

1 Introduction

In nearly all areas of high-dimensional statistics distinguishing between the important and unimportant predictors is of interest. In genetics many candidate snips and their association with some outcome is of interest. Many studies take measurements on a large number of biomarkers to predict medical outcomes such as the presence or progression of a disease. While there has been much focus on finding important predictors, there has been less interest in knowing if any predictors are actually important. The vast number of potential predictors creates new statistical challenges to answering this less obvious question.

In the univariate case, once a measure of association is selected, standard approaches exist to construct asymptotically valid (and sometimes optimal) tests (?). In this setting, a test is considered optimal if it achieves the largest power for every alternative compared to every other test with the same type one error. However, scientists are often interested if any of a multitude of covariate are associated with the outcome of interest. Even in the case with two covariates, a test with the above optimal property no longer exists and the alternative plays a larger role in which tests are more or less powerful.

To illustrate this, consider two tests with equal type one error rates. Both test the null that both of two association measures ψ_1, ψ_2 are zero. The first test focuses only on the ψ_1 and the second test considers both ψ_1 and ψ_2 . When the ψ_1 is non-zero ψ_2 is zero the first test can outperform all other tests. Conversely when only ψ_2 is non-zero the first test will be outperformed by the second test. Additionally in this setting, the power of these tests will be tied to the covariance between the covariates (a problem that did not exist for a single covariate). These two factors make it difficult to define a test that performs well across all possible alternatives.

Work with **simultaneous hypothesis testing** began with Tukey in 1953 (?). Previous work by Bonferroni was used by (??) to come up with some of the first multiple hypothesis testing procedures. Further improvements were proposed by (???). Bonferroni-based correction procedures have the advantage of being easy to apply to already existing tests, while guaranteeing family-wise error control. However, these tests suffer from low power, especially in cases where the probability of rejecting each hypothesis is highly correlated.

Newer procedures (?) and **add other papers here** attain improved power compared to Bonferroni-based methods, but often rely on asymptotics to obtain these results, and don't account for the irregularity of the estimator on which their test is based. Subsequent work (???) addressed these concerns by accounting for the adaptive nature of the considered tests. Still, these newer tests are restricted to testing a particular hypothesis, must make assumptions about the data-generating mechanism to obtain theoretical guarantees.

In this article a test is proposed that works across a wide variety of data generating mechanisms and parameters of interest, but also achieves comparable power to tailor made procedures. Section ?? describes the data generating mechanisms that are considered to evaluate the performance of the test and the competing test made for the given data generating mechanism. Section ?? proposes the testing procedure. **Add other sections here**

2 Working Examples

Let X_1, \dots, X_n be independent identically distributed draws from some distribution P , and let $\mathbf{X} = \{X_1, \dots, X_n\}$. Let $X_i = (Y_i, W_{1i}, \dots, W_{di}), i \in \{1, \dots, n\}$ where Y is the outcome of interest, and each W is a covariate. Let $\psi_1 = \Psi_1(P), \dots, \psi_a = \Psi_a(P)$ be measures of association between Y and some combination of the W_j 's. While the results found in this article are valid for any integer a , for the remainder of this article assume $a = d$, the number of covariates. Also, let ψ_j correspond to a measure of association between Y and W_j . The null hypothesis for our test will be the strong null:

$$H_0 : \psi_1 = \psi_2 = \dots = \psi_a = 0 \text{ versus } H_1 : \psi_j \neq 0 \text{ for some } j \in \{1, \dots, a\}.$$

Last, let \mathcal{M} denote the set of all possible distributions. Define $\mathcal{M}_0 = \{P \in \mathcal{M} : H_0 \text{ holds}\}$

2.1 Correlation Parameter

We will compare our method to both a simple bonferroni correction method, and the method described in (?). The settings considered will be the same as the first setting in (?). The parameter of interest, $\psi_j(P)$ will be the correlation between the outcome of interest and the j 'th covariate.

The vector of covariates in this setting will be generated from a normal distribution with mean zero and a variance covariance of Σ with Σ_{ij} equal to ρ when $i \neq j$ and equal to 1 when $i = j$, and . Three different models are considered. Three different models for the outcome of interest (Y) will be considered. Letting $\varepsilon \sim N(0, 1)$ and be independent of X , the first model let $Y = \varepsilon$, the second has $Y = X_1/4$, and third has $Y = \sum_{k=1}^{10} \beta_k X_k + \varepsilon$ where $\beta_k = 0.15$ for $k = \{1, \dots, 5\}$, and $\beta_k = -0.1$ for $k = \{6 \dots 10\}$. in which Sample sizes of 100 and 200, dimensions of 10, 50, 100, 150, and 200, and ρ of 0, 0.5, or 0.8 will be considered,

2.2 Missing Data Example

In the second example, Y is binary, Δ is a missingness indicator, and each W_j is a covariate of interest. When $\Delta = 0$ we don't observe Y . The identifying assumption is $\Delta \perp\!\!\!\perp Y|W$. The parameter of interest is the risk ratio,

$$\Psi_j(P^{\text{full}}) = \frac{\text{Cov}(\log(Pr(Y = 1|W_j)), W_j)}{\text{Var}(W_j)}.$$

Using the identifying assumption, the observed data parameter is:

$$\tilde{\Psi}_j(P^{\text{obs}}) = \frac{\text{Cov}(\log(E[Pr(\Delta Y = 1|\Delta = 1, W = W)|W_j]), W_j)}{\text{Var}(W_j)}$$

Right now, it seems that the time for a single test is quite a bit larger for this test than it is for the correlation test (which is reasonable), though it may speed up if Superlearner can be improved. I am thinking of using similar simulation settings as was used before. There will be a single mechanism for the missingness, but it will be variable if the mechanism is known. There will be three settings for between W correlation. All W will be equally correlated, with $\rho = 0, 0.3$, or 0.7 . Dimension will (depending on computation time) be 10, 50, 100 and possibly 200.

There will be three settings for the β 's in the data generating model:

$$\log(Pr(Y = 1|W)) = \beta_0 + W^\top \beta$$

In one setting $\beta = 0$. In the second only one or two entries of β will be non-zero. In the last setting, many of the entries of β will be non-zero (possibly 70 or 80 percent of the entries).

2.3 Marginal Structural Model

For this marginal structural model, we are interested in if the average treatment is modified by any covariates. The marginal structural model for each W_j is defined by

$$\text{logit}(Pr(Y^{(a)} = 1|w)) = \beta_0 + \beta_1 a + \beta_2 w_j + \beta_3 w_j a$$

$$(\beta_0^*, \beta_1^*, \beta_2^*, \beta_3^*) = \text{argmin}_{\beta_0, \beta_1, \beta_2, \beta_3} \int \left(\text{logit}(Pr(Y^{(a)} = 1|w)) - (\beta_0 + \beta_1 a + \beta_2 w_j + \beta_3 w_j a) \right)^2 dP(w, a)$$

The parameter of interest is β_3^* . The other W_j 's are included in the analysis and are marginalized over, but are not part of the defined model. I am planning on using the same simulation settings as above, but also including non-changing effects for a , but am open to other ideas.

3 proposed testing procedure

As with all tests, our test will be defined by an acceptance region. Any test of $H_0 : P \in \mathcal{M}_0$ versus $H_1 : P \notin \mathcal{M}_0$ can be characterized by an acceptance region $\Theta_0(P) \subset \mathbb{R}^d$. This region $\Theta_0(P)$ can be chosen so the probability of rejection under the null is controlled asymptotically:

$$PR_Q\{Z \notin \Theta_0(p)\} = 1 - \alpha \text{ for every } P \in \mathcal{M}_0 \text{ where } Z \sim Q(P). \quad (1)$$

While there are infinitely many regions satisfying (??), we focus for now on a particular class of regions defined using ℓ_p norms which will naturally lead to a straightforward testing procedure. For simplicity, first consider regions defined using an ℓ_2 norm:

$$\Theta_0(r) = \{\omega : \|\omega\|_2 \leq r\}. \quad (2)$$

a region satisfying (??) and (??) has a radius defined by:

$$r_\alpha(P) = \min \{r : PR_Q(\|Z\|_2 \leq r) \geq 1 - \alpha\}.$$

By constraining the possible regions we consider, we can now define our test. Let $\hat{\Psi} : \mathbb{R}^d \rightarrow \mathbb{R}$ be an estimator of ψ , and let $\hat{\psi} \equiv \hat{\Psi}(x)$ be an estimate of ψ . Suppose for now that $\sqrt{n}\hat{\psi}$ converges in law to a normal distribution $Q(p)$ when p is contained in the model space \mathcal{M}_0 . Our test can be defined by

$$\text{reject } H_0 \text{ if } \|\sqrt{n}\hat{\psi}\|_2 \geq r_\alpha, \quad (3)$$

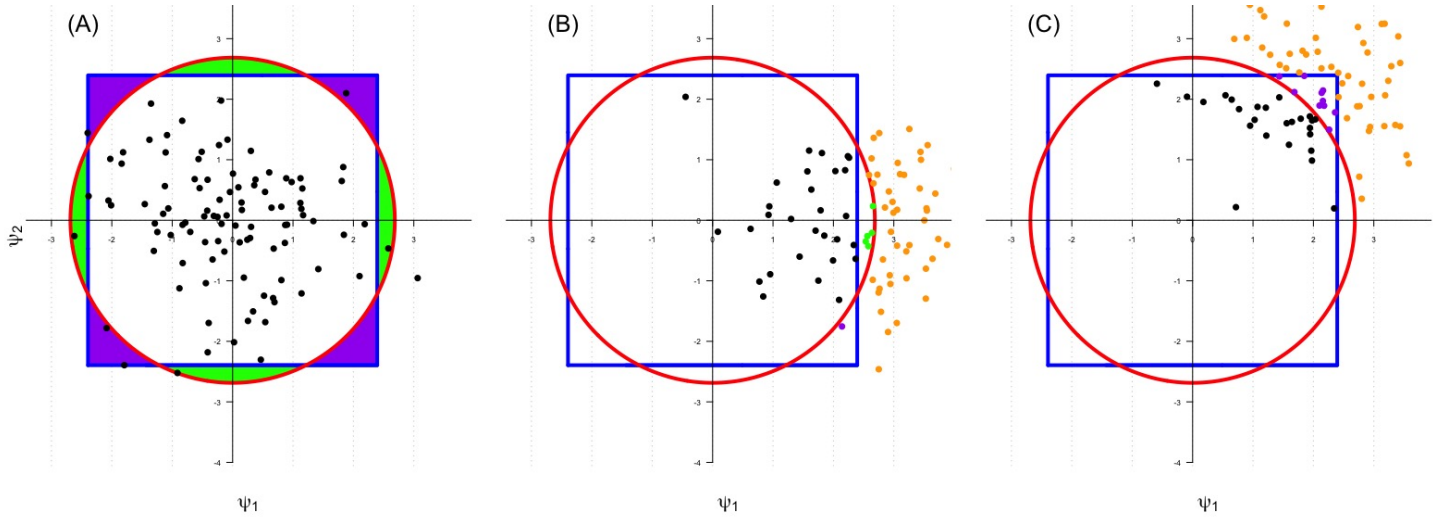


Figure 1: Plots of 100 observations from a limiting distribution of a hypothetical vector of parameter estimators in \mathbb{R}^2 (A) under the null, (B) under an alternative with $\psi_1 = 0, \psi_2 \neq 0$, and (C) under an alternative with $\psi_1, \psi_2 \neq 0$. The 95% quantiles for the data based on the max (blue) and ℓ_2 (red) norms under the null are given in all three panels. If a test statistic fell within the purple regions the test would fail to reject H_0 if the ℓ_∞ norm was used, but would reject H_0 if the ℓ_2 norm was used. The converse is true for the green regions. Depending on the alternative, the ℓ_∞ norm (B) or the ℓ_2 norm (C) will achieve higher power.

and p-values can now be defined by $\Pr_Q(\|Z\|_2 \geq \|\sqrt{n}\hat{\psi}\|_2)$. This example illustrates how functions mapping into \mathbb{R} (the ℓ_2 norm in this case) allow us to define a test and p-value for H_0 . Next consider this procedure in more generality.

Consider a function Γ , of a vector in \mathbb{R}^d (corresponding to $\sqrt{n}\hat{\psi}$) and a d dimensional distribution function (corresponding to $\sqrt{n}\hat{\psi}$'s limiting distribution). Because the limiting distribution of $\sqrt{n}\hat{\psi}$ under H_0 is always a mean zero multivariate normal, we will index our general functions Γ by the d by d variance covariance matrix instead of the infinite dimensional distribution function. An example of a more complicated function Γ_Σ is given below:

$$\Gamma_\Sigma(x) = \Pr_Q(\|\tilde{Z} + x\| > c_{0.8}) \text{ where } c_{\Sigma,0.8} \equiv \min\{c : \Pr_Q(\|\tilde{Z}\| < c) \geq 0.8\} \text{ and } \tilde{Z} \sim N(0, \Sigma). \quad (4)$$

While we have outlined a way of defining different tests of H_0 , the natural next question is which test will perform the best. To explore this question, consider a simple example comparing the test described in equation (??) with a test that is identical except it uses the maximum absolute value instead of ℓ_2 norm. Figure ??, panel (A) illustrates these two tests in \mathbb{R}^2 . One hundred draws are taken from a bivariate normal distribution with mean zero and identity covariance matrix. All observations in panel (A) except the five with the largest ℓ_2 norm are contained within the red circle. The blue squares in each panel contains all observations except the five with the largest ℓ_∞ norm. The circle and square represent the cutoffs of tests using empirical estimates of the 95th percentile of $\ell_2(Z)$ and $\ell_\infty(Z)$ respectively. Observations that fall within the purple region would result in a rejected null hypothesis if the ℓ_2 norm was used to define the test, but not if the ℓ_∞ norm was used. The converse is true of the green region. Panels B shows draws from an alternative in which $\psi_2 \neq 0$ and $\psi_1 = 0$ and the ℓ_∞ norm performs better (achieves higher power). Panel C shows draws from an alternative in which ψ_1 and $\psi_2 \neq 0$ and the ℓ_2 norm performs better.

While both acceptance regions are created so to achieve type 1 error control, depending on the alternative one test will outperform the other. Panel B shows an alternative in which only ψ_2 is non-zero. Because the max norm only considers the largest coordinate, shifting each observations in only a single direction will have larger impact on the max norm of the observations compared to the ℓ_2 norm. This trend is shown by the numerous green observations outside of the blue box (equivalent to rejecting H_0) and inside the red circle (equivalent to failing to reject H_0). In contrast, there is only a single observation that is outside the red circle, and inside the blue box. The converse trend is shown in panel C. Here, the ℓ_2 norm performs better, because it takes into account both coordinates of the shift whereas the ℓ_∞ norm can only take into account one of these coordinate shift.

The efficiency that can be gained from using the correct norm can be quite large, especially when dimension grows. Work done by (?) states that even for dimensions as small as 3 the gains in asymptotic efficiency can become arbitrarily large between two potential norms.

3.1 Adaptive selection of a norm

In the previous section we showed that a test can be defined with a summarizing function and norm. However, the choice of norm can influence the power of the test. In most scenarios it will not be clear a priori which test will have maximal power

for the true alternative. The procedure proposed in this section adaptively selects a norm and it will be shown that this procedure achieves greater power than a test with a fixed norm, while maintaining type 1 error control.

To achieve this, first define $\Gamma_{1,\Sigma}(x), \dots, \Gamma_{p,\Sigma}(x)$ as collection of functions in which the functions only differ by the norm used in their definition. Next, define

$$\Gamma_{\Sigma}^*(x) = \max \{ \Gamma_{1,\Sigma}(x), \Gamma_{2,\Sigma}(x), \dots, \Gamma_{p,\Sigma}(x) \}.$$

This could be the min if smaller values indicate larger distances away from the norm (such as p-values).

While this function is more complicated than before, $\Gamma_{\Sigma}^*(Z)$ can still be compared to $\Gamma_{\Sigma}^*(\sqrt{n}\hat{\psi})$ to obtain a p-value. Also, while it may be difficult to obtain the exact distribution of $\Gamma_{\Sigma}^*(Z)$, the distribution is a function of Σ , so obtaining good approximations of $\Gamma_{\Sigma}^*(Z)$ is possible by taking many draws from Z .

3.2 Obtaining the null distribution

The above procedure requires knowledge of the limiting distribution of $\sqrt{n}\hat{\psi}$ when $P \in \mathcal{M}_0$. To obtain an estimate of this limiting distribution, assume the vector of parameter estimates $\hat{\psi}$, converges to a normal distribution with an estimable variance covariance matrix when properly centered and normalized. Also assume each of the estimators $\hat{\psi}_1, \dots, \hat{\psi}_d$ of ψ_1, \dots, ψ_d is asymptotically linear. That is for each $j \in \{1, \dots, d\}$:

$$\hat{\psi}_j = \psi_j + \frac{1}{n} \sum_{i=1}^n D_j(\mathbf{x}_i) + o_p(1/\sqrt{n}) \text{ for some function } D_j$$

When there is a fixed number of covariates, the Cramer-Wold device can be used to show that the vector of parameter estimates is asymptotically normal with mean zero, and variance covariance matrix given by $\Sigma = E_{P_0} [D(X)D(X)^{\top}]$:

$$\sqrt{n}(\hat{\psi} - \psi) \xrightarrow{d} Z \sim N(0, \Sigma)$$

Under H_0 , $\sqrt{n}\hat{\psi}$ converges to Z , and Σ can be approximated with $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n D(\mathbf{x}_i)D(\mathbf{x}_i)^{\top}$. In practice, a consistent estimator of Σ , $\hat{\Sigma}$ will be used in place of Σ . Similarly, the random variable \hat{Z} with distribution \hat{Q} will be used in place of Z . Thus, the test statistic will be $\Gamma_{\hat{\Sigma}}(\sqrt{n}\hat{\psi})$ and will be compared to $\Gamma_{\hat{\Sigma}}(\hat{Z})$.

3.3 Using a permutation test for the test statistic

While the above approach works asymptotically, there can be issues for small sample sizes. To avoid inflated type one error, a permutation based test can be used. Here, $\hat{Q}^{\#}$ is used to define Γ , and the test will compare $\Gamma_{\hat{Q}}(\sqrt{n}\hat{\psi})$ to $\Gamma_{\hat{Q}}(Z^{\#})$. To determine $\hat{Q}^{\#}$, the Y 's from the observed data are permuted before calculating $\hat{\Sigma}$. Draws from $Z^{\#}$ are taken by permuting all of the Y 's of the observed data. [There are a few more complications here that need to be figured out.](#)

4 Simulation Study

4.1 Correlation

In the first set of experiments, we compare our method to both a Bonferroni based method and a tailor made method. The Bonferroni based method computes a p-value for each covariate using a marginal linear regression. If any p-value is smaller than the Bonferroni adjusted cutoff value, the null hypothesis is rejected. The method outlined by (?) is similar to our method, but uses a parametric estimate of the correlation instead of using an influence function based estimator. Four versions of our test are displayed in these plots. Our test that uses the ℓ_2 norm, our test that uses the ℓ_{∞} or max norm, and our test that selects over five possible ℓ_p norms. The last version of our test selects over many different versions of the following norm : $\|x\|_k = \sum_{i=1}^k x_{(i)}^2$ where $x_{(i)}^2$ are the ordered statistics of x^2 .

In each setting, each test is run 1000 times. The proportion of tests that reject the null is plotted for each simulation setting and for each test.

We find that our procedure uses obtains relatively good power across a variety of settings. While in certain settings our test is beaten by that of Zhang and Laber, this is to be expected because their procedure is tailor made for the setting in which there is normal data that follows a linear model. We also find that the adaptive version of the test consistently outperforms the Bonferroni based method in all settings except those in which there is no between-covariate correlation.

1. Two-Phase sampling Scheme of correlation
2. Marginal structural model example.

All three examples will come with simulation results in the form of figures. Not sure how many figures to use or what sample size(s) to use.

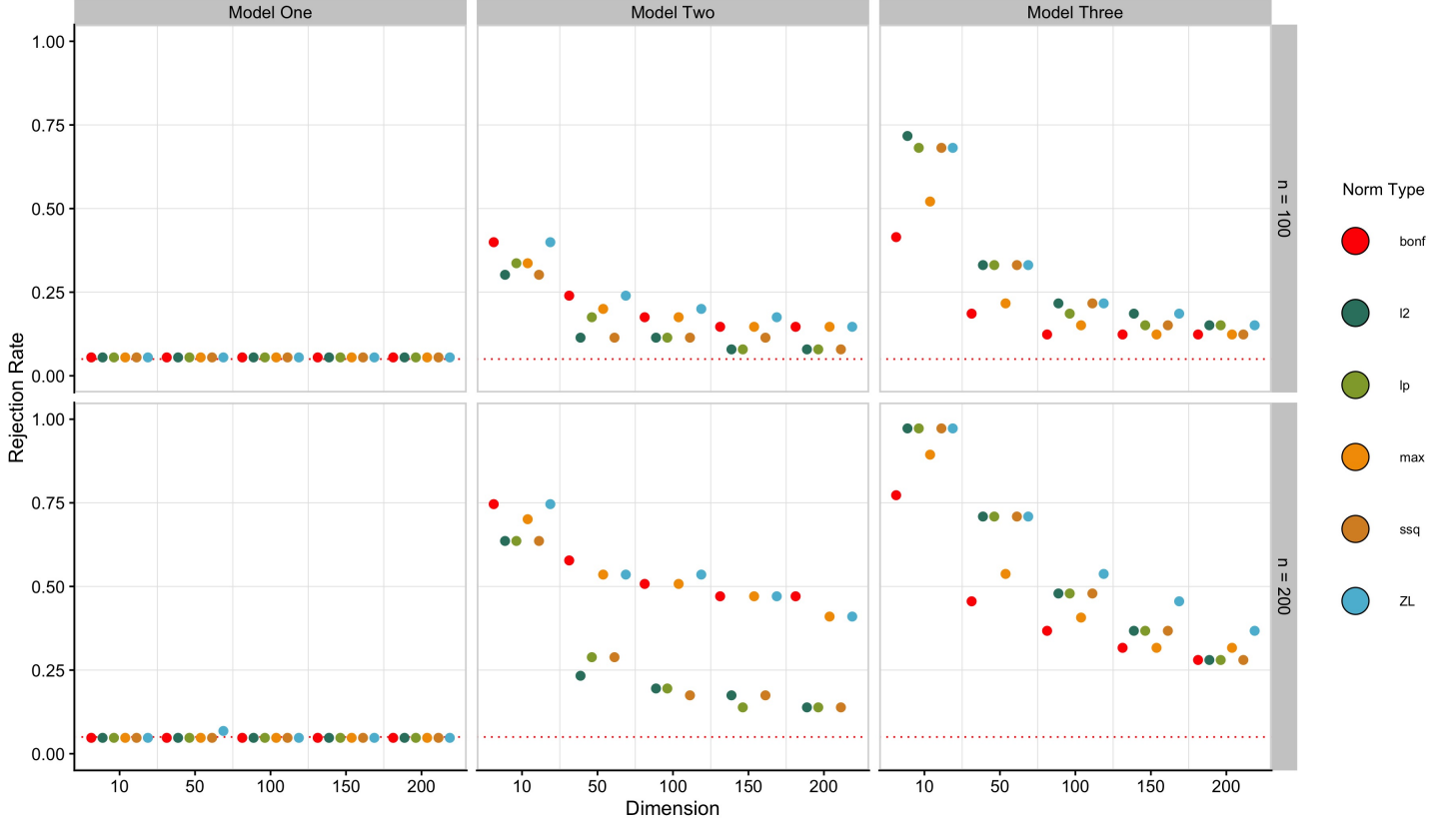


Figure 2: Display of simulations for vector of covariates in this setting will be generated from a normal distribution with mean zero and a variance covariance of Σ with Σ_{ij} equal to 0 when $i \neq j$ and equal to 1 when $i = j$. Three different models for the outcome of interest (Y) will be considered. Letting $\varepsilon \sim N(0, 1)$ and be independent of X , in the first model $Y = \varepsilon$, in the second $Y = X_1/4$, and in the third $Y = \sum_{k=1}^{10} \beta_k X_k + \varepsilon$ where $\beta_k = 0.15$ for $k = \{1, \dots, 5\}$, and $\beta_k = -0.1$ for $k = \{6 \dots 10\}$. Sample sizes of 100 and 200, dimensions of 10, 50, 100, 150, and 200 are considered.

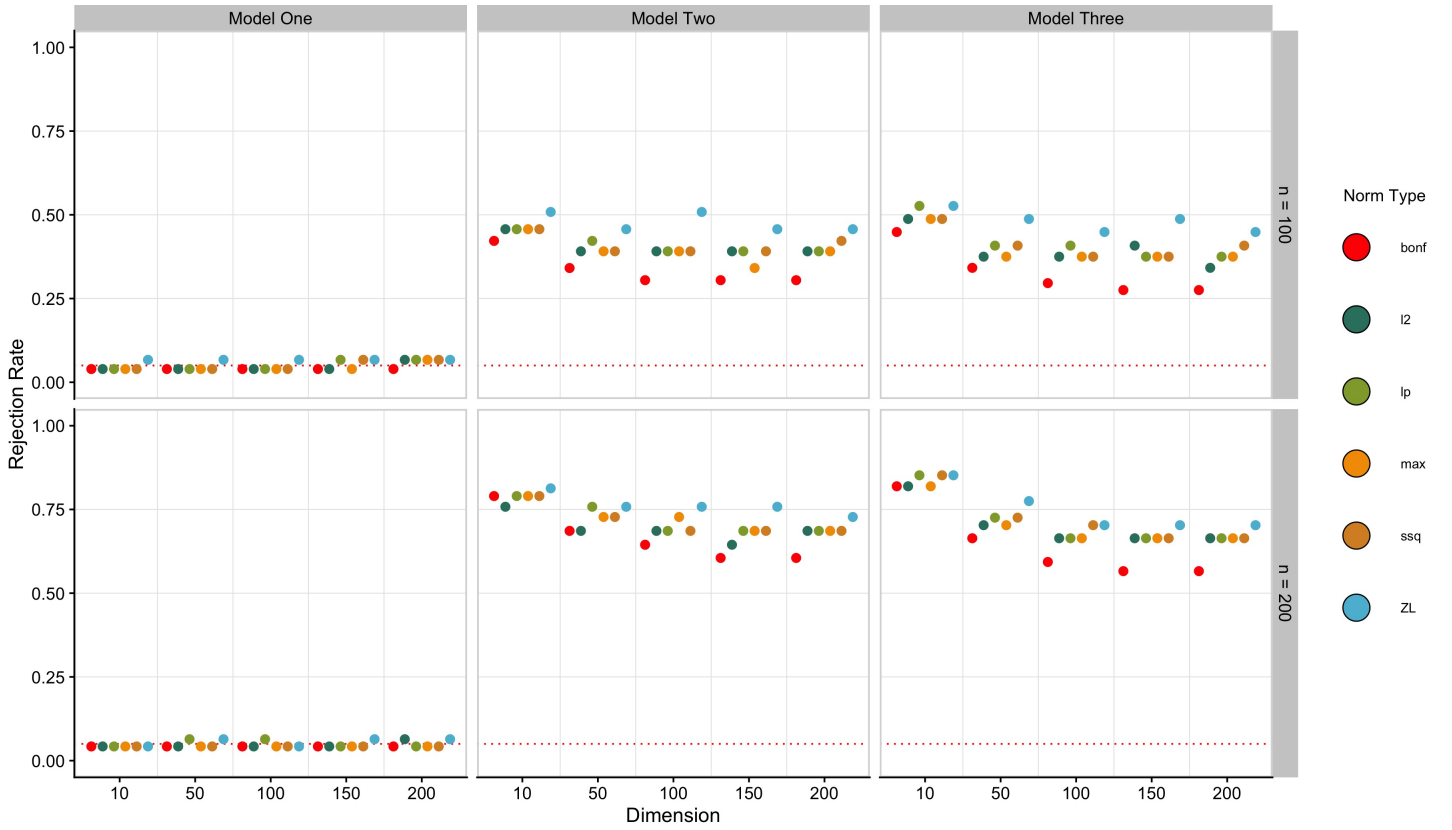


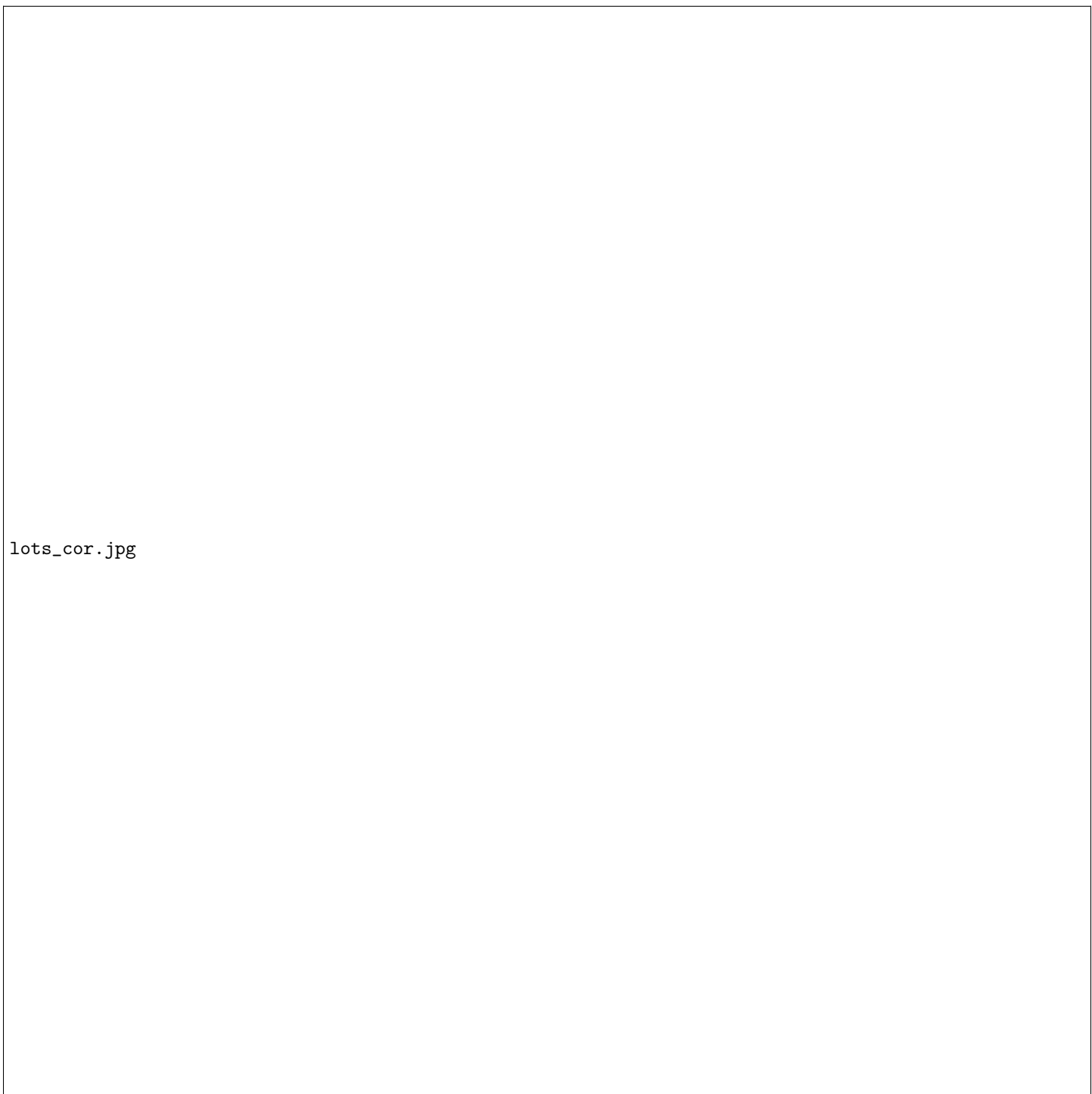
Figure 3: The same simulation settings as those used in Figure ??, but $\Sigma_{ij} = 0.5$ for $i \neq j$.

5 Data Application

1. Data from Peter Gilbert

6 Discussion

- The most generalizable methods for hypothesis testing suffer from poor power, and ignore correlation between covariates.
- Tests with high power are often difficult to generalize, working only for a single parameter, or relying on parametric assumptions.
- We find a happy medium in this article.
- Our method has comparable power to tailor made methods, without require modeling assumptions or parameters of interest.
- We showed the methods flexibility through the use of three different data examples that include missing data and a causal model.
- This method can also be expanded by using different Γ functions, potentially optimizing of entire classes of Γ functions with machine learning algorithms.
- Also look at multiplying $\sqrt{n}\hat{\psi}$ by inverse of influence function.
- There are issues with using parametric bootstrap, which can be avoided in certain cases by using a permutation test.



lots_cor.jpg

Figure 4: The same simulation settings as those used in Figure ??, but $\Sigma_{ij} = 0.8$ for $i \neq j$.

7 Conclusion

8 Appendix

8.1 Test Consistency

Theorem 1. Assume the performance metric of choice is $\hat{r}_{\alpha,p}(a)$, and each the norms g_a considered have the following properties:

$$\text{for } x \in \mathbb{R}^d, \text{ and } s, l > 0, s \cdot \max(\mathbf{x}) \leq g_a(\mathbf{x}) \leq l \cdot d \cdot \max(\mathbf{x}) \quad (5)$$

$$\text{for } s \leq 1, l \geq 1, g_a(s \cdot \mathbf{x}) \leq g_a(\mathbf{x}) \leq g_a(l \cdot \mathbf{x}) \quad (6)$$

Then for all $P \notin \mathcal{M}_0$

$$\Pr \left(\frac{1}{B} \sum_{k=1}^B I \left\{ \hat{T} \leq \hat{T}_k^\# \right\} < 0.05 \right) \xrightarrow{P} 1 \text{ as } n \rightarrow \infty$$

or Equivalently

$$\Pr \left(\hat{T} \leq F_{\hat{T}^\#}^{-1}(0.05) \right) \xrightarrow{P} 1 \text{ as } n \rightarrow \infty$$

Proof. This proof will consist of three parts. We will first show that as $n \rightarrow \infty$, $\hat{T}_a^\#$ becomes bounded away from 0 for any valid norm. Next we will show \hat{T}_a converges to 0 in probability for any valid norm. Last we will show the two previous findings imply theorem ??.

Let $P_X^\#$ denote the distribution of the randomly permuted observations. Because $\Pr(Y_i^\# \perp\!\!\!\perp \mathbf{W}_i^\#) \rightarrow 0$, $\Psi(P_X^\#) = \mathbf{0}$. Additionally, $\sqrt{n}(\hat{\psi}^\# - \psi^\#) \xrightarrow{d} Z^\# \sim N(\mathbf{0}, \Sigma_{\text{perm}})$. Define $\hat{T}_a^\# = \min_s \{s : g_a(s \cdot \sqrt{n}\hat{\psi}^\#) \geq C_{0.95,a}^\#\}$ and $C_{0.95,a}^\#$ is $F_{Z^\#}^{-1}(0.95)$. Since we know that $\psi^\# = \mathbf{0}$, it follows that $\sqrt{n}\hat{\psi}^\# \overset{d}{\approx} Z^\#$.

Now, consider:

$$\begin{aligned} \Pr \left(\hat{T}_a^\# > \epsilon \right) &= \Pr \left(g_a \left(\epsilon \cdot \sqrt{n}\hat{\psi}^\# \right) \leq C_{0.95,a} \right) \\ &\geq \Pr \left(\epsilon \cdot d \cdot \max \left(\sqrt{n}\hat{\psi}^\# \right) \leq C_{0.95,a} \right) \\ &= \Pr \left(\max \left(\sqrt{n}\hat{\psi}^\# \right) \leq C_{0.95,a} / (\epsilon \cdot d) \right) \end{aligned}$$

Because $\max \left(\left| \sqrt{n}\hat{\psi}^\# \right| \right)$ converges to a well defined, positive distribution as a result of the continuous mapping theorem, for each constant $c < 1$, we know there exists an ϵ_c such that $\Pr \left(\hat{T}_a^\# > \epsilon_c \right) \geq c$.

Now, shifting our focus to \hat{T}_a , under alternatives, $\psi \neq \mathbf{0}$. Define $\psi_{\max} = \max(\psi_1, \dots, \psi_d)$. Using this knowledge, and (??), note that

$$\begin{aligned} \Pr \left(\hat{T}_a < \epsilon \right) &= \Pr \left(g_a \left(\epsilon \cdot \sqrt{n}\hat{\psi} \right) \geq C_{0.95,a} \right) \\ &\geq \Pr \left(\epsilon \cdot \sqrt{n} \max(\hat{\psi}_1, \dots, \hat{\psi}_d) \geq C_{0.95,a} \right) \\ &= \Pr \left(\max(\hat{\psi}_1, \dots, \hat{\psi}_d) \geq C_{0.95,a} / (\epsilon \cdot \sqrt{n}) \right) \\ &\geq \Pr \left(\max(\hat{\psi}_1, \dots, \hat{\psi}_d) \geq \psi_{\max}/2 \right) \Pr \left(\psi_{\max}/2 \geq C_{0.95,a} / (\epsilon \cdot \sqrt{n}) \right) \end{aligned} \quad (7)$$

The first factor of the product in (??) will converge to 1 as $n \rightarrow \infty$ from the consistency of $\hat{\psi}$. The second quantity will be equal to 1 for sufficiently large n . Thus $\hat{T}_a \xrightarrow{P} 0$ under any alternative.

It was shown that for each a that $\hat{T}_a \xrightarrow{P} 0$. This means that our adaptive estimator $\hat{T} \xrightarrow{P} 0$ as well. Now, let $c = 0.05/k$ and ϵ_c be small enough that $\Pr \left(\hat{T}_a^\# > \epsilon_c \right) \geq 1 - (0.05/k)$. The permutation version of the adaptive estimator $\hat{T}^\#$ has the property that

$$\Pr \left(\hat{T}^\# < \epsilon_c \right) \leq \Pr(\hat{T}_1^\# < \epsilon_c) + \dots + \Pr(\hat{T}_k^\# < \epsilon_c) \leq 0.05,$$

and the theorem's conclusion follows. \square

8.2 Unbiasedness at local alternatives

Consider a local alternative in which the true value of ψ is shrinking towards zero at a root n rate: $\psi = h/\sqrt{n}$. We assume that each potential norm is convex. This assumption can be relaxed to what was described by Eaton and Perlman (1991) (Concentration inequalities for multivariate distributions: mv normal). Under this local alternative, we will have $\sqrt{n}\hat{\psi} \xrightarrow{d} N(\underline{h}, \Sigma)$. Show that the test will reject the null with a probability greater than α for an α level test

Theorem 2. *Under local alternatives described above,*

$$Pr_P \left(\Gamma_{\hat{\Sigma}}^*(\sqrt{n}\hat{\psi}) \leq F_{\Gamma_{\hat{\Sigma}}^*(\hat{Z})}^{-1}(\alpha) \right) > \alpha,$$

Where $\hat{Z} \sim \hat{Q}$ and $Z \sim Q$. Here (unlike other parts of the paper) small values of Γ provide evidence against the null. Values of Γ can be thought of as similar to p -values.

Lemma 3. *The function:*

$$\Gamma_{\Sigma}(t) = Pr_Q(\|\tilde{Z} + t\| > c_{0.8}) \text{ where } c_{\Sigma,0.8} \equiv \min_c \{c : Pr_Q(\|\tilde{Z}\| < c) \geq 0.8\} \text{ and } \tilde{Z} \sim N(0, \Sigma)$$

is continuous with respect to Σ and t .

This lemma will follow because Γ is the integral of a composition of bounded, continuous functions.

Proof. The multivariate normal probability density function

$$\phi(x, \mu, \Sigma) = (2\pi)^{-k} \det(\Sigma)^{-\frac{1}{2}} \exp \left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu) \right)$$

is continuous with respect to both Σ and t .

Since the exponential function, determinate, matrix inverses, and linear operators are all differentiable, they are also all continuous. Because the multivariate normal pdf is the composition of continuous functions, it is also continuous.

It follows that if $\mu_n \rightarrow \mu$, and $\Sigma_n \rightarrow \Sigma$, then for each x , $\phi(x, \mu_n, \Sigma_n) \rightarrow \phi(x, \mu, \Sigma)$. This result and the dominated convergence theorem imply the corresponding CDF's are also continuous ($\phi(x, \hat{\mu}, \hat{\Sigma}) + \phi(x, \mu, \Sigma)$ can be used as the dominating measure): (dominating function shouldn't depend on n) To find a dominating measure, assume that Σ_n^{-1} is close enough to Σ^{-1} so the smallest eigenvalue is within ε then bound the pdf using this fact and that $(x - \mu)^\top \Sigma^{-1}(x - \mu) = (x - \mu)^\top A D A^{-1}(x - \mu) \leq \|(x - \mu)\|^2 c$ where c is the largest eigen value (can be done to provide opposite direction as well).

$$\Phi_{\hat{\Sigma}, \hat{\mu}}(t) = \int_{\|x\| < t} \phi(x, \hat{\mu}, \hat{\Sigma}) dx \rightarrow \int_{\|x\| < t} \phi(x, \mu, \Sigma) dx = \Phi_{\Sigma, \mu}(t).$$

Show (or find reference showing that) $\Phi_{\Sigma_n, 0}^{-1}(0.8) \rightarrow \Phi_{\Sigma, 0}^{-1}(0.8)$ when $\Sigma_n \rightarrow \Sigma$. E.g., implicit function theorem will do this

Because the normal cdf is continuous, we know the normal quantile function will be continuous as well. Potentially we can prove the existence of the inverse using the fact the CDF has a bounded derivative on all of \mathbb{R} . Let $\Sigma_n \rightarrow \Sigma$ and $h_n \rightarrow h$. Also, let Q_n be a normal distribution with mean zero and variance-covariance Σ_n . These findings imply the following:

$$\begin{aligned} |\Gamma_{\Sigma_n}(h) - \Gamma_{\Sigma}(h)| &= \left| Pr_{Q_n}(\|\tilde{Z} + h\| > \Phi_{\Sigma_n, 0}^{-1}(0.8)) - Pr_Q(\|\tilde{Z} + h\| > \Phi_{\Sigma, 0}^{-1}(0.8)) \right| \\ &= \left| \int_{\mathbb{R}^d} \phi(t, 0, \Sigma_n) I\{\|t + h\| > \Phi_{\Sigma_n, 0}^{-1}(0.8)\} - \phi(t, 0, \Sigma) I\{\|t + h\| > \Phi_{\Sigma, 0}^{-1}(0.8)\} dt \right| \\ &= \left| \int_{\mathbb{R}^d \setminus \{t: \|t\| = \Phi_{\Sigma, 0}^{-1}(0.8)\}} \phi(t, h, \Sigma_n) I\{\|t\| > \Phi_{\Sigma_n, 0}^{-1}(0.8)\} - \phi(t, h, \Sigma) I\{\|t\| > \Phi_{\Sigma, 0}^{-1}(0.8)\} dt \right| \\ &\leq \int_{\mathbb{R}^d \setminus \{t: \|t\| = \Phi_{\Sigma, 0}^{-1}(0.8)\}} |\phi(t, h, \Sigma_n) I\{\|t\| > \Phi_{\Sigma_n, 0}^{-1}(0.8)\} - \phi(t, h, \Sigma) I\{\|t\| > \Phi_{\Sigma, 0}^{-1}(0.8)\}| dt \end{aligned}$$

The above quantity converges to zero by the dominated convergence theorem.

$$\begin{aligned}
&= \int_{\|t\| > \Phi_{\Sigma_n,0}^{-1}(0.8), \Phi_{\Sigma,0}^{-1}(0.8)} |\phi(t, h, \Sigma_n) - \phi(t, h, \Sigma)| dt + \\
&\quad \int_{\Phi_{\Sigma_n,0}^{-1}(0.8) \geq \|t\| > \Phi_{\Sigma_n,0}^{-1}(0.8)} |\phi(t, h, \Sigma_n)| dt + \int_{\Phi_{\Sigma_n,0}^{-1}(0.8) \geq \|t\| > \Phi_{\Sigma,0}^{-1}(0.8)} |\phi(t, h, \Sigma)| dt \\
&\leq \int_{\|t\| > \Phi_{\Sigma_n,0}^{-1}(0.8)} |\phi(t, h, \Sigma_n) - \phi(t, h, \Sigma)| dt + \\
&\quad \int_{\Phi_{\Sigma_n,0}^{-1}(0.8) \geq \|t\| > \Phi_{\Sigma_n,0}^{-1}(0.8)} \phi(t, h, \Sigma_n) dt + \int_{\Phi_{\Sigma_n,0}^{-1}(0.8) \geq \|t\| > \Phi_{\Sigma,0}^{-1}(0.8)} \phi(t, h, \Sigma) dt \\
&= \int_{\|t\| > \Phi_{\Sigma_n,0}^{-1}(0.8)} |\phi(t, h, \Sigma_n) - \phi(t, h, \Sigma)| dt + \\
&\quad \Phi_{\Sigma_n,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) - \Phi_{\Sigma_n,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) + \Phi_{\Sigma,h}(\Phi_{\Sigma,0}^{-1}(0.8)) - \Phi_{\Sigma,h}(\Phi_{\Sigma,0}^{-1}(0.8)) \\
&= \int_{\|t\| > \Phi_{\Sigma_n,0}^{-1}(0.8)} |\phi(t, h, \Sigma_n) - \phi(t, h, \Sigma)| dt + \\
&\quad \Phi_{\Sigma_n,h}(\Phi_{\Sigma,0}^{-1}(0.8)) - \Phi_{\Sigma_n,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) - \left(\Phi_{\Sigma,h}(\Phi_{\Sigma,0}^{-1}(0.8)) - \Phi_{\Sigma,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) \right) + \\
&\quad \left(\Phi_{\Sigma,h}(\Phi_{\Sigma,0}^{-1}(0.8)) - \Phi_{\Sigma,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) \right) + \Phi_{\Sigma,h}(\Phi_{\Sigma,0}^{-1}(0.8)) - \Phi_{\Sigma,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) \\
&= \int_{\|t\| > \Phi_{\Sigma_n,0}^{-1}(0.8)} |\phi(t, h, \Sigma_n) - \phi(t, h, \Sigma)| dt + \\
&\quad \left(\Phi_{\Sigma_n,h}(\Phi_{\Sigma,0}^{-1}(0.8)) - \Phi_{\Sigma,h}(\Phi_{\Sigma,0}^{-1}(0.8)) \right) - \left(\Phi_{\Sigma_n,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) - \Phi_{\Sigma,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) \right) + \\
&\quad \left(\Phi_{\Sigma,h}(\Phi_{\Sigma,0}^{-1}(0.8)) - \Phi_{\Sigma,h}(\Phi_{\Sigma_n,0}^{-1}(0.8)) \right) + \Phi_{\Sigma,h}(\Phi_{\Sigma,0}^{-1}(0.8)) - \Phi_{\Sigma,h}(\Phi_{\Sigma_n,0}^{-1}(0.8))
\end{aligned}$$

Taking the limit as $n \rightarrow \infty$ of the quantity above, we find that the first term is zero by dominated convergence theorem. The other four terms are also zero because of the uniform continuity of $\Phi_{\Sigma,x}$ and continuity of $\Phi_{\Sigma,x}^{-1}$ \square

Lemma 4. *Under local alternatives,*

$$\Gamma_{\hat{\Sigma}}^*(\sqrt{n}\hat{\psi}) \xrightarrow{d} \Gamma_{\Sigma}^*(Z + \underline{h})$$

Proof. It was shown in Lemma ?? that for each norm $\|\cdot\|$, $\Gamma_{\Sigma}(x)$ is continuous with respect to Σ and x . This finding, the continuity of the max function, and because a composition of continuous functions is also continuous it follows that

$$\Gamma_{\Sigma}^* \equiv \max \{ \Gamma_{1,\Sigma}, \dots, \Gamma_{d,\Sigma} \}$$

is also continuous with respect to Σ and x . Under local alternatives, $\sqrt{n}\hat{\psi} \xrightarrow{d} Z + \underline{h}$ and $\hat{\Sigma} \xrightarrow{p} \Sigma$. It follows from the continuous mapping theorem that $\Gamma_{\hat{\Sigma}}^*(\sqrt{n}\hat{\psi}) \xrightarrow{d} \Gamma_{\Sigma}^*(Z + \underline{h})$ \square

It has now been established that under local alternatives the distribution of our estimate converges to $\Gamma_{\Sigma}^*(Z + \underline{h})$. Denote the CDF of $\Gamma_{\Sigma}^*(Z + t)$ this distribution by $F_{\Sigma,t}$ and the corresponding quantile function of the distribution by $F_{\Sigma,t}^{-1}$.

To prove unbiasedness at local alternatives, we show $>?$

$$F_{\Sigma,t}(F_{\Sigma,0}^{-1}(1 - \alpha)) \geq F_{\Sigma,0}(F_{\Sigma,0}^{-1}(1 - \alpha)) = \alpha$$

using results from the (?) manuscript.

A result of (?) is that for two centrally symmetric, unimodal functions, $f_1(x)$ and $f_2(x)$, the convolution,

$$g(\theta) = \int f_1(x)f_2(x - \theta)$$

is centrally symmetric and ray decreasing. A function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is centrally symmetric if $f(-x) = f(x)$ for every x . A function is unimodal if for every k , the set $\{x : f(x) \geq k\}$ is convex. A function f on \mathbb{R}^p is ray decreasing if for every x in \mathbb{R}^p , the function $g(\beta) = f(\beta x), \beta \in \mathbb{R}$ is a decreasing function of β .

For now, we will assume that. We assume that our tests will always be define our tests in such a way that we reject the null when $\Gamma_{\Sigma}^*(\sqrt{n}\hat{\psi}) > c$ for some c . It is expected that $\Gamma_{\Sigma}^*(x)$ increases as x moves away from the origin. For the purposes of this proof it will be useful to define new functions

$$\Upsilon_{i,\Sigma} \equiv (\Upsilon_{i,\Sigma})^{-1} \text{ and } \Upsilon_{\Sigma}^* = \min \{\Upsilon_{1,\Sigma}, \dots, \Upsilon_{d,\Sigma}\}$$

that decrease as x moves away from the origin.

Lemma 5. *Let $\Upsilon_1, \dots, \Upsilon_d$ all be centrally symmetric, unimodal functions. Then $\Upsilon^* = \min \{\Upsilon_1, \dots, \Upsilon_d\}$ is also centrally symmetric and unimodal.*

Proof. Because each Υ_i is a centrally symmetric and Υ^* is a function of x only through the Υ_i 's, Υ^* is also centrally symmetric:

$$\Upsilon^*(x) = \min \{\Upsilon_1(x), \dots, \Upsilon_d(x)\} = \min \{\Upsilon_1(-x), \dots, \Upsilon_d(-x)\} = \Upsilon^*(-x)$$

The set $M^* \equiv \{x : \Upsilon^*(x) \geq k\}$ contains all of the x for which $\Upsilon_i(x) \geq k$ for all $i \in \{1, \dots, d\}$. Thus M^* is the intersection of all sets $M_i \equiv \{x : \Upsilon_i(x) \geq k\}$. Each M_i is convex because each Υ_i is unimodal, and the intersection of a countable number of convex sets is convex. Thus M^* is convex and Υ^* is unimodal. \square

Lemma 6. *Let $f(x)$ be a centrally symmetric, unimodal function. Then $g(x) = I\{f(x) \geq c\}$ where $c \in \mathbb{R}$ is also centrally symmetric and unimodal.*

Proof. Because f is centrally symmetric and g is a function of x only through g , g is also centrally symmetric:

$$g(x) = I\{f(x) \geq c\} = I\{f(-x) \geq c\} = g(-x)$$

If $k < 0$ or $k > 1$, then $\{x : g(x) \geq k\}$ will be the empty set or all of \mathbb{R}^d respectively, and both sets are convex.

Otherwise, the set $\{x : g(x) \geq k\} = \{x : I\{f(x) \geq c\} \geq k\}$. The indicator function will be greater than or equal to k whenever $f(x) \geq c$, so $\{x : g(x) \geq k\} = \{x : f(x) \geq c\}$ which is convex because f is unimodal.

Thus g is unimodal. \square

Because each $\Upsilon_{\Sigma,i}$ is unimodal and centrally symmetric, lemma ?? implies $\Upsilon_{\Sigma,i}^*$ is also unimodal and centrally symmetric. It follows from lemma ?? and the previous finding that $I\{\Upsilon_Q^*(x) \geq c_{0.95}\}$ is centrally symmetric and unimodal.

$$\begin{aligned} \int I\{\Gamma_Q^*(x) > c_{0.95}\} \phi(x - \mu) dx &= \int I\{\Upsilon_Q^*(x) < c_{0.95}^{-1}\} \phi(x - \mu) dx \\ &= \int (1 - I\{\Upsilon_Q^*(x) \geq c_{0.95}^{-1}\}) \phi(x - \mu) dx \\ &= 1 - \int I\{\Upsilon_Q^*(x) \geq c_{0.95}^{-1}\} \phi(x - \mu) dx \end{aligned}$$

Since the subtracted quantity is decreasing by (?), the quantity as a whole will be increasing. Thus local power is obtained.

Proof Outline:

- Under local alternatives, $\sqrt{n}\hat{\psi} \xrightarrow{d} N(\mathbf{c}, \Sigma)$
- With a large enough sample size, and enough MC draws, we have that for each norm:
- Think of Gamma as a function indexed by Σ to allow it to be similar. Also make some assumptions about how smooth Υ is with respect to this parameters. For any given value of the input (t) think about
- Right now we have given up on proving things for the permutation test, but we may try to do it again at some point. The permutation test statistic will converge in distribution to a standard normal. This paper should help: (?)
- For each norm selected, (?) states that the power will be non-decreasing as long as the rejection region is convex (This should be true most of our rejection regions), and the probability density is decreasing away from the mean (which is true of a normal distribution).
- Show that for n large, the norm is selected to give the best power. Thus since each norm obtains local power, the adaptive test will also obtain local power for n large.

$$\mathcal{L}\left(\Gamma_{\hat{\Sigma}}\left(\sqrt{n}\hat{\Psi}(\mathbf{X})\right), \Gamma_{\Sigma}\left(\sqrt{n}\hat{\Psi}(\mathbf{X})\right)\right) \rightarrow 0$$

8.3 Consistency of norm selection

This proof would likely be difficult and require some slower than \sqrt{n} convergence rates of the local alternative considered. Under fixed alternatives, all norms perform equally perfectly

8.4 Type 1 error control

Proof is so short, I am not sure if it is worth including Under the null, $Y \perp\!\!\!\perp X$. Thus the test statistic will be taken from the same distribution as all of the permutation based test statistics used to estimate the distribution of the test statistic under the null. Therefore, as B grows, $Pr(T_n \geq F_{T_{k,n}^\#|X_n}(0.95)) \Rightarrow 0.05$