_m = \frac{1}{M_m} \sum_{i=1}^{M_m} \frac{p_{mi}(X_m)}{p_{m0}(X_m)},$$

where p_{mi} , $i \in [M_m] \cup \{0\}$ are, respectively, the densities of P_{mi} , $i \in [M_m] \cup \{0\}$ with respect to a common dominating sigma-finite measure on V_m . Then, it is well known (see, e.g., Ingster and Suslina (2012), and Section III.B of Banks et al. (2018) for a very clear statement) that to achieve consistency, that is, to have $R_m \to 0$, we must have $\lim_{m \to \infty} \operatorname{Var}_{P_{m0}}[L_m] = \infty$.

A further key result holds when the null distribution P_{m0} is $\mathcal{N}(0, I_{p_m})$ and the alternative H_{m1} contains distributions of the form $\mathcal{N}(s_m, I_{p_m})$, for $s_m \in \Theta_{m1}$. Consider a prior Π_m on Θ_{m1} . Then we have—see, for example, Ingster and Suslina (2012) or Lemma 1 of Banks et al. (2018)—for two independent copies $S, S' \sim \Pi_m$,

(4)
$$\operatorname{Var}_{P_{m0}}[L_m] = \mathbb{E}_{S,S' \sim \Pi_m} \exp(S^{\top}S').$$

2.3.2. Tail bounds of random variables. We recall some well-known tail bounds for random variables. Suppose that for all $m \ge 1$ and $i \in [m]$, Z_i are i.i.d. random variables with a probability distribution π . Let $F_{\pi}(t,n) = P(|n^{-1}\sum_{i=1}^{n} Z_i| > t)$, with $Z_i \sim \pi$ i.i.d. for all $i \in [n]$.

There is a vast number of well-known results on tail bounds of sums of i.i.d. random variables under a variety of conditions; see, for example, Boucheron, Lugosi and Massart (2013), Petrov (2012), Vershynin (2018), etc. Each of these can be used together with our framework to obtain consistency results. In a very rough order of increasing generality:

- 1. The tail of sums of *subexponential random variables* (including *sub-Gaussian and bounded variables*) can be controlled via Bernstein-type inequalities, which lead to $F_{\pi}(t,n) \leq C \exp(-cn \min\{t,t^2\})$ for some C, c depending only on π (Vershynin (2018)). Bernstein-Orlicz random variables interpolate between sub-Gaussian and subexponential random variables (van de Geer and Lederer (2013)).
- 2. There are various *Orlicz norms* for random variables, and corresponding tail bounds, for instance, for random variables with tail decay of order roughly $\exp(-x^{\alpha})$, $\alpha > 0$ (which have all polynomial moments but for $\alpha < 1$ have no moment generating function) (Chamakh, Gobet and Szabó (2020)), or of order roughly $\exp(-\ln[x+1]^{\kappa})$ for $\kappa > 0$ (which have all polynomial moments but no moment $\mathbb{E}_{X \sim \pi} \exp(|X|^c)$, c > 0 (Chamakh, Gobet and Liu (2021)).

For instance, the results of Chamakh, Gobet and Liu (2021) imply the following. Consider $\Psi: \mathbb{R}^+ \mapsto \mathbb{R}^+$, $\Psi(x) = \exp(\ln[x+1]^\kappa) - 1$, and for a random vector Z, the Ψ -Orlicz "norm" $\|Z\|_{\Psi} = \inf\{c > 0 : \mathbb{E}\Psi(\|Z\|/c) \le 1\}$. Then, for i.i.d. random variables $Z_1, \ldots, Z_n \sim \pi$, with finite Ψ -Orlicz norm and finite variance, $F_\pi(t,n) \le 2\exp(-\ln[Cn_m^{1/2}t+1]^\kappa)$ for some C depending on π , see the remark after Corollary 2.3 of Chamakh, Gobet and Liu (2021).

- 3. For random variables with *finitely many polynomial moments*, one has Khintchine-type inequalities (Boucheron, Lugosi and Massart (2013), Petrov (2012)), as well as Rosenthal- and Fuk-Nagaev-type inequalities (Marchina (2019), Rio (2017)).
- 4. For more heavy-tailed random variables with only a variance, Chebyshev's inequality applies to the sample mean, but there are tighter tail bounds for other mean estimators, see, for example, Catoni (2012), Lugosi and Mendelson (2019), Lugosi and Mendelson (2021).
- 2.3.3. Bernoulli processes. Here we review the definition of Bernoulli processes, which we will use later in our consistency results. For any positive integer q, a subset T of \mathbb{R}^q , and a vector $b = (b_1, \ldots, b_p)$ of independent Rademacher random variables, the map $t \mapsto t^\top b$, with $t \in T$, is referred to as a Bernoulli process (also called a Rademacher process, especially in learning theory) see, for example, Boucheron, Lugosi and Massart (2013), Talagrand (2014).

In this case, for any function class $F_m = \{f_m^* = (f_{m,1}, \ldots, f_{m,n_m})\}$, such that each $f_{m,i}$: $\mathbb{R}^{p_m} \mapsto \mathbb{R}$ is an odd function, and any random vectors $N_m = (N_{m,1}, \ldots, N_{m,n_m})$ that are mutually independent and sign-symmetric, that is, $N_{m,j} =_d - N_{m,j}$ for all $j \in [n_m]$, for i.i.d. signflips b_1, \ldots, b_{n_m} , conditional on $N_{m,i} \in \{\pm N_{m,i}^0\}$ for fixed $N_{m,i}^0$, $i \in [n_m]$, the randomization distribution $(b_1 N_{m,1}, \ldots, b_{n_m} N_{m,n_m})$ for test statistics of the form

$$f_m(N_m) = \sup_{f_m^* \in F_m} \sum_{i=1}^{n_m} f_{m,i}(N_{m,i})$$

is a Bernoulli process. Indeed, one can take $q = n_m$, and the index set $T = \{(f_{m,1}(N_{m,1}^0), \ldots, f_{m,n_m}(N_{m,n_m}^0)) : f_m^* \in F_m\}.$

The fundamental result for bounding expectations of suprema of Bernoulli processes is the Bednorz-Latala theorem (Bednorz and Latała (2013)), see also Proposition 5.14 & Theorem 5.1.5 in Talagrand (2014) for an expository presentation. Consider a subset T of \mathbb{R}^q for some q > 0 and a vector $b = (b_1, \ldots, b_q)$ of i.i.d. Rademacher random variables. Then, for $Z \sim \mathcal{N}(0, I_q)$, the Bernoulli complexity of T is characterized as

$$b(T) := \mathbb{E} \sup_{t \in T} t^{\top} b \sim \inf \Big\{ \mathbb{E} \sup_{t \in T_1} t^{\top} Z + \sup_{t \in T_2} \|t\|_1 : T \subset T_1 + T_2 \Big\}.$$

²This may nor may not satisfy the triangle inequality, see Chamakh, Gobet and Liu (2021) for discussion.

In turn, the Gaussian complexity $\mathbb{E} \sup_{t \in T_1} t^\top Z$ is characterized up to constants by the generic chaining (Talagrand (2014)).

Further, Bernoulli processes concentrate around their mean with a sub-Gaussian tail: assuming $T \subset B(t_0, \sigma)$ (where B(x, r) is the ℓ_2 ball of radius r centered at x), for any u > 0,

$$P\left(\left|\sup_{t\in T} t^{\top} b - b(T)\right| \ge u\right) \le c \exp\left(-cu^2/\sigma^2\right),$$

for a universal constant c, see Theorem 5.3.2 in Talagrand (2014). We define the infimum of the radii of all ℓ_2 balls containing the set T as the radius r(T) of T. Further, for any scalar l, we denote

(5)
$$U^{+}(T,l) := b(T) + l \cdot r(T).$$

The above results imply that, for any sequence of positive integers $(q_m)_{m\geq 1}$, any sequence of sets $(T_m)_{m\geq 1}$ with $T_m \subset \mathbb{R}^{q_m}$, and any sequence $(l_m)_{m\geq 1}$ such that $l_m > 0$ for all m and $l_m \to \infty$ as $m \to \infty$, $P(\sup_{t \in T_m} t^\top b \leq U^+(T_m, l_m)) \to 1$. In principle, these results provide basic tools to control the tails of Bernoulli processes. However, they can require some work to use in specific cases; thus, more specific results (which we will discuss later) are of interest.

- **3. Examples.** In this section, we apply our theory to several important statistical problems. Our results allow us to determine consistency conditions in a broad range of settings. Beyond the examples below, we also provide results for two-sample testing in Section 2 of the Supplementary Material (Dobriban (2022b)).
- 3.1. Detecting sparse vectors. Our first example is the fundamental statistical problem of sparse vector detection. We make n_m noisy observations $X_{m,i}$, $i=1,\ldots,n_m$ of a signal vector s_m . We assume that the signal vector is either zero, or "sparse" in the sense that it has only a few nonzero coordinates. We are interested to detect—or test—if there is indeed a nonzero signal buried in the noisy observations. This is challenging due to the potentially large and unknown level of noise. Randomization tests can be useful, because they do not require the user to know the level of noise. Indeed, they only require one to know some symmetries of the noise, and automatically adapt to the other nuisance parameters such as the noise level.

Formally, we observe n_m vectors $X_{m,i} = s_m + N_{m,i}$, $i = 1, ..., n_m$ of dimension p_m , which are sampled from a signal-plus-noise model. We arrange them into an $n_m \times p_m$ matrix X_m , which has the form $X_m = 1_{n_m} s_m^\top + N_m$. We are interested to detect "sparse" vectors s_m ; more specifically, we are interested to test against s_m with a large ℓ_∞ norm $||s_m||_\infty$. We use the test statistic $f_m(X_m) = n_m^{-1} ||1_{n_m}^\top X_m||_\infty$.

3.1.1. Sign-symmetric noise. Based on specific assumptions on the noise, various different randomization tests are valid. To illustrate our theory, we will make the relatively weak nonparametric assumption that the noise vectors $(N_{m,i})_{i \in [n_m]}$ are mutually independent, and the distribution of each noise vector N_m is sign-symmetric, independently of all other noise vectors, that is, for any vector $b \in \{\pm 1\}^{n_m}$, $(N_{m,1}, \ldots, N_{m,n_m}) =_d (b_1 N_{m,1}, \ldots, b_{n_m} N_{m,n_m})$.

We consider the randomization test from equation (1), where we randomly flip the sign of the datapoints K times using diagonal matrices $B_{m,i}$, i = 1, ..., K, with i.i.d. Rademacher entries on the diagonal. We have the following result.

PROPOSITION 3.1 (Consistency of randomization test for sparse vector detection). Let $X_{m,i} = s_m + N_{m,i}$, $i = 1, ..., n_m$, where s_m are p_m -dimensional signal vectors and $N_{m,i}$, $i = 1, ..., n_m$, are mutually independent vectors such that $N_{m,i} = d - N_{m,i}$. As $m \to \infty$, the

sequence of randomization tests (1) of the sequence of null hypotheses $s_m = 0$, with statistics $f_m(X_m) = n_m^{-1} \| 1_{n_m}^\top X_m \|_{\infty}$ and randomization distribution uniform over $n_m \times n_m$ diagonal matrices with independent Bernoulli entries is consistent against the sequence of alternatives with $s_m \in \Theta_{m1}$, if there is a sequence $(t_m)_{m\geq 1}$ such that with probability tending to unity, $\|n_m^{-1}\sum_{i=1}^{n_m} N_{m,i}\|_{\infty} \le t_m$, and for any sequence $(s_m)_{m\geq 1}$ such that for all $m \ge 1$, $s_m \in \Theta_{m1}$,

(6)
$$\liminf_{m \to \infty} \frac{\|s_m\|_{\infty}}{2t_m} > 1.$$

See Section 5 of the Supplementary Material (Dobriban (2022b)) for the proof. Roughly speaking, this result shows the consistency of the signflip-based randomization test when the signal strength is at least "twice above the noise level," as formalized in equation (6). Intriguingly, Proposition 2.3 leads to the same condition; thus suggesting that the additional noise created by randomization may be small in this case. However, of course, the two results provide only sufficient conditions; not necessary ones.

Obtaining specific consistency results. Therefore, obtaining specific consistency results boils down to controlling $\|n_m^{-1}\sum_{i=1}^{n_m}N_{m,i}\|_{\infty}$, the ℓ_{∞} norm of a mean of potentially non-i.i.d. random vectors. This can be accomplished under a variety of conditions, and has been widely studied in the areas of concentration inequalities and empirical processes. We need to find t_m such that $\|n_m^{-1}\sum_{i=1}^{n_m}N_{m,i}\|_{\infty} \leq t_m$ holds with probability tending to unity. Consider first the simplest setting: for all $m \geq 1$ and $i \in [m]$, $N_{m,i}$ are i.i.d. and have

Consider first the simplest setting: for all $m \ge 1$ and $i \in [m]$, $N_{m,i}$ are i.i.d. and have p_m i.i.d. coordinates sampled from a probability distribution π . Then by a union bound, the required condition holds with probability least $1 - p_m F_\pi(t_m; n_m)$, where $F_\pi(t, n) = P(|n^{-1}\sum_{i=1}^n Z_i| > t)$, with $Z_i \sim \pi$ i.i.d. for all $i \in [n]$. To ensure consistency, it is thus enough if t_m is such that $\lim_{m\to\infty} p_m F_\pi(t_m; n_m) = 0$. The tail bounds from Section 2.3.2 imply the following:

- 1. For subexponential random variables (including sub-Gaussian and bounded variables), Bernstein-type inequalities imply $\lim_{m\to\infty} p_m F_{\pi}(t_m; n_m) = 0$ if $t_m \sim \sqrt{(\log p_m)/n_m}$, assuming $t_m \leq 1$.
- 2. For random variables with a finite Ψ -Orlicz norm and finite variance, where Ψ : $\mathbb{R}^+ \mapsto \mathbb{R}^+$, $\Psi(x) = \exp(\ln[x+1]^k) 1$, the results of Chamakh, Gobet and Liu (2021) imply that $\lim_{m\to\infty} p_m F_\pi(t_m; n_m) = 0$ if $t_m \sim \exp[(\log p_m)^{1/\kappa}]/\sqrt{n_m}$.

Non-i.i.d. noise vectors with possibly dependent entries. Beyond the simplest setting of i.i.d. noise vectors with i.i.d. entries, one can consider more general, nonidentically distributed noise vectors with possibly dependent entries. The sign-symmetry requirement $N_m := (N_{m,1}, \ldots, N_{m,n_m}) =_d (b_1 N_{m,1}, \ldots, b_{n_m} N_{m,n_m})$ for the validity of the randomization test is equivalent to taking an arbitrary random vector $N_m^0 = (N_{m,1}^0, \ldots, N_{m,n_m}^0)$, and then multiplying each $N_{m,j}^0$, $j \in [n_m]$, by an independent Rademacher random variable.

To bound the tail of such a test statistic $f_m(N_m)$ for an arbitrary noise distribution, one general approach is to first condition on the "orbit" of N_m under the signflip group, $G(N_m) = \{(v_1N_{m,1}, \ldots, v_{n_m}N_{m,n_m}), v \in \{\pm 1\}^{n_m}\}$, apply a bound accounting for the random signflips (possibly using bounds on Bernoulli processes), and finally control the resulting tail bound over the unconditional distribution of N_m .

Rate-optimality. Next, using tools from Section 2.3.1, we discuss certain rate-optimality results for the randomization tests discussed in this section. In the setting of Proposition 3.1, consider P_{m0} specifying the distribution of the noise N_m , and $P_{mi} \in H_{m1}$, $i = 1, ..., M_m$. Then

$$L_m = \frac{1}{M_m} \sum_{j=1}^{M_m} \frac{p_{mj}(X_m)}{p_{m0}(X_m)} = \frac{1}{M_m} \sum_{j=1}^{M_m} \frac{p_{m0}(X_m - 1_{n_m} S_{mj}^{\top})}{p_{m0}(X_m)} = \frac{1}{M_m} \sum_{j=1}^{M_m} \prod_{i=1}^{n_m} \frac{p_{m,i,0}(X_{m,i} - S_{mj})}{p_{m,i,0}(X_{m,i})}.$$

Suppose all $N_{m,i}$ have equal distribution, with p_m i.i.d. coordinates with density π . Let $M_m = p_m$, and $S_{mj} = \tau_m \cdot e_j$, where e_j is the *j*th standard basis vector, and $\tau_m > 0$ will be chosen below. Then

$$L_m = \frac{1}{p_m} \sum_{j=1}^{p_m} \prod_{i=1}^{n_m} \frac{\pi(X_{m,i,j} - \tau_m)}{\pi(X_{m,i,j})}.$$

Thus,

$$\operatorname{Var} L_{m} = \frac{1}{p_{m}} \operatorname{Var} \left[\prod_{i=1}^{n_{m}} \frac{\pi(X_{m,i,1} - \tau_{m})}{\pi(X_{m,i,1})} \right] = \frac{1}{p_{m}} \left\{ \left(\operatorname{Var}_{Z \sim \pi} \left[\frac{\pi(Z - \tau_{m})}{\pi(Z)} \right] + 1 \right)^{n_{m}} - 1 \right\}.$$

Under appropriate regularity conditions in parametric statistical models

$$\operatorname{Var}_{Z \sim \pi} \left[\frac{\pi(Z - \tau_m)}{\pi(Z)} \right] = \chi^2 \left(\pi(\cdot - \tau_m), \pi \right) = I_{\pi} \cdot \tau_m^2 + o(\tau_m^2),$$

where $I_{\pi} = \int \pi'(x)^2/\pi(x) dx$ is the Fisher information of π (see, e.g., Polyanskiy (2019), Theorem 7.12.). Consistency requires that $\lim_{m\to\infty} \operatorname{Var}_{P_{m0}}[L_m] = \infty$, so that for any C > 0, $\lim_{m\to\infty} (1 + I_{\pi}\tau_m^2)/\log(Cp_m + 1) \ge 1$. Thus, the minimal signal strength required for detection is at least $\sim \sqrt{\log(p_m)/n_m}$. For subexponential random variables, this shows that the signflip randomization test is rate-optimal in this case.

To summarize this discussion, we can formulate the following result:

PROPOSITION 3.2 (Rate-optimality of signflip test for sparse vector detection). Under the assumptions of Proposition 3.1, suppose that $N_{m,i}$, $i=1,\ldots,n_m$, have i.i.d. entries from a distribution π that is subexponential and symmetric about zero, with a finite Fisher information. Let $\Theta_{m1}(\tau_m) = \{s_m \in \mathbb{R}^{p_m} : \|s_m\|_{\infty} \geq \tau_m\}$. The sequence of signflip-based randomization tests (1) of the sequence of null hypotheses $s_m = 0$ from Proposition 3.1 is consistent against the sequence of alternatives with $s_m \in \Theta_{m1}(\tau_m)$ when $\tau_m = C\sqrt{\log(p_m)/n_m}$ for a sufficiently large constant C > 0. Moreover, when $\tau_m = o(\sqrt{\log(p_m)/n_m})$, there is no consistent sequence of tests of $s_m = 0$ against $s_m \in \Theta_{m1}(\tau_m)$.

3.1.2. Spherical noise. We also study the case of spherical noise. Since the symmetry group of the noise is larger, it turns out that is enough to have a single observation $X_m = s_m + N_m \in \mathbb{R}^{p_m}$ to obtain a consistent test for a reasonable signal strength. We consider the randomization test from equation (1), with a randomization distribution that rotates the data K times using uniformly chosen rotation matrices $O_{m,i} \in O(p_m)$, i = 1, ..., K. Thus, in contrast to the previous section where each data point was transformed individually via a sign flip, here the single data vector X_m is transformed via an arbitrary rotation.

PROPOSITION 3.3 (Consistency of orthogonal randomization test for sparse vector detection). Let $X_m = s_m + N_m$, where X_m , s_m , N_m are p_m -dimensional vectors and N_m has a spherical distribution. As $m \to \infty$, the sequence of randomization tests (1) with statistics $\|X_m\|_{\infty}$ and randomization distributions uniform over $O(p_m)$ is consistent against the sequence of alternatives with $s_m \in \Theta_{m1}$, if there is a sequence $(t_{m,2})_{m \ge 1}$ such that with probability tending to unity, $\|N_m\|_2 \le t_{m,2}$, and for any sequence $(s_m)_{m \ge 1}$ such that for all $m \ge 1$, $s_m \in \Theta_{m1}$,

(7)
$$\liminf_{m \to \infty} \frac{\|s_m\|_{\infty}/(2\log p_m)^{1/2}}{(\|s_m\|_2 + 2t_{m,2})/p_m^{1/2}} > 1.$$

See Section 5 of the Supplementary Material (Dobriban (2022b)) for the proof.

This condition is a form of *relative sparsity*: the maximal absolute coordinate $||s_m||_{\infty}$ is large compared to the ℓ_2 norm $||s_m||_2$ and to the noise level $||N_m||_2$. Proposition 2.3 leads to the condition $\lim \inf_{m\to\infty} ||s_m||_{\infty}/||N_m||_{\infty} \ge 2$. Now, one can check (and we do in the proof) that $||N_m||_{\infty} =_d ||N_m||_2 \cdot ||Z_m||_{\infty}/||Z_m||_2$, where $Z_m \sim \mathcal{N}(0, I_{p_m})$. Moreover, as we also check in the proof, $||Z_m||_{\infty} \sim (2\log p_m)^{1/2}$ and $||Z_m||_2 \sim p_m^{1/2}$. Hence, the condition for the deterministic test is (roughly)

$$\liminf_{m \to \infty} \frac{\|s_m\|_{\infty}/(2\log p_m)^{1/2}}{2t_{m,2}/p_m^{1/2}} > 1.$$

We can see that the condition is milder that (7) (compare the denominators); but may be asymptotically equivalent if $||s_m||_2 = o(t_{m,2})$.

Obtaining consistency results. Therefore, obtaining specific consistency results boils down to controlling $||N_m||_2$, the ℓ_2 norm of a spherically invariant random vector. This distribution can be completely arbitrary. We give a few examples of such random vectors in Table 1, including normal, multivariate t, and multivariate Cauchy distributions. See Fang, Kotz and Ng (2018), Chapter 3, for more examples.

1. For $Z_m \sim \mathcal{N}(0, I_{p_m})$, we have $\|Z_m\|_2^2 \sim \chi_{p_m}^2$. By the chi-squared tail bound in Lemma 8.1 of Birgé (2001), when $\Gamma_m \sim \chi_{p_m}^2$

$$\mathbb{P}\left(\Gamma_m/p_m \ge 1 + 2\sqrt{\frac{x}{p_m}} + \frac{2x}{p_m}\right) \le e^{-x}.$$

Hence, for any sequence $(l_m)_{m\geq 1}$ such that $l_m>0$ for all m and $l_m\to\infty$ as $m\to\infty$, $\Gamma_m^{1/2}\leq p_m^{1/2}(l_m^{1/2}\wedge[1+O((l_m/p_m)^{1/2})])$ with probability tending to unity. Thus, we can take

$$t_{m,2} = p_m^{1/2} (l_m^{1/2} \wedge [1 + O((l_m/p_m)^{1/2})]).$$

2. For a multivariate Cauchy distribution (and more generally a multivariate t distribution with $d_m \geq 1$ degrees of freedom), by the chi-squared tail bound in Lemma 8.1 of Birgé (2001), when $\Gamma_m \sim \chi^2_{d_m}$, $\Gamma_m/d_m \geq 1 - 2\sqrt{x/d_m}$ with probability at most $\exp(-x)$. Hence, for any sequence $(l_m)_{m\geq 1}$ such that $l_m > 0$ for all m and $l_m \to \infty$ as $m \to \infty$, $1/\Gamma_m^{1/2} \leq d_m^{1/2}(l_m^{1/2} \wedge [1 + O((l_m/d_m)^{1/4})])$ with probability tending to unity. Thus, we can take

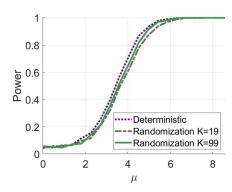
$$t_{m,2} = \frac{p_m^{1/2}(l_m^{1/2} \wedge [1 + O((l_m/p_m)^{1/2})])}{d_m^{1/2}(l_m^{1/2} \wedge [1 + O((l_m/d_m)^{1/4})])}.$$

See Section 5 of the Supplementary material (Dobriban (2022b)) for a discussion of rate-optimality.

TABLE 1

Classical examples of spherical distributions, for random vectors $Z \in \mathbb{R}^p$, for p > 0. The densities are given up to constants independent of the argument $z \in \mathbb{R}^p$, and the distribution of $\|Z\|_2^2$ is given in terms of classical distributions such as the chi-squared distribution with p degrees of freedom (χ_p^2) , and the F-distribution with p and d > 0 degrees of freedom $(F_{p,d})$

Distribution	Density	Distribution of $ Z ^2$
Normal Multivar. Cauchy	$\sim \exp(-\ z\ _2^2/2)$	χ_p^2
Multivar. t with d d.o.f.	$\sim (1 + z _2^2)^{-(p+1)/2}$ $\sim (1 + z _2^2/d)^{-(p+d)/2}$	$p \cdot F_{p,1} \ p \cdot F_{p,d}$



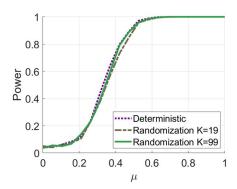


FIG. 2. Evaluating the power of the randomization test in comparison with the deterministic test as a function of signal strength in sparse vector detection. Left plot: rotation test; right plot: signflip test. See the text for details.

Numerical example. We support our theoretical result by a numerical example. We generate data from the signal-plus-noise model $X_m = s_m + N_m$, where $N_m \sim \mathcal{N}(0, I_{p_m})$, with $p_m = 100$ and $s_m = (\mu, 0, 0, \dots, 0)^{\top}$ with the signal strength parameter μ taking values over a grid of size 20 spaced equally between 0 and $4 \cdot \sqrt{\log p_m}$. We evaluate the power of the deterministic test based on $\|X_m\|_{\infty}$, tuned to have level equal to $\alpha = 0.05$. The critical value t_{α} is set so that $P_{H_{m0}}(\|X_m\|_{\infty} \ge t_{\alpha}) = 0.05$, and thus equals $t_{\alpha} = \Phi^{-1}([(1-\alpha)^{1/p_m} + 1]/2)$, where Φ^{-1} is the standard normal quantile function, that is, the inverse of the standard normal cumulative distribution function. In this case, the noise has rotational symmetry. This model can also be viewed as having sign symmetric noise (Section 3.1.1), with $n_m = p_m$, and with each $N_{m,i}$, $i = 1, \dots, n_m$ being a one-dimensional standard normal random variable.

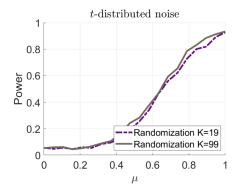
We evaluate the power of the randomization test based on K=19 and K=99 random orthogonal rotations as well as the same number of random signflips, with $\alpha=0.05$. Since the two randomization tests are appropriate under different observation models, we emphasize that this experiment does not aim to compare the two types of randomization methods. We repeat the experiment 1000 times and plot the average frequency of rejections.

On Figure 2, we observe that, as expected, the randomization tests correctly controls the level (under the null when $\mu=0$). Moreover, the power of all tests increases to unity over the range of signals considered, and the deterministic test has only slightly higher power than the randomization tests. In particular, the randomization tests achieve power almost equal to unity at almost the same point as the deterministic test. This is aligned with our results, and supports our claims that the randomization tests are near-optimal. Further, we also observe that the power with K=99 random transforms is slightly higher.

Heavy tailed example. One of the the strengths of randomization tests is that they seamlessly apply to heavy tailed noise. To illustrate this, we repeat the above experiment with t-distributed noise entries (with three and five degrees of freedom, respectively) instead of normal noise, and using the signflip randomization test. On Figure 3, we observe that the power of the randomization test increases over the range studied; but since the t distribution has heavier tails than the normal, the power increases at a slower rate than in our previous experiment, especially for the t distribution with three degrees of freedom.

3.2. Detecting spikes/low-rank matrices. A second example is the important problem of detecting low-rank matrices, which is fundamental in multivariate statistical analysis, including in PCA and factor analysis; see, for example, Anderson (1958), Dobriban (2020), Hong, Sheng and Dobriban (2020), Johnstone (2001), Johnstone and Onatski (2020), Johnstone and Paul (2018), Muirhead (2009).

Here the data X_m is represented as an $n_m \times p_m$ matrix, where often n_m is the number of samples/datapoints, and p_m is the number of features. We are interested to detect if there is



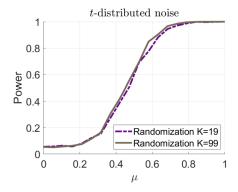


FIG. 3. Evaluating the power of the randomization test for t-distributed noise. Left plot: t distribution with three degrees of freedom; right plot: t distribution with five degrees of freedom. See the text for details.

a latent signal in the highly noisy observation matrix; and we model this by a matrix with a large operator norm. Formally, $X_m = s_m + N_m$, where s_m , N_m are $n_m \times p_m$ matrices, and we use the operator norm test statistic $f_m(X_m) = \|X_m\|_{\text{op}} = \sigma_{\max}(X_m)$. This is just one of the many possibilities. One could consider other ψ -subadditive test statistics; and in particular norms, such as the maximum absolute entry, $\max_{i,j} |X_{m,ij}|$, or generalized Ky Fan norms of the form $X \mapsto (\sum_{i=1}^{\kappa} \sigma_i(X_m)^{\zeta})^{1/\zeta}$, where $\sigma_1(X_m) \ge \dots \sigma_{n_m \wedge p_m}(X_m) \ge 0$ are the singular values of X_m , $\kappa \ge 1$, and $\zeta \ge 1$ (Li and Tsing (1988)).

As in the previous sections, there are many possible models for the structure of the noise and its corresponding group of invariances. For illustration, we only study one of them here. We consider a model where the columns of N_m are independent, and each has a spherical distribution. As in the general theory, we consider a sequence of such signal-plus-noise matrices, for a sequence of signals s_m . We can then randomize via independent uniform rotations of the columns. Recall that $||s_m||_{2,\infty}$ is the maximum of the column ℓ_2 norms of s_m .

PROPOSITION 3.4. Let the observations follow the matrix signal-plus-noise model $X_m = s_m + N_m$, where X_m , s_m , N_m are $n_m \times p_m$ -dimensional matrices and each column of N_m is independent, with a spherical distribution. As n_m , $p_m \to \infty$ such that $c_0 \le n_m/p_m \le c_1$ for arbitrary fixed $0 < c_0 < c_1$, the sequence of randomization tests (1) with test statistics $\|X_m\|_{op}$ and randomization distributions uniform over the direct product of orthogonal groups $G_m = O(n_m) \otimes O(n_m) \dots \otimes O(n_m)$ rotating the columns of the data is consistent against the sequence of alternatives with $s_m \in \Theta_{m1}$, if there is a sequence $(t_{m,2})_{m \ge 1}$ such that with probability tending to unity, $\|N_m\|_{2,\infty} \le t_{m,2}$, and for any sequence $(s_m)_{m \ge 1}$ such that for all $m \ge 1$, $s_m \in \Theta_{m1}$,

$$\liminf_{m \to \infty} \frac{\|s_m\|_{\text{op}}/(n_m^{1/2} + p_m^{1/2})}{(\|s_m\|_{2,\infty} + 2t_{m,2})/n_m^{1/2}} > 2.$$

See Section 5 of the Supplementary Material (Dobriban (2022b)) for the proof. One can verify that Proposition 2.3 implies that the deterministic test based on $\|\hat{\beta}_m\|_{\infty}$ is consistent when

$$\liminf_{m\to\infty} \frac{\|s_m\|_{\text{op}}/(n_m^{1/2}+p_m^{1/2})}{2t_{m,2}/n_m^{1/2}} > 2.$$

When $N_m \sim \mathcal{N}(0, I_{n_m} \otimes I_{p_m})$, one can verify that we can take $t_{m,2} = n_m^{1/2}(1 + o_P(1))$, thus the condition in Proposition 3.4 can be verified to simplify to

$$\liminf_{m\to\infty} [\|s_m\|_{\text{op}}/(n_m^{1/2}+p_m^{1/2})-\|s_m\|_{2,\infty}/n_m^{1/2}] > 1.$$

More generally, suppose that $N_m = [\nu_{m,1} O_{m,1}; \dots, \nu_{m,p_m} O_{m,p_m}]$, where ν_{i,m_i} , $i \in [p_m]$ are i.i.d. from a distribution with cdf F_m , and O_{i,m_i} , $i \in [p_m]$ are i.i.d. according to the Haar measure on the orthogonal group $O(p_m)$. Then the condition on $t_{m,2}$ is that $P(\max_{i=1}^{p_m} \nu_{i,m_i} \le t_{m,2}) = F_m(t_{m,2})^{p_m} \to 1$. Consider any sequence $(l_m)_{m \ge 1}$ such that $l_m > 0$ for all m and $l_m \to 0$ as $m \to \infty$. Then, we can take $t_{m,2} = F_m^{-1}(1 - l_m/p_m)$.

See Section 5 of the Supplementary Material (Dobriban (2022b)) for a rate-optimality result.

3.3. Sparse detection in linear regression. We consider the fundamental linear regression problem $Y_m = X_m \beta_m + \varepsilon_m$, where ε_m is random. The null hypothesis is that $\beta_m = 0$, and we are interested to detect "sparse" alternatives in the same way as in Section 3.1, that is, vectors β_m with a large ℓ_∞ norm.

We can directly view this as a signal plus noise model, where $s_m = X_m \beta_m$. However, the most direct approach of using a test statistic such as $f_m(Y_m) = \|Y_m\|_{\infty}$ leads to a condition for consistency that depends on the ℓ_{∞} norm $X_m \beta_m$ as opposed to β_m only. Instead, we write the ordinary least squares (OLS) estimator $\hat{\beta}_m$ as

$$\hat{\beta}_m = X_m^{\dagger} Y_m = P_{X_m} \beta_m + X_m^{\dagger} \varepsilon_m,$$

where X_m^{\dagger} is the pseudo-inverse of X_m , and P_{X_m} is the projection into the row space of X_m . Formally, this is the OLS estimator if $n_m \geq p_m$ and X_m has full rank; otherwise it is the minimum ℓ_2 norm interpolator of the normal equations $X_m^{\top}(Y_m - X_m \hat{\beta}_m) = 0$. We can view this as a signal-plus-noise model with observation $X_m' = \hat{\beta}_m$, signal $s_m = P_{X_m} \beta_m$, and noise $N_m = X_m^{\dagger} \varepsilon_m$. If $n_m \geq p_m$ and X_m has full rank, $s_m = \beta_m$, but in general this approach only provides information about the projection of β_m into the row span of X_m . We are interested to detect sparse signals using the test statistic $f_m(\hat{\beta}_m) = \|\hat{\beta}_m\|_{\infty}$.

As before, there are many possibilities for the structure of the noise. As in Section 3.1.1, we consider coordinatewise sign-symmetric noise, assuming that for any vector $b \in \{\pm 1\}^{n_m}$, $(\varepsilon_{m,1},\ldots,\varepsilon_{m,n_m})=_d(b_1\varepsilon_{m,1},\ldots,b_{n_m}\varepsilon_{m,n_m})$. We consider the randomization test from equation (1), where we randomly flip the sign of the data K times using diagonal matrices $B_{m,i}$, $i=1,\ldots,K$, with i.i.d. Rademacher entries on the diagonal. For any n_m -dimensional vector v, define the matrix

(8)
$$\mathcal{X}_m(v) = \left[X_m^{\dagger} \operatorname{diag}(v); -X_m^{\dagger} \operatorname{diag}(v) \right].$$

For $j = 1, ..., p_m$, let $[X_m^{\dagger}]_{j,.}$ be the jth row of X_m^{\dagger} . Let

(9)
$$T(X_m) = \{ \operatorname{diag}([X_m^{\dagger}]_{j,\cdot}) X_m w : w \in \mathbb{R}^{p_m}, ||w||_{\infty} \le 1, j \in [p_m] \}.$$

Define the vector $|\varepsilon_m| = (|\varepsilon_{m,1}|, \dots, |\varepsilon_{m,n_m}|)^{\top}$. Recall U^+ from (5). Below, $||M||_{\infty,\infty} = \sup_{\|v\|_{\infty} \le 1} ||Mv||_{\infty}$ is the induced matrix norm, which is also the maximum of the ℓ_1 norms of the rows of M.

PROPOSITION 3.5. Let the data (X_m, Y_m) follow the linear regression model $Y_m = X_m \beta_m + \varepsilon_m$, where Y_m is an n_m -dimensional vector of outcomes, X_m is and $n_m \times p_m$ -dimensional observation matrix, and β_m is an unknown p_m -dimensional vector of regression parameters. Let ε_m have independent entries $\varepsilon_{m,i}$, $i=1,\ldots,n_m$, such that $\varepsilon_{m,i}=_d-\varepsilon_{m,i}$. The sequence of randomization tests (1) of the null hypothesis $P_{X_m}\beta_m=0$ with test statistics $\|\hat{\beta}_m\|_{\infty}$, where $\hat{\beta}_m=X_m^{\dagger}Y_m$, and randomization distributions uniform over $n_m\times n_m$ diagonal matrices with independent Bernoulli entries is consistent against the sequence of alternatives with $P_{X_m}\beta_m\in\Theta_{m1}$, if there are two sequences $(l_m)_{m\geq 1}$ and $(t_m)_{m\geq 1}$ such that the following hold:

- 1. $l_m > 0$ for all m and $l_m \to \infty$ as $m \to \infty$,
- 2. $t_m > 0$ for all m and, with U^+ from (5) and \mathcal{X}_m from (8),

$$P(U^+(\mathcal{X}_m(|\varepsilon_m|), l_m) \leq t_m) \to 1,$$

3. for any sequence $(P_{X_m}\beta_m)_{m\geq 1}$ such that for all $m\geq 1$, $P_{X_m}\beta_m\in\Theta_{m1}$, with $T(X_m)$ from (9),

$$\liminf_{m\to\infty} \left(\|P_{X_m}\beta_m\|_{\infty} \frac{1 - U^+(T(X_m), l_m)}{2t_m} \right) > 1.$$

See Section 5 of the Supplementary Material (Dobriban (2022b)) for the proof. This result bounds $\|X_m^{\dagger}\varepsilon_m\|_{\infty}$ by an "asymmetrization" argument first, by conditioning on $|\varepsilon_m|$ and using the Bernoulli/Rademacher randomness over the signs of the entries of ε_m . However, in specific cases when more is known about the distribution of ε_m , one may obtain simpler results by directly bounding this quantity. For instance, when $\varepsilon_m \sim \mathcal{N}(0, I_{p_m})$, $X_m^{\dagger}\varepsilon_m \sim \mathcal{N}(0, X_m^{\dagger}(X_m^{\dagger})^{\top})$, and under certain structural conditions on X_m , one may be able to derive sharp bounds for the required maximum $\|X_m^{\dagger}\varepsilon_m\|_{\infty}$ of a correlated multivariate Gaussian random vector.

For comparison, one can verify that Proposition 2.3 implies that the deterministic test based on $\|\hat{\beta}_m\|_{\infty}$ is consistent when the (at least as liberal) condition $\lim \inf_{m\to\infty} \|P_{X_m}\beta_m\|_{\infty}/(2t_m) > 1$ holds. See Section 5 of the Supplementary Material (Dobriban (2022b)) for a discussion of rate-optimality.

4. Discussion. We developed a set of results on the consistency of randomization tests. While we think that our results are quite powerful, they also have limitations to be addressed in future work. A limitation is the restriction to signal plus noise models. This is needed in the current proof technique; in fact our entire approach is based on this structure. However, to broaden the scope of our results, it would be important to extend to more general statistical models.

Acknowledgments. We thank Peter Bickel, Edward I. George, Jesse Hemerik, Ilmun Kim, Panos Toulis, and Larry Wasserman for valuable discussions. We are grateful to the Associate Editor and the referees for many helpful comments that have helped to significantly improve the paper.

Funding. This work was supported in part by NSF BIGDATA grant IIS 1837992 and NSF CAREER award DMS-2046874.

SUPPLEMENTARY MATERIAL

Supplementary technical material (DOI: 10.1214/22-AOS2200SUPPA; .pdf). The supplementary technical material contains the proofs of all technical results, additional results on two-sample testing, further related work and discussion.

Code supplement (DOI: 10.1214/22-AOS2200SUPPB; .zip). The code supplement contains the scripts used to produce the numerical experiments from the paper.

REFERENCES

ALBERS, W., BICKEL, P. J. and VAN ZWET, W. R. (1976). Asymptotic expansions for the power of distribution free tests in the one-sample problem. *Ann. Statist.* **4** 108–156. MR0391373

ANDERSON, T. W. (1958). An Introduction to Multivariate Statistical Analysis. Wiley Publications in Statistics. Wiley, New York. MR0091588

- ANDERSON, M. J. and ROBINSON, J. (2001). Permutation tests for linear models. *Aust. N. Z. J. Stat.* **43** 75–88. MR1837497 https://doi.org/10.1111/1467-842X.00156
- BANKS, J., MOORE, C., VERSHYNIN, R., VERZELEN, N. and XU, J. (2018). Information-theoretic bounds and phase transitions in clustering, sparse PCA, and submatrix localization. *IEEE Trans. Inf. Theory* **64** 4872–4994. MR3819345 https://doi.org/10.1109/tit.2018.2810020
- BARRY, W. T., NOBEL, A. B. and WRIGHT, F. A. (2005). Significance analysis of functional categories in gene expression studies: A structured permutation approach. *Bioinformatics* **21** 1943–1949.
- BEDNORZ, W. and LATAŁA, R. (2013). On the suprema of Bernoulli processes. C. R. Math. Acad. Sci. Paris 351 131–134. MR3038002 https://doi.org/10.1016/j.crma.2013.02.013
- BERRY, K. J., JOHNSTON, J. E. and MIELKE, P. W. JR. (2014). A Chronicle of Permutation Statistical Methods. Springer, Cham. MR3289450 https://doi.org/10.1007/978-3-319-02744-9
- BICKEL, P. J. and VAN ZWET, W. R. (1978). Asymptotic expansions for the power of distribution free tests in the two-sample problem. *Ann. Statist.* **6** 937–1004. MR0499567
- BICKEL, P. J., KLAASSEN, C. A. J., RITOV, Y. and WELLNER, J. A. (1993). Efficient and Adaptive Estimation for Semiparametric Models. Johns Hopkins Series in the Mathematical Sciences. Johns Hopkins Univ. Press, Baltimore, MD. MR1245941
- BIRGÉ, L. (2001). An alternative point of view on Lepski's method. In State of the Art in Probability and Statistics (Leiden, 1999). Institute of Mathematical Statistics Lecture Notes—Monograph Series 36 113–133. IMS, Beachwood, OH. MR1836557 https://doi.org/10.1214/lnms/1215090065
- BOUCHERON, S., LUGOSI, G. and MASSART, P. (2013). Concentration Inequalities: A Nonasymptotic Theory of Independence. Oxford Univ. Press, Oxford. MR3185193 https://doi.org/10.1093/acprof:oso/9780199535255. 001.0001
- CATONI, O. (2012). Challenging the empirical mean and empirical variance: A deviation study. *Ann. Inst. Henri Poincaré Probab. Stat.* **48** 1148–1185. MR3052407 https://doi.org/10.1214/11-AIHP454
- CHAMAKH, L., GOBET, E. and LIU, W. (2021). Orlicz norms and concentration inequalities for β -heavy tailed random variables.
- CHAMAKH, L., GOBET, E. and SZABÓ, Z. (2020). Orlicz random Fourier features. J. Mach. Learn. Res. 21 145. MR4138129
- CHURCHILL, G. A. and DOERGE, R. W. (1994). Empirical threshold values for quantitative trait mapping. *Genetics* **138** 963–971.
- DAVID, H. A. (2008). The beginnings of randomization tests. Amer. Statist. 62 70–72. MR2416900 https://doi.org/10.1198/000313008X269576
- DIACONIS, P. (1988). Group Representations in Probability and Statistics. Institute of Mathematical Statistics Lecture Notes—Monograph Series 11. IMS, Hayward, CA. MR0964069
- DOBRIBAN, E. (2020). Permutation methods for factor analysis and PCA. Ann. Statist. 48 2824–2847. MR4152122 https://doi.org/10.1214/19-AOS1907
- DOBRIBAN, E. (2022a). Code supplement for "Consistency of invariance-based randomization tests." https://doi.org/10.1214/22-AOS2200SUPPB.
- DOBRIBAN, E. (2022b). Supplementary technical material for "Consistency of invariance-based randomization tests." https://doi.org/10.1214/22-AOS2200SUPPA.
- DWASS, M. (1957). Modified randomization tests for nonparametric hypotheses. Ann. Math. Stat. 28 181–187. MR0087280 https://doi.org/10.1214/aoms/1177707045
- EATON, M. L. (1989). Group Invariance Applications in Statistics. NSF-CBMS Regional Conference Series in Probability and Statistics. IMS, Hayward, CA. MR1089423
- EDEN, T. and YATES, F. (1933). On the validity of Fisher's z test when applied to an actual example of non-normal data. *J. Agric. Sci.* **23** 6–17.
- EFRON, B. and TIBSHIRANI, R. (2007). On testing the significance of sets of genes. *Ann. Appl. Stat.* 1 107–129. MR2393843 https://doi.org/10.1214/07-AOAS101
- EKLUND, A., NICHOLS, T. E. and KNUTSSON, H. (2016). Cluster failure: Why fMRI inferences for spatial extent have inflated false-positive rates. *Proc. Natl. Acad. Sci. USA* 113 7900–7905.
- EPSTEIN, M. P., DUNCAN, R., JIANG, Y., CONNEELY, K. N., ALLEN, A. S. and SATTEN, G. A. (2012). A permutation procedure to correct for confounders in case-control studies, including tests of rare variation. *Am. J. Hum. Genet.* **91** 215–223.
- ERNST, M. D. (2004). Permutation methods: A basis for exact inference. Statist. Sci. 19 676–685. MR2185589 https://doi.org/10.1214/088342304000000396
- FANG, K.-T., KOTZ, S. and NG, K. W. (2018). Symmetric Multivariate and Related Distributions. CRC Press/CRC, Boca Raton.
- FANG, K. T. and ZHANG, Y. T. (1990). Generalized Multivariate Analysis. Springer, Berlin. MR1079542
- FISHER, R. A. (1935). The Design of Experiments. Oliver & Boyd, Edinburgh.

- FOLLAND, G. B. (2016). A Course in Abstract Harmonic Analysis, 2nd ed. Textbooks in Mathematics. CRC Press, Boca Raton, FL. MR3444405
- FULTON, W. and HARRIS, J. (2013). Representation Theory: A First Course. Springer, Berlin.
- GIRI, N. C. (1996). Group Invariance in Statistical Inference. World Scientific, River Edge, NJ. MR1626154 https://doi.org/10.1142/9789812831705
- GOOD, P. I. (2006). Permutation, Parametric, and Bootstrap Tests of Hypotheses. Springer, Berlin.
- GUPTA, A. K. and VARGA, T. (2012). Elliptically Contoured Models in Statistics. Springer, Berlin.
- HALL, B. (2015). Lie Groups, Lie Algebras, and Representations: An Elementary Introduction, 2nd ed. Graduate Texts in Mathematics 222. Springer, Cham. MR3331229 https://doi.org/10.1007/978-3-319-13467-3
- HEMERIK, J. and GOEMAN, J. (2018a). Exact testing with random permutations. TEST 27 811–825. MR3878362 https://doi.org/10.1007/s11749-017-0571-1
- HEMERIK, J. and GOEMAN, J. J. (2018b). False discovery proportion estimation by permutations: Confidence for significance analysis of microarrays. J. R. Stat. Soc. Ser. B. Stat. Methodol. 80 137–155. MR3744715 https://doi.org/10.1111/rssb.12238
- HEMERIK, J. and GOEMAN, J. J. (2020). Another look at the lady tasting tea and differences between permutation tests and randomisation tests. *Int. Stat. Rev.*. **89**. 367-381.
- HEMERIK, J., GOEMAN, J. J. and FINOS, L. (2020). Robust testing in generalized linear models by sign flipping score contributions. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **82** 841–864. MR4112787 https://doi.org/10.1111/rssb.12369
- HEMERIK, J., SOLARI, A. and GOEMAN, J. J. (2019). Permutation-based simultaneous confidence bounds for the false discovery proportion. *Biometrika* **106** 635–649. MR3992394 https://doi.org/10.1093/biomet/asz021
- HOEFFDING, W. (1952). The large-sample power of tests based on permutations of observations. *Ann. Math. Stat.* **23** 169–192. MR0057521 https://doi.org/10.1214/aoms/1177729436
- HONG, D., SHENG, Y. and DOBRIBAN, E. (2020). Selecting the number of components in PCA via random signflips. ArXiv preprint. Available at arXiv:2012.02985.
- INGSTER, Y. I. (1987). Minimax testing of nonparametric hypotheses on a distribution density in the L_p metrics. Theory Probab. Appl. 31 333–337.
- INGSTER, Y. and SUSLINA, I. A. (2012). Nonparametric Goodness-of-Fit Testing Under Gaussian Models Springer, Berlin.
- JAMES, G. and LIEBECK, M. (2001). Representations and Characters of Groups, 2nd ed. Cambridge Univ. Press, New York. MR1864147 https://doi.org/10.1017/CBO9780511814532
- JOHNSTONE, I. M. (2001). On the distribution of the largest eigenvalue in principal components analysis. *Ann. Statist.* **29** 295–327. MR1863961 https://doi.org/10.1214/aos/1009210544
- JOHNSTONE, I. M. and ONATSKI, A. (2015). Testing in high-dimensional spiked models. *Annals of Statistics*. **48**. 1231–1254.
- JOHNSTONE, I. M. and PAUL, D. (2018). PCA in high dimensions: An orientation. *Proc. IEEE* **106** 1277–1292. https://doi.org/10.1109/JPROC.2018.2846730
- KALLENBERG, O. (2006). Probabilistic Symmetries and Invariance Principles. Springer, Berlin.
- KENNEDY, P. E. (1995). Randomization tests in econometrics. J. Bus. Econom. Statist. 13 85–94. MR1323048 https://doi.org/10.2307/1392523
- KIM, I., BALAKRISHNAN, S. and WASSERMAN, L. (2020). Robust multivariate nonparametric tests via projection averaging. Ann. Statist. 48 3417–3441. MR4185814 https://doi.org/10.1214/19-AOS1936
- KIM, I., BALAKRISHNAN, S. and WASSERMAN, L. (2020b). Minimax optimality of permutation tests. ArXiv preprint. Available at arXiv:2003.13208.
- KNAPP, A. W. (2013). Lie Groups Beyond an Introduction Springer, Berlin.
- KONING, N. W. and HEMERIK, J. (2022). Faster exact permutation testing: Using a representative subgroup. ArXiv preprint. Available at arXiv:2202.00967.
- LEHMANN, E. L. and ROMANO, J. P. (2005). *Testing Statistical Hypotheses*, 3rd ed. *Springer Texts in Statistics*. Springer, New York. MR2135927
- LEHMANN, E. L. and STEIN, C. (1949). On the theory of some nonparametric hypotheses. *Ann. Math. Stat.* **20** 28–45. MR0030168
- LI, C.-K. and TSING, N.-K. (1988). Some isometries of rectangular complex matrices. *Linear Multilinear Algebra* 23 47–53. MR0943768 https://doi.org/10.1080/03081088808817855
- LUGOSI, G. and MENDELSON, S. (2019). Mean estimation and regression under heavy-tailed distributions: A survey. *Found. Comput. Math.* 19 1145–1190. MR4017683 https://doi.org/10.1007/s10208-019-09427-x
- LUGOSI, G. and MENDELSON, S. (2021). Robust multivariate mean estimation: The optimality of trimmed mean. Ann. Statist. 49 393–410. MR4206683 https://doi.org/10.1214/20-AOS1961
- MARCHINA, A. (2019). About the rate function in concentration inequalities for suprema of bounded empirical processes. *Stochastic Process. Appl.* **129** 3967–3980. MR3997668 https://doi.org/10.1016/j.spa.2018.11.010

- MUIRHEAD, R. J. (2009). Aspects of Multivariate Statistical Theory Wiley, New York.
- NARICI, L. and BECKENSTEIN, E. (2010). Topological Vector Spaces. CRC Press, Boca Raton.
- ONGHENA, P. (2018). Randomization, Masking, and Allocation Concealment. Chapman and Hall/CRC.
- PESARIN, F. (2001a). Multivariate Permutation Tests: With Applications in Biostatistics. Wiley, Chichester. MR1855501
- PESARIN, F. (2015). Some elementary theory of permutation tests. *Comm. Statist. Theory Methods* **44** 4880–4892. MR3424814 https://doi.org/10.1080/03610926.2013.802350
- PESARIN, F. and SALMASO, L. (2010a). Permutation Tests for Complex Data: Theory, Applications and Software. Wiley, New York.
- PESARIN, F. and SALMASO, L. (2010). Finite-sample consistency of combination-based permutation tests with application to repeated measures designs. *J. Nonparametr. Stat.* **22** 669–684. MR2682214 https://doi.org/10.1080/10485250902807407
- PESARIN, F. and SALMASO, L. (2012). A review and some new results on permutation testing for multivariate problems. *Stat. Comput.* **22** 639–646. MR2865041 https://doi.org/10.1007/s11222-011-9261-0
- PESARIN, F. and SALMASO, L. (2013). On the weak consistency of permutation tests. *Comm. Statist. Simulation Comput.* **42** 1368–1379.
- PETROV, V. (2012). Sums of Independent Random Variables Springer, Berlin.
- PITMAN, E. J. (1937). Significance tests which may be applied to samples from any populations. *Suppl. J. R. Stat. Soc.* 4 119–130.
- PITMAN, E. J. G. (1939). Tests of hypotheses concerning location and scale parameters. *Biometrika* **31** 200–215. MR0000382 https://doi.org/10.1093/biomet/31.1-2.200
- POLYANSKIY, Y. (2019). Information Theoretic Methods in Statistics and Computer Science.
- Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D., Maller, J., Sklar, P., De Bakker, P. I. et al. (2007). PLINK: A tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81** 559–575.
- RIO, E. (2017). About the constants in the Fuk-Nagaev inequalities. *Electron. Commun. Probab.* **22** 28. MR3652041 https://doi.org/10.1214/17-ECP57
- SERRE, J.-P. (1977). Linear Representations of Finite Groups. Graduate Texts in Mathematics Springer, New York. MR0450380
- SIMPSON, S. L., LYDAY, R. G., HAYASAKA, S., MARSH, A. P. and LAURIENTI, P. J. (2013). A permutation testing framework to compare groups of brain networks. *Front. Comput. Neurosci.* **7** 171. https://doi.org/10.3389/fncom.2013.00171
- SOUTHWORTH, L. K., KIM, S. K. and OWEN, A. B. (2009). Properties of balanced permutations. *J. Comput. Biol.* **16** 625–638. MR2511824 https://doi.org/10.1089/cmb.2008.0144
- SUBRAMANIAN, A., TAMAYO, P., MOOTHA, V. K., MUKHERJEE, S., EBERT, B. L., GILLETTE, M. A., PAULOVICH, A., POMEROY, S. L., GOLUB, T. R. et al. (2005). Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA* **102** 15545–15550.
- TALAGRAND, M. (2014). Upper and Lower Bounds for Stochastic Processes: Modern Methods and Classical Problems. Ergebnisse der Mathematik und Ihrer Grenzgebiete. 3. Folge. A Series of Modern Surveys in Mathematics Springer, Heidelberg. MR3184689 https://doi.org/10.1007/978-3-642-54075-2
- TUSHER, V. G., TIBSHIRANI, R. and CHU, G. (2001). Significance analysis of microarrays applied to the ionizing radiation response. *Proc. Natl. Acad. Sci. USA* **98** 5116–5121.
- VAN DE GEER, S. and LEDERER, J. (2013). The Bernstein–Orlicz norm and deviation inequalities. *Probab. Theory Related Fields* **157** 225–250. MR3101846 https://doi.org/10.1007/s00440-012-0455-y
- VERSHYNIN, R. (2018). High-Dimensional Probability. Cambridge Series in Statistical and Probabilistic Mathematics Cambridge Univ. Press, Cambridge. MR3837109 https://doi.org/10.1017/9781108231596
- WIJSMAN, R. A. (1990). Invariant Measures on Groups and Their Use in Statistics. Institute of Mathematical Statistics Lecture Notes—Monograph Series 14. IMS, Hayward, CA. MR1218397
- WINKLER, A. M., RIDGWAY, G. R., WEBSTER, M. A., SMITH, S. M. and NICHOLS, T. E. (2014). Permutation inference for the general linear model. *NeuroImage* **92** 381–397.

SUPPLEMENTARY TECHNICAL MATERIAL FOR "CONSISTENCY OF INVARIANCE-BASED RANDOMIZATION TESTS"

By Edgar Dobriban¹

¹Department of Statistics and Data Science, The Wharton School, University of Pennsylvania, dobriban@wharton.upenn.edu

SUPPLEMENTARY MATERIAL

1. Additional related works. Here we discuss additional related works, which were not included in the main text due to space limitations.

A number of invariance-based randomization based tests have been developed for linear and generalized linear models (Freedman and Lane, 1983; Perry and Owen, 2010; Winkler et al., 2014; Hemerik, Goeman and Finos, 2020; Lei and Bickel, 2021). The works by Anderson and Legendre (1999); Winkler et al. (2014) review and compare a number of previously proposed permutation methods for inference in linear models with nuisance parameters. Hemerik, Thoresen and Finos (2020) show empirically that permutation tests can control type I error even in certain high dimensional linear models. Hemerik, Goeman and Finos (2020) develop tests for potentially mis-specified generalized linear models by randomly flipping signs of score contributions.

Other specific problems where invariance-based randomization tests have been developed include independence tests (Pitman, 1937), location and scale problems (Pitman, 1939), parallel analysis type methods for PCA and factor analysis (Horn, 1965; Buja and Eyuboglu, 1992; Dobriban, 2020; Dobriban and Owen, 2019), and time series data, where Jentsch and Pauly (2015) randomly permute entries between periodograms to test for equality of spectral densities. In addition, randomization based inference has been useful to study factorial designs (Pauly, Brunner and Konietschke, 2015), regression kink designs (Ganong and Jäger, 2018), and linear mixed-effects models (Rao, Drikvandi and Saville, 2019).

Most works assume exact invariance of the distribution. Romano (1990) studies the behavior of invariance-based randomization tests beyond the exact group invariance framework. This work shows that asymptotic validity holds in certain cases, and fails in others. Canay, Romano and Shaikh (2017) relax assumptions to only require a form of limiting invariance in distribution. They show that the group randomization test has an asymptotically correct level.

Chung and Romano (2013) develop general permutation tests with finite-sample error control based on studentization. Further studies include discussions of conditioning on sufficient statistics (Welch, 1990), combination methods (Pesarin, 1990), and others (Janssen and Pauls, 2003; Kim, Balakrishnan and Wasserman, 2020).

Beyond permutation tests, flipping signs is considered in many works, see e.g., Pesarin and Salmaso (2010). Following Wedderburn (1975), Langsrud (2005) discusses rotation tests in Gaussian linear regression. This approach assumes data $X_m \sim \mathcal{N}(0, I_{p_m} \otimes \Sigma_m)$, and computes the values of test statistics on $X_R = R_m X_m$, where R_m are uniformly distributed orthogonal matrices over the symmetric group $O(p_m)$. This is applied to testing independence of two random vectors, as well as to more general tests in multivariate linear regression. Perry and Owen (2010) extends the method to verify latent structure. Solari, Finos and Goeman (2014) argues for the importance of this method in multiple testing adjusting for confounding. The theoretical aspects of rotation tests for sphericity testing of densities are discussed briefly by Romano (1989), Proposition 3.2.

Toulis (2019) develops residual invariance-based randomization methods for inference in regression. This work considers a general invariance assumption $\varepsilon_m =_d g_m \varepsilon_m$ for the noise

 ε_m , for all group elements $g_m \in \mathcal{G}_m$. For ordinary least squares (OLS), it considers the test statistic $t(\hat{\varepsilon}_m) = a^\top (X_m^\top X_m)^{-1} X_m^\top g_m \hat{\varepsilon}_m$, where $\hat{\varepsilon}_m$ are the OLS residuals, and a is a vector. This work discusses many examples, including clustered observations such that the noise is correlated within clusters, proposing to flip the signs of the cluster residuals.

2. Two-sample testing. We study a two-sample testing problem, which is a classical and fundamental problem of exceeding importance in statistics, see e.g., (Lehmann and Casella, 1998; Lehmann and Romano, 2005). We study this for illustration purposes only, as there are well-established tests. We do not claim that randomization tests are better, merely that they are applicable, and it is of interest to understand what they lead to.

We consider permutation based randomization tests, valid when the entries of the noise are exchangeable. For a given integer $m \geqslant 1$ and dimension p_m , let $(f_\mu)_{\mu \in \mathbb{R}^{p_m}}$ be a location family of densities on \mathbb{R}^{p_m} . Let $\|\cdot\|_{\mathbb{R}^{p_m}}$ be a norm on \mathbb{R}^{p_m} . Let $\varepsilon_{m,i} \sim f_{0_m}$ sampled from the location family at the all-zero vector be iid for $i \in [n_m]$, and $\varepsilon'_{m,i} \sim f_{0_m}$ also be iid for $i \in [n'_m]$.

PROPOSITION 2.1. Let $Z_{m,1},\ldots,Z_{m,n_m}\sim f_{\mu_m}$, $Y_{m,1},\ldots,Y_{m,n'_m}\sim f_{\mu'_m}$ be independent observations, and test the null hypothesis that $\mu_m=\mu'_m$ against the alternative that $\mu_m\neq \mu'_m$. Consider the randomization test (1) with test statistic $\|\bar{Z}_m-\bar{Y}_{m'}\|_{\mathbb{R}^{p_m}}$, where $\bar{Z}_m=n_m^{-1}\sum_{i=1}^{n_m}Z_{m,i}$ and $\bar{Y}'_m=(n'_m)^{-1}\sum_{i=1}^{n'_m}Y_{m,i}$. For a randomization distribution uniform over the symmetric group of all permutations

For a randomization distribution uniform over the symmetric group of all permutations $S_{n_m+n'_m}$, the sequence of randomization tests (1) of the sequence of null hypotheses $\mu_m = \mu'_m$ is consistent against the sequence of alternatives with $(\mu_m, \mu'_m) \in \Theta_{m1}$, if

- 1. as $m \to \infty$, $n_m + n'_m \to \infty$,
- 2. there is a sequence $(t_m)_{m\geqslant 1}$ such that

$$P\left(\left\|\frac{1}{n_m'}\sum_{i=1}^{n_m'}\varepsilon_{m,i}'-\frac{1}{n_m}\sum_{i=1}^{n_m}\varepsilon_{m,i}\right\|_{\mathbb{R}^{p_m}}\leqslant t_m\right)\to 1,$$

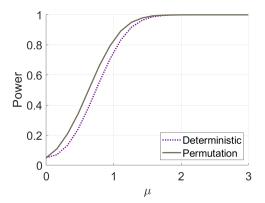
where $\varepsilon_{m,i}, \varepsilon'_{m,j} \sim f_{0_m}$, $i \in [n_m]$, $j \in [n'_m]$ are iid.

3. for any sequence $(\mu_m, \mu'_m)_{m\geqslant 1}$ such that for all $m\geqslant 1$, $(\mu_m, \mu'_m)\in\Theta_{m1}$,

$$\liminf_{m \to \infty} \frac{\|\mu'_m - \mu_m\|_{\mathbb{R}^{p_m}}}{t_m} > 2.$$

See Supplement 5.5 for the proof. As for the one-sample test for sparse detection, Proposition 2.3 leads to the same condition; thus suggesting that the additional noise due to randomization is small. The condition looks similar to the one we obtained for the one-sample test; however this concerns a different randomization distribution (permutations), and thus requires a different analysis. Bounding t_m depends on the conditions we impose on the location family, on the growth of the dimension and sample sizes, and on the specific norm used. For instance, in certain cases one may use Orlicz-norm based concentration inequalities (see e.g., Section 3.1.1 for examples), which can be adapted to the norm $\|\cdot\|_{\mathbb{R}^{p_m}}$.

Following the approach from Section 3.1.1, for $\|\cdot\|_{\mathbb{R}^{p_m}} = \|\cdot\|_{\infty}$, the same results stated there apply by assuming the same conditions on the noise vectors for both samples, and by bounding the noise vectors of the two samples separately. For instance, if the entries of $\varepsilon_{m,i}$, $i \in [n_m]$, $\varepsilon'_{m,i}$, $i \in [n'_m]$ are iid sub-exponential, then we can take $t_m \sim \sqrt{(\log p_m)/\min(n_m, n'_m)}$.



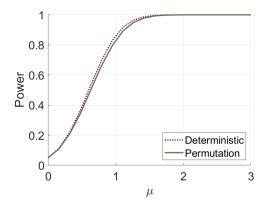


Fig 1: Evaluating the power of the permutation and deterministic versions of the t-statistic (left) and the difference-in-means statistic (right) as a function of signal strength in two-sample testing. See the text for details.

Rate-optimality. It is straightforward to see that the lower bound technique from Section 3.1.1 generalizes, and leads to a bound of the order $\tau_m = \sqrt{(\log p_m)/\min(n_m, n_m')}$. Indeed, when $n_m \leqslant n_m'$, one can take $\mu_{m',j} = 0$ and $\mu_{m,j} = \tau_m \cdot e_j$, for $j \in [n_m]$ in the construction of the alternatives in Ingster's method, and it is straightforward to see that the desired conclusion holds by the same calculation as in Section 3.1.1. This shows that for noise with iid sub-exponential entries, the signflip based randomization test is rate-optimal. To summarize:

PROPOSITION 2.2 (Rate-optimality of permutation test for sparse two-sample testing). Under the assumptions of Proposition 2.1, suppose that $\varepsilon_{m,i} \sim f_{0_m}$ for $i \in [n_m]$, and $\varepsilon'_{m,i} \sim f_{0_m}$ for $i \in [n'_m]$, each with iid entries following a sub-exponential distribution π . Let $\Theta_{m1}(\tau_m) = \{(\mu_m, \mu'_m) \in \mathbb{R}^{p_m} \times \mathbb{R}^{p_m} : \|\mu_m - \mu'_m\|_{\infty} \geqslant \tau_m\}$. The permutation test of the sequence of null hypotheses $\mu_m = \mu'_m$ from Proposition 2.1 is consistent against the sequence of alternatives with $(\mu_m, \mu'_m) \in \Theta_{m1}(\tau_m)$ when $\tau_m = C\sqrt{\log(p_m)/\min(n_m, n'_m)}$ for a sufficiently large constant C > 0. Moreover, when

$$\tau_m = o\left(\sqrt{\log(p_m)/\min(n_m, n'_m)}\right),$$

there is no consistent sequence of tests of $\mu_m = \mu'_m$ against $(\mu_m, \mu'_m) \in \Theta_{m1}(\tau_m)$, $m \ge 1$.

Numerical example. We support our theoretical result by a numerical example, using tests based on the two-sample t-statistic and the difference in means for two-sample testing. We generate data from the Gaussian signal-plus-noise model $Z_{m,i} \sim \mathcal{N}(s_m,1)$, for $i \in [n_m]$, and $Y_{m,i} \sim \mathcal{N}(0,1)$, for $i \in [n'_m]$, where $s_m = \mu$, with the signal strength parameter μ taking values over a grid of size 20 spaced equally between 0 and 3. We take $n_m = n'_m = 15$. We evaluate the power of the deterministic tests based on the two-sample t-statistic, and the mean difference of the two samples, tuned to have levels equal to $\alpha = 0.05$. We also evaluate the power of their randomization versions based on K = 99 random permutations. We repeat the experiment 100,000 times and plot the average frequency of rejections.

On Figure 1, we observe similar phenomena to those mentioned before: the randomization tests correctly control the level, and the power of all tests increases to unity over the range of signals considered. The powers of the tests are very close.² In this experiment,

¹We thank a referee for suggesting this experiment.

²We note that similar observations have been made by Lehmann (2009).

the permutation version of the two-sample t-test has a slightly higher power than the deterministic t-test, however this is reversed for the difference in means.³

3. Practical Considerations. When are invariance based tests applicable in practice? When can one invoke the group invariance hypothesis? We think that this is a challenging applied statistics problem, and we provide some discussion here. When a data analyst is performing a hypothesis test, and they have reason to think that under the null hypothesis the distribution of the data is (nearly) unchanged under some operation, then one can invoke a group invariance condition. Suppose for instance that the data analyst thinks that under the null hypothesis, the data is equally likely to have come in any order — then one can invoke permutation invariance. However, suppose that the data comes in predefined clusters (such as strata, or classes based on some key distinguishing property), and under the null hypothesis it is only reasonable to think that that data is equally likely to appear in any order in some specific clusters. Then one can use permutation invariance only over the permutations within those clusters.

This type of reasoning is more readily justifiable when testing a point null. In that case, since we only consider one distribution, assumptions can be justified with greater ease. However, if we consider composite null hypotheses, such as those in two-sample testing, then it becomes much more challenging to justify invariance assumptions.

However, one difficulty is that formally testing (evaluating) invariance assumptions can be very difficult, especially if the invariance groups are large (for instance suppose that we only have one observation; then it is impossible to test that its density is symmetric around zero).⁵ In our view these type of decisions can be quite application-specific. Further, there are a number of books and reviews on group invariance and permutation tests in statistics, and the interested statistical data analyst can study them for additional insights (see e.g., Pesarin, 2001; Ernst, 2004; Pesarin and Salmaso, 2010, 2012; Good, 2006; Kennedy, 1995; Eaton, 1989; Wijsman, 1990; Giri, 1996, etc.).

4. Proofs for the general theory.

4.1. Proof of ψ -sub-additivity in Section 2.2. Let $x,y\in\mathbb{R}$ and suppose first that $0\leqslant x< y$. Then, by concavity, $c(y)=c(x[x/y]+[x+y][1-x/y])\geqslant c(x)(x/y)+c(x+y)(1-x/y)$, or equivalently, $c(x+y)\leqslant [yc(y)-xc(x)]/[y-x]$. Thus, $c(x+y)\leqslant c(x)+c(y)$ follows if $xc(y)\leqslant yc(x)$. By concavity again, and also using that $c(0)\geqslant 0$, we have $c(x)=c(y[x/y]+0[1-x/y])\geqslant c(y)(x/y)+c(0)(1-x/y)\geqslant c(y)(x/y)$, as required. Next, if $0\leqslant x=y$, then the above argument used for y=2x shows that $c(x)\geqslant c(2x)/2$, thus $c(x+y)=c(2x)\leqslant 2c(x)=c(x)+c(y)$. This finishes the argument when $x,y\geqslant 0$. The same argument applies when $x,y\leqslant 0$.

The remaining case is when x,y have opposite signs. We can assume without loss of generality that x < 0 < y and that $|y| \ge |x|$ (otherwise we can consider (-x,-y)). Then $f(x+y) = c(|x+y|) = c(x+y) = c(y-|x|) \le c(y) + c(|x|)$, where the last inequality follows because c is non-decreasing, and also as $0 \le c(0) \le c(|x|)$.

³This may appear to be a bit surprising, and we have double checked it experimentally with an alternative implementation.

⁴We thank a reviewer for raising this question.

⁵This does not contradict that the randomization tests considered in the paper can be viewed as tests of invariance. Indeed, our claim here is that generally ascertaining invariance is difficult. The tests considered in this paper can have power to detect certain deviations from invariance, but in general they may not universally detect all deviations from invariance.

4.2. *Proof of Theorem 2.1.* **Control of type I error.** The first claim, about the level/Type I error control, is discussed at various levels of generality in many works. The textbook result, e.g., Problem 15.3 in Lehmann and Romano (2005) considers finite groups, and for infinite groups (e.g., Problem 15.1 in the same reference), assumes that we average over the full group. See also more general statements in theorem 2 in Hemerik and Goeman (2018a) and theorem 2 in Hemerik and Goeman (2018b). We provide a simple argument to show a key required exchangeability claim, which extends the above results allowing for compact topological groups at a full level of generality, and applies to random sampling of a finite number of group elements. This is crucial for our results, because we use continuous groups such as orthogonal groups in many of our examples.

Let $T_0 = f_m(N_m)$, and $T_i = f_m(G_{mi}N_m)$ for i = 1, ..., K. Note that due to noise invariance, T_i , i = 0, ..., K are exchangeable when $N_m, G_{m1}, ..., G_{mK}$ are all considered random: the random variables in the vector $L = (N_m, G_{m1}N_m, ..., G_{mK}N_m)$ are exchangeable.

LEMMA 4.1. The random vectors $\{N_m, G_{m1}N_m, ..., G_{mK}N_m\}$ are mutually exchangeable.

PROOF. To see this, we will show that $L=(N_m,G_{m1}N_m,\ldots,G_{mK}N_m)$ has the same distribution as $B=(G_mN_m,G_{m1}N_m,\ldots,G_{mK}N_m)$, where $G_m\sim Q_m$ is independent of N_m,G_{m1},\ldots,G_{mK} . Denote $G_mN_m=N_m'$. Then this is equivalent to the statement that L has the same distribution as $(N_m',G_{m1}G_m^{-1}N_m',\ldots,G_{mK}G_m^{-1}N_m')$.

Let $G'_{mi} = G_{mi}G_m^{-1}$, for i = 1, ..., K. Since $N_m =_d N'_m$, the above claim follows from because the vectors $(G_{m1}, ..., G_{mK})$ and $(G'_{m1}, ..., G'_{mK})$ have an identical distribution. For simplicity, we show this for K = 2. The proof for the more general case is very similar.

We can write for $i \neq j$, $Q_m(G'_{mi} \in M_i, G'_{mj} \in M_j) = Q_m(G_{mi}G_m^{-1} \in M_i, G_{mj}G_m^{-1} \in M_j) = Q_m(G_{mi} \in G_m M_i, G_{mj} \in G_m M_j)$. Now, let us condition on G_m . Then, we can write using the independence of G_{mi}, G_{mj} that $Q_m(G_{mi} \in G_m M_i, G_{mj} \in G_m M_j | G_m) = Q_m(G_{mi} \in G_m M_i | G_m)Q_m(G_{mj} \in G_m M_j | G_m)$. Recall that $G_{mi} \sim Q_m$ are iid from the Haar/uniform probability measure on G_m . Using the left-invariance of the Haar measure, we have $Q_m(G_{mi} \in G_m M_i | G_m) = Q_m(G_{mi} \in M_i | G_m) = Q_m(G_{mi} \in M_i)$, and similarly for j. Hence, we find, using again the independence of G_{mi}, G_{mj} that

$$Q_m(G'_{mi} \in M_i, G'_{mj} \in M_j) = Q_m(G_{mi} \in M_i)Q_m(G_{mj} \in M_j)$$

= $Q_m(G_{mi} \in M_i, G_{mj} \in M_j).$

This shows that the joint distribution of (G_{m1}, \ldots, G_{mK}) and $(G'_{m1}, \ldots, G'_{mK})$ is the same for K = 2. The same argument works for K > 2. This finishes the proof.

One can then finish the proof of type I error control as in the proof of theorem 2 in Hemerik and Goeman (2018b).

Consistency. Now we move to the part about consistency. We will consider a slight variant of the invariance-based randomization test, where for a fixed $K \ge 1$ we reject the null when

(1)
$$f_m(X_m) > \max(f_m(G_{m1}X_m), \dots, f_m(G_{mK}X_m)),$$

and where each G_{mi} , i = 1, ..., K is chosen uniformly at random over \mathcal{G}_m . The type I error probability over the random X_m and G_{mi} of this test is at most 1/(K+1), see Theorem 2.1. The consistency of this test implies the consistency of the quantile-based test. Specifically,

given any $\alpha \in (0,1)$, choose any positive integer K such that $1/(K+1) \leqslant \alpha$. Let $R_{m,K}$ denote the event (1) and let $R_{m,\alpha}$ denote the event (1). Then, $R_{m,K} \subset R_{m,\alpha}$, and hence $P_{H_{m1}}(R_{m,K}) \leqslant P_{H_{m1}}(R_{m,\alpha})$. We will show that $P_{H_{m1}}(R_{m,K}) \to 1$. Thus, it will follow that $P_{H_{m1}}(R_{m,\alpha}) \to 1$. Therefore, it is enough to study the test (1). A simplification is given by the following lemma.

LEMMA 4.2. Suppose K is fixed. Then we have

$$P(f_m(X_m) > \max_{i=1}^K f_m(G_{mi}X_m)) \to 1$$

if and only if we have $P(f_m(X_m) > f_m(G_mX_m)) \to 1$ for a single $G_m \sim Q_m$.

PROOF OF LEMMA 4.2. Consider the events $A_i = \{f_m(X_m) \leq f_m(G_{mi}X_m)\}$. By taking complements, it is enough to show that $P(\bigcup_{i=1}^K A_i) \to 0$ if and only if $P(A_1) \to 0$.

Since G_{mi} have the same distribution for all $i \in [k]$, we have $P(A_i) = P(A_j)$ for all i, j. Moreover, since $A_1 \subset \bigcup_{i=1}^K A_i$, we have by the union bound that

(2)
$$P(A_1) \leqslant P(\bigcup_{i=1}^K A_i) \leqslant \sum_{i=1}^K P(A_i) = K \cdot P(A_1).$$

Hence, as K is bounded, we have $P(\bigcup_{i=1}^K A_i) \to 0$ iff $P(A_1) \to 0$.

Thus, for consistency to hold, it is enough to show that with probability tending to unity,

$$f_m(X_m) > f_m(G_m X_m).$$

Now, $f_m(G_mX_m) = f_m(G_ms_m + G_mN_m)$. We have the following:

LEMMA 4.3 (Independence Lemma). If $g_m N_m =_d N_m$ for any fixed $g_m \in \mathcal{G}_m$, then $G_m \perp \!\!\! \perp G_m N_m$ when $G_m \sim Q_m$.

PROOF OF LEMMA 4.3. We can write, for a measurable set A

$$P(G_m N_m \in A | G_m = g_{m0}) = P(g_{m0} N_m \in A | G_m = g_{m0})$$
$$= P(g_{m0} N_m \in A) = P(N_m \in A).$$

Since this expression does not depend on g_{m0} , the distribution of $G_m N_m$ does not depend on the value of G_m ; thus $G_m N_m$ is independent of G_m .

This implies that for G_m, N_m sampled independently, $G_m s_m + G_m N_m$ has the same distribution as $G_m s_m + N_m$. Therefore, $f_m(G_m X_m) =_d f_m(G_m s_m + N_m)$, and it is enough to give conditions for the potentially stronger condition that there is a deterministic sequence of critical values t_m' such that

(3)
$$P_{H_{m1}}(f_m(G_m s_m + N_m) \leqslant t'_m) + P_{H_{m1}}(f_m(X_m) > t'_m) \to 2.$$

By ψ -subadditivity, we can write

(4)
$$f_m(X_m) = f_m(s_m + N_m) \geqslant \psi f_m(s_m) - f_m(-N_m).$$

Since t_m is such that $P(f_m(-N_m) \leqslant t_m) \to 1$, we conclude that $P(f_m(X_m) \geqslant \psi f_m(s_m) - t_m) \to 1$. Hence, if $f_m(s_m) > \psi^{-1}[t_m'(s_m) + t_m]$, then the desired condition $P_{H_{m1}}(f_m(X_m) > t_m)$

 $t_m') \to 1$ holds, provided that $P_{H_{m1}} (f_m(G_m s_m + N_m) \leqslant t_m') \to 1$. By ψ -subadditivity again, we can write

$$f_m(G_m s_m + N_m) \leq \psi^{-1} [f_m(G_m s_m) + f_m(N_m)].$$

Now, with probability tending to unity, $f_m(G_m s_m) + f_m(N_m) \leqslant \tilde{t}_m + t_m$. Hence, taking $t_m' = \psi^{-1}[\tilde{t}_m + t_m]$ finishes the proof.

Proof of Proposition 2.2. The proof of this result consists of the proof of Theorem 2.1, until equation (3). The assumption about the separating sequence ensures precisely that this condition holds, and thus finishes the proof.

Comments on the proof. The sequence $(t_m')_{m\geqslant 1}$ can be viewed as a "separating sequence", which deterministically separates the values of the original test statistic $f_m(X_m)$ from the randomized test statistic $f_m(G_mX_m)$. The current proof technique relies crucially on the existence of this sequence. However, randomization tests may be consistent even if such a sequence does not exist; thus, this step is not always sharp.

Consider for instance a sequence of observations $(X_m)_{m\geqslant 1}$, and test statistics defined by a sequence of norms $\|\cdot\|$ defined on their respective sample spaces (where the dependence of the norm on m is suppressed for readability). Suppose that randomization tests are consistent, and that there is a separating sequence $(t'_m)_{m\geqslant 1}$ such that $\|X_m\| > t'_m$ and $\|G_mX_m\| \leqslant t'_m$ both hold with probability tending to unity as $m\to\infty$.

Consider now a new observation model, where the observation X_m' equals X_m with probability 1/2, and equals $A_m X_m$ with probability 1/2; where $(A_m)_{m\geqslant 1}$ is a certain deterministic sequence of positive scalars. We choose A_m such that there is no separating sequence for X_m' . This can be accomplished by first choosing two sequences $(A_{m1})_{m\geqslant 1}$, $(A_{m2})_{m\geqslant 1}$ of positive scalars, such that $\|X_m\| \leqslant A_{m1}$ and $\|G_m X_m\| > A_{m2}$ with probability tending to unity as $m\to\infty$. The existence of A_{m1} is clear, while for A_{m2} , we only need that $G_m X_m \neq 0$ with probability tending to unity, which is a mild condition that holds in all examples we have considered. Then, we can take $A_m = A_{m2}/A_{m1}$, and it follows that, with probability tending to unity, $\|G_m (A_m X_m)\| > \|X_m\|$. Hence, there is a deterministic sequence $(t_m'')_{m\geqslant 1}$, specifically $t_m'' = A_{m1}$, such that with positive probability, $\|G_m X_m'\| > t_m''$ and $t_m'' \geqslant \|X_m'\|$, and thus there is no deterministic separating sequence for X_m' . However, the random separating sequence that equals t_m' when $X_m = X_m'$ and $A_m t_m'$ otherwise, shows that the randomization test is consistent when the observation is X_m' .

This shows that the current separating sequence approach is only sufficient and not necessary for the consistency of randomization tests.

4.3. Proof of Proposition 2.3. As in the proof of Theorem 2.1, it is enough to give conditions for the analogue of (3), i.e., that there is a deterministic sequence of critical values t'_m such that

$$P_{H_{m0}}(f_m(N_m) \leq t'_m) + P_{H_{m1}}(f_m(X_m) > t'_m) \to 2.$$

By condition 2(a) of Theorem 2.1, we can take $t_m' = t_m$, and $P_{H_{m0}}$ $(f_m(N_m) \leq t_m') \to 1$. By ψ -subadditivity, we have (4). Thus, we only need that $\psi f_m(s_m) - t_m > t_m$, which is true by (3). This shows that we can take $\tilde{c}_m \leq t_m$ and finishes the proof.

5. Proofs and discussion for the examples.

5.1. Proof of Proposition 3.1. Since $\|\cdot\|_{\infty}$ is a norm, it is 1-subadditive. Thus, the condition from Theorem 2.1 reads $n_m^{-1}\|1_{n_m}^{\top}s_m\|_{\infty} > \tilde{t}_m(s_m) + 2t_m$. Moreover, $n_m^{-1}\|1_{n_m}^{\top}s_m\|_{\infty} = \|s_m\|_{\infty}$. The requirement on t_m, \tilde{t}_m is that with probability tending to unity, one has

$$\begin{split} \|n_m^{-1} \sum_{i=1}^{n_m} N_{m,i}\|_{\infty} &\leqslant t_m, \text{ and for Rademacher random variables } b_{m,i}, \ i \in [n_m], \text{ with probability tending to unity, } \|n_m^{-1} \sum_{i=1}^{n_m} b_{m,i} s_m\|_{\infty} = |n_m^{-1} \sum_{i=1}^{n_m} b_{m,i}| \cdot \|s_m\|_{\infty} \leqslant \tilde{t}_m. \\ \text{By Hoeffding's inequality, for any } C > 0, \ P(|n_m^{-1} \sum_{i=1}^{n_m} b_{m,i}| \geqslant C) \leqslant 2 \exp(-2n_m C^2). \end{split}$$

By Hoeffding's inequality, for any C>0, $P(|n_m^{-1}\sum_{i=1}^{n_m}b_{m,i}|\geqslant C)\leqslant 2\exp(-2n_mC^2)$. Hence, we can take $\tilde{t}_m=(a_m/[2n_m])^{1/2}\cdot\|s_m\|_\infty$, for any sequence $(a_m)_{m\geqslant 1}$ with $a_m\to\infty$. Thus, the condition is that for all m large enough,

$$||s_m||_{\infty} > (a_m/[2n_m])^{1/2} \cdot ||s_m||_{\infty} + 2t_m.$$

This requires that $a_m/[2n_m] < 1$, which we can ensure holds for all large enough n_m by taking a_m to grow sufficiently slowly. For such large n_m , the condition is

$$||s_m||_{\infty} > \frac{2t_m}{1 - (a_m/[2n_m])^{1/2}}.$$

Clearly, this holds when a_m grows sufficiently slowly, for instance when $a_m = \log n_m$, if $\lim \inf_{m \to \infty} \frac{\|s_m\|_{\infty}}{2t_m} > 1$.

- 5.2. Proof and discussion of Proposition 3.3.
- 5.2.1. Proof of Proposition 3.1. Since $\|\cdot\|_{\infty}$ is a norm, it is 1-subadditive. Thus, the condition from Theorem 2.1 reads $\|s_m\|_{\infty} > \tilde{t}_m(s_m) + 2t_m$. The requirement on t_m, \tilde{t}_m is that with probability tending to unity, $\|N_m\|_{\infty} \leqslant t_m$, and for $O_m \sim O(p_m)$, with probability tending to unity, $\|O_m s_m\|_{\infty} \leqslant \tilde{t}_m$.

Now, for a normal random vector $Z_m \sim \mathcal{N}(0, I_{p_m})$, we have $\|O_m s_m\|_{\infty} =_d \|Z_m\|_{\infty}/\|Z_m\|_2$ $\|s_m\|_2$. For $Z_m \sim \mathcal{N}(0, I_{p_m})$, using standard chi-squared concentration of measure (Boucheron, Lugosi and Massart, 2013), we have $\|Z_m\|_2 = p_m^{1/2}(1+o_P(1))$. Moreover, $\|Z_m\|_{\infty} \leq (1+o_P(1))\sqrt{2\log p_m}$ with probability tending to unity. Hence, we can take $\tilde{t}_m = (1+o_P(1))(2[\log p_m]/p_m)^{1/2} \cdot \|s_m\|_2$. Similarly, $\|N_m\|_{\infty} = \|O_m N_m\|_{\infty} =_d \|Z_m\|_{\infty}/\|Z_m\|_2 \cdot \|N_m\|_2 = (1+o_P(1))(2[\log p_m]/p_m)^{1/2} \cdot \|N_m\|_2$.

Thus, the condition is that there is a sequence $t_{m,2}$ such that $P(\|N_m\|_2 \le t_{m,2}) \to 1$ and for all m large enough,

$$||s_m||_{\infty} > (1 + o_P(1))(2[\log p_m]/p_m)^{1/2} \cdot (||s_m||_2 + 2t_{m,2}).$$

This holds when

$$\liminf_{m \to \infty} \frac{\|s_m\|_{\infty}/(2\log p_m)^{1/2}}{(\|s_m\|_2 + 2t_{m,2})/p_m^{1/2}} > 1.$$

5.2.2. Discussion of rate-optimality. In this case, obtaining explicit lower bounds on detection thresholds is much more difficult. We are not aware of any results in this direction under the full level of generality of our model, and thus we discuss the difficulties here. Suppose that the noise distribution has density \tilde{p}_m with respect to the Lebesgue measure; since the distribution is rotationally invariant, we have $\tilde{p}_m(N_m) = \pi_m(\|N_m\|_2)$ for some density π_m on $[0,\infty)$. The chi-squared method shows that to achieve consistency, one must have

$$\lim_{m \to \infty} \int_{x_m \in \mathbb{R}^{p_m}} \frac{\pi_m(\|x_m - s_m\|_2)^2}{\pi_m(\|x_m\|_2)} dx_m = \infty.$$

For instance, if the noise is distributed as a multivariate t distribution with d_m degrees of freedom, with density $c_m(1+\|z\|_2^2)^{-(p_m+d_m)/2}$, where $c_m=\Gamma[(p_m+d_m)/2]/[\Gamma(d_m/2)(\pi d_m)^{p_m/2}]$, then we must show that, with $e_m=(p_m+d_m)/2$,

$$\lim_{m \to \infty} c_m \int_{x_m \in \mathbb{R}^{p_m}} \left(\frac{1 + \|x_m\|_2^2 / d_m}{(1 + \|x_m - s_m\|_2^2 / d_m)^2} \right)^{e_m} dx_m = \infty.$$

By changing variables to $x_m - s_m$, using the rotational invariance of the density, denoting $\nu_m = \|s_m\|$, we can express the integral as an expectation with respect to X_m distributed as a multivariate t distribution with d_m degrees of freedom as

$$\mathbb{E}\left(1 + \frac{\nu_m(2X_{m,p_m} + \nu_m)}{d_m + ||X_m||_2^2}\right)^{e_m}.$$

However, there does not appear to be a simple way to evaluate, or obtain sharp bounds on, this expectation, showing the difficulty of obtaining lower bounds for this problem.

5.3. Proof and discussion of Proposition 3.4.

5.3.1. Proof of Proposition 3.4. Since the maximal singular value is a norm, it is 1-subadditive. Thus, the condition from Theorem 2.1 reads $\|s_m\|_{\text{op}} > \tilde{t}_m(s_m) + 2t_m$. The requirement on t_m, \tilde{t}_m is that with probability tending to unity, $\|N_m\|_{\text{op}} \leqslant t_m$, and for $O_{m,1},\ldots,O_{m,p_m} \sim O(n_m)$, with probability tending to unity, $\|[O_{m,1}s_{m,1};\ldots;O_{m,p_m}s_{m,p_m}]\|_{\text{op}} \leqslant \tilde{t}_m$.

Now, for iid normal random vectors $Z_{m,i} \sim \mathcal{N}(0, I_{n_m})$, $i \in [p_m]$, we have $O_{m,i}s_{m,i} =_d Z_{m,i}/\|Z_{m,i}\|_2 \cdot \|s_{m,i}\|_2$. Thus,

$$\begin{split} &\|[O_{m,1}s_{m,1};\ldots;O_{m,p_m}s_{m,p_m}]\|_{\text{op}} \\ &=_{d}\|[Z_{m,1}/\|Z_{m,1}\|_{2}\cdot\|s_{m,1}\|_{2};\ldots;Z_{m,p_m}/\|Z_{m,p_m}\|_{2}\cdot\|s_{m,p_m}\|_{2}]\|_{\text{op}}. \end{split}$$

Further, for any matrix $M = [m_1; m_2; ...; m_{p_m}]$ and scalars $d_i, i \in [p_m]$,

$$||[d_1m_1; d_2m_2; \dots; d_mm_{p_m}]||_{op} \leqslant \max_i |d_i| \cdot ||M||_{op}.$$

Now, from standard concentration inequalities we have $P(|\|Z_{m,i}\|/n_m^{1/2}-1| \geqslant \delta+1/\sqrt{n_m}) \leqslant 2\exp(-n_m\delta^2/2)$. This follows from the Lipschitz concentration of Gaussian random variables, see e.g., Example 2.28 in Wainwright (2019), and from the fact that the mean of the $\chi(n_m)$ random variable $\|Z_{m,i}\|$ is bounded as $\sqrt{n_m}-1\leqslant \mathbb{E}\|Z_{m,i}\|\leqslant \sqrt{n_m}$, see exercise 3.1 in Boucheron, Lugosi and Massart (2013).

Taking a union bound, we find that $P(\max_{i=1,\dots,p_m}|\|Z_{m,i}\|_2/n_m^{1/2}-1| \geqslant \delta+1/\sqrt{n_m}) \leqslant 2\exp(\log p_m-n_m\delta^2/2)$. So, $\max_{i=1,\dots,p_m}|\|Z_{m,i}\|_2/n_m^{1/2}-1| \to_P 0$ as long as there is a sequence $\delta=\delta_m$ such that $\delta_m\to 0$ and $n_m\delta_m^2-2\log p_m\to\infty$. This holds if $\log p_m=o(n_m)$. Then, we also have that $\max_{i=1,\dots,p_m}|n_m^{1/2}/\|Z_{m,i}\|_2-1|\to_P 0$.

Thus denoting $Z_m = [Z_{m,1}; \ldots; Z_{m,p_m}]$, with probability tending to unity,

$$\begin{split} & \|[Z_{m,1}/\|Z_{m,1}\|_2 \cdot \|s_{m,1}\|_2; \dots; Z_{m,p_m}/\|Z_{m,p_m}\|_2 \cdot \|s_{m,p_m}\|_2]\|_{\text{op}} \\ & \leqslant (1+o_P(1))\|s_m\|_{2,\infty}/n_m^{1/2} \cdot \|Z_m\|_{\text{op}}. \end{split}$$

It is well known that as $n_m, p_m \to \infty$ such that $c_0 \leqslant n_m/p_m \leqslant c_1$ for some $0 < c_0 < c_1$, we have almost surely that $\|Z_m\|_{\text{op}} \leqslant (1+o_P(1))(\sqrt{n_m}+\sqrt{p_m})$. This follows from (Davidson and Szarek, 2001, Theorem 2.13). Hence, we can take $\tilde{t}_m = (1+o_P(1))\|s_m\|_{2,\infty}(\sqrt{n_m}+\sqrt{p_m})/n_m^{1/2}$.

Now, due to the distributional invariance of N_m , we have

$$||N_m||_{\text{op}} =_d ||[Z_{m,1}/||Z_{m,1}||_2 \cdot ||N_{m,1}||_2; \dots; Z_{m,p_m}/||Z_{m,p_m}||_2 \cdot ||N_{m,p_m}||_2]||_{\text{op}}$$

Hence, using the same argument as above, for any sequence $t_{m,2}$ such that $||N_m||_{2,\infty} \le t_{m,2}$ with probability tending to unity, we can take $t_m = (1 + o_P(1))(\sqrt{n_m} + \sqrt{p_m})$.

 $t_{m,2}/n_m^{1/2}$. Thus, a sufficient condition is that there is a sequence $t_{m,2}$ such that $P(\|N_m\|_{2,\infty} \le t_{m,2}) \to 1$ and

$$||s_m||_{\text{op}} > (1 + o_P(1))[1 + (p_m/n_m)^{1/2}] \cdot (||s_m||_{2,\infty} + 2t_{m,2}).$$

This holds when

$$\liminf_{m \to \infty} \frac{\|s_m\|_{\text{op}} / \left(n_m^{1/2} + p_m^{1/2}\right)}{\left(\|s_m\|_{2,\infty} + 2t_{m,2}\right) / n_m^{1/2}} > 1.$$

This finishes the proof.

5.3.2. Rate-optimality. Suppose that $N_m \sim \mathcal{N}(0, I_{n_m} \otimes I_{p_m})$, and let $\Theta_{m1} = \{\sqrt{n_m/2} \cdot v \cdot uv^\top, v \in \mathbb{R}^{n_m}, u \in \mathbb{R}^{p_m}, \|u\| = \|v\| = 1\}$. Suppose without loss of generality that $n_m \leq p_m$; otherwise flip the roles of n_m and p_m . Consider a prior Π_m on Θ_{m1} such that $u = [v; 0_{p_m - n_m}]$, and v follows a distribution Π'_m . Based on (4), we have

$$\begin{aligned} \operatorname{Var}_{P_{m0}}[L_m] &= \mathbb{E}_{S,S' \sim \Pi_m} \exp(S^{\top}S') \\ &= \mathbb{E}_{uv^{\top},u'(v')^{\top} \sim \Pi_m} \exp(n_m \tau^2 / 2 \cdot u^{\top} u' v^{\top} v') \\ &= \mathbb{E}_{v,v' \sim \Pi'_m} \exp\left(n_m \tau^2 / 2 \cdot (v^{\top} v')^2\right). \end{aligned}$$

This has the exact same form as the expression studied in Theorem 1 of Banks et al. (2018). From that result, it follows that, if Π'_m is uniform over $\{\pm 1\}^{n_m}/\sqrt{n_m}$ and $\tau < 1$, then $\mathrm{Var}_{P_{m0}}[L_m] \leqslant C$ for a constant $C < \infty$ not depending on n_m . This shows a lower bound of order $\tau \gtrsim n_m^{1/2}$. Meanwhile, our upper bound simplifies to $\tau \lesssim n_m^{1/2}$, showing that randomization tests are rate-optimal in this case.

To summarize:

PROPOSITION 5.1 (Rate-optimality of rotation test for low-rank matrix detection). Under the assumptions of Proposition 3.4, suppose that $N_m \sim \mathcal{N}(0, I_{n_m} \otimes I_{p_m})$, and let $\Theta_{m1}(\tau_m) = \left\{s_m = \sqrt{\min(n_m, p_m)/2} \cdot \tau_m \cdot uv^\top, v \in \mathbb{R}^{n_m}, u \in \mathbb{R}^{p_m}, \|u\| = \|v\| = 1\right\}$. The sequence of rotation tests (1) of the sequence of null hypotheses $s_m = 0$ from Proposition 3.4 is consistent against the sequence of alternatives with $s_m \in \Theta_{m1}(\tau_m)$ when $\tau_m = C\sqrt{\min(n_m, p_m)}$ for a sufficiently large constant C > 0. Moreover, when $\tau_m = o\left(\sqrt{\min(n_m, p_m)}\right)$, there is no consistent sequence of tests of $s_m = 0$ against $s_m \in \Theta_{m1}(\tau_m)$.

- 5.4. Proof and discussion of Proposition 3.5.
- 5.4.1. Proof of Proposition 3.5. Since the map $Y_m \mapsto \|X_m^\dagger Y_m\|_{\infty}$ is a quasi-norm, it is 1-subadditive. Thus, the condition from Theorem 2.1 reads $\|P_{X_m}\beta_m\|_{\infty} > \tilde{t}_m + 2t_m$. The requirement on t_m, \tilde{t}_m is that with probability tending to unity, $\|X_m^\dagger \varepsilon_m\|_{\infty} \leqslant t_m$, and for $B_m = \mathrm{diag}(b_{m,1},\ldots,b_{m,p_m})$ with iid Rademacher entries $b_{m,i}$, $i \in [p_m]$, with probability tending to unity, $\|X_m^\dagger B_m X_m \beta_m\|_{\infty} \leqslant \tilde{t}_m$.

Let $(l_m)_{m\geqslant 1}$ be any sequence such that $l_m>0$ for all m and $l_m\to\infty$ as $m\to\infty$. Now, conditional on the vector $|\varepsilon_m|=(|\varepsilon_{m,1}|,\ldots,|\varepsilon_{m,n_m}|)$, $X_m^\dagger\varepsilon_m$ is an n_m -dimensional Bernoulli process over the rows of the matrix $\mathcal{X}_m(|\varepsilon_m|)$. Thus, conditional on $|\varepsilon_m|$, we have $\|X_m^\dagger\varepsilon_m\|_\infty\leqslant U^+(\mathcal{X}_m(|\varepsilon_m|),l_m)$ with probability going to unity, see (5). Thus, it is enough to take t_m to be an upper bound of this quantity with probability tending to unity.

Next, writing $B_m = \text{diag}(b_m)$,

$$\begin{split} \|X_m^{\dagger} B_m X_m \beta_m \|_{\infty} &\leqslant \|X_m^{\dagger} B_m X_m \|_{\infty,\infty} \cdot \|\beta_m \|_{\infty} \\ &= \max_{j \in [p_m]} |[X_m^{\dagger}]_{j,\cdot}^{\top} \cdot B_m X_m \|_1 \cdot \|\beta_m \|_{\infty} \\ &= \max_{j \in [p_m]} \|X_m^{\top} \operatorname{diag}([X_m^{\dagger}]_{j,\cdot}) \cdot b_m \|_1 \cdot \|\beta_m \|_{\infty} \\ &= \|\beta_m \|_{\infty} \cdot \sup_{v \in T(X_m)} v^{\top} b_m. \end{split}$$

Thus, it is enough if $\tilde{t}_m = U^+(T(X_m), l_m)$. Thus, a sufficient condition is that there is a sequence $(l_m)_{m\geqslant 1}$ such that $l_m>0$ for all m and $l_m\to\infty$ as $m\to\infty$, and a sequence $(t_m)_{m\geqslant 1}$ such that $P(U^+(\mathcal{X}_m(|\varepsilon_m|), l_m)\leqslant t_m)\to 1$ and

$$\liminf_{m\to\infty} \|P_{X_m}\beta_m\|_{\infty} \frac{1-U^+(T(X_m),l_m)}{2t_m} > 1.$$

This finishes the proof.

- 5.4.2. Discussion of rate-optimality. There is a large literature on optimal hypothesis testing for linear regression, see for instance Ingster, Tsybakov and Verzelen (2010); Arias-Castro, Candès and Plan (2011); Mukherjee and Sen (2020); Carpentier and Verzelen (2021) and references therein. These works essentially only study iid Gaussian (or sub-Gaussian) noise, and make varying assumptions on the design matrix and signal strength. In general it appears quite difficult to make a direct comparison to our assumptions. For instance the work of Arias-Castro, Candès and Plan (2011) (their Theorem 2) implies that if $[X_m]_{j,\cdot}$ is the j-th row of X_m , and $(c_m)_{m\geqslant 1}$ is a sequence such that $c_m>0$ for all m and $c_m\to 0$ as $m\to \infty$, then if $X_m^{\top}X_m$ is normalized to have unit diagonal entries, if for all $i \in [p_m]$, $|\{j \in [p_m]:$ $|[X_m]_{i,\cdot}^{\top}[X_m]_{i,\cdot}| \ge c_m(\log p_m)^{-4}\}| = O(p_m^{\delta})$ for all $\delta > 0$, and if the regression coefficient β_m can be any 1-sparse vector, then it is required that $\liminf_{m\to\infty} \|\beta_m\|_{\infty}/\sqrt{2\log p_m} \geqslant 1$ in order for any test to have non-vanishing detection power. The main assumption is that for any feature, the number of other features with correlation above the level $c_m(\log p_m)^{-4}$ is smaller than any positive power of p_m . This assumption does not appear to be easily comparable to our conditions. Indeed, our conditions require (among others) to bound $||X_m^{\dagger}\varepsilon_m||_{\infty}$, where $X_m^{\dagger} \varepsilon_m \sim \mathcal{N}(0, X_m^{\dagger}(X_m^{\dagger})^{\top})$, which does not appear to be directly related to the conditions from Arias-Castro, Candès and Plan (2011). Thus, our conditions under which the randomization test works appear to be different from the ones that have been studied before for rate optimality in this problem. Since our main goal in this paper was to develop a general framework that enables proving consistency results for randomization tests, we view it as beyond our scope to fully elucidate the relationships between our conditions and those variously proposed in the literature. We would like to emphasize that our consistency results cover settings where the noise for every observation is assumed to be merely independent and symmetrically distributed, potentially heteroskedastic and heavy-tailed. This goes beyond the settings in which lower bounds have been proved for this problem.
- 5.5. Proof of Proposition 2.1. We can write $Z_{m,i} = \mu_m + \varepsilon_{m,i}$, for $i \in [n_m]$, where $\varepsilon_{m,i} \sim f_{0_m}$ are iid. Similarly, we can write $Y_{m,i} = \mu_m + \varepsilon'_{m,i}$, for $i \in [n'_m]$, where $\varepsilon'_{m,i} \sim f_{0_m}$ are also iid. We can arrange the datapoints as the rows of a matrix. This model has a signal-plus-noise form with nuisance $\mu_{m,*} = 1_{n_m + n'_m} \cdot \mu_m^\top$ and signal $S = [0_{n_m}; 1_{n'_m}] \cdot \Delta_m^\top$, where $\Delta_m = \mu'_m \mu_m$.

We can follow our general approach for problems with nuisance parameters, see Section 2. Let P_m be the projection in the orthogonal complement of the span of the nuisance. We project

$$X_m = [Z_{m,1}; \dots; Z_{m,n_m}; Y_{m,1}; \dots; Y_{m,n'_m}]$$

to $\tilde{X}_m = P_m X_m$, and we obtain a standard signal-plus-noise model $\tilde{X}_m = \tilde{s}_m + \tilde{N}_n$. Since $P_m = I_{n_m + n'_m} - 1_{n_m + n'_m} 1_{n_m + n'}^{\top} / (n_m + n'_m)$, we have

$$\tilde{X}_m = [I_{n_m + n_m'} - 1_{n_m + n_m'} \mathbf{1}_{n_m + n_m'}^\top / (n_m + n_m')] \tilde{X}_m = \tilde{X}_m - 1_{n_m + n_m'} \bar{X}_m^\top.$$

Also

$$\tilde{s}_m = P_m s_m = s_m - 1_{n_m + n_m'} \bar{s}_m^\top = [-n_m' \cdot 1_{n_m}; n_m \cdot 1_{n_m'}] / (n_m + n_m') \cdot \Delta_m^\top.$$

We can write the test statistic $\|\bar{Y}_m - \bar{Z}_m\|_{\mathbb{R}^{p_m}}$ as $\|w^\top \tilde{X}_m\|_{\mathbb{R}^{p_m}}$, where $w = [-1_{n_m}/n_m; 1_{n_m'}/n_m']$. Note that $P_m w = w$.

The test statistic is clearly 1-subadditive. Thus, the condition from Theorem 2.1 reads $\|\Delta_m\|_{\mathbb{R}^{p_m}} > \tilde{t}_m + 2t_m$. The requirement on t_m, \tilde{t}_m is that with probability tending to unity, $\|w^\top \tilde{N}_m\|_{\mathbb{R}^{p_m}} \leqslant t_m$, and for a uniformly random permutation matrix Π_m of $n_m + n'_m$ entries, with probability tending to unity, $\|w^\top \Pi_m \tilde{s}_m\|_{\mathbb{R}^{p_m}} \leqslant \tilde{t}_m$.

Now,

$$||w^{\top} \tilde{N}_{m}||_{\mathbb{R}^{p_{m}}} = ||w^{\top} N_{m}||_{\mathbb{R}^{p_{m}}} = ||\bar{Y}'_{m} - \bar{Z}_{m}||_{\mathbb{R}^{p_{m}}}$$
$$= ||(n'_{m})^{-1} \sum_{i=1}^{n'_{m}} \varepsilon'_{m,i} - n_{m}^{-1} \sum_{i=1}^{n_{m}} \varepsilon_{m,i}||_{\mathbb{R}^{p_{m}}}.$$

Also,

$$\|w^{\top}\Pi_{m}\tilde{s}_{m}\|_{\mathbb{R}^{p_{m}}} = w^{\top}\Pi_{m}w \cdot \frac{n'_{m}n_{m}}{n_{m} + n'_{m}}\|\Delta_{m}\|_{\mathbb{R}^{p_{m}}}.$$

Consider the random variable $U = w^{\top} \Pi_m w$, where the randomness is due to the random permutation matrix Π_m . Let $d = n_m + n'_m$ be the dimension of w. Now, if $\pi_m : [d] \mapsto [d]$ denotes the permutation represented by Π_m ,

$$\mathbb{E}U^2 = \mathbb{E}w^{\top} \Pi_m w \cdot w^{\top} \Pi_m w = \mathbb{E}\sum_{ij} w_i w_{\pi_m(i)} w_j w_{\pi_m(j)}$$
$$= \sum_{ij} w_i w_j \mathbb{E}w_{\pi_m(i)} w_{\pi_m(j)}.$$

If i=j, then $\mathbb{E} w_{\pi_m(i)}w_{\pi_m(j)}=\mathbb{E} w_{\pi_m(i)}^2=\|w\|^2/d$. If $i\neq j$, then, since $\sum_k w_k=0$,

$$\mathbb{E} w_{\pi_m(i)} w_{\pi_m(j)} = \frac{1}{d(d-1)} \sum_{k \neq l} w_k w_l = -\frac{\|w\|^2}{d(d-1)}.$$

Thus,

$$\mathbb{E}U^{2} = \sum_{i} w_{i}^{2} \|w\|^{2} / d + \sum_{i \neq j} w_{i} w_{j} \left(-\frac{\|w\|^{2}}{d(d-1)}\right)$$
$$= \|w\|^{4} \left(\frac{1}{d} + \frac{1}{d(d-1)}\right) = \frac{\|w\|^{4}}{d-1}.$$

Now, we can check that $||w||^2 = \frac{n_m + n'_m}{n'_m n_m}$. Therefore, by Chebyshev's inequality,

$$P(w^{\top}\Pi_{m}w \cdot \frac{n'_{m}n_{m}}{n_{m} + n'_{m}} \geqslant l_{m}) = P(U/\|w\|^{2} \geqslant l_{m})$$

$$\leq \frac{\mathbb{E}(U^{2}/\|w\|^{4})}{l_{m}^{2}} = \frac{1}{(n_{m} + n'_{m} - 1)l_{m}^{2}}.$$

Thus, if $l_m\to\infty$, we can take $\tilde{t}_m=l_m\cdot\frac{\|\Delta_m\|_{\mathbb{R}^{p_m}}}{(n_m+n'_m-1)^{1/2}}$. Thus, a sufficient condition is that there is a sequence $(l_m)_{m\geqslant 1}$ such that $l_m>0$ for all m and $l_m\to\infty$ as $m\to\infty$, and a sequence $(t_m)_{m\geqslant 1}$ such that $P(\|(n'_m)^{-1}\sum_{i=1}^{n'_m}\varepsilon'_{m,i}-n_m^{-1}\sum_{i=1}^{n_m}\varepsilon_{m,i}\|_{\mathbb{R}^{p_m}}\leqslant t_m)\to 1$ and

$$\|\Delta_m\|_{\mathbb{R}^{p_m}} > l_m \cdot \frac{\|\Delta_m\|_{\mathbb{R}^{p_m}}}{(n_m + n'_m - 1)^{1/2}} + 2t_m.$$

This requires that $n_m + n'_m \to \infty$. Then, we can take l_m to grow sufficiently slowly, and the above condition holds if

$$\liminf_{m \to \infty} \frac{\|\Delta_m\|_{\mathbb{R}^{p_m}}}{t_m} > 2.$$

This finishes the proof.

REFERENCES

- ANDERSON, M. J. and LEGENDRE, P. (1999). An empirical comparison of permutation methods for tests of partial regression coefficients in a linear model. *Journal of statistical computation and simulation* **62** 271–303.
- ARIAS-CASTRO, E., CANDÈS, E. J. and PLAN, Y. (2011). Global testing under sparse alternatives: ANOVA, multiple comparisons and the higher criticism. *The Annals of Statistics* 2533–2556.
- BANKS, J., MOORE, C., VERSHYNIN, R., VERZELEN, N. and XU, J. (2018). Information-theoretic bounds and phase transitions in clustering, sparse PCA, and submatrix localization. *IEEE Transactions on Information Theory* **64** 4872–4894.
- BOUCHERON, S., LUGOSI, G. and MASSART, P. (2013). *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press.
- BUJA, A. and EYUBOGLU, N. (1992). Remarks on parallel analysis. *Multivariate behavioral research* **27** 509-540.
- CANAY, I. A., ROMANO, J. P. and SHAIKH, A. M. (2017). Randomization tests under an approximate symmetry assumption. *Econometrica* **85** 1013–1030.
- CARPENTIER, A. and VERZELEN, N. (2021). Optimal sparsity testing in linear regression model. *Bernoulli* 27 727–750.
- CHUNG, E. and ROMANO, J. P. (2013). Exact and asymptotically robust permutation tests. *Annals of Statistics* **41** 484–507.
- DAVIDSON, K. R. and SZAREK, S. J. (2001). Local operator theory, random matrices and Banach spaces. *Handbook of the geometry of Banach spaces* 1 131.
- DOBRIBAN, E. (2020). Permutation methods for factor analysis and PCA. The Annals of Statistics 48 2824–2847.
- DOBRIBAN, E. and OWEN, A. B. (2019). Deterministic parallel analysis: an improved method for selecting factors and principal components. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **81** 163–183.
- EATON, M. L. (1989). Group invariance applications in statistics. In *Regional conference series in Probability* and *Statistics*.
- ERNST, M. D. (2004). Permutation methods: a basis for exact inference. Statistical Science 19 676-685.
- FREEDMAN, D. and LANE, D. (1983). A nonstochastic interpretation of reported significance levels. *Journal of Business & Economic Statistics* 1 292–298.
- GANONG, P. and JÄGER, S. (2018). A permutation test for the regression kink design. *Journal of the American Statistical Association* **113** 494–504.
- GIRI, N. C. (1996). Group invariance in statistical inference. World Scientific.

- GOOD, P. I. (2006). Permutation, parametric, and bootstrap tests of hypotheses. Springer Science & Business Media
- HEMERIK, J. and GOEMAN, J. J. (2018a). False discovery proportion estimation by permutations: confidence for significance analysis of microarrays. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80 137–155.
- HEMERIK, J. and GOEMAN, J. (2018b). Exact testing with random permutations. Test 27 811-825.
- HEMERIK, J., GOEMAN, J. J. and FINOS, L. (2020). Robust testing in generalized linear models by sign flipping score contributions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **82** 841–864.
- HEMERIK, J., THORESEN, M. and FINOS, L. (2020). Permutation testing in high-dimensional linear models: an empirical investigation. *Journal of Statistical Computation and Simulation* 1–18.
- HORN, J. L. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika* 30 179–185.
- INGSTER, Y. I., TSYBAKOV, A. B. and VERZELEN, N. (2010). Detection boundary in sparse regression. *Electronic Journal of Statistics* **4** 1476–1526.
- JANSSEN, A. and PAULS, T. (2003). How do bootstrap and permutation tests work? *Annals of statistics* **31** 768–806
- JENTSCH, C. and PAULY, M. (2015). Testing equality of spectral densities using randomization techniques. Bernoulli 21 697–739.
- KENNEDY, F. E. (1995). Randomization tests in econometrics. *Journal of Business & Economic Statistics* 13 85–94
- KIM, I., BALAKRISHNAN, S. and WASSERMAN, L. (2020). Minimax optimality of permutation tests. arXiv preprint arXiv:2003.13208.
- LANGSRUD, O. (2005). Rotation tests. Statistics and computing 15 53-60.
- LEHMANN, E. L. (2009). Parametric versus nonparametrics: two alternative methodologies. *Journal of Nonparametric Statistics* 21 397–405.
- LEHMANN, E. and CASELLA, G. (1998). Theory of Point Estimation. Springer Texts in Statistics.
- LEHMANN, E. L. and ROMANO, J. P. (2005). Testing statistical hypotheses. Springer Science & Business Media.
- LEI, L. and BICKEL, P. J. (2021). An assumption-free exact test for fixed-design linear models with exchangeable errors. *Biometrika* **108** 397–412.
- MUKHERJEE, R. and SEN, S. (2020). On minimax exponents of sparse testing. arXiv preprint arXiv:2003.00570.
- PAULY, M., BRUNNER, E. and KONIETSCHKE, F. (2015). Asymptotic permutation tests in general factorial designs. *Journal of the Royal Statistical Society: Series B: Statistical Methodology* 461–473.
- PERRY, P. O. and OWEN, A. B. (2010). A rotation test to verify latent structure. *Journal of Machine Learning Research* 11.
- PESARIN, F. (1990). On a nonparametric combination method for dependent permutation tests with applications. *Psychotherapy and Psychosomatics* **54** 172–179.
- PESARIN, F. (2001). Multivariate permutation tests: with applications in biostatistics. Wiley Chichester.
- PESARIN, F. and SALMASO, L. (2010). Permutation tests for complex data: theory, applications and software. John Wiley & Sons.
- PESARIN, F. and SALMASO, L. (2012). A review and some new results on permutation testing for multivariate problems. *Statistics and Computing* **22** 639–646.
- PITMAN, E. J. G. (1937). Significance tests which may be applied to samples from any populations. II. The correlation coefficient test. Supplement to the Journal of the Royal Statistical Society 4 225–232.
- PITMAN, E. (1939). Tests of hypotheses concerning location and scale parameters. Biometrika 31 200-215.
- RAO, K., DRIKVANDI, R. and SAVILLE, B. (2019). Permutation and Bayesian tests for testing random effects in linear mixed-effects models. *Statistics in medicine* **38** 5034–5047.
- ROMANO, J. P. (1989). Bootstrap and randomization tests of some nonparametric hypotheses. *The Annals of Statistics* 141–159.
- ROMANO, J. P. (1990). On the behavior of randomization tests without a group invariance assumption. *Journal of the American Statistical Association* **85** 686–692.
- SOLARI, A., FINOS, L. and GOEMAN, J. J. (2014). Rotation-based multiple testing in the multivariate linear model. *Biometrics* 70 954–961.
- Toulis, P. (2019). Life after bootstrap: Residual randomization inference in regression models. *arXiv preprint* arXiv:1908.04218.
- WAINWRIGHT, M. J. (2019). High-Dimensional Statistics: A Non-Asymptotic Viewpoint. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- WEDDERBURN, R. W. M. (1975). Random Rotations and Multivariate Normal Simulation. *Research Report, Rothamsted Experimental Station*.
- WELCH, W. J. (1990). Construction of permutation tests. *Journal of the American Statistical Association* 85 693–698.

WIJSMAN, R. A. (1990). Invariant measures on groups and their use in statistics. IMS. WINKLER, A. M., RIDGWAY, G. R., WEBSTER, M. A., SMITH, S. M. and NICHOLS, T. E. (2014). Permutation inference for the general linear model. *Neuroimage* **92** 381–397.