

# Forward screening and post-screening inference in factorial designs

Lei Shi\*

Jingshen Wang<sup>†</sup>

Peng Ding <sup>‡</sup>

## Abstract

Ever since the seminal work of F. Yates, factorial designs have been an important experimental tool to simultaneously estimate the treatment effects of multiple factors. In factorial designs, the number of treatment levels may grow exponentially with the number of factors, which motivates the forward screening strategy based on the sparsity, hierarchy, and heredity principles for factorial effects. Although this strategy is intuitive and has been widely used in practice, its rigorous statistical theory has not been formally established. To fill this gap, we establish a design-based theory for forward factor screening in factorial designs based on the potential outcome framework. We not only prove its consistency property but also discuss statistical inference after factor screening. In particular, with perfect screening, we quantify the advantages of forward screening based on asymptotic efficiency gain in estimating factorial effects. With imperfect screening in higher-order interactions, we propose two novel strategies and investigate their impact on subsequent inference. Our formulation differs from the existing literature on variable selection and post-selection inference because our theory is based solely on the physical randomization of the factorial design and does not rely on a correctly-specified outcome model.

**Keywords:** Causal inference; Design-based inference; Forward selection; Post-selection inference.

---

\*Division of Biostatistics, University of California, Berkeley. leishi@berkeley.edu

<sup>†</sup>Division of Biostatistics, University of California Berkeley. jingshenwang@berkeley.edu. Corresponding author.

<sup>‡</sup>Peng Ding, Department of Statistics, University of California, Berkeley. pengdingpku@berkeley.edu

# 1 Introduction

## 1.1 Factorial experiments: opportunities and challenges

Ever since the seminal work of Yates (1937), factorial designs have been widely used in many fields, including agricultural, industrial, and biomedical sciences (e.g., Box et al., 2005; Wu and Hamada, 2011; Gerber and Green, 2012). For example, in social science, one government funded research by Zhang (2022) studied the social construction of hate crime in the U.S. using factorial experiments based on three factors: race, sexual orientation, and religious affiliation. As another example, in ecology, Rillig et al. (2019) studied multiple global change factors in driving soil functions and microbial biodiversity with factorial designs involving up to ten factors, such as **XXXXX**. Factorial experiments are popular partially because they can simultaneously accommodate multiple factors and offer opportunities to estimate not only the main causal effects of factors but also their interactions.

We focus on the  $2^K$  factorial design in which  $K$  binary factors are randomly assigned to  $N$  experimental units. With a small  $K$ , we can simultaneously estimate the  $2^K - 1$  main effects and interactions. Nevertheless, when  $K$  is large, the number of factorial effects grows exponentially with  $K$ . This motivates us to conduct factor screening based on sparsity, hierarchy, and heredity principles for factorial effects. More precisely, Wu and Hamada (2011) summarized these three principles as below:

- (a) (sparsity) The number of important factorial effects is small.
- (b) (hierarchy) Lower-order effects are more important than higher-order effects, and effects of the same order are equally important.
- (c) (heredity) Higher-order effects are significant only if their corresponding lower-order effects are significant.

The sparsity principle motivates conducting factor screening in factorial designs. The hierarchy principle motivates the forward screening strategy that starts from lower-order effects and then moves on to higher-order effects. The heredity principle motivates using structural restrictions on higher-order effects based on the selected lower-order effects. Due to its simplicity and computational efficiency, while the forward screening strategy has been widely used in data analysis (Wu and Hamada, 2011; Espinosa et al., 2016), its design-based theory under the potential outcome framework has not been formally established. Moreover, it is often challenging to understand the

impact of factor screening on the subsequent statistical inference. The overarching goal of this manuscript is to fill these gaps.

## 1.2 Our contributions and literature review

We summarize our contribution from three perspectives:

First, our study adds to the growing literature of factorial designs with growing number of factors under the potential outcome framework (Dasgupta et al., 2015; Branson et al., 2016; Lu, 2016b; Espinosa et al., 2016; Egami and Imai, 2019; Blackwell and Pashley, 2021; Zhao and Ding, 2021; Pashley and Bind, 2023; Wu et al., 2022). To deal with a large number of factors, Espinosa et al. (2016) and Egami and Imai (2019) informally used factor screening without studying its statistical properties, whereas Zhao and Ding (2021) discussed parsimonious model specifications that are chosen a priori and independent of data. The rigorous theory for factor screening is generally missing in this literature, let alone the theory for statistical inference after factor screening. At a high level, our contributions fill the gaps.

Second, we formalize forward factor screening and establish its consistency under the design-based framework under few outcome modeling assumptions; see Section 3. Factor screening in factorial design sounds like a familiar statistical task if we formulate it as a variable selection problem in a linear model. Thus, forward screening is reminiscent of the vast literature on forward selection. Wang (2009) and Wieczorek and Lei (2022) proved the consistency of forward selection for the main effects in a linear model, whereas Hao and Zhang (2014) and Hao et al. (2018) moved further to allow for second-order interactions. Other researchers proposed various penalized regressions to encode the sparsity, hierarchy, and heredity principles (e.g., Yuan et al., 2007; Zhao et al., 2009; Bickel et al., 2010; Bien et al., 2013; Lim and Hastie, 2015; Haris et al., 2016), without formally studying the statistical properties of the selected model. Our design-based framework departs from the literature without assuming a correctly-specified linear outcome model. This framework is classic in experimental design and causal inference with randomness coming solely from the design of experiments rather than the error terms in a linear model (Neyman, 1923/1990; Kempthorne, 1952; Freedman, 2008; Lin, 2013; Dasgupta et al., 2015). This framework invokes fewer outcome modeling assumptions but consequently imposes technical challenges for developing the theory. Bloniarz et al. (2016) discussed the design-based theory for covariate selection in treatment-control experiments, but the corresponding theory for factorial designs is largely unexplored.

Third, we discuss statistical inference after forward factor screening with (Sections 4 and 6) or without perfect screening (Section 5). On the one hand, we prove the screening consistency of the

forward screening procedure, which ensures that the selected factorial effects are the true, non-zero ones. With this perfect screening property, we can then proceed as if the selected working model is the true model. This allows us to ignore the impact of forward screening on the subsequent inference, which is similar to the proposal of Zhao et al. (2021) for statistical inference after Lasso Tibshirani (1996). In particular, we quantify the advantages of conducting forward screening based on the asymptotic efficiency gain for estimating factorial effects. As an application under perfect screening, we discuss statistical inference for the mean outcome under the best factorial combination (Andrews et al., 2019; Guo et al., 2021; Wei et al., 2022). On the other hand, we acknowledge that perfect screening can be too much to hope for in practice as it requires strong regularity conditions on factorial effects. As a remedy, we propose two strategies to deal with imperfect screening in higher-order interactions, and we study their impacts on post-screening inference. A key motivation for our strategies is to ensure that the parameters of interest after forward factorial screening are not data-dependent, avoiding philosophical debates in the current literature of post-selection inference (Fithian et al., 2014; Kuchibhotla et al., 2022).

### 1.3 Notation

We will use the following notation throughout. For asymptotic analyses,  $a_N = O(b_N)$  denotes that there exists a positive constant  $C > 0$  such that  $a_N \leq Cb_N$ ;  $a_N = o(b_N)$  denotes that  $a_N/b_N \rightarrow 0$  as  $N$  goes to infinity;  $a_N = \Theta(b_N)$  denotes that there exists positive constants  $c$  and  $C$  such that  $cb_N \leq a_N \leq Cb_N$ .

For matrix  $V$ , define  $\varrho_{\max}(V)$  and  $\varrho_{\min}(V)$  as the largest and smallest eigen-values, respectively, and define  $\kappa(V) = \varrho_{\max}(V)/\varrho_{\min}(V)$  as its condition number. For two positive semi-definite matrices  $V_1$  and  $V_2$ , we write  $V_1 \preceq V_2$  or  $V_2 \succeq V_1$  if  $V_2 - V_1$  is positive semi-definite.

We will use different levels of sets. For an integer  $K$ , let  $[K] = \{1, \dots, K\}$ . We use  $\mathcal{K}$  in calligraphic to denote a subset of  $[K]$ . Let  $\mathbb{K} = \{\mathcal{K} \mid \mathcal{K} \subset [K]\}$  denote the power set of  $[K]$ . We also use blackboard bold font to denote subsets of  $\mathbb{K}$ . For example,  $\mathbb{M} \subset \mathbb{K}$  denotes that  $\mathbb{M}$  is a subset of  $\mathbb{K}$ .

We will use  $A_i \sim B_i$  to denote the least-squares fit of  $A_i$ 's on  $B_i$ 's, which is purely a numerical procedure without assuming a linear model. Let  $\xrightarrow{\mathbb{P}}$  denote convergence in probability and  $\rightsquigarrow$  denote convergence in distribution.

## 2 Setup of factorial designs

This section introduces the key mathematical components of factorial experiments. Section 2.1 introduces the notation of potential outcomes and the definitions of the factorial effects. Section 2.2 introduces the treatment assignment mechanism, the observed data, and regression analysis of the data. Section 2.3 uses a concrete example of a  $2^3$  factorial experiment to illustrate the key concepts.

### 2.1 Potential outcomes and factorial effects

We first introduce the general framework of a  $2^K$  factorial design, with  $K \geq 2$  being an integer. This design has  $K$  binary factors, and factor  $k$  can take value  $z_k \in \{0, 1\}$  for  $k = 1, \dots, K$ . Let  $\mathbf{z} = (z_1, \dots, z_K)$  denote the treatment combining all  $K$  factors. The  $K$  factors in total define  $Q = 2^K$  treatment combinations, collected in the set below:

$$\mathcal{T} = \{\mathbf{z} = (z_1, \dots, z_K) \mid z_k \in \{0, 1\} \text{ for } k = 1, \dots, K\} \quad \text{with} \quad |\mathcal{T}| = Q.$$

We follow the potential outcome notation of Dasgupta et al. (2015) for  $2^K$  factorial designs. Unit  $i$  has potential outcome  $Y_i(\mathbf{z})$  under each treatment level  $\mathbf{z}$ . Corresponding to the  $Q = 2^K$  treatment levels, each unit  $i$  has  $Q$  potential outcomes, vectorized as  $\mathbf{Y}_i = \{Y_i(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}$  using the lexicographic order. Over units  $i = 1, \dots, N$ , the potential outcomes have finite-population mean vector  $\bar{\mathbf{Y}} = (\bar{Y}(\mathbf{z}))_{\mathbf{z} \in \mathcal{T}}$  and covariance matrix  $S = (S(\mathbf{z}, \mathbf{z}'))_{\mathbf{z}, \mathbf{z}' \in \mathcal{T}}$ , with elements defined as follows:

$$\bar{Y}(\mathbf{z}) = \frac{1}{N} \sum_{i=1}^N Y_i(\mathbf{z}), \quad S(\mathbf{z}, \mathbf{z}') = \frac{1}{N-1} \sum_{i=1}^N (Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z}))(Y_i(\mathbf{z}') - \bar{Y}(\mathbf{z}')).$$

We then use the potential outcomes to define factorial effects. For a subset  $\mathcal{K} \subset [K]$  of the  $K$  factors, we introduce the following “contrast vector” notation to facilitate the presentation. To start with, we define the main causal effect for factor  $k$ . For a treatment level  $\mathbf{z} = (z_1, \dots, z_K) \in \mathcal{T}$ , we use  $g_{\{k\}}(\mathbf{z}) = 2z_k - 1$  to denote the “centered” treatment indicator  $z_k$ . We then define a  $Q$ -dimensional contrast vector  $g_{\{k\}}$  by aggregating these centered treatment variables into a vector using the lexicographic order, that is

$$g_{\{k\}} = \{g_{\{k\}}(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}, \text{ where } g_{\{k\}}(\mathbf{z}) = 2z_k - 1. \quad (2.1)$$

Next, for the interactions of multiple factors with  $|\mathcal{K}| \geq 2$ , we define the contrast vector  $g_{\mathcal{K}} \in \mathbb{R}^Q$  as

$$g_{\mathcal{K}} = \{g_{\mathcal{K}}(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}, \text{ where } g_{\mathcal{K}}(\mathbf{z}) = \prod_{k \in \mathcal{K}} g_{\{k\}}(\mathbf{z}). \quad (2.2)$$

As a special case, when no factor is considered, we define  $g_\emptyset = \mathbf{1}_Q$ . Stack the contrast vectors into a  $Q \times Q$  matrix

$$G = (g_\emptyset, g_{\{1\}}, \dots, g_{\{K\}}, g_{\{1,2\}}, \dots, g_{\{K-1,K\}}, \dots, g_{[K]}),$$

which has orthogonal columns with  $G^\top G = Q \cdot I_Q$ . We refer to  $G$  as the contrast matrix.

Equipped with the contrast vector notation, we are ready to introduce the main effects and interactions. More concretely, define the main causal effect of a single factor and the  $k$ -way interaction causal effect of multiple factors ( $k \geq 2$ ) as the inner product of the contrast vector  $g_K$ , and the averaged potential outcome  $\bar{Y}$ , that is

$$\tau_K = Q^{-1} \cdot g_K^\top \bar{Y} \quad \text{for } K \subset [K].$$

For convenience in description, we use  $\tau_\emptyset = Q^{-1} g_\emptyset^\top \bar{Y}$  to denote the average of potential outcomes. We call the effect  $\tau_K$  a *parent* of  $\tau_{K'}$  if  $K \subset K'$  and  $|K| = |K'| - 1$ . More compactly, we summarize the entire collection of causal parameters in factorial experiments as

$$\tau = (\tau_K)_{K \subset [K]} = Q^{-1} \cdot G^\top \bar{Y}.$$

## 2.2 Treatment assignment, observed data, and regression analysis

Under the design-based framework, the treatment assignment mechanism characterizes the completely randomized factorial design. In other words, the experimenter randomly assigns  $N(\mathbf{z})$  units to treatment level  $\mathbf{z} \in \mathcal{T}$ , with  $\sum_{\mathbf{z} \in \mathcal{T}} N(\mathbf{z}) = N$ . Assume  $N(\mathbf{z}) \geq 2$  to allow for variance estimation within each treatment level. Let  $Z_i \in \mathcal{T}$  denote the treatment level for unit  $i$ . The treatment vector  $(Z_1, \dots, Z_N)$  is a random permutation of a vector with prespecified number  $N(\mathbf{z})$  of the corresponding treatment level  $\mathbf{z}$ , for  $\mathbf{z} \in \mathcal{T}$ .

For each unit  $i$ , the treatment level  $Z_i$  only reveals one potential outcome. We use  $Y_i = Y_i(Z_i) = \sum_{\mathbf{z} \in \mathcal{T}} Y_i(\mathbf{z}) \mathbf{1}\{Z_i = \mathbf{z}\}$  to denote the observed outcome. We also use  $N_i = N(Z_i)$  to denote the number of units for the treatment group in which unit  $i$  is assigned to. The central task of causal inference in factorial designs is to use the observed data  $(Z_i, Y_i)_{i=1}^N$  to estimate factorial effects. Define

$$\hat{Y}(\mathbf{z}) = N(\mathbf{z})^{-1} \sum_{i=1}^N \mathbf{1}\{Z_i = \mathbf{z}\} Y_i, \quad \hat{S}(\mathbf{z}, \mathbf{z}) = \{N(\mathbf{z}) - 1\}^{-1} \sum_{i=1}^N \mathbf{1}\{Z_i = \mathbf{z}\} (Y_i - \hat{Y}(\mathbf{z}))^2$$

as the sample mean and variance of the observed outcomes under treatment  $\mathbf{z}$ . Vectorize the sample means as  $\hat{Y} = (\hat{Y}(\mathbf{z}))_{\mathbf{z} \in \mathcal{T}}$ , which has mean  $\bar{Y}$  and covariance matrix  $V_{\hat{Y}} = D_{\hat{Y}} - N^{-1}S$  (Li and Ding, 2017), where

$$D_{\hat{Y}} = \text{Diag} \{N(\mathbf{z})^{-1} S(\mathbf{z}, \mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}.$$

An unbiased estimator for  $D_{\hat{Y}}$  is

$$\hat{V}_{\hat{Y}} = \text{Diag} \left\{ N(\mathbf{z})^{-1} \hat{S}(\mathbf{z}, \mathbf{z}) \right\}_{\mathbf{z} \in \mathcal{T}},$$

whereas  $S$  does not have an unbiased sample analogue because the potential outcomes across treatment levels are never jointly observed for the same units. Therefore,  $\hat{V}_{\hat{Y}}$  is a conservative estimator of the covariance matrix in the sense that  $\mathbb{E}\{\hat{V}_{\hat{Y}}\} = D_{\hat{Y}} \succcurlyeq V_{\hat{Y}}$ .

A dominant approach to estimate factorial effects from factorial designs is through estimating least-squares coefficients based on appropriate model specifications. Let  $g_i$  denote the row vector in the contrast matrix  $G$  corresponding to unit  $i$ 's treatment level  $Z_i$ , that is,  $g_i = \{g_{\mathcal{K}}(Z_i)\}_{\mathcal{K} \subset [K]} \in \mathbb{R}^Q$  with  $g_{\mathcal{K}}(\mathbf{z})$  defined in (2.2). For a set of target effects  $\{\tau_{\mathcal{K}}\}_{\mathcal{K} \in \mathbb{M}}$  indexed by  $\mathbb{M}$ , we can run weighted least squares (WLS) to obtain unbiased estimates:

$$\hat{\tau} = \arg \min_{\tau} \sum_{i=1}^N w_i (Y_i - g_i^{\top} \tau)^2 \text{ with } w_i = 1/N_i.$$

With a small  $K$ , we can simply fit the *saturated regression* by regressing the observed outcome  $Y_i$  on the regressor  $g_i$ . The saturated regression involves  $Q = 2^K$  coefficients without any restrictions on the targeted factorial effects.

In contrast, an *unsaturated regression* involves fewer coefficients by regressing the observed outcome  $Y_i$  on  $g_{i,\mathbb{M}}$ , a subvector of  $g_i$ , where  $\mathbb{M} \subset \mathbb{K}$  is a subset of the power set of all factors. That is,

$$\hat{\tau} = \arg \min_{\tau} \sum_{i=1}^N w_i (Y_i - g_{i,\mathbb{M}}^{\top} \tau)^2 \text{ with } w_i = 1/N_i. \quad (2.3)$$

For the convenience of description, we will call  $\mathbb{M}$  a *working model*. We use a working model to generate estimates based on least squares without assuming its correctness. When  $\mathbb{M} = \mathbb{K}$ , (2.3) incorporates the saturated regression 2.2. Based on the unsaturated regression with working model  $\mathbb{M}$ , let

$$\hat{\tau}(\mathbb{M}) = \{\hat{\tau}_{\mathcal{K}}\}_{\mathcal{K} \in \mathbb{M}} \quad \text{and} \quad \tau(\mathbb{M}) = \{\tau_{\mathcal{K}}\}_{\mathcal{K} \in \mathbb{M}}$$

denote the vectors of estimated and true coefficients, respectively. Zhao and Ding (2021) showed that if we run unsaturated regressions with weights  $1/N_i$  for unit  $i$ , then the obtained estimated coefficients are unbiased for the true factorial effects within the working model  $\mathbb{M}$ . More precisely,  $\hat{\tau}(\mathbb{M}) = Q^{-1} G(\cdot, \mathbb{M})^{\top} \hat{Y}$ , where  $G(\cdot, \mathbb{M})$  to denote the columns in  $G$  indexed by  $\mathbb{M}$ . Because  $\hat{\tau}(\mathbb{M})$  is a linear transformation of  $\hat{Y}$ , we can use the following estimator for its covariance matrix:

$$\hat{\Sigma}(\mathbb{M}) = \frac{1}{Q^2} G(\cdot, \mathbb{M})^{\top} \hat{V}_{\hat{Y}} G(\cdot, \mathbb{M}). \quad (2.4)$$

See Lemma S1 in Section A.1 of the supplementary material for more discussions on the above algebraic results for unsaturated regressions.

### 2.3 An illustrating example of a $2^3$ factorial design

We realize that the above notation can be rather abstract. In what follows, we provide an illustrative Example 1 below with  $K = 3$  factors.

**Example 1** ( $2^3$  factorial design). *Suppose we have three binary factors  $z_1$ ,  $z_2$ , and  $z_3$ . These three factors generate 8 treatment combinations, indexed by a triplet  $(z_1 z_2 z_3)$  with  $z_1, z_2, z_3 \in \{0, 1\}$ , in the set*

$$\mathcal{T} = \{(000), (001), (010), (011), (100), (101), (110), (111)\}.$$

*Each unit  $i$  has a potential outcome vector  $\mathbf{Y}_i = \{Y_i(z_1 z_2 z_3)\}_{z_1, z_2, z_3=0,1}^\top$ . The vector of causal parameters in this factorial experiment is*

$$\tau = \frac{1}{2^3} G^\top \bar{Y} \triangleq (\tau_\emptyset, \tau_{\{1\}}, \tau_{\{2\}}, \tau_{\{3\}}, \tau_{\{1,2\}}, \tau_{\{1,3\}}, \tau_{\{2,3\}}, \tau_{\{1,2,3\}})^\top,$$

where  $G$  is the contrast matrix

$$G = \begin{matrix} & \tau_\emptyset & \tau_{\{1\}} & \tau_{\{2\}} & \tau_{\{3\}} & \tau_{\{1,2\}} & \tau_{\{1,3\}} & \tau_{\{2,3\}} & \tau_{\{1,2,3\}} \\ \begin{matrix} (000) \\ (001) \\ (010) \\ (011) \\ (100) \\ (101) \\ (110) \\ (111) \end{matrix} & \begin{pmatrix} 1 & -1 & -1 & -1 & 1 & 1 & 1 & -1 \\ 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 \\ 1 & -1 & 1 & 1 & -1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 \\ 1 & 1 & -1 & 1 & -1 & 1 & -1 & -1 \\ 1 & 1 & 1 & -1 & 1 & -1 & -1 & -1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \end{pmatrix}.$$

We observe the pair  $(Y_i, Z_i)$  for unit  $i$ , where  $Z_i = (z_{i,1}, z_{i,2}, z_{i,3})$  is the observed treatment combinations. Let  $g_{\{k\}}(Z_i) = 2z_{i,k} - 1$  be the centered version of  $z_{i,k}$ . For the factor-based regression, the regressor  $g_i$  corresponding to the treatment level  $Z_i$  equals

$$t_i = \left[ 1, g_{\{1\}}(Z_i), g_{\{2\}}(Z_i), g_{\{3\}}(Z_i), g_{\{2,3\}}(Z_i), g_{\{1,3\}}(Z_i), g_{\{1,2\}}(Z_i), g_{\{1,2,3\}}(Z_i) \right].$$

For instance, when  $Z_i = (101)$ , the regressor  $g_i$  corresponds to the row (101) of the contrast matrix  $G$ . Then, a saturated regression is to regress  $Y_i$  on  $g_i$ . For the unsaturated regression, if we only



include indices  $\emptyset$  (the intercept),  $\{1\}$ ,  $\{1, 2\}$ ,  $\{1, 3\}$ ,  $\{123\}$  in our regression, we can form the working model  $\mathbb{M} = \{\emptyset, \{1\}, \{1, 2\}, \{1, 3\}, \{1, 2, 3\}\}$  and perform weighted least squares  $Y_i \sim t_{i,\mathbb{M}}$ , where

$$t_{i,\mathbb{M}} = \begin{bmatrix} 1, & g_{\{1\}}(Z_i), & g_{\{1,3\}}(Z_i), & g_{\{1,2\}}(Z_i), & g_{\{1,2,3\}}(Z_i) \end{bmatrix}$$

and the weight for unit  $i$  equals  $1/N_i = 1/N(Z_i)$ .

### 3 Forward screening in factorial experiments

In factorial designs with small  $K$ , we can simply run the saturated regression to estimate all factorial effects simultaneously. However, when  $K$  is large, saturated regression can be computationally unwieldy and scientifically unreasonable by delivering potentially noisy estimates of all higher-order interactions. As a remedy, forward screening is a popular strategy frequently adopted in practice to analyze data collected from factorial experiments, due to its clear benefits in screening out a large number of zero nuisance factorial effects. In this section, we formalize forward screening as a principled procedure to carefully decide an unsaturated working model  $\hat{\mathbb{M}}$ . We first present a formal version of forward screening and then demonstrate its consistency property.

#### 3.1 A formal forward screening procedure

In this subsection, we introduce a principled forward screening procedure that not only fully respects the effect hierarchy, sparsity, and heredity principles but also results in an interpretable parsimonious model with statistical guarantees. More concretely, the algorithm starts by performing factor screening over lower-order effects, then move forward to select the significant higher order effects following the heredity principle. Algorithm 1 summarizes the forward screening procedure.

In what follows, we illustrate why the proposed procedure in Algorithm 1 respects the three fundamental principles in factorial experiments.

First, Algorithm 1 obeys the hierarchy principle as it performs factor screening in a forward style (coded in the global loop from  $d = 1$  to  $d = D$ , Step 2 in particular). More concretely, we begin with an empty working model. We then select relevant main effects (Steps 4 and 8) and add them into the working model. Once the working model is updated, we continue to select relevant higher-order interaction effects in a forward fashion. Such a forward screening procedure is again motivated by the hierarchy principle that lower-order effects are more important than higher-order ones.

---

**Algorithm 1:** Forward factorial screening

---

**Input:** Factorial data  $\{(Y_i, Z_i)\}_{i=1}^N$ ; predetermined integer  $D \leq K$ ; initial working model

$\widehat{\mathbb{M}} = \{\emptyset\}$ ; significance levels  $\{\alpha_d\}_{d=1}^D$ .

**Output:** Selected working model  $\widehat{\mathbb{M}}$ .

1 Define an intermediate working model  $\widehat{\mathbb{M}}' = \widehat{\mathbb{M}}$  for convenience.

2 **for**  $d = 1, \dots, D$  **do**

3     Update the intermediate working model to include all the  $d$ -order (interaction) terms:

$$\widehat{\mathbb{M}}' = \widehat{\mathbb{M}} \cup \{\mathcal{K} \mid |\mathcal{K}| = d\} \triangleq \widehat{\mathbb{M}} \cup \mathbb{K}_d.$$

4     Screen out indices in  $\widehat{\mathbb{M}}'$  according to either the weak/strong heredity principles, and renew the screened working model as  $\widehat{\mathbb{M}}'$ .

5     Run the unsaturated regression with the working model  $\widehat{\mathbb{M}}'$ :

$$Y_i \sim g_{i, \widehat{\mathbb{M}}'}, \text{ with weights } w_i = N/N_i.$$

6     Obtain coefficients  $\widehat{\tau}(\widehat{\mathbb{M}}')$  and robust covariance estimation  $\widehat{\Sigma}(\widehat{\mathbb{M}}')$  defined in (2.4).

7     Extract  $\widehat{\tau}_{\mathcal{K}}(\widehat{\mathbb{M}}')$  and  $\widehat{\sigma}_{\mathcal{K}}(\widehat{\mathbb{M}}')$  for all  $\mathcal{K} \in \widehat{\mathbb{M}}'$  with  $|\mathcal{K}| = d$ .

8     Run marginal t-tests using the above  $\widehat{\tau}_{\mathcal{K}}(\widehat{\mathbb{M}}')$  and  $\widehat{\sigma}_{\mathcal{K}}(\widehat{\mathbb{M}}')$  under the significance level  $\min\{\alpha_d/(|\widehat{\mathbb{M}}'| - |\widehat{\mathbb{M}}|), 1\}$  and remove the non-significant terms from  $\widehat{\mathbb{M}}' \setminus \widehat{\mathbb{M}}$ .

9     Set  $\widehat{\mathbb{M}} = \widehat{\mathbb{M}}'$ .

10 **return**  $\widehat{\mathbb{M}}$ 

---

Second, Algorithm 1 operates under the sparsity principle as it removes potentially unimportant effects using marginal t-tests with the Bonferroni correction (Step 8). This step induces a sparse working model and helps us to identify essential factorial effects. The sparsity-inducing step can incorporate many popular selection frameworks, such as marginal t-tests, Lasso (Tibshirani, 1996), sure independence screening (Fan and Lv, 2008), etc. For simplicity, we present Algorithm 1 with marginal t-tests and relegate more general discussions to Section B of the supplementary material.

Third, Algorithm 1 incorporates the heredity principle as it screens out the interaction effects (Wu and Hamada, 2011; Hao and Zhang, 2014; Lim and Hastie, 2015) when either none of their parent effects is included (weak heredity) or some of their parent effects are excluded (strong heredity) in the previous working model (Step 4).

Lastly, we note that our forward screening procedure enhances the interpretability of the selected working model by iterating between the ‘‘Sparsity-screening’’ step (called the S-step in the rest of the manuscript), captured by a data-dependent operator  $\widehat{\mathbf{S}} = \widehat{\mathbf{S}}(\cdot; \{Y_i, Z_i\}_{i=1}^N)$ , and the ‘‘Heredity-

screening” step (called the H-step in the rest of the manuscript), captured by a deterministic operator  $\mathbf{H} = \mathbf{H}(\cdot)$ . Because the working model is updated in an iterative fashion,

$$\widehat{\mathbb{M}}_1 \xrightarrow{\mathbf{H}} \widehat{\mathbb{M}}_{2,+} \xrightarrow{\widehat{\mathbf{S}}} \widehat{\mathbb{M}}_2 \cdots \xrightarrow{\widehat{\mathbf{S}}} \widehat{\mathbb{M}}_{d-1} \xrightarrow{\mathbf{H}} \widehat{\mathbb{M}}_{d,+} \xrightarrow{\widehat{\mathbf{S}}} \widehat{\mathbb{M}}_d \rightarrow \cdots \xrightarrow{\widehat{\mathbf{S}}} \widehat{\mathbb{M}}_D. \quad (3.5)$$

the final working model includes a small number of statistically significant effects that fully respect the heredity principle.

### 3.2 Consistency of forward screening

We are now ready to analyze the screening consistency property of Algorithm 1. We shall show that the proposed algorithm selects the targeted working model up to level  $D$  with probability tending to one as the sample size goes to infinity. Here, the targeted working model at level  $k \in [K]$ , denoted as  $\mathbb{M}_k^*$ , is the collection of  $\mathcal{K}$ ’s where  $|\mathcal{K}| = k$  and  $\tau_{\mathcal{K}} \neq 0$ . Define the full targeted working model up to level  $D$  as

$$\mathbb{M}_{1:D}^* = \bigcup_{d=1}^D \mathbb{M}_d^*.$$

In particular, when  $D = K$ , we omit the subscript and simply denote  $\mathbb{M}^* = \mathbb{M}_{1:K}^*$ .

We start by introducing the following condition on *nearly uniform designs*:

**Condition 1** (Nearly uniform design). *There exists a positive integer  $N_0$  and absolute constants  $\underline{c} \leq \bar{c}$ , such that*

$$N(\mathbf{z}) = c(\mathbf{z})N_0 \geq 2, \text{ where } \underline{c} \leq c(\mathbf{z}) \leq \bar{c}.$$

Condition 1 allows for diverging  $Q$  and bounded  $N(\mathbf{z})$ ’s across all treatment levels (Shi and Ding, 2022). It generalizes the classical assumption in the fixed  $Q$  regime where  $Q$  is fixed, and each treatment arm contains a sufficiently large number of replications (Li and Ding, 2017).

Next, we quantify the order of the true factorial effect sizes  $\tau_{\mathcal{K}}$ ’s and the tuning parameters  $\alpha_d$ ’s adopted in the Bonferroni correction. We allow these parameters to change with the sample size  $N$ :

**Condition 2** (Order of parameters). *The true factorial effects  $\tau_{\mathcal{K}}$ ’s and tuning parameters  $\alpha_d$ ’s have the following orders:*

- (i) *True nonzero factorial effects:  $|\tau_{\mathcal{K}}| = \Theta(N^\delta)$  for some  $-1/2 < \delta \leq 0$  and all  $\mathcal{K} \in \mathbb{M}_{1:D}^*$ .*

- (ii) *Tuning parameters in Bonferroni correction:*  $\alpha_d = \Theta(N^{-\delta'})$  for all  $d \in [D]$  with some  $\delta' > 0$ .
- (iii) *Size of the targeted working model:*  $\sum_{d=1}^D |\mathbb{M}_d^*| = \Theta(N^{\delta''})$  for some  $0 \leq \delta'' < 1/3$ .

Condition 2(i) specifies the allowable order of the true factorial effects. If this condition fails, the effect size is of the same order as the statistical error and thus is too small to be detected by marginal t-test. Similar conditions are also adopted in model selection literature, including Zhao and Yu (2006) and Wieczorek and Lei (2022). Condition 2(ii) requires the tuning parameter  $\alpha_d$  to converge to zero, which ensures that there is no Type I error in our procedure as  $N$  goes to infinity and hence guarantees the selection consistency. Wasserman and Roeder (2009, Theorems 4.1 and 4.2) assumed similar conditions in high-dimensional model selection settings for linear models. Condition (iii) restricts the size of the targeted working model. The rate is due to our technical analysis. Similar conditions also appeared in Zhao and Yu (2006), Wieczorek and Lei (2022) and Wasserman and Roeder (2009).

The next condition specifies a set of regularity assumptions on the potential outcomes.

**Condition 3** (Regularity conditions on the potential outcomes). *The potential outcomes satisfy the following conditions:*

- (i) *Nondegenerate correlation matrix.* Let  $V^*$  be the correlation matrix of  $\hat{Y}$ . There exists  $\sigma > 0$  such that the condition number of  $V^*$  is smaller than or equal to  $\sigma^2$ .
- (ii) *Bounded fourth central moments.* There exists a universal constant  $\Delta > 0$  such that

$$\max_{z \in [Q]} \frac{1}{N} \sum_{i=1}^N \{Y_i(z) - \bar{Y}(z)\}^4 \leq \Delta^4.$$

- (iii) *Bounded standardization scales.* There exists a constant  $M > 0$  such that  $M_N \leq M$  where

$$M_N = \frac{\max_{i \in [N], q \in [Q]} |Y_i(q) - \bar{Y}(q)|}{\{\min_{q \in [Q]} S(q, q)\}^{1/2}}.$$

Condition 3(i) requires the correlation matrix of  $\hat{Y}$  to be well-behaved. Condition 3(ii) controls the moments of the potential outcomes. Condition 3(iii) imposes a universal bound on the standardization of potential outcomes, which is required by Shi and Ding (2022) to prove the Berry–Esseen bound based on Stein’s method.

Lastly, we impose the following structural conditions on the factorial effects:

**Condition 4** (Hierarchical structure in factorial effects). *The nonzero true factorial effects align with the effect heredity principle:*

- *Weak heredity:*  $\tau_K \neq 0$  only if there exists  $K' \subset K$  with  $|K'| = |K| - 1$  such that  $\tau_{K'} \neq 0$ .
- *Strong heredity:*  $\tau_K \neq 0$  only if  $\tau_{K'} \neq 0$  for all  $K' \subset K$  with  $|K'| = |K| - 1$ .

Finally, we present the screening consistency property of Algorithm 1:

**Theorem 1** (Perfect screening property). *Under Conditions 1-4, the working model selected by Algorithm 1 converges to the targeted working model with probability one as the sample size goes to infinity:*

$$\lim_{N \rightarrow \infty} \mathbb{P} \left( \widehat{\mathbb{M}} = \mathbb{M}_{1:D}^* \right) = 1.$$

## 4 Inference under perfect screening

Statistical inference is relatively straightforward under the perfect screening of the factorial effects. If forward screening correctly identifies the true, nonzero factorial effects with probability approaching one, we can proceed as if the selected working model is predetermined. In Section 4.1, we present the point estimators and confidence intervals for general causal parameters. In Section 4.2, we study the advantages of forward screening in terms of asymptotic efficiency in estimating general causal parameters, compared with the corresponding estimators without forward screening. We relegate the extensions to vector parameters to Section A.2 of the supplementary material since it is conceptually straightforward.

### 4.1 Post-screening inference for general causal parameters

Define a general causal parameter of interest as a weighted combination of average potential outcomes:

$$\gamma = \sum_{z \in \mathcal{T}} \mathbf{f}(z) \bar{Y}(z) \triangleq \mathbf{f}^\top \bar{Y},$$

where  $\mathbf{f} = \{\mathbf{f}(z)\}_{z \in \mathcal{T}}$  is a pre-specified weighting vector. For example, if one is interested in estimating the main factorial effects,  $\mathbf{f}$  can be taken as the contrast vectors  $g_{\{k\}}$  given in (2.1). If one wants to estimate interaction effects, then  $\mathbf{f}$  can be constructed from (2.2). However, we allow  $\mathbf{f}$  to be different from the contrast vectors  $g_K$ . For instance, if one wants to focus on the first two arms in factorial experiments and estimate the average treatment effect, we shall choose

$$\mathbf{f} = (1, -1, 0, \dots, 0)^\top.$$

In general, researchers may tailor the choice of  $\mathbf{f}$  to the specific research questions of interest.

Without factor screening, a well-studied plug-in estimator of  $\gamma$  in the existing literature is to replace  $\bar{Y}$  with its sample analogue (Li and Ding, 2017; Zhao and Ding, 2021; Shi and Ding, 2022):

$$\hat{\gamma} = \mathbf{f}^\top \hat{Y} = \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z}) \hat{Y}(\mathbf{z}). \quad (4.6)$$

Under regularity conditions in Shi and Ding (2022), the plug-in estimator  $\hat{\gamma}$  satisfies a central limit theorem  $(\hat{\gamma} - \gamma)/v \rightsquigarrow \mathcal{N}(0, 1)$  with the variance  $v^2 = \mathbf{f}^\top V_{\hat{Y}} \mathbf{f}$ . When  $N(\mathbf{z}) \geq 2$ , its variance can be estimated by:

$$\hat{v}^2 = \mathbf{f}^\top \hat{V}_{\hat{Y}} \mathbf{f} = \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z})^2 N(\mathbf{z})^{-1} \hat{S}(\mathbf{z}, \mathbf{z}).$$

With the help of factor screening, based on the selected working model  $\hat{\mathbb{M}}$ , we consider a potentially more efficient estimator of  $\bar{Y}$  via the restricted least squares (RLS)

$$\hat{Y}_R = \arg \min_{\mu \in \mathbb{R}^Q} \left\{ \|\hat{Y} - \mu\|_2^2 : G(\cdot, \hat{\mathbb{M}}^c)^\top \mu = 0 \right\}, \quad (4.7)$$

which leverages the information that the nuisance effects  $G(\cdot, \hat{\mathbb{M}}^c)^\top \bar{Y}$  are all zero. The  $\hat{Y}_R$  in (4.7) has a closed form solution (see Lemma S6 in the supplementary material):

$$\hat{Y}_R = Q^{-1} G(\cdot, \hat{\mathbb{M}}) G(\cdot, \hat{\mathbb{M}})^\top \hat{Y}.$$

Under perfect screening,  $\hat{Y}_R$  is also a consistent estimator for  $\bar{Y}$ , so  $\hat{\gamma}_R = \mathbf{f}^\top \hat{Y}_R$  is also consistent for  $\gamma$ . Introduce the following notation

$$\mathbf{f}[\mathbb{M}] = Q^{-1} G(\cdot, \mathbb{M}) G(\cdot, \mathbb{M})^\top \mathbf{f} \quad (4.8)$$

to simplify  $\hat{\gamma}_R$  and its variance estimator as

$$\hat{\gamma}_R = \mathbf{f}[\hat{\mathbb{M}}]^\top \hat{Y} \quad \text{and} \quad \hat{v}_R^2 = \mathbf{f}[\hat{\mathbb{M}}]^\top \hat{V}_{\hat{Y}} \mathbf{f}[\hat{\mathbb{M}}].$$

Construct a Wald-type level- $(1 - \alpha)$  confidence interval for  $\gamma$ :

$$\left[ \hat{\gamma}_R \pm z_{1-\alpha/2} \times \hat{v}_R \right], \quad (4.9)$$

where  $z_{1-\alpha/2}$  is  $(1 - \alpha/2)$ th quantile of a standard normal distribution. We can also obtain point estimates and confidence intervals handily from WLS regression of  $Y_i$  on  $g_{i, \hat{\mathbb{M}}}$  with weights  $1/N_i$ . See Section A.1 in the supplementary material for more details.

In the following subsection, we provide the theoretical properties of  $\hat{\gamma}_R$  and  $\hat{v}_R^2$ , and compare their asymptotic behaviors with the plug-in estimators  $\hat{\gamma}$  and  $\hat{v}^2$  in various settings.

## 4.2 Theoretical properties under perfect screening

In this section, we first present the asymptotic normality result for  $\hat{\gamma}_R$ . To simplify discussion, we denote  $\mathbf{f}^\star = \mathbf{f}[\mathbb{M}^\star]$ . Given  $\mathbb{M}^\star$  is the true working model, we have  $(\mathbf{f}^\star)^\top \bar{Y} = \mathbf{f}^\top \bar{Y}$ , for all  $\mathbf{f} \in \mathbb{R}^Q$ . This identity holds for the true working model, not a general model, suggested by the following algebraic facts:

$$\begin{aligned} \mathbf{f}^\top \bar{Y} &= \mathbf{f}^\top \{Q^{-1}G(\cdot, \mathbb{M}^\star)G(\cdot, \mathbb{M}^\star)^\top + Q^{-1}G(\cdot, \mathbb{M}^{\star c})G(\cdot, \mathbb{M}^{\star c})^\top\} \bar{Y} \text{ (orthogonality of } G) \\ &= (\mathbf{f}^\star)^\top \bar{Y} + G(\cdot, \mathbb{M}^{\star c})\tau(\mathbb{M}^{\star c}) \text{ (definition of } \mathbf{f}^\star \text{ based on (4.8))} \\ &= (\mathbf{f}^\star)^\top \bar{Y}. \text{ (using } \tau(\mathbb{M}^{\star c}) = 0) \end{aligned}$$

We are now ready to present the asymptotic properties of  $\hat{\gamma}_R$  and  $\hat{v}_R^2$ :

**Theorem 2** (Statistical properties of  $\hat{\gamma}_R$  and  $\hat{v}_R^2$ ). *Let  $N \rightarrow \infty$ . Assume Conditions 1-4. We have*

$$\frac{\hat{\gamma}_R - \gamma}{v_R} \rightsquigarrow \mathcal{N}(0, 1)$$

where  $v_R^2 = \mathbf{f}^{\star\top} V_{\hat{Y}} \mathbf{f}^\star$ . Further assume  $\|\mathbf{f}^\star\|_\infty = O(Q^{-1})$ . The variance estimator  $\hat{v}_R^2$  is conservative in the sense that:

$$N(\hat{v}_R^2 - v_{R,\text{lim}}^2) \xrightarrow{\mathbb{P}} 0, \quad v_{R,\text{lim}}^2 \geq v_R^2,$$

where  $v_{R,\text{lim}}^2 = \mathbf{f}^{\star\top} D_{\hat{Y}} \mathbf{f}^\star$  is the limiting value of  $\hat{v}_R^2$ .

Theorem 2 above guarantees that the proposed confidence interval in (4.9) for  $\gamma$  attains the nominal coverage probability asymptotically. Furthermore, it allows us to compare the conditions for reaching asymptotic normality of  $\hat{\gamma}$ , which we formalize in the following remark:

**Remark 1** (Comparison of conditions for asymptotic normality). *Without factor screening, the simple plug-in estimator  $\hat{\gamma}$  in (4.6) satisfies a central limit theorem if*

$$N_0^{-1/2} \cdot \frac{\|\mathbf{f}\|_\infty}{\|\mathbf{f}\|_2} \rightarrow 0 \tag{4.10}$$

recalling the definition of  $N_0$  in Condition 1 (Shi and Ding, 2022, Theorem 1). Condition (4.10) fails when  $N_0$  is small and  $\mathbf{f}$  is sparse. Besides, it does not incorporate the sparsity information in the structure of factorial effects. With factor screening, however, we can borrow the benefit of a sparse working model and overcome the above drawbacks. Therefore, factor screening broadens the applicability of our proposed estimator  $\hat{\gamma}_R$  by weakening the assumptions for the Wald-type inference.

To elaborate the benefits of conducting forward factorial screening in terms of asymptotic efficiency, we make a simple comparison of the asymptotic variances of  $\hat{\gamma}$  and  $\hat{\gamma}_R$  in Proposition 1 below. In the most general setup, there is no ordering relationship between  $v_R^2$  and  $v^2$ . That is, the RLS based estimator may have higher variance than the unrestricted OLS estimator. This is a known fact due to heteroskedasticity and the use of sandwich variance estimators (Meng and Xie, 2014; Zhao and Ding, 2021). Nevertheless, in many interesting scenarios, we can prove an improvement of efficiency by factor screening. Two conditions are summarized in Proposition 1:

**Proposition 1** (Asymptotic relative efficiency comparison between  $\hat{\gamma}$  and  $\hat{\gamma}_R$ ). *Assume that both  $\hat{\gamma}$  and  $\hat{\gamma}_R$  converge to a normal distribution as  $N \rightarrow \infty$ .*

(i) *If the condition number of  $V_{\hat{Y}}$  satisfies  $\kappa(V_{\hat{Y}}) = 1$ , then*

$$\frac{v_R^2}{v^2} \leq 1.$$

(ii) *Let  $s^*$  denote the number of nonzero elements in  $\mathbf{f}$ . Then the asymptotic relative efficiency between  $\hat{\gamma}$  and  $\hat{\gamma}_R$  is upper bounded by*

$$\frac{v_R^2}{v^2} \leq \kappa(V_{\hat{Y}}) \cdot \frac{s^*|\mathbb{M}^*|}{Q}.$$

Now we add some interpretation for Proposition 1. The condition  $\kappa(V_{\hat{Y}}) = 1$  in Part (i) holds when the variance of  $\hat{Y}(\mathbf{z})$  does not change with its treatment group membership  $\mathbf{z}$ . One concrete problem of interest is testing the *sharp null hypothesis* of constant effects in uniform factorial designs (with  $N_0$  replications in each arm), i.e.,

$$H_{0F} : Y_i(\mathbf{z}) = Y_i \text{ for all } i \in [N] \text{ and } \mathbf{z} \in \mathcal{T}.$$

Under  $H_{0F}$ , we have

$$V_{\hat{Y}} = N_0^{-1} \sigma \cdot I_Q, \text{ where } \sigma = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2 \text{ and } \bar{Y} = \frac{1}{N} \sum_{i=1}^N Y_i.$$

The proposed RLS-based estimator  $\hat{\gamma}_R$  is in general more efficient than the plug-in estimator  $\hat{\gamma}$ . Part (ii) studies a general heteroskedastic setting with sparse weighting vector  $\mathbf{f}$  and working model size  $|\mathbb{M}^*|$ . The condition number  $\kappa(V_{\hat{Y}})$  captures the variability of the variances of  $\hat{Y}(\mathbf{z}) = N(\mathbf{z})^{-1} \sum_{Z_i=\mathbf{z}} Y_i$  across multiple treatment combination groups in  $\mathcal{T}$ . When the variability of such changes is limited in the sense that  $\kappa(V_{\hat{Y}}) < Q/(s^*|\mathbb{M}^*|)$ , the RLS-based estimator is more efficient than  $\hat{\gamma}$ . Moreover, the above result can be extend to compare the length of the confidence intervals as well. The conclusion is similar. See Proposition S1 in the supplementary material for the details.



## 5 Post-screening inference under imperfect screening

Similar to many other consistency results for variable selection, the perfect screening property can be too much to hope for in practical data analysis in factorial designs. This is because the consistency of forward screening is a theoretical property under regularity conditions. Motivated by the hierarchy principle for factorial effects, the main factorial effects and lower-order factorial effects are more likely to be non-negligible than the higher-order factorial effects. With smaller effect size for the higher-order interactions, the perfect screening property more less likely to hold even asymptotically. Technically, when Condition 2(i) is violated, Algorithm 1 may no longer enjoy the perfect screening property.

Statistical inference without perfect screening is a non-trivial problem in factorial designs. If we do not put any restriction on the imperfect selection, the selected model can be anything, even without a stable limit. Classical strategies for post-selection inference (Kuchibhotla et al., 2022) will encounter various drawbacks in our current setup. For example, data splitting (Wasserman and Roeder, 2009) is a widely used strategy to validate inference after variable selection due to its simplicity. However, it highly relies on the independent sampling assumptions, which is violated under complete randomization. On the other hand, selective inference (Fithian et al., 2014) is another widely studied strategy. As another example, selective inference delivers valid inference for data-dependent parameters. However, it cannot be directly applied to analyze data collected in factorial designs. This is because selective inference strategy tends to rely on specific selection methods and parametric modelling assumptions on the outcome variables.

Rather than directly generalizing classical post selection inference methods to factorial experiments, in this section, we shall discuss two alternative strategies leveraging the special data structures in factorial experiments, along with with their statistical inference results (summarized in Figure 1).

### 5.1 Two alternative strategies for imperfect screening and statistical inference

The two proposed strategies are built on a belief that perfect screening is more plausible for selecting the main factorial effects and lower-order factorial effects up to level  $d^*$  than the high order effects. We will add more discussion on  $d^*$  after presenting the two strategies.

For Strategy 1, when the higher-order factorial effects are not necessary in the study, we may stop our forward screening procedure in Algorithm 1 at  $d = d^*$  (instead of  $d = D$ ). Such a strategy

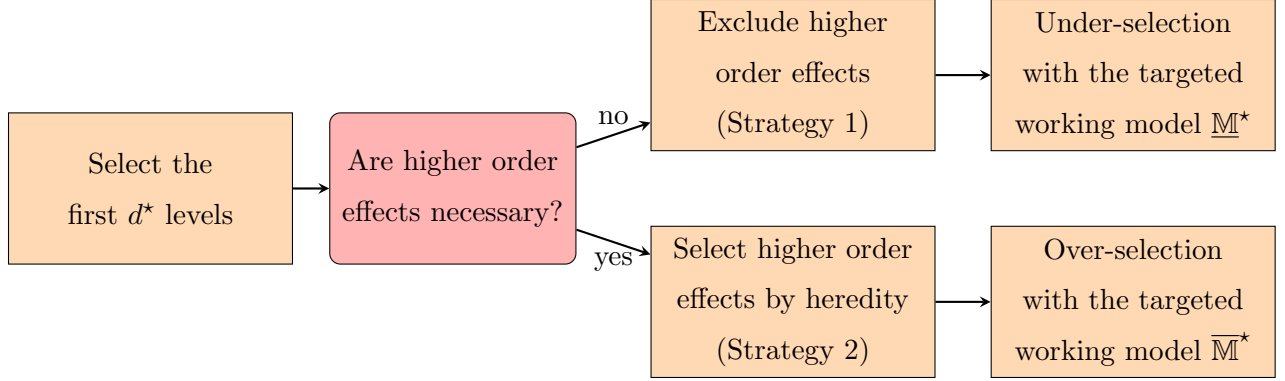


Figure 1: Two strategies for factorial screening: Strategy 1 under selects whereas Strategy 2 over selects.

focus on recovering a targeted working model  $\underline{\mathbb{M}}^*$  up to level  $d^*$ , that is,

$$\underline{\mathbb{M}}^* = \cup_{d=1}^{d^*} \mathbb{M}_d^* \subseteq \mathbb{M}^*,$$

which leads to an under-selected parsimonious working model. We summarize this strategy below.

**Strategy 1** (Under selection by excluding high-order interactions). *In Algorithm 1, we stop the screening procedure at  $d = d^*$ . Or equivalently, we set  $\alpha_d = \infty$  for  $d \geq d^* + 1$  so that no effects beyond level  $d^*$  will be selected and  $\underline{\hat{\mathbb{M}}} = \cup_{d=1}^{d^*} \hat{\mathbb{M}}_d$ .*

Given the selected working model  $\underline{\hat{\mathbb{M}}}$ , we can again construct an estimator of  $\gamma = \mathbf{f}^\top \bar{\mathbf{Y}}$  (defined in Section 4.1) based on RLS:

$$\hat{\gamma}_{\text{RU}} = \mathbf{f}[\underline{\hat{\mathbb{M}}}]^\top \hat{\mathbf{Y}}, \quad \text{and} \quad \hat{v}_{\text{RU}}^2 = \mathbf{f}[\underline{\hat{\mathbb{M}}}]^\top \hat{\mathbf{V}}_{\hat{\mathbf{Y}}} \mathbf{f}[\underline{\hat{\mathbb{M}}}]. \quad (5.11)$$

For Strategy 2, rather than excluding all higher-order interactions with negligible effects, we may further leverage the heredity principle and continue our screening procedure beyond level  $d^*$ . This means that instead of selecting the higher order interactions via marginal t-test and Bonferroni correction, we select the higher order interaction terms whenever either all of their parent effects are selected (strong heredity) or one of their parent effects is selected (weak heredity). While such a strategy takes higher order factorial effects into account, it often targets a working model  $\overline{\mathbb{M}}^*$  that includes the true model  $\mathbb{M}^*$ , that is,

$$\mathbb{M}^* \subseteq \overline{\mathbb{M}}^* = \bigcup_{d=1}^D \overline{\mathbb{M}}_d^*, \quad \text{where } \overline{\mathbb{M}}_d^* = \begin{cases} \mathbb{M}_d^*, & d \leq d^*; \\ \mathbf{H}^{(d-d^*)}(\mathbb{M}_{d^*}^*), & d^* + 1 \leq d \leq D. \end{cases}$$

The selected model by this strategy is expected to introduce an over-selected model that includes  $\mathbb{M}^*$  as well. We summarize this strategy as follows:

**Strategy 2** (Over selection by including higher-order interactions through the heredity principle). *In Algorithm 1, set  $\alpha_d = 0, d \geq d^* + 1$  and apply a heredity principle (either weak or strong, depending on people's knowledge on the structure of the effects). Then the high order effects beyond level  $d^*$  are selected merely by heredity principle and*

$$\widehat{\mathbb{M}} = \cup_{d=1}^D \widehat{\mathbb{M}}_d \text{ where } \widehat{\mathbb{M}}_d = \begin{cases} \text{Algorithm 1 Output,} & d \leq d^*; \\ H^{(d-d^*)}(\widehat{\mathbb{M}}_{d^*}), & d^* + 1 \leq d \leq D. \end{cases}$$

Here  $H^{(d-d^*)}$  is the  $(d - d^*)$ -order composition of  $H$ , meaning applying  $H$  for  $(d - d^*)$  times.

Given the selected working model  $\widehat{\mathbb{M}}$ , similarly, we can construct an estimator of  $\gamma = \mathbf{f}^\top \bar{Y}$  based on RLS:

$$\hat{\gamma}_{\text{RO}} = \mathbf{f}[\widehat{\mathbb{M}}]^\top \hat{Y}, \quad \text{and} \quad \hat{v}_{\text{RO}}^2 = \mathbf{f}[\widehat{\mathbb{M}}]^\top \hat{V}_{\hat{Y}} \mathbf{f}[\widehat{\mathbb{M}}]. \quad (5.12)$$

In terms of implementation, one can use WLS to conveniently obtain the point estimators in (5.11) and (5.12) and slightly more conservative variance estimators. Due to the orthogonality of the contrast matrix  $G$ , perfect screening is not required for computation. See Section A.1 in the supplementary material for more detailed discussions.

In real-world factorial experiments, how should practitioners decide which strategy to work with? This relies on the domain knowledge and the research question of interest. Strategy 1 is more suitable when there are domain-specific messages indicating that higher order interactions are small, or the research question only involves lower order factorial effects. Moreover, Strategy 1 is helpful when the number of active lower order interaction is large and Strategy 2 cannot be applied. Meanwhile, Strategy 2 works better when domain knowledge suggests non-negligible higher order interactions or the research question targets a more general parameter beyond factorial effects themselves. It also works well when the number of active lower order interaction is small, and we can include a small set of high order terms according to the heredity principle.

In the following subsection, we study the statistical properties of  $\hat{\gamma}_{\text{RO}}$  and  $\hat{\gamma}_{\text{RU}}$  and demonstrate the trade-offs between the two strategies for statistical inference from a theoretical perspective.

## 5.2 Theoretical properties under imperfect screening

Throughout this subsection, we discuss the scenario where perfect screening is hard to achieve. We work under a relaxed condition of Condition 2 defined as follows:

**Condition 5** (Order of parameters up to level  $d^*$ ). *Condition 2 holds with  $D = d^*$ .*

Condition 5 no longer imposes any restriction on the order of the parameters beyond level  $d^*$ . By Theorem 1, Condition 5 guarantees that Algorithm 1 perfectly screens the first  $d^*$  levels of factorial effects in the sense that

$$\mathbb{P} \left\{ \widehat{\mathbb{M}}_d = \mathbb{M}_d^* \text{ for } d = 1, \dots, d^* \right\} \rightarrow 1.$$

We start by analyzing the statistical property of  $\widehat{\gamma}_{\text{RU}}$  with  $\widehat{\mathbb{M}}$  obtained from the under selection Strategy 1. Because the selected working model might deviate from the truth beyond level  $d^*$ ,  $\widehat{\gamma}_{\text{RU}}$  may not be an consistent estimator of  $\gamma$ . Therefore, we focus on weighting vectors  $\mathbf{f}$  that satisfy certain orthogonality conditions as introduced in Theorem 3 below:

**Theorem 3** (Guarantee for Strategy 1). *Recall the equation (4.8) and define  $\underline{\mathbf{f}}^* = \mathbf{f}[\underline{\mathbb{M}}^*] = Q^{-1}G(\cdot, \underline{\mathbb{M}}^*)G(\cdot, \underline{\mathbb{M}}^*)^\top \mathbf{f}$ . Assume Conditions 1, 3, 4, 5, and  $\mathbf{f}$  satisfies the following orthogonality condition:*

$$G(\cdot, \mathbb{M}_d^*)^\top \mathbf{f} = 0 \text{ for } d^* + 1 \leq d \leq K. \quad (5.13)$$

Let  $N \rightarrow \infty$ . We have

$$\frac{\widehat{\gamma}_{\text{RU}} - \gamma}{v_{\text{RU}}} \rightsquigarrow \mathcal{N}(0, 1),$$

where  $v_{\text{RU}}^2 = \underline{\mathbf{f}}^{*\top} V_{\widehat{\gamma}} \underline{\mathbf{f}}^*$ . Further assume  $\|\underline{\mathbf{f}}^*\|_\infty = O(Q^{-1})$ . The variance estimator  $\widehat{v}_{\text{RU}}^2$  is conservative in the sense that:

$$N(\widehat{v}_{\text{RU}}^2 - v_{\text{RU}, \text{lim}}^2) \xrightarrow{\mathbb{P}} 0, \quad v_{\text{RU}, \text{lim}}^2 \geq v_{\text{RU}}^2,$$

where  $v_{\text{RU}, \text{lim}}^2 = \underline{\mathbf{f}}^{*\top} D_{\widehat{\gamma}} \underline{\mathbf{f}}^*$  is the limiting value of  $\widehat{v}_{\text{RU}}^2$ .

Now we add some discussion on Theorem 3. The orthogonality condition presented in (5.13) restricts the weighting vector  $\mathbf{f}$  to be orthogonal to the higher order contrasts. Intuitively, because the higher-order interactions are excluded from the model, making inference on a weighted combination of those excluded interactions is infeasible. One set of weighting vectors satisfying (5.13) is the contrast vectors of nonzero canonical lower-order interactions, given by  $\mathbf{f} = G(\cdot, \cup_{d=1}^{d^*} \mathbb{M}_d^*)$ . In large  $K$  settings, the lower order interactions can also grow polynomially fast in  $K$  and add difficulty for interpretation. As an example, when  $K = 10$ , for the first two levels of factorial effects without screening, there are a total of more than 50 estimates. It can still greatly benefit the analysis and interpretation to filter out the insignificant ones and obtain a parsimonious, structured working model.

As for Strategy 2, similarly, we have the following results:

**Theorem 4** (Guarantee for Strategy 2). Recall the equation (4.8) and define  $\bar{\mathbf{f}}^\star = f[\bar{\mathbf{M}}^\star] = Q^{-1}G(\cdot, \bar{\mathbf{M}}^\star)G(\cdot, \bar{\mathbf{M}}^\star)^\top \mathbf{f}$ . Assume Conditions 1, 3, 4 and 5. Let  $N \rightarrow \infty$ . If  $|\bar{\mathbf{M}}^\star|/N \rightarrow 0$ , then

$$\frac{\hat{\gamma}_{\text{RO}} - \gamma}{v_{\text{RO}}} \rightsquigarrow \mathcal{N}(0, 1),$$

where  $v_{\text{RO}}^2 = \bar{\mathbf{f}}^{\star\top} V_{\hat{\gamma}} \bar{\mathbf{f}}^\star$ . Further assume  $\|\bar{\mathbf{f}}^\star\|_\infty = O(Q^{-1})$ . The variance estimator  $\hat{v}_{\text{RO}}^2$  is conservative in the sense that:

$$N(\hat{v}_{\text{RO}}^2 - v_{\text{RO,lim}}^2) \xrightarrow{\mathbb{P}} 0, \quad v_{\text{RO,lim}}^2 \geq v_{\text{RO}}^2,$$

where  $v_{\text{RO,lim}}^2 = \bar{\mathbf{f}}^{\star\top} D_{\hat{\gamma}} \bar{\mathbf{f}}^\star$  is the limiting value of  $\hat{v}_{\text{RO}}^2$ .

We comment that there is an additional technical requirement in Theorem 4 for over-selection: we assume  $|\bar{\mathbf{M}}^\star|/N \rightarrow 0$ . This equation mainly serves as a sufficient condition for CLT. The reason is that we need to control the size of the target model  $\bar{\mathbf{M}}^\star$  compared to the sample size  $N$  in order to infer a general causal parameter.

When analyzing Strategies 1 and 2, Algorithm 1 recovers a targeted model with high probability. Both strategies have advantages and disadvantages. Under-selection reflects bias-variance trade-off: it can induce more bias for certain weighting vectors, but the constructed estimator typically enjoys smaller variance. Over-selection can reduce bias for estimation, but may not be feasible if there are too many lower order terms which can result in many redundant terms into the selected model. In practice, if higher order interactions are not crucial for study, Strategy 1 should be applied. If high order interactions are of interest and hard to select, one could pursue Strategy 2 as a practically useful and interpretable solution.

**Remark 2.** Under equal-sized designs with  $N(\mathbf{z}) = N_0$  and homoskedasticity with  $S(\mathbf{z}, \mathbf{z}) = S_0$  for all  $\mathbf{z} \in \mathcal{T}$ , we can prove  $v_{\text{RU}}^2 \leq v_{\text{RO}}^2$ . Therefore, by excluding higher order terms and pursuing under-selection, we can obtain an equal or smaller asymptotic variance compared with over-selection. In general, due to heteroskedasticity, the order of  $v_{\text{RU}}^2$  and  $v_{\text{RO}}^2$  depends on the choice of target weighing vector  $\mathbf{f}$ . Here we take a sparse  $\mathbf{f} = \mathbf{e}_1 = (1, 0, \dots, 0)^\top$  as an example. We can show that

$$\frac{v_{\text{RU}}^2}{v_{\text{RO}}^2} \leq \kappa(V_{\hat{\gamma}}) \cdot \frac{|\bar{\mathbf{M}}^\star|}{|\underline{\mathbf{M}}^\star|}.$$

When the variability of  $V_{\hat{\gamma}}$  between treatment arms is small in the sense that  $\kappa(V_{\hat{\gamma}}) < |\bar{\mathbf{M}}^\star|/|\underline{\mathbf{M}}^\star|$ , under-selection leads to smaller asymptotic variance for inferring  $\mathbf{e}_1^\top \bar{\mathbf{Y}}$ .

## 6 Application to inference on the best arm in factorial experiments

In the previous sections, we consider the problem of making inference on a single factorial causal effect  $\gamma = \mathbf{f}^\top \bar{Y}$ . As an application of the proposed framework, we study the problem of inference on best causal effects. Section 6.1 introduces the basic problem and an inferential procedure. Section 6.2 presents the theoretical guarantees.

### 6.1 Best arms, tie set and statistical inference

Suppose we have a set of causal effects  $\Gamma$  defined by pre-specified weighting vectors  $\mathbf{f}_1, \dots, \mathbf{f}_L$  ( $L$  is potentially large), that is

$$\Gamma = \{\gamma_1, \dots, \gamma_L\}, \quad \gamma_l = \mathbf{f}_l^\top \bar{Y}.$$

We aim to perform statistical inference on their ordered values

$$\gamma_{(1)} \geq \dots \geq \gamma_{(l_0)}$$

with  $l_0 < L$  being a fixed positive integer. As a concrete example, if we choose  $\{\mathbf{f}_l\}_{l \in [L]} = \{\mathbf{e}(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}$  to be the set of the canonical bases  $\{\mathbf{e}(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}$ , then our inferential target is the maximal potential outcome means:

$$\bar{Y}_{(1)} = \max_{\mathbf{z} \in \mathcal{T}} \bar{Y}(\mathbf{z}). \quad (6.14)$$

A more practically appealing consideration in factorial experiments is to incorporate structural constraints into the choices of  $\{\mathbf{f}_l\}_{l \in [L]}$ . For example, when inferring (6.14) in a real-world factorial experiment, it might be infeasible to consider all treatment levels  $\mathcal{T}$  due to budget concerns or resource limitation. Instead, we may only be interested in factor combinations  $\mathbf{z} = (z_1, \dots, z_K)$  with at most  $K_0 (\leq K)$  1's. In other words, we replace  $\mathcal{T}$  with the following  $\mathcal{T}'$  in (6.14) and obtain:

$$\mathcal{T}' = \left\{ \mathbf{z} = (z_1, \dots, z_K) \mid \sum_{k=1}^K z_k \leq K_0 \right\}, \quad \bar{Y}_{(1)} = \max_{\mathbf{z} \in \mathcal{T}'} \bar{Y}(\mathbf{z}). \quad (6.15)$$

Such constraints are unique and important to factorial experiments, especially when  $K$  is large. By focusing on  $\{\mathbf{f}_l\}_{l \in [L]}$  that are most pertinent to factorial designs, the problem allows us to use the available data to decide if the best causal parameter among those practically interesting ones has a non-zero causal effect.

Two challenges exist in delivering valid statistical inference on  $\gamma_{(1)}, \dots, \gamma_{(l_0)}$  in factorial experiments. On the one hand, sample analogs of the order parameters,  $(\hat{\gamma}_{(1)}, \dots, \hat{\gamma}_{(l_0)})$ , are often biased

estimates of  $(\gamma_{(1)}, \dots, \gamma_{(l_0)})$  due to the well-known winner's curse phenomenon (Andrews et al., 2019; Guo et al., 2021; Wei et al., 2022). On the other hand, although one might argue that existing approaches can be applied to remove the winner's curse bias in  $\hat{\gamma}_{(l)}$ , these approaches do not account for the special structural constraint in factorial experiments. Moreover, rigorous statistical guarantees are lacking in the context due to the unique presence of both large  $L$  and large  $Q$  in factorial designs.

To simultaneously address the above challenges, we propose a procedure that tailors the tie-set identification approach proposed in Claggett et al. (2014) and Wei et al. (2022) to our current problem setup. We focus on making inference on the first ordered value  $\gamma_{(1)}$  to simplify discussion, and our approach extends naturally to other ordered values. The proposed procedure is provided in Algorithm 2.

Algorithm 2 consists of three major components. First, we need to construct  $\hat{\gamma}_l = \mathbf{f}_l^\top \hat{\mathbf{Y}}_R$  with feature screening (Step 1-2). These RLS based estimators enjoy great benefits for large  $Q$  and small  $N_0$  regimes based on our previous discussion. Second, we construct  $\hat{\mathcal{L}}_1$  to include the estimates that are close to  $\hat{\gamma}_{(1)}$  (Step 3). Intuitively, these collected estimates are different due to random error. We will show that with proper tuning, this procedure will include all the  $l$  for which  $\gamma_l$  are statistically indistinguishable from  $\gamma_{(1)}$  with high probability. Third, we construct estimators by averaging over  $\hat{\mathcal{L}}_1$  (Step 4). By averaging the estimates over the selected  $\hat{\mathcal{L}}_1$  we alleviate the impact of randomness and obtain accurate estimates for the maximal effect.

## 6.2 Theoretical guarantees

In the following, we present theoretical guarantees for Algorithm 2. We introduce the following notation  $\mathcal{L}_1$  to include all effects that stay in a local neighbourhood of  $\gamma_{(1)}$ :

$$\mathcal{L}_1 = \left\{ l \in [L] \mid |\gamma_l - \gamma_{(1)}| = O(N^{-\delta_3}) \right\}, \text{ for some } \delta_3 > 0.$$

A well-known fact is that the naive estimator  $\max_{\mathbf{z} \in [Q]} \hat{Y}(\mathbf{z})$  can be an overly optimistic estimate for  $\gamma_{(1)}$  when  $\mathcal{L}_1$  contains more than one element (Andrews et al., 2019; Wei et al., 2022). Define

$$d_h = \max_{\mathbf{z} \in \mathcal{L}_1} |\gamma_l - \gamma_{(1)}|, \quad d_h^* = \min_{\mathbf{z} \notin \mathcal{L}_1} |\gamma_l - \gamma_{(1)}|.$$

as within- and between-group distances, respectively. We work under the following condition:

**Condition 6** (Order of  $d_h$ ,  $d_h^*$  and  $\eta_N$ ). *Assume the within and between group distances satisfy:*

$$d_h^* = \Theta(N^{\delta_1}), \quad \eta_N = \Theta(N^{\delta_2}), \quad d_h = \Theta(N^{\delta_3}).$$

---

**Algorithm 2:** Inference on best causal effect(s)

---

**Input:** Factorial data  $(Y_i, Z_i)$ ; predetermined integer  $D$ ; initial model for factorial effects

$\widehat{\mathbb{M}} = \{\emptyset\}$ ; significance level  $\{\alpha_d\}_{d=1}^D$ ; set of weighting vectors  $\{\mathbf{f}_l\}_{l \in [L]}$ ; thresholds  $\eta_N$ .

**Output:** Selected working model  $\widehat{\mathbb{M}}$ .

- 1 Perform forward effects screening with Algorithm 1 and obtain working model  $\widehat{\mathbb{M}}$ .
- 2 Obtain RLS based estimates: use Equation (4.8) and definition of  $\widehat{Y}_R$  (4.7) to compute

$$\mathbf{f}_l[\widehat{\mathbb{M}}] = Q^{-1}G(\cdot, \widehat{\mathbb{M}})G(\cdot, \widehat{\mathbb{M}})^\top \mathbf{f}_l, \quad \widehat{\gamma}_l = \mathbf{f}_l^\top \widehat{Y}_R = \mathbf{f}_l[\widehat{\mathbb{M}}]^\top \widehat{Y}, \quad l \in [L].$$

- 3 Record the set of effects close to  $\widehat{\gamma}_{(1)}$ :

$$\widehat{\mathcal{L}}_1 = \{l \in [L] \mid |\widehat{\gamma}_l - \widehat{\gamma}_{(1)}| \leq \eta_N\}.$$

Here,  $\eta_N$  is a tuning parameter which can be selected using the algorithm provided in Wei et al. (2022), Appendix C.1.

- 4 Define

$$\mathbf{f}_{(1)} = (Q|\widehat{\mathcal{L}}_1|)^{-1} \sum_{l \in \widehat{\mathcal{L}}_1} G(\cdot, \widehat{\mathbb{M}})G(\cdot, \widehat{\mathbb{M}})^\top \mathbf{f}_l.$$

Generate point estimates and variance estimator for  $\gamma_{(1)}$ :

$$\widehat{Y}_{(1)} = \frac{1}{|\widehat{\mathcal{L}}_1|} \sum_{l \in \widehat{\mathcal{L}}_1} \widehat{\gamma}_l = \mathbf{f}_{(1)}^\top \widehat{Y}, \quad \widehat{v}_{(1)}^2 = \mathbf{f}_{(1)}^\top \widehat{V}_Y \mathbf{f}_{(1)}.$$

**5 return**  $\widehat{\mathcal{L}}_1, \widehat{Y}_{(1)}, \widehat{v}_{(1)}^2$

---

with  $\delta_3 \leq -1/2 < \delta_2 < \delta_1 \leq 0$ .

Define the population counterpart of  $\mathbf{f}_{(1)}$  as

$$\mathbf{f}_{(1)}^* = (Q|\mathcal{L}_1|)^{-1} \sum_{l \in \mathcal{L}_1} G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \mathbf{f}_l.$$

We establish the following result for the procedure provided in Algorithm 2. Recall  $\delta_2$  from Condition 6 and  $\delta''$  from Condition 2(iii), which characterizes the magnitude of the within/between group distances and the size of the true working model, respectively.

**Theorem 5** (Asymptotic results on the estimated effects using Algorithm 2). *Assume Condition*



1-4 and 6. Let  $N \rightarrow \infty$ . If

$$N^{-(1+2\delta_2-\delta'')} \rightarrow 0, \quad (6.16)$$

$$L \cdot |\mathcal{L}_1| \cdot N^{-\frac{1-\delta''}{2}} \rightarrow 0, \quad (6.17)$$

then

$$\frac{\hat{\gamma}_{(1)} - \gamma_{(1)}}{v_{(1)}} \rightsquigarrow \mathcal{N}(0, 1),$$

where  $v_{(1)}^2 = \mathbf{f}_{(1)}^{\star\top} V_Y \mathbf{f}_{(1)}^{\star}$ . Moreover,  $\hat{v}_{(1)}^2$  is conservative in the sense that

$$N(\hat{v}_{(1)}^2 - v_{(1),\lim}^2) \xrightarrow{\mathbb{P}} 0, \quad v_{(1),\lim}^2 \geq v_{(1)}^2,$$

where  $v_{(1),\lim}^2 = \mathbf{f}_{(1)}^{\star\top} D_{\hat{Y}} \mathbf{f}_{(1)}^{\star}$  is the limiting value of  $v_{(1)}^2$ .

The conditions in Theorem 5 are mild and reveals an interesting trade-off between some mathematical quantities. For the first asymptotic condition (6.16), when the size of the targeted working model is small compared to  $N$ , say  $\delta'' = 0$  (meaning  $|\mathbb{M}^{\star}|$  does not grow with  $N$ ), this condition always holds. More generally, (6.16) is easier to satisfy with larger gap size (larger  $\delta_2$ ) and smaller true working model size (smaller  $\delta''$ ). The second condition (6.17) reflects the tradeoff among the total number of interested parameters ( $L$ , which is also  $|\mathcal{T}'|$ ), the size of the neighborhood of  $\gamma_{(1)}$  ( $|\mathcal{L}_1|$ ), and the size of the true working model ( $\delta''$ ). The smaller these quantities are compared to  $N$ , the easier inference will be. Moreover, (6.17) is easily justifiable. Going back to the previous example (6.15), (6.17) translates into

$$\sum_{K=0}^{K_0} \binom{K}{k} \cdot |\mathcal{L}_1| \cdot \left( \frac{|\mathbb{M}^{\star}|}{N} \right)^{1/2} \rightarrow 0. \quad (6.18)$$

(6.18) accommodates a variety of interesting regimes with different specifications of  $K_0$ ,  $|\mathcal{L}_1|$  and  $|\mathbb{M}^{\star}|$ . We omit the discussion here.

Theorem 5 also suggests the benefits of factor screening compared to procedures where no screening is involved. Recall our discussion in Remark 1. Without screening, one requires  $Q$  to be small compared to  $N$  or  $\{\mathbf{f}_l\}_{l \in [L]}$  are dense, which is violated in large  $Q$  setups and many practical scenarios such as (6.14). On the contrary, factor screening can capture the shared information across treatment arms and still deliver valid inference.

As a final comment, the result of our Theorem 5 relies on the perfect screening property (Theorem 1), which are ensured by Conditions 1 - 4. Without perfect screening, there might be additional sources of bias due to the uncertainty induced by the screening step and possible under-selection results. Nevertheless, one can consider applying the over-selection strategy (Strategy 2 in Section 5.1) to facilitate inference on the best factorial effects.

## 7 Simulation

In this section, we use simulation studies to demonstrate the finite-sample performance of the proposed forward screening framework and the inferential properties of the RLS-based estimator. More concretely, our simulation results verify the following properties of the proposed procedure and estimators:

**(G1)** The RLS-based estimator  $\hat{\gamma}_R$  demonstrates efficiency gain (in terms of improved power and shortened confidence interval) compared to the simple moment estimator  $\hat{\gamma}$  for general causal parameters defined by sparse weighting vectors.

**(G2)** The factorial forward screening procedure provided in Algorithm 1 can improve the performance of effect screening compared to naive procedure (i.e., screening without leveraging the heredity principle).

**(G1)** echoes our discussion on the comparison of CLT conditions and asymptotic variance in Remark 1 and Proposition 1. **(G2)** verifies the results in Theorem 1 and 2 and checks the finite sample behaviors of the proposed procedures. For both goals, we will vary the sample size and effect size to provide a comprehensive understanding of their performance.

### 7.1 Simulation setup

We set up a  $2^8$  factorial experiment ( $K = 8$ ). There are  $N_0$  units in each treatment arm where  $N_0$  is set to be a varying number. We generate independent potential outcomes from a shifted exponential distribution:

$$Y_i(\mathbf{z}) \sim \text{EXP}(1) - 1 + \mu(\mathbf{z}).$$

Here  $\mu(\mathbf{z})$  are super population means of potential outcomes under treatment  $\mathbf{z}$ . We choose  $\mu(\mathbf{z})$  such that the factorial effects satisfy the following structure:

- Main effects: the main effects corresponding to the first five factors,  $\tau_{\{1\}}, \dots, \tau_{\{5\}}$ , are nonzero; the rest three main effects,  $\tau_{\{6\}}, \dots, \tau_{\{8\}}$ , are zero.
- Two-way interactions: the two-way interactions associated with the first five factors are nonzero, i.e.,  $\tau_{\{kl\}} \neq 0$  for  $k \neq l, k, l \in [5]$ . All the rest of the two-way interactions are zero.
- Higher order interactions: all the higher-order interactions  $\tau_{\mathcal{K}}$  are zero if  $|\mathcal{K}| \geq 3$ .

The above setup of factorial effects guarantees that they are sparse and follow the strong heredity principle. In the provided simulation results, we will vary the number of units in each treatment arm and the size of the nonzero factorial effects. More details can be found in the R code attached to the support materials.

## 7.2 Simulation results supporting (G1)

In this subsection, we evaluate the performance of the RLS-based estimators  $(\hat{\gamma}_R, \hat{v}_R)$  compared to  $(\hat{\gamma}, \hat{v})$  for testing a causal effect  $\gamma_{\text{target}} = \mathbf{f}^\top \bar{Y}$  specified by a sparse vector:  $\mathbf{f} = (0, \dots, 0, 1)^\top \in \mathbb{R}^Q$ . Intuitively,  $\gamma_{\text{target}}$  measures the average of potential outcomes in the last level. For each estimator, we report: (i) power for testing  $H_0 : \gamma_{\text{target}} = 0$ . (ii) coverage probability of the confidence intervals for  $\gamma_{\text{target}}$  at level 0.95. Figure 2 summarizes the results.

Figure 2 shows that the RLS-based estimator  $\hat{\gamma}_R$  has power one and coverage probability 0.95 asymptotically when the effect size or the number of replications is large enough. The simple moment estimator  $\hat{\gamma}$ , on the other hand, suffers from deficiency in power in finite sample. While the coverage probability appears high for small  $N_0$  and effect size, it is largely due to an overly conservative confidence interval given the low power. This echoes our conclusion that effect screening can incorporate information across treatment levels and improve the power for testing general causal effects.

## 7.3 Simulation results for (G2)

In this subsection, we compare the performance of forward and naive screening procedures. We use four candidate effect screening methods:

- *Forward Bonferroni*. Forward screening based on Bonferroni corrected marginal t tests;
- *Forward Lasso*. Forward screening based on Lasso;
- *Naive Bonferroni*. Screening with the full working model based on Bonferroni corrected marginal t tests;
- *Naive Lasso*. Screening with the full working model based on Lasso.

For each of the screening methods, we evaluate their performance from several perspectives: (i) perfect screening probability  $\mathbb{P}\{\hat{\mathbb{M}} = \mathbb{M}^*\}$ . (ii) power of  $\hat{\gamma}_R$  for testing  $H_0 : \gamma_{\text{target}} = 0$  for the

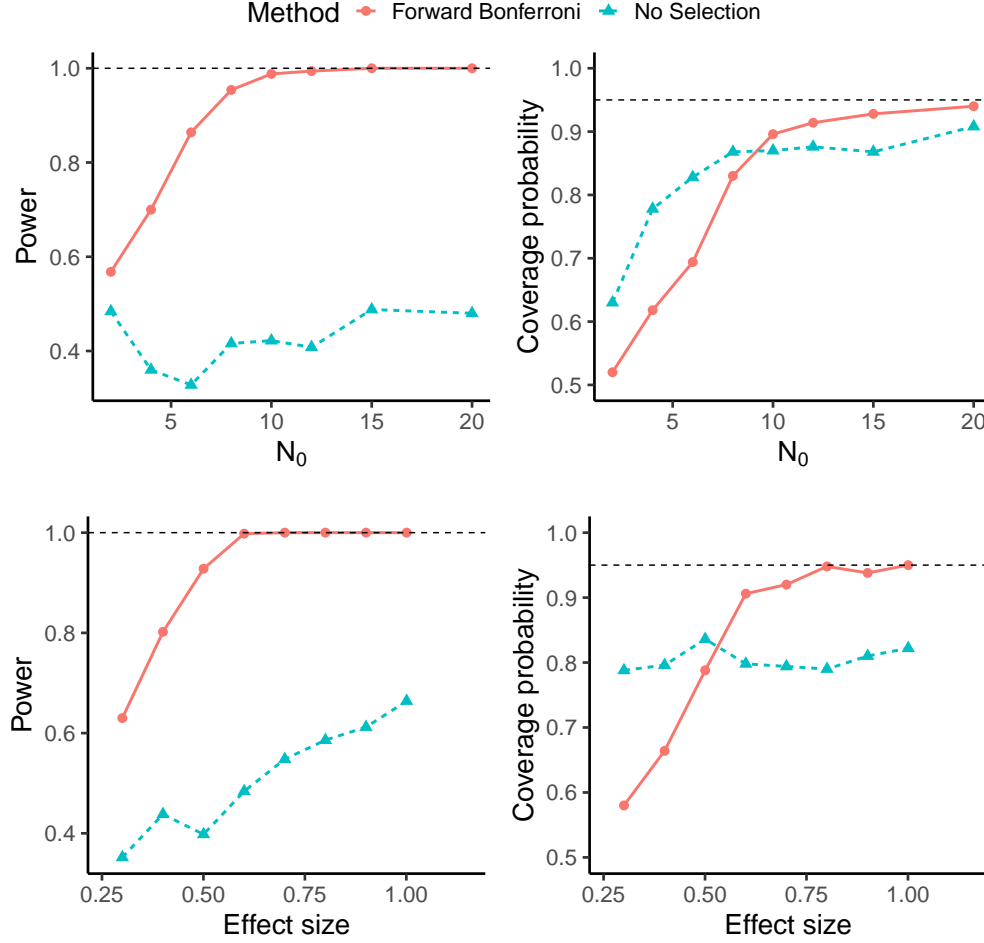


Figure 2: Simulation results on (G1). (i) Top left panel: power curve with varying  $N_0$ ; (ii) Top right panel: coverage probability with varying  $N_0$ ; (iii) Bottom left panel: power curve with varying effect size; (iv) Bottom right panel: coverage probability with varying effect size.

same  $\gamma_{\text{target}}$  defined in the previous section. (iii) coverage probability of the RLS based confidence interval for  $\gamma_{\text{target}}$  at level 0.95. The results are summarized in Figure 3.

From Figure 3, all four effect screening methods lead to perfect selection with high probability as the number of replications  $N_0$  or effect size increases. Nevertheless, with the forward screening procedure, the probability of perfect screening are higher than naive screening procedure. Besides, forward screening keeps the heredity structure and demonstrates higher interpretability than the naive screening methods. In terms of the power of  $\hat{\gamma}_R$  and  $\hat{v}_R$  for testing  $H_0 : \gamma_{\text{target}} = 0$ , all four methods have asymptotic power 1, while forward screening possesses higher power in finite samples. Analogously, we can see an improvement in coverage probability of the RLS based confidence

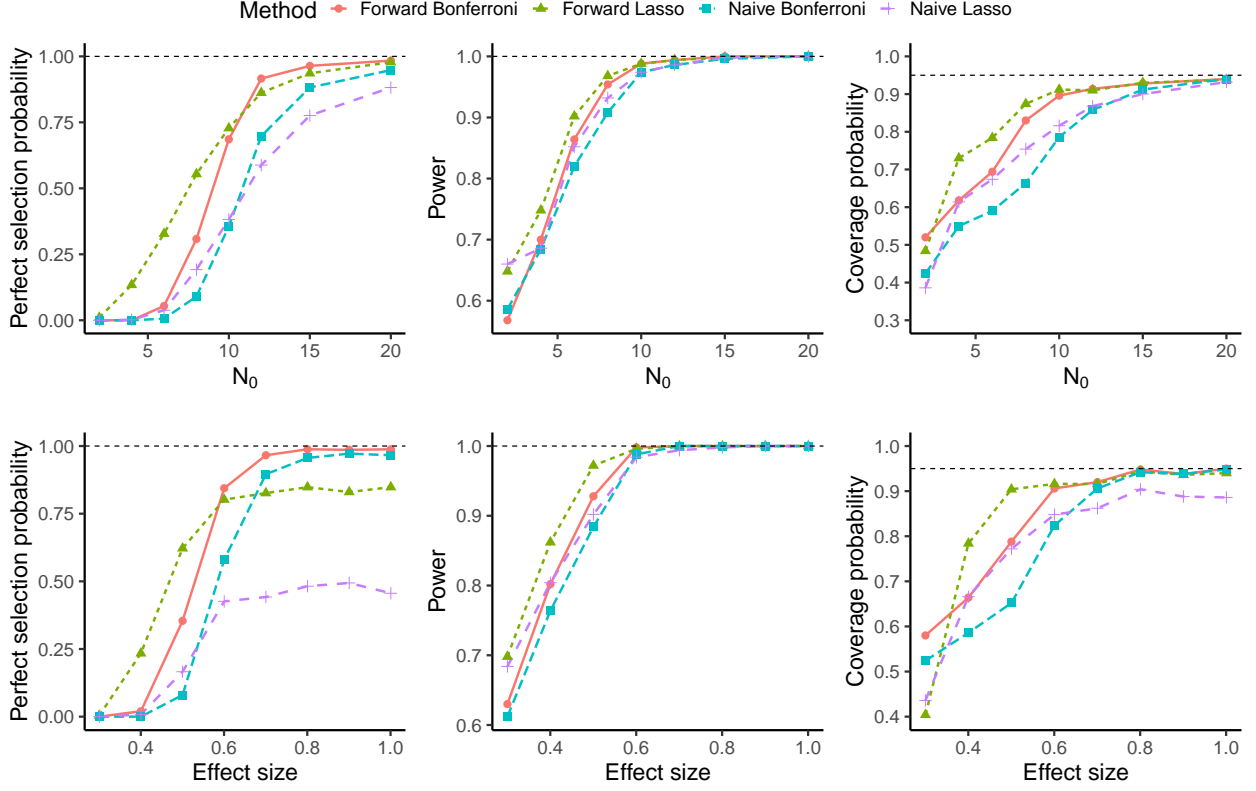


Figure 3: Simulation results on **(G2)**. (i) Top left panel: perfect screening probability with a small fixed effect size 0.20 and varying  $N_0$ ; (ii) Top middle panel: power curve with a small fixed effect size 0.20 and varying  $N_0$ ; (iii) Top right panel: coverage probability with a small fixed effect size 0.20 and varying  $N_0$ ; (iv) Bottom left panel: perfect screening probability with a small fixed replication  $N_0 = 2$  and varying effect size; (v) Bottom middle panel: power curve with a small fixed replication  $N_0 = 2$  and varying effect size; (vi) Bottom right panel: coverage probability with a small fixed replication  $N_0 = 2$  and varying effect size.

intervals with the forward screening procedure.

## 8 Discussion

We have discussed the formal theory for forward screening and post-screening inference in  $2^K$  factorial designs. It is conceptually straightforward to extend the theory to general factorial designs with multi-valued factors. We omit the technical details. Another important direction is covariate adjustment in factorial experiments. Lin (2013), Lu (2016a) and Liu et al. (2022) demonstrated

the efficiency gain of covariate adjustment with small  $K$ . Zhao and Ding (2023) discussed covariate adjustment in factorial experiment with factors and covariates selected independent of data. We leave it to future research to establish the theory for factor screening and covariate selection in factorial designs.

## References

- Andrews, I., Kitagawa, T., and McCloskey, A. (2019), “Inference on winners,” Tech. rep., National Bureau of Economic Research.
- Angrist, J. D. and Pischke, J.-S. (2009), *Mostly Harmless Econometrics: An Empiricist’s Companion*, Princeton: Princeton University Press.
- Bai, Z., Choi, K. P., Fujikoshi, Y., and Hu, J. (2022), “Asymptotics of AIC, BIC and Cp model selection rules in high-dimensional regression,” *Bernoulli*, 28, 2375–2403.
- Bickel, P. J., Ritov, Y., Tsybakov, A. B., et al. (2010), “Hierarchical selection of variables in sparse high-dimensional regression,” *IMS Collections*, 6, 28.
- Bien, J., Taylor, J., and Tibshirani, R. (2013), “A lasso for hierarchical interactions,” *Annals of Statistics*, 41, 1111.
- Blackwell, M. and Pashley, N. E. (2021), “Noncompliance and Instrumental Variables for  $2^K$  Factorial Experiments,” *Journal of the American Statistical Association*, in press.
- Bloniarz, A., Liu, H., Zhang, C.-H., Sekhon, J. S., and Yu, B. (2016), “Lasso adjustments of treatment effect estimates in randomized experiments,” *Proceedings of the National Academy of Sciences*, 113, 7383–7390.
- Box, G., Hunter, J., and Hunter, W. (2005), *Statistics for Experimenters: Design, Innovation, and Discovery*, Hoboken, NJ: Wiley.
- Branson, Z., Dasgupta, T., and Rubin, D. B. (2016), “Improving covariate balance in  $2^K$  factorial designs via rerandomization with an application to a New York City Department of Education High School Study,” *Annals of Applied Statistics*, 10, 1958–1976.
- Claggett, B., Xie, M., and Tian, L. (2014), “Meta-analysis with fixed, unknown, study-specific parameters,” *Journal of the American Statistical Association*, 109, 1660–1671.

- Dasgupta, T., Pillai, N. S., and Rubin, D. B. (2015), “Causal inference from  $2^K$  factorial designs by using potential outcomes,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77, 727–753.
- Egami, N. and Imai, K. (2019), “Causal interaction in factorial experiments: application to conjoint analysis,” *Journal of the American Statistical Association*, 114, 529–540.
- Espinosa, V., Dasgupta, T., and Rubin, D. B. (2016), “A Bayesian perspective on the analysis of unreplicated factorial experiments using potential outcomes,” *Technometrics*, 58, 62–73.
- Fan, J. and Lv, J. (2008), “Sure independence screening for ultrahigh dimensional feature space,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70, 849–911.
- Fithian, W., Sun, D., and Taylor, J. (2014), “Optimal inference after model selection,” *arXiv preprint arXiv:1410.2597*.
- Freedman, D. A. (2008), “On regression adjustments to experimental data,” *Advances in Applied Mathematics*, 40, 180–193.
- Gerber, A. S. and Green, D. P. (2012), *Field Experiments: Design, Analysis, and Interpretation*, New York, NY: Norton.
- Guo, X., Wei, L., Wu, C., and Wang, J. (2021), “Sharp Inference on Selected Subgroups in Observational Studies,” *arXiv preprint arXiv:2102.11338*.
- Hao, N., Feng, Y., and Zhang, H. H. (2018), “Model selection for high-dimensional quadratic regression via regularization,” *Journal of the American Statistical Association*, 113, 615–625.
- Hao, N. and Zhang, H. H. (2014), “Interaction screening for ultrahigh-dimensional data,” *Journal of the American Statistical Association*, 109, 1285–1301.
- Haris, A., Witten, D., and Simon, N. (2016), “Convex modeling of interactions with strong heredity,” *Journal of Computational and Graphical Statistics*, 25, 981–1004.
- Hastie, T., Tibshirani, R., Friedman, J. H., and Friedman, J. H. (2009), *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, vol. 2, New York: Springer.
- Kempthorne, O. (1952), *The Design and Analysis of Experiments*, New York: Wiley.
- Kuchibhotla, A. K., Kolassa, J. E., and Kuffner, T. A. (2022), “Post-selection inference,” *Annual Review of Statistics and Its Application*, 9, 505–527.

- Li, X. and Ding, P. (2017), “General forms of finite population central limit theorems with applications to causal inference,” *Journal of the American Statistical Association*, 112, 1759–1769.
- Lim, M. and Hastie, T. (2015), “Learning interactions via hierarchical group-lasso regularization,” *Journal of Computational and Graphical Statistics*, 24, 627–654.
- Lin, W. (2013), “Agnostic notes on regression adjustments to experimental data: Reexamining Freedman’s critique,” *Annals of Applied Statistics*, 7, 295–318.
- Liu, H., Ren, J., and Yang, Y. (2022), “Randomization-based Joint Central Limit Theorem and Efficient Covariate Adjustment in Randomized Block  $2^K$  Factorial Experiments,” *Journal of the American Statistical Association*, in press.
- Lu, J. (2016a), “Covariate adjustment in randomization-based causal inference for  $2^K$  factorial designs,” *Statistics and Probability Letters*, 119, 11–20.
- (2016b), “On randomization-based and regression-based inferences for  $2^K$  factorial designs,” *Statistics and Probability Letters*, 112, 72–78.
- Meng, X.-L. and Xie, X. (2014), “I got more data, my model is more refined, but my estimator is getting worse! Am I just dumb?” *Econometric Reviews*, 33, 218–250.
- Neyman, J. (1923/1990), “On the application of probability theory to agricultural experiments. Essay on principles. Section 9.” *Statistical Science*, 465–472.
- Pashley, N. E. and Bind, M.-A. C. (2023), “Causal Inference for Multiple Non-Randomized Treatments using Fractional Factorial Designs,” *Canadian Journal of Statistics*, in press.
- Rillig, M. C., Ryo, M., Lehmann, A., Aguilar-Trigueros, C. A., Buchert, S., Wulf, A., Iwasaki, A., Roy, J., and Yang, G. (2019), “The role of multiple global change factors in driving soil functions and microbial biodiversity,” *Science*, 366, 886–890.
- Shi, L. and Ding, P. (2022), “Berry–Esseen bounds for design-based causal inference with possibly diverging treatment levels and varying group sizes,” *arXiv preprint arXiv:2209.12345*.
- Tibshirani, R. (1996), “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society: Series B (Methodological)*, 58, 267–288.
- Wainwright, M. J. (2019), *High-dimensional Statistics: A Non-asymptotic Viewpoint*, vol. 48, Cambridge: Cambridge University Press.



- Wang, H. (2009), “Forward regression for ultra-high dimensional variable screening,” *Journal of the American Statistical Association*, 104, 1512–1524.
- Wasserman, L. and Roeder, K. (2009), “High dimensional variable selection,” *Annals of Statistics*, 37, 2178.
- Wei, W., Zhou, Y., Zheng, Z., and Wang, J. (2022), “Inference on the Best Policies with Many Covariates,” *arXiv preprint arXiv:2206.11868*.
- Wieczorek, J. and Lei, J. (2022), “Model selection properties of forward selection and sequential cross-validation for high-dimensional regression,” *Canadian Journal of Statistics*, 50, 454–470.
- Wu, C. J. and Hamada, M. S. (2011), *Experiments: Planning, Analysis, and Optimization*, vol. 552, Hoboken, NJ: John Wiley & Sons.
- Wu, Y., Zheng, Z., Zhang, G., Zhang, Z., and Wang, C. (2022), “Non-stationary a/b tests: Optimal variance reduction, bias correction, and valid inference,” *Bias Correction, and Valid Inference (May 20, 2022)*.
- Yates, F. (1937), “The design and analysis of factorial experiments,” Tech. Rep. Technical Communication 35, Imperial Bureau of Soil Science, London, U. K.
- Yuan, M., Joseph, V. R., and Lin, Y. (2007), “An efficient variable selection approach for analyzing designed experiments,” *Technometrics*, 49, 430–439.
- Zhang, C. (2022), “Social construction of hate crimes in the U.S.: A factorial survey experiment,” *Theses and Dissertations–Sociology*, 49.
- Zhao, A. and Ding, P. (2021), “Regression-based causal inference with factorial experiments: estimands, model specifications and design-based properties,” *Biometrika*, 109, 799–815.
- (2023), “Covariate adjustment in multi-armed, possibly factorial experiments,” *Journal of the Royal Statistical Society, Series B (Statistical Methodology)*, in press.
- Zhao, P., Rocha, G., and Yu, B. (2009), “The composite absolute penalties family for grouped and hierarchical variable selection,” *The Annals of Statistics*, 37, 3468–3497.
- Zhao, P. and Yu, B. (2006), “On model selection consistency of Lasso,” *The Journal of Machine Learning Research*, 7, 2541–2563.
- Zhao, S., Witten, D., and Shojaie, A. (2021), “In defense of the indefensible: A very naive approach to high-dimensional inference,” *Statistical Science*, 36, 562–577.

# Supplementary material

Section A provides more discussions/extensions to the results introduced in the main paper. More concretely, Section A.1 presents detailed discussion of the use of weight least squares in factorial experiments. Section A.2 extends the inference results in Section 4 to a vector of causal effects.

Section B presents general results on consistency of forward factor screening. Theorem 1 is a corollary of the results in Section B.

Section C gives the technical proofs of the results in the main paper and the Appendix.

## A Additional results

This section provides more extensions to the results in the main paper. Section A.1 discusses the use of WLS in analyzing factorial experiments. Section A.2 extends the inference results under perfect screening (Section 4) to a vector of causal effects.

### A.1 Weighted least squares for estimating factorial effects

In this subsection, we briefly state and prove some useful facts about weighted least squares in estimating factorial effects. More discussions can be found in Zhao and Ding (2021). Denote the design matrix as  $X = (g_{1,\mathbb{M}}, \dots, g_{N,\mathbb{M}})^\top$ . Let  $W = \text{Diag}\{w_i\}$ . The problem (2.3) has closed-form solution:

$$\begin{aligned}\hat{\tau} &= (X^\top W X)^{-1} (X^\top W Y) \text{ (closed form solution of WLS)} \\ &= \{G(\cdot, \mathbb{M})^\top G(\cdot, \mathbb{M})\}^{-1} \{G(\cdot, \mathbb{M})^\top \hat{Y}\} \\ &\quad \text{(units under the same treatment arm share the same regressor)} \\ &= Q^{-1} G(\cdot, \mathbb{M})^\top \hat{Y}.\end{aligned}\tag{S1}$$

The closed form (S1) motivates the variance estimation:

$$\hat{V}_{\hat{\tau}} = Q^{-2} G(\cdot, \mathbb{M})^\top \hat{V}_{\hat{Y}} G(\cdot, \mathbb{M}).\tag{S2}$$

Alternatively, one can use the Eicker–Huber–White (EHW) variance estimation with the HC2 correction (Angrist and Pischke, 2009):

$$\hat{V}_{\text{EHW}} = (X^\top W X)^{-1} X^\top W \text{Diag} \left\{ \frac{\hat{\epsilon}_i^2}{1 - N_i^{-1}} \right\} W X (X^\top W X)^{-1}, \quad \hat{\epsilon}_i = Y_i - g_{i,\mathbb{M}}^\top \hat{\tau}.\tag{S3}$$

Again, because units under the same treatment arm share the same regressor,  $\widehat{V}_{\text{EHW}}$  simplifies to

$$\widehat{V}_{\text{EHW}} = Q^{-2} G(\cdot, \mathbb{M})^\top \widehat{V}'_{\widehat{\tau}} G(\cdot, \mathbb{M}), \quad (\text{S4})$$

where

$$\widehat{V}'_{\widehat{\tau}} = \text{Diag} \left\{ N(\mathbf{z})^{-1} \widehat{S}'(\mathbf{z}, \mathbf{z}) \right\}_{\mathbf{z} \in \mathcal{T}} \text{ with } \widehat{S}'(\mathbf{z}, \mathbf{z}) = \frac{1}{N(\mathbf{z}) - 1} \sum_{Z_i = \mathbf{z}} (Y_i - g_{i, \mathbb{M}}^\top \widehat{\tau})^2.$$

Following some algebra, we can show

$$\begin{aligned} \widehat{S}'(\mathbf{z}, \mathbf{z}) &= \frac{1}{N(\mathbf{z}) - 1} \sum_{Z_i = \mathbf{z}} (Y_i - \widehat{Y}(\mathbf{z}))^2 + \frac{N(\mathbf{z})}{N(\mathbf{z}) - 1} \{\widehat{Y}(\mathbf{z}) - G(\mathbf{z}, \mathbb{M})\widehat{\tau}\}^2 \\ &= \widehat{S}(\mathbf{z}, \mathbf{z}) + \frac{N(\mathbf{z})}{N(\mathbf{z}) - 1} \{\widehat{Y}(\mathbf{z}) - G(\mathbf{z}, \mathbb{M})\widehat{\tau}\}^2. \end{aligned}$$

Hence  $\widehat{S}'(\mathbf{z}, \mathbf{z}) \geq \widehat{S}(\mathbf{z}, \mathbf{z})$ . In general  $\widehat{Y}(\mathbf{z}) \neq G(\mathbf{z}, \mathbb{M})\widehat{\tau}$ , so the difference is not negligible. The following Lemma S1 formally summarizes the statistical property of  $\widehat{\tau}$  and its two variance estimators,  $\widehat{V}_{\widehat{\tau}}$  and  $\widehat{V}_{\text{EHW}}$ . The proof can be done by utilizing the moment facts from Section C.2 and C.3 of Shi and Ding (2022), which we omit here.

**Lemma S1.** *Assume Conditions 1 and 3. For the WLS in (2.3), we have*

1.  $\widehat{\tau} = Q^{-1} G(\cdot, \mathbb{M})^\top \widehat{Y}$  is unbiased for the true factorial effects  $\tau(\mathbb{M})$ ; i.e.,  $\mathbb{E}\{\widehat{\tau}\} = \tau(\mathbb{M})$ .
2. Both variance estimators are consistent and robust:  $N(\widehat{V}_{\widehat{\tau}} - V_{\widehat{\tau}, \text{lim}}) = o_{\mathbb{P}}(1)$ ,  $N(\widehat{V}_{\text{EHW}} - V_{\text{EHW}, \text{lim}}) = o_{\mathbb{P}}(1)$ , with  $V_{\widehat{\tau}, \text{lim}} \succcurlyeq V_{\widehat{\tau}}$  and  $V_{\text{EHW}} \succcurlyeq V_{\widehat{\tau}}$ , where

$$V_{\widehat{\tau}, \text{lim}} = Q^{-2} G(\cdot, \mathbb{M})^\top D_{\widehat{\tau}} G(\cdot, \mathbb{M}),$$

and

$$V_{\text{EHW}, \text{lim}} = Q^{-2} G(\cdot, \mathbb{M})^\top \text{Diag} \left\{ \frac{1 - N^{-1}}{N(\mathbf{z}) - 1} S(\mathbf{z}, \mathbf{z}) + \frac{1}{N(\mathbf{z}) - 1} \{\bar{Y}(\mathbf{z}) - G(\mathbf{z}, \mathbb{M})\tau(\mathbb{M})\}^2 \right\} G(\cdot, \mathbb{M}).$$

3. EHW variance estimator is more conservative than the direct variance estimator:  $\widehat{V}_{\text{EHW}} \succcurlyeq \widehat{V}_{\widehat{\tau}}$ .

It is worthy of mentioning that in the fixed  $Q$  setting, if we assume that the factorial effects that are not included in  $\mathbb{M}$  are all zero, Lemma S1 implies EHW variance estimator (S3) or (S4) has the same asymptotic statistical property as the direct variance estimator (S2), which agrees with the conclusion of Zhao and Ding (2021).

## A.2 Extension of post-screening inference to vector parameters

In this subsection we present an extension of Theorem 2 to a vector of causal parameters:

$$\Gamma = (\gamma_1, \dots, \gamma_L)^\top, \quad \text{where } \gamma_l = \mathbf{f}_l^\top \bar{Y}.$$

For convenience we can stack  $\mathbf{f}_1, \dots, \mathbf{f}_L$  into a weighting matrix  $F = (\mathbf{f}_1, \dots, \mathbf{f}_L)$  and write

$$\Gamma = F^\top \bar{Y}.$$

We will focus on linear projections of  $\Gamma$ , defined as  $\gamma_b = b^\top \Gamma$  for a given  $b \in \mathbb{R}^L$ . Naturally, we can apply forward screening and construct RLS-based estimators for  $\Gamma$ :

$$\hat{\Gamma}_R = (\hat{\gamma}_{1,R}, \dots, \hat{\gamma}_{L,R})^\top, \quad \hat{V}_{\hat{\Gamma},R} = F[\hat{\mathbb{M}}]^\top \hat{V}_{\hat{Y}} F[\hat{\mathbb{M}}], \quad (\text{S5})$$

where

$$F[\hat{\mathbb{M}}] = Q^{-1} G(\cdot, \hat{\mathbb{M}}) G(\cdot, \hat{\mathbb{M}})^\top F.$$

For  $\gamma_b$ , an estimator based on (S5) is

$$\hat{\gamma}_{b,R} = b^\top \hat{\Gamma}_R, \quad \hat{v}_{b,R}^2 = b^\top \hat{V}_{\hat{\Gamma},R} b.$$

For standard factorial effects, we can use WLS to obtain the robust covariance matrix (Section A.1).

For one single  $b$ , we can actually apply Theorem 2 with

$$\mathbf{f}_b = Fb = \sum_{l=1}^L b_l \mathbf{f}_l.$$

Define  $\mathbf{f}_b^* = F[\mathbb{M}^*]b$ . We then get the following theorem:

**Theorem S1** (Statistical properties linear projections of  $\Gamma$ ). *Assume Conditions 1-4. Let  $N \rightarrow \infty$ . Then*

$$\frac{\hat{\gamma}_{b,R} - \gamma}{v_{b,R}} \rightsquigarrow \mathcal{N}(0, 1)$$

where  $v_{b,R}^2 = \mathbf{f}_b^{*\top} \hat{V}_{\hat{Y}} \mathbf{f}_b^*$ . Further assume  $\|\mathbf{f}_b^*\|_\infty = O(Q^{-1})$ . The variance estimator  $\hat{v}_{b,R}^2$  is conservative in the sense

$$N(\hat{v}_{b,R}^2 - v_{b,R,\text{lim}}^2) \xrightarrow{\mathbb{P}} 0, \quad v_{b,R,\text{lim}}^2 \geq v_{b,R}^2,$$

where  $v_{b,R,\text{lim}}^2 = \mathbf{f}_b^{*\top} D_{\hat{Y}} \mathbf{f}_b^*$  is the limiting value of  $\hat{v}_{b,R}^2$ .

The proof of Theorem S1 is similar to that of Theorem 2, which is mainly based on Lemma S5 and thus omitted here. Moreover, for a fixed integer  $L$ , Theorem S1 implies joint normality of  $\widehat{\Gamma}_R$ , a result due to the Cramér-Wold theorem. We summarize the result as the following corollary and omit the proof:

**Corollary S1.** *Assume a fixed  $L$ . Assume Conditions 1-4. We have*

$$V_{\widehat{\Gamma},R}^{-1/2}(\widehat{\Gamma}_R - \Gamma) \rightsquigarrow \mathcal{N}(0, I_L),$$

where  $V_{\widehat{\Gamma},R} = F[\mathbb{M}^\star]^\top V_{\widehat{Y}} F[\mathbb{M}^\star]$ . Further assume  $\max_{\|b\|_2=1} \|\mathbf{f}_b^\star\|_\infty = O(Q^{-1})$ . The variance estimator  $\widehat{v}_{b,R}^2$  is conservative in the sense that

$$N(\widehat{V}_{\widehat{\Gamma},R} - V_{\widehat{\Gamma},R,\text{lim}}) \xrightarrow{\mathbb{P}} 0, \quad V_{\widehat{\Gamma},R,\text{lim}} \succcurlyeq V_{\widehat{\Gamma},R},$$

where  $V_{\widehat{\Gamma},R,\text{lim}} = F[\mathbb{M}^\star]^\top D_{\widehat{Y}} F[\mathbb{M}^\star]$  is the limiting value of  $\widehat{V}_{\widehat{\Gamma},R}$ .

## B General results on consistency of forward screening

In this section we provide some theoretical insights into the forward factor screening algorithm (Algorithm 1). The discussion in this section starts from a more broad discussion where we allow the S-step to be general procedures that satisfy certain conditions. We will show Bonferroni corrected marginal t-test is a special case of these procedures.

We start with some regularization conditions to characterize a “good” layer-wise S-step, and ensure the P-step is compatible with the structure of the true factorial effects. In light of this, we use  $\mathbb{M}_{d,+}^\star$  to denote the pruned set of effects on the  $d$ -th layer based on the true model  $\mathbb{M}_{d-1}^\star$  on the previous layer; that is,

$$\mathbb{M}_{d,+}^\star = \mathbb{H}(\mathbb{M}_{d-1}^\star).$$

These discussions motivate the following assumption on the layer-wise selection procedure  $\widehat{\mathbb{S}}(\cdot)$ :

**Assumption 1** (Validity and consistency of the selection operator). *We denote*

$$\widetilde{\mathbb{M}}_d = \widehat{\mathbb{S}}(\mathbb{M}_{d,+}^\star; \{Y_i, Z_i\}_{i=1}^N),$$

where  $\mathbb{M}_{d,+}^\star = \mathbb{H}(\mathbb{M}_{d-1}^\star)$  is defined as above. Let  $\{\alpha_d\}_{d=1}^D$  be a sequence of significance levels in  $(0, 1)$ .

We assume that the following validity and consistency property hold for  $\mathbb{S}_N(\cdot)$ :

$$\text{Validity: } \limsup_{N \rightarrow \infty} \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d \cap \mathbb{M}_d^{\star c} \neq \emptyset \right\} \leq \alpha_d,$$

$$\text{Consistency: } \limsup_{N \rightarrow \infty} D \sum_{d=1}^D \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^\star \neq \emptyset \right\} = 0.$$

This assumption can be verified for many screening procedures. In Theorem 1 we will show it holds for the layer-wise Bonferroni corrected marginal testing procedure in Algorithm 1. Moreover, in the high dimensional super population study, a combination of data splitting, adaptation of  $\ell_1$  regularization and marginal t tests can also fulfill such a requirement (Wasserman and Roeder, 2009).

Besides, we assume the  $H(\cdot)$  operator respects the structure of the nonzero factorial effects:

**Assumption 2** (H-heredity). *For  $d = 1, \dots, D - 1$ , it holds*

$$\mathbb{M}_{d+1}^* \subset \mathbb{P}(\mathbb{M}_d^*).$$

One special case of  $H(\cdot)$  operator satisfying Assumption 2 is naively adding all the higher order interactions regardless of the lower-order screening results. Besides, if we have evidence that the effects have particular hierarchical structure, applying the heredity principles can improve screening accuracy as well as interpretability of the screening results.

**Theorem S2** (Screening consistency). *Assume Assumption 1 and 2. Then the forward screening procedure (3.5) has the following properties:*

(i) Type I error control. *Forward screening controls the Type I error rate, in the sense that*

$$\limsup_{N \rightarrow \infty} \mathbb{P} \left( \widehat{\mathbb{M}}_d \cap \mathbb{M}_d^{*c} \neq \emptyset \text{ for some } d \in [D] \right) \leq \alpha = \sum_{d=1}^D \alpha_d.$$

(ii) Screening consistency. *Further assume  $\alpha = \alpha_N \rightarrow 0$ . The forward procedure consistently selects all the nonzero effects up to  $D$  levels with probability tending to 1:*

$$\limsup_{N \rightarrow \infty} \mathbb{P} \left( \widehat{\mathbb{M}}_d = \mathbb{M}_d^* \text{ for all } d \in [D] \right) = 1.$$

Theorem S2 consists of two parts. First, one can control the type I error rate, which is defined as the probability of over-selects at least one zero effect. The definition is introduced and elaborated detailedly in Wasserman and Roeder (2009) for model selection. Second, if the tuning parameter  $\alpha = \sum_{d=1}^D \alpha_d$  vanish asymptotically, one can actually achieve perfect screening up to  $D$  levels of effects. To apply Theorem S2 to specific procedures, the key step is to verify Assumption 1 and justify Assumption 2, which we will do for Bonferroni corrected marginal t tests as an example in the next section.

Moreover, the scaling of  $\alpha_N$  plays an important role in theoretical discussion. To achieve perfect selection, we hope  $\alpha_N$  decays as fast as possible; ideally if  $\alpha_N$  equals zero then we do not commit

any type I error (or equivalently, we will never select redundant effects). However, for many data-dependent selection procedure  $\alpha$  can only decay at certain rates, because a fast decaying  $\alpha$  means higher possibility of rejection, thus can lead to strict under-selection. Therefore, in the tuning process,  $\alpha_d$  should be scaled properly if one wants to pursue perfect selection. Nevertheless, even if the tuning is hard and perfect model selection can not be achieved, we still have many strategies to exploit the advantage of the forward screening procedure. We will have more discussions in later sections.

Lastly, as we have commented earlier, in practice people have many alternative methods for the S-step. They are attractive in factorial experiments because many lead to simple form solutions due to the orthogonality of factorial designs. For example, Lasso is a commonly adopted strategy for variable selection in linear models (Zhao and Yu, 2006). It solves the following penalized WLS problem in factorial settings:

$$\hat{\mathbb{M}}_L = \{\mathcal{K} : \hat{\tau}_{L,\mathcal{K}} \neq 0\}, \quad \hat{\tau}_{L,\mathcal{K}} = \min_{\tau' \in \mathbb{R}^H} \frac{1}{2} \sum_{z \in \mathcal{T}} w_i (Y_i - g_i^\top \tau')^2 + \lambda_L \|\tau'\|_1.$$

Due to the orthogonality of  $G$ , the resulting  $\hat{\mathbb{M}}$  has a closed-form solution (Hastie et al., 2009):

$$\hat{\mathbb{M}}_L = \{\mathcal{K} : |\hat{\tau}_{\mathcal{K}}| \geq \lambda_L\}.$$

Other methods, such as AIC/BIC (Bai et al., 2022), sure independence screening (Fan and Lv, 2008), etc., are also applicable. With more delicate assumptions and tuning parameter choices, these methods can also be justified theoretically for screening consistency and post-screening inference. We omit the details.

## C Technical proofs

In this section we present the technical proofs for the results across the whole paper. Section C.1 presents some preliminary probabilistic results that are useful in randomized experiments which are mainly attributed to Shi and Ding (2022). The main proof starts from Section C.2.

## C.1 Preliminaries: some important probabilistic results in randomized experiments

In this subsection we present some preliminary probability results that are crucial for our theoretical discussion. Consider an estimator of the form

$$\hat{\gamma} = Q^{-1} \sum_{\mathbf{z} \in \mathcal{T}} w(\mathbf{z}) \hat{Y}(\mathbf{z}),$$

with variance estimator

$$\hat{v}^2 = Q^{-2} \sum_{\mathbf{z} \in \mathcal{T}} w(\mathbf{z})^2 \hat{S}(\mathbf{z}, \mathbf{z}).$$

Li and Ding (2017) showed that

$$\mathbb{E}\{\hat{Y}\} = \bar{Y}, \quad V_{\hat{Y}} = \text{Var}\{\hat{Y}\} = D_{\hat{Y}} - N^{-1}S. \quad (\text{S6})$$

Then (S6) further leads to the following facts:

$$\begin{aligned} \mathbb{E}\{\hat{\gamma}\} &= \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z}) \bar{Y}(\mathbf{z}) = \gamma, \\ \text{Var}\{\hat{\gamma}\} &= \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z})^2 N(\mathbf{z})^{-1} S(\mathbf{z}, \mathbf{z}) - N^{-1} \mathbf{f}^\top S \mathbf{f}, \\ \mathbb{E}\{\hat{v}^2\} &= \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z})^2 N(\mathbf{z})^{-1} S(\mathbf{z}, \mathbf{z}). \end{aligned} \quad (\text{S7})$$

We have the following variance estimation results and Berry–Esseen bounds:

**Lemma S2** (Variance concentration and Berry–Esseen bounds in finite population). *Define  $\gamma = \mathbb{E}\{\hat{\gamma}\}$ ,  $v^2 = \text{Var}(\hat{\gamma})$  and  $v_{\text{lim}}^2 = \mathbb{E}\{\hat{v}^2\}$ . Suppose the following conditions hold:*

- *Nondegenerate variance. There exists a  $\sigma_w > 0$ , such that*

$$Q^{-2} \sum_{\mathbf{z}=1}^Q w(\mathbf{z})^2 N_{\mathbf{z}}^{-1} S(\mathbf{z}, \mathbf{z}) \leq \sigma_w^2 v^2. \quad (\text{S8})$$

- *Bounded fourth moments. There exists a  $\delta > 0$  such that*

$$\max_{\mathbf{z} \in [Q]} \frac{1}{N} \sum_{i=1}^N \{Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})\}^4 \leq \Delta^4. \quad (\text{S9})$$

*Then we have the following conclusions:*



1. The variance estimator is conservative for the true variance:  $v_{\text{lim}}^2 \geq v^2$ . Besides, the following tail bound holds:

$$\mathbb{P} \{N|\hat{v}^2 - v_{\text{lim}}^2| > t\} \leq \frac{C\bar{c}^3 \underline{c}^{-4} \|w\|_\infty^2 \Delta^4}{QN_0} \cdot \frac{1}{t^2}.$$

2. We have a Berry–Esseen bound with the true variance:

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{\gamma} - \gamma}{v} \leq t \right\} - \Phi(t) \right| \leq 2C\sigma_w \frac{\underline{c}^{-1} \|w\|_\infty \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\|w\|_2 \sqrt{\bar{c}^{-1} \min_{z \in [Q]} S(z, z)} \cdot \sqrt{N_0}}.$$

3. We have a Berry–Esseen bound with the estimated variance: for any  $\epsilon_N \in (0, 1/2]$ ,

$$\begin{aligned} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{\gamma} - \gamma}{\hat{v}} \leq t \right\} - \Phi \left( \frac{v_{\text{lim}}}{v} t \right) \right| &\leq \epsilon_N + \frac{C\bar{c}^3 \underline{c}^{-4} \|w\|_\infty^2 \Delta^4}{QN_0} \cdot \frac{1}{(Nv^2\epsilon_N)^2} \\ &\quad + 2C\sigma_w \frac{\underline{c}^{-1} \|w\|_\infty \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\|w\|_2 \sqrt{\bar{c}^{-1} \min_{z \in [Q]} S(z, z)} \cdot \sqrt{N_0}}. \end{aligned}$$

*Proof of Lemma S2.* 1. See Lemma S13 of Shi and Ding (2022).

2. See Theorem 1 of Shi and Ding (2022).

3. First we show a useful result: for  $|a| \leq 1/2$  and any  $b \in \mathbb{R}$ ,

$$\sup_{t \in \mathbb{R}} |\Phi\{(1+a)t+b\} - \Phi\{t\}| \leq |a| + |b|. \quad (\text{S10})$$

(S10) is particularly useful for small choices of  $a$  and  $b$ . Intuitively, it evaluates the change of  $\Phi$  under a small affine perturbation of  $t$ .

The proof of (S10) is based on a simple step of the mean value theorem: for any  $t \in \mathbb{R}$ ,

$$\begin{aligned} &|\Phi\{(1+a)t+b\} - \Phi\{t\}| \\ &= |\phi(\xi_{t,(1+a)t}) \cdot (at+b)| \\ &= |\phi(\xi_{t,(1+a)t}) \cdot at| + |\phi(\xi_{t,(1+a)t}) \cdot b| \\ &= |a| \cdot |\phi(\xi_{t,(1+a)t}) \cdot t| \cdot \mathbf{1}\{|t| \leq 1\} + |a| \cdot |\phi(\xi_{t,(1+a)t}) \cdot t| \cdot \mathbf{1}\{|t| > 1\} + |\phi(\xi_{t,(1+a)t}) \cdot b| \\ &\leq \frac{1}{\sqrt{2\pi}} |a| \cdot \mathbf{1}\{|t| \leq 1\} + \frac{1}{\sqrt{2\pi}} |a||t| \cdot \exp(-t^2/8) \cdot \mathbf{1}\{|t| > 1\} + \frac{1}{\sqrt{2\pi}} |b| \\ &\leq |a| + |b|. \end{aligned}$$

We consider  $t \geq 0$  because  $t < 0$  can be handled similarly. For any  $\epsilon_N > 0$ , We have

$$\begin{aligned} \mathbb{P} \left\{ \frac{\hat{\gamma} - \gamma}{\hat{v}} \leq t \right\} &= \mathbb{P} \left\{ \frac{\hat{\gamma} - \gamma}{v} \leq \frac{\hat{v}}{v} t \right\} \\ &= \mathbb{P} \left\{ \frac{\hat{\gamma} - \gamma}{v} \leq \frac{\hat{v}}{v} t, \left| \frac{\hat{v} - v_{\text{lim}}}{v} \right| \leq \epsilon_N \right\} + \mathbb{P} \left\{ \frac{\hat{\gamma} - \gamma}{v} \leq \frac{\hat{v}}{v} t, \left| \frac{\hat{v} - v_{\text{lim}}}{v} \right| > \epsilon_N \right\}. \end{aligned}$$

Then we can show that

$$\begin{aligned}\mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{\hat{v}} \leq t\right\} &\leq \mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{v} \leq \frac{\hat{v}}{v}t, \left|\frac{\hat{v}-v_{\text{lim}}}{v}\right| \leq \epsilon_N\right\} + \mathbb{P}\left\{\left|\frac{\hat{v}-v_{\text{lim}}}{v}\right| > \epsilon_N\right\} \\ &\leq \mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{v} \leq \left(\frac{v}{v} + \epsilon_N\right)t\right\} + \mathbb{P}\left\{\left|\frac{\hat{v}-v_{\text{lim}}}{v}\right| > \epsilon_N\right\}.\end{aligned}$$

For the first term, we have

$$\begin{aligned}\sup_{t \geq 0} \left| \mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{v} \leq \left(\frac{v_{\text{lim}}}{v} + \epsilon_N\right)t\right\} - \Phi\left\{\left(\frac{v_{\text{lim}}}{v} + \epsilon_N\right)t\right\} \right| \\ \leq 2C\sigma_w \frac{\underline{c}^{-1}\|w\|_{\infty} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\|w\|_2 \sqrt{\bar{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})} \cdot \sqrt{N_0}}.\end{aligned}$$

For the second term, using the variance estimation results in Part 1 we have

$$\begin{aligned}\mathbb{P}\left\{\left|\frac{\hat{v}-v_{\text{lim}}}{v}\right| \geq \epsilon_N\right\} &\leq \mathbb{P}\left\{\left|\frac{\hat{v}-v_{\text{lim}}}{v}\right| \cdot \left|\frac{\hat{v}+v_{\text{lim}}}{v}\right| \geq \epsilon_N\right\} \quad (\text{because } v_{\text{lim}} \text{ is conservative}) \\ &= \mathbb{P}\left\{\left|\frac{N\hat{v}^2 - Nv_{\text{lim}}^2}{Nv^2}\right| \geq \epsilon_N\right\} \\ &\leq \frac{C\bar{c}^3 \underline{c}^{-4} \|w\|_{\infty}^2 \Delta^4}{QN_0} \cdot \frac{1}{(Nv^2\epsilon_N)^2}.\end{aligned}$$

Besides, by (S10), when  $\epsilon_N \leq 1/2$ , we also have

$$\sup_{t \in \mathbb{R}} \left| \Phi\left\{\left(\frac{v_{\text{lim}}}{v} + \epsilon_N\right)t\right\} - \Phi\left(\frac{v_{\text{lim}}}{v}t\right) \right| \leq \frac{v\epsilon_N}{v_{\text{lim}}} \leq \epsilon_N.$$

Aggregating all the parts above, we can show that for any  $t \geq 0$ ,

$$\begin{aligned}\mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{\hat{v}} \leq t\right\} &\leq \Phi\left(\frac{v_{\text{lim}}}{v}t\right) + \epsilon_N + \frac{C\bar{c}^3 \underline{c}^{-4} \|w\|_{\infty}^2 \Delta^4}{QN_0} \cdot \frac{1}{(Nv^2\epsilon_N)^2} \\ &\quad + 2C\sigma_w \frac{\underline{c}^{-1}\|w\|_{\infty} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\|w\|_2 \sqrt{\bar{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})} \cdot \sqrt{N_0}}.\end{aligned}$$

On the other hand, we can show that

$$\begin{aligned}\mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{\hat{v}} \leq t\right\} &\geq \mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{v} \leq \frac{\hat{v}}{v}t, \left|\frac{\hat{v}-v_{\text{lim}}}{v}\right| \leq \epsilon_N\right\} \\ &\geq \mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{v} \leq \left(\frac{v_{\text{lim}}}{v} - \epsilon_N\right)t\right\} - \mathbb{P}\left\{\left|\frac{\hat{v}-v_{\text{lim}}}{v}\right| \geq \epsilon_N\right\}.\end{aligned} \tag{S11}$$

By (S10), when  $\epsilon_N \leq 1/2$ , we also have

$$\sup_{t \in \mathbb{R}} \left| \Phi\left\{\left(\frac{v_{\text{lim}}}{v} - \epsilon_N\right)t\right\} - \Phi\left(\frac{v_{\text{lim}}}{v}t\right) \right| \leq \epsilon_N.$$

So we can derive a lower bound analogous to (S11). Note that the results can be analogously generalized to  $t \leq 0$ . Putting pieces together, we can show that for any  $t \geq 0$ ,

$$\begin{aligned} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{\gamma} - \gamma}{\hat{v}} \leq t \right\} - \Phi \left( \frac{v_{\text{lim}}}{v} t \right) \right| &\leq \epsilon_N + \frac{C \bar{c}^3 \underline{c}^{-4} \|w\|_\infty^2 \Delta^4}{Q N_0} \cdot \frac{1}{(N v^2 \epsilon_N)^2} \\ &\quad + 2C \sigma_w \frac{\underline{c}^{-1} \|w\|_\infty \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\|w\|_2 \sqrt{\bar{c}^{-1} \min_{z \in [Q]} S(z, z)} \cdot \sqrt{N_0}}. \end{aligned}$$

□

The following corollary shows a Berry–Esseen bound for the studentized statistic in the special case where  $w = (w(z))_{z \in [Q]}$  is a contrast vector for factorial effects. That is,  $w = g_{\mathcal{K}}$  for some  $\mathcal{K} \in \mathbb{K}$ .

**Corollary S2.** *Assume Condition (S8) and (S9) hold. Let  $w = g_{\mathcal{K}}$  for some  $\mathcal{K} \in \mathbb{K}$ . Then we have a Berry–Esseen bound with the estimated variance:*

$$\begin{aligned} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{\tau}_{\mathcal{K}} - \tau_{\mathcal{K}}}{\hat{v}} \leq t \right\} - \Phi \left( \frac{v_{\text{lim}}}{v} t \right) \right| &\leq 2 \left( \frac{C \sigma_w^4 \bar{c}^5 \underline{c}^{-6} \Delta^4}{\{\min_{z \in \mathcal{T}} S(z, z)\}^2} \right)^{1/3} \cdot \frac{1}{(Q N_0)^{1/3}} \\ &\quad + 2C \sigma_w \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\sqrt{\bar{c}^{-1} \min_{z \in [Q]} S(z, z)}} \cdot \frac{1}{(Q N_0)^{1/2}}. \end{aligned}$$

*Proof of Corollary S2. Lower bound for  $N v^2$ .* Note that  $\|w\|_2^2 = Q$  and  $\|w\|_\infty = 1$ . Using Condition (S8), we have

$$\begin{aligned} N v^2 &\geq N \sigma_w^{-2} Q^{-2} \sum_{z=1}^Q w(z)^2 N_z^{-1} S(z, z) \\ &\geq (\underline{c} Q N_0) \cdot \sigma_w^{-2} \bar{c}^{-1} Q^{-1} N_0^{-1} \min_{z \in \mathcal{T}} S(z, z) \cdot (Q^{-1} \|w\|_2^2) \\ &= \sigma_w^{-2} \underline{c} \bar{c}^{-1} \min_{z \in \mathcal{T}} S(z, z). \end{aligned}$$

Therefore, the Berry–Esseen bound becomes

$$\begin{aligned} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{\tau}_{\mathcal{K}} - \tau_{\mathcal{K}}}{\hat{v}} \leq t \right\} - \Phi \left( \frac{v_{\text{lim}}}{v} t \right) \right| &\leq \epsilon_N + \frac{C \sigma_w^4 \bar{c}^5 \underline{c}^{-6} \Delta^4}{(Q N_0) \{\min_{z \in \mathcal{T}} S(z, z)\}^2} \cdot \frac{1}{\epsilon_N^2} \\ &\quad + 2C \sigma_w \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\sqrt{\bar{c}^{-1} \min_{z \in [Q]} S(z, z)} \cdot \sqrt{Q N_0}}. \end{aligned}$$

**Optimize the summation of the first and second term.** By taking derivative over  $\epsilon_N$  on the upper bound and solving for the zero point, we know that when

$$\epsilon_N = \left( \frac{2C \sigma_w^4 \bar{c}^5 \underline{c}^{-6} \Delta^4}{(Q N_0) \{\min_{z \in \mathcal{T}} S(z, z)\}^2} \right)^{1/3},$$

the upper bound is minimized and

$$\begin{aligned} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\tau}_{\mathcal{K}} - \tau_{\mathcal{K}}}{\widehat{v}} \leq t \right\} - \Phi \left( \frac{v_{\text{lim}}}{v} t \right) \right| &\leq 2 \left( \frac{C \sigma_w^4 \bar{c}^5 \underline{c}^{-6} \Delta^4}{\{\min_{\mathbf{z} \in \mathcal{T}} S(\mathbf{z}, \mathbf{z})\}^2} \right)^{1/3} \cdot \frac{1}{(QN_0)^{1/3}} \\ &\quad + 2C \sigma_w \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\sqrt{\bar{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})}} \cdot \frac{1}{(QN_0)^{1/2}}. \end{aligned}$$

□

Additionally, we have a Berry–Esseen bounds after screening the effects:

**Lemma S3** (Berry Esseen bound with screening). *Assume there exists  $\sigma_w > 0$  such that*

$$\sum_{\mathbf{z}=1}^Q \mathbf{f}[\mathbb{M}](\mathbf{z})^2 N_{\mathbf{z}}^{-1} S(\mathbf{z}, \mathbf{z}) \leq \sigma_w^2 v^2(\mathbb{M}). \quad (\text{S12})$$

Then

$$\begin{aligned} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}[\widehat{\mathbb{M}}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t \right\} - \Phi(t) \right| \\ \leq 2\mathbb{P} \left\{ \widehat{\mathbb{M}} \neq \mathbb{M} \right\} + 2C \sigma_w \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\sqrt{\bar{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})} \cdot \sqrt{N_0}} \cdot \frac{\|\mathbf{f}[\mathbb{M}]\|_{\infty}}{\|\mathbf{f}[\mathbb{M}]\|_2}. \end{aligned}$$

*Proof of Lemma S3.* With the selected working model we have

$$\begin{aligned} &\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}[\widehat{\mathbb{M}}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t \right\} - \Phi(t) \right| \\ &= \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}[\widehat{\mathbb{M}}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t, \widehat{\mathbb{M}} = \mathbb{M} \right\} - \Phi(t) + \mathbb{P} \left\{ \frac{\widehat{\gamma}[\widehat{\mathbb{M}}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t, \widehat{\mathbb{M}} \neq \mathbb{M} \right\} \right| \\ &\leq \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}[\widehat{\mathbb{M}}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t, \widehat{\mathbb{M}} = \mathbb{M} \right\} - \Phi(t) \right| + \mathbb{P} \left\{ \frac{\widehat{\gamma}[\widehat{\mathbb{M}}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t, \widehat{\mathbb{M}} \neq \mathbb{M} \right\} \\ &= \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}[\mathbb{M}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t, \widehat{\mathbb{M}} = \mathbb{M} \right\} - \Phi(t) \right| + \mathbb{P} \left\{ \frac{\widehat{\gamma}[\widehat{\mathbb{M}}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t, \widehat{\mathbb{M}} \neq \mathbb{M} \right\} \\ &\leq \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}[\mathbb{M}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t \right\} - \Phi(t) \right| + 2\mathbb{P} \left\{ \widehat{\mathbb{M}} \neq \mathbb{M} \right\}. \end{aligned}$$

Now we have

$$\begin{aligned} \widehat{\gamma}(\mathbb{M}) &= \mathbf{f}^{\top} G(\cdot, \mathbb{M}) \widehat{\tau}(\mathbb{M}) \\ &= \mathbf{f}^{\top} G(\cdot, \mathbb{M}) G(\cdot, \mathbb{M})^{\top} \widehat{Y} \\ &= \mathbf{f}[\mathbb{M}]^{\top} \widehat{Y}. \end{aligned}$$

By Theorem 1 of Shi and Ding (2022), we have a Berry–Esseen bound with the true variance:

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{\gamma}(\mathbb{M}) - \gamma[\mathbb{M}]}{v} \leq t \right\} - \Phi(t) \right| \leq 2C\sigma_w \frac{\|\mathbf{f}[\mathbb{M}]\|_\infty \bar{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\|\mathbf{f}[\mathbb{M}]\|_2 \sqrt{\bar{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z}) \cdot \sqrt{N_0}}}.$$

□

A crucial quantity that appeared in Lemma S3 is the ratio of norms:

$$\frac{\|\mathbf{f}[\mathbb{M}]\|_\infty}{\|\mathbf{f}[\mathbb{M}]\|_2}. \quad (\text{S13})$$

The following Lemma S4 provides an explicit bound on (S13) which reveals how the ratio is controlled with respect to the size of the working model.

**Lemma S4.** *For  $\mathbf{f}[\mathbb{M}] \neq 0$ , we have*

$$\frac{\|\mathbf{f}[\mathbb{M}]\|_\infty}{\|\mathbf{f}[\mathbb{M}]\|_2} \leq \left( \frac{|\mathbb{M}|}{Q} \right)^{1/2}. \quad (\text{S14})$$

*Proof of Lemma S4.* Because the LHS of (S14) is a ratio, based on the definition of  $\mathbf{f}^\star$  (4.8) we can assume  $\|\mathbf{f}\|_2 = 1$  without loss of generality. Due to the orthogonality of  $G$ , we can use the columns of  $G$  as bases and express  $\mathbf{f}$  as

$$\mathbf{f} = \frac{1}{\sqrt{Q}} G(\cdot, \mathbb{M}) b_1 + \frac{1}{\sqrt{Q}} G(\cdot, \mathbb{M}^c) b_2,$$

where  $b_1 \in \mathbb{R}^{|\mathbb{M}|}$  and  $b_2 \in \mathbb{R}^{|\mathbb{M}^c|}$  and  $\|(b_1^\top, b_2^\top)^\top\|_2 = 1$ . Then

$$\mathbf{f}[\mathbb{M}] = Q^{-1} G(\cdot, \mathbb{M}) G(\cdot, \mathbb{M})^\top \mathbf{f} = \frac{1}{\sqrt{Q}} G(\cdot, \mathbb{M}) b_1.$$

Hence

$$\|\mathbf{f}[\mathbb{M}]\|_\infty \leq \frac{1}{\sqrt{Q}} \|b_1\|_1, \quad \|\mathbf{f}[\mathbb{M}]\|_2 = \|b_1\|_2, \quad \frac{\|\mathbf{f}[\mathbb{M}]\|_\infty}{\|\mathbf{f}[\mathbb{M}]\|_2} \leq \frac{1}{\sqrt{Q}} \cdot \frac{\|b_1\|_1}{\|b_1\|_2} \leq \left( \frac{|\mathbf{f}[\mathbb{M}]|}{Q} \right)^{1/2}.$$

□

## C.2 Proof of Theorem S2

*Proof of Theorem S2.* According to the orthogonality of designs, the signs for all terms in the studied unsaturated population regressions are consistent with those of saturated regressions, which saves the effort of differentiating true models for partial and full regression. We introduce several key events that will play a crucial role in the proof: for  $D_0 \in [D]$ , define

$$\text{Under-selection: } \mathcal{E}_u(D_0) = \{\widehat{\mathbb{M}}_d \subset \mathbb{M}_d^\star, d \in [D_0]\},$$

$$\text{Strict under-selection: } \mathcal{E}_{\text{su}}(D_0) = \{\widehat{\mathbb{M}}_d \subset \mathbb{M}_d^\star, d \in [D_0]; \text{ there exists } d \in [D_0], \widehat{\mathbb{M}}_d \subsetneq \mathbb{M}_d^\star\}.$$

**High level idea of the proof.** To prove screening consistency, we will prove two facts:

$$\mathbb{P}\{\mathcal{E}_U(D) \text{ holds}\} \rightarrow 1, \quad \mathbb{P}\{\mathcal{E}_{\text{SU}}(D) \text{ holds}\} \rightarrow 0.$$

Combining these two results together, we can conclude asymptotic screening consistency.

We start from the strict under-selection probability.

**Step 1: Prove that asymptotically, there is no strict under-selection.**

By definition,

$$\mathbb{P}\{\mathcal{E}_{\text{SU}}(1) \text{ holds}\} = \mathbb{P}\left\{\tilde{\mathbb{M}}_1 \subsetneq \mathbb{M}_1^*\right\} \leq \mathbb{P}\left\{\tilde{\mathbb{M}}_1^c \cap \mathbb{M}_1^* \neq \emptyset\right\}.$$

We now derive a recursive bound for  $\mathbb{P}\{\mathcal{E}_{\text{SU}}(D_0 + 1) \text{ holds}\}$  where  $1 \leq D_0 \leq D - 1$ . We have decomposition

$$\begin{aligned} \mathcal{E}_{\text{SU}}(D_0 + 1) &= \left\{\widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*, d \leq D_0 + 1\right\} - \left\{\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0 + 1\right\} \\ &= \mathcal{E}_{\text{SU},1}(D_0 + 1) \cup \mathcal{E}_{\text{SU},2}(D_0 + 1), \end{aligned}$$

where

$$\begin{aligned} \mathcal{E}_{\text{SU},1}(D_0 + 1) &= \left\{\widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*, d \leq D_0 + 1\right\} - \left\{\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subset \mathbb{M}_{D_0+1}^*\right\}, \\ \mathcal{E}_{\text{SU},2}(D_0 + 1) &= \left\{\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subset \mathbb{M}_{D_0+1}^*\right\} - \left\{\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0 + 1\right\}. \end{aligned}$$

For  $\mathcal{E}_{\text{SU},1}(D_0 + 1)$ , we have

$$\begin{aligned} \mathbb{P}\{\mathcal{E}_{\text{SU},1}(D_0 + 1) \text{ holds}\} &= \mathbb{P}\left(\left\{\widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*, d \leq D_0 + 1\right\} - \left\{\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subset \mathbb{M}_{D_0+1}^*\right\}\right) \\ &\leq \mathbb{P}\left(\forall d \in [D_0 + 1], \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*; \exists d \in [D_0], \widehat{\mathbb{M}}_d \subsetneq \mathbb{M}_d^*\right) \\ &\leq \mathbb{P}\left(\forall d \in [D_0], \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*; \exists d \in [D_0], \widehat{\mathbb{M}}_d \subsetneq \mathbb{M}_d^*\right) \\ &= \mathbb{P}\{\mathcal{E}_{\text{SU}}(D_0) \text{ holds}\}. \end{aligned} \tag{S15}$$

For  $\mathcal{E}_{\text{SU},2}(D_0 + 1)$ , we notice that  $\widehat{\mathbb{M}}_{D_0+1}$  is generated based on  $\widehat{\mathbb{M}}_{D_0}$  and the set of estimates over the prescreened effect set  $\widehat{\mathbb{M}}_{D_0+1,+}$ . Under Assumption 2, on the event  $\widehat{\mathbb{M}}_d = \mathbb{M}_d^*$  we have

$$\widehat{\mathbb{M}}_{d+1} = \tilde{\mathbb{M}}_{d+1}.$$

Hence we can compute

$$\begin{aligned} \mathbb{P}\{\mathcal{E}_{\text{SU},2}(D_0 + 1) \text{ holds}\} &= \mathbb{P}\left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subsetneq \mathbb{M}_{D_0+1}^*\right) \\ &= \mathbb{P}\left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \tilde{\mathbb{M}}_{D_0+1} \subsetneq \mathbb{M}_{D_0+1}^*\right) \\ &\leq \mathbb{P}\left(\tilde{\mathbb{M}}_{D_0+1}^c \cap \mathbb{M}_{D_0+1}^* \neq \emptyset\right). \end{aligned} \tag{S16}$$

Now (S15) and (S16) together suggest that

$$\begin{aligned}
& \mathbb{P} \{ \mathcal{E}_{\text{SU}}(D_0 + 1) \text{ holds} \} \\
& \leq \mathbb{P} \{ \mathcal{E}_{\text{SU}}(D_0) \text{ holds} \} + \mathbb{P} \left\{ \tilde{\mathbb{M}}_{D_0+1}^c \cap \mathbb{M}_{D_0+1}^* \neq \emptyset \right\} \\
& \leq \dots \leq \sum_{d=1}^{D_0+1} \mathbb{P} \left\{ \tilde{\mathbb{M}}_{D_0+1}^c \cap \mathbb{M}_{D_0+1}^* \neq \emptyset \right\}.
\end{aligned} \tag{S17}$$

Taking  $D_0 = D - 1$  in (S17) and apply Assumption 1, we conclude

$$\mathbb{P} \{ \mathcal{E}_{\text{SU}}(D) \text{ holds} \} \rightarrow 0.$$

**Step 2: Prove the first part of Theorem S2 and give a probability bound for under-selection.** We compute the probability for under-selection:

$$\begin{aligned}
& \mathbb{P} \{ \mathcal{E}_{\text{U}}(D) \text{ fails} \} \\
& = \mathbb{P} \{ \mathcal{E}_{\text{U}}(1) \text{ fails} \} + \sum_{D_0=2}^D \mathbb{P} \{ \mathcal{E}_{\text{U}}(D_0 - 1) \text{ holds}; \mathcal{E}_{\text{U}}(D_0) \text{ fails} \} \\
& = \mathbb{P} \{ \mathcal{E}_{\text{U}}(1) \text{ fails} \} (\triangleq \otimes_1) + \sum_{D_0=2}^D \mathbb{P} \{ \mathcal{E}_{\text{P}}(D_0 - 1) \text{ holds}; \mathcal{E}_{\text{U}}(D_0) \text{ fails} \} (\triangleq \otimes_2) \\
& \quad + \sum_{D_0=2}^D \mathbb{P} \{ \mathcal{E}_{\text{SU}}(D_0 - 1) \text{ holds}; \mathcal{E}_{\text{U}}(D_0) \text{ fails} \} (\triangleq \otimes_3).
\end{aligned}$$

For  $\otimes_1$ , by definition of  $\mathcal{E}_{\text{U}}(1)$  we have

$$\otimes_1 = \mathbb{P} \{ \mathcal{E}_{\text{U}}(1) \text{ fails} \} = \mathbb{P} \left\{ \widehat{\mathbb{M}}_1 \cap \mathbb{M}_1^{*c} \neq \emptyset \right\} = \mathbb{P} \left\{ \tilde{\mathbb{M}}_1 \cap \mathbb{M}_1^{*c} \neq \emptyset \right\}. \tag{S18}$$

For  $\otimes_2$ , we have

$$\otimes_2 \leq \sum_{D_0=2}^D \mathbb{P} \left( \widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \in [D_0 - 1]; \tilde{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{*c} \neq \emptyset \right) \leq \sum_{D_0=2}^D \mathbb{P} \left\{ \tilde{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{*c} \neq \emptyset \right\}, \tag{S19}$$

which is because on the given event,  $\widehat{\mathbb{M}}_{D_0,+} = \text{H}(\widehat{\mathbb{M}}_{D_0-1}) = \text{H}(\mathbb{M}_{D_0-1}^*) = \mathbb{M}_{D_0,+}^*$  and  $\widehat{\mathbb{M}}_{D_0} = \widehat{\text{S}}(\widehat{\mathbb{M}}_{D_0,+}) = \tilde{\mathbb{M}}_{D_0}$ .

From (S18) and (S19),

$$\limsup_{N \rightarrow \infty} (\otimes_1 + \otimes_2) = \sum_{D_0=1}^D \mathbb{P} \left\{ \tilde{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{*c} \neq \emptyset \right\} \leq \sum_{D_0=1}^D \alpha_{D_0} = \alpha. \text{ (by Assumption 1)} \tag{S20}$$

For  $\otimes_3$ , we have

$$\begin{aligned}
\otimes_3 &\leq \sum_{D_0=2}^D \mathbb{P} \{ \mathcal{E}_{\text{SU}}(D_0 - 1) \text{ holds} \} \\
&\leq \sum_{D_0=2}^D \sum_{d=1}^{D_0-1} \mathbb{P} \left\{ \tilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^* \neq \emptyset \right\} \text{ (using (S17))} \\
&= \sum_{d=1}^{D-1} (D-d) \mathbb{P} \left\{ \tilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^* \neq \emptyset \right\} \rightarrow 0. \text{ (using Assumption 1)} \tag{S21}
\end{aligned}$$

Therefore, by (S20) and (S21), the probability of failure of under-selection gets controlled under  $\alpha$  asymptotically.

As a side product, we obtain finite sample bounds:

$$\mathbb{P} \{ \mathcal{E}_{\text{U}}(D) \text{ fails} \} \leq \sum_{D_0=1}^D \mathbb{P} \left\{ \tilde{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{*c} \neq \emptyset \right\} + \sum_{d=1}^{D-1} (D-d) \mathbb{P} \left\{ \tilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^* \neq \emptyset \right\}.$$

**Step 3. Prove of the second part of Theorem S2 and conclude screening consistency.**

Under  $\alpha = \alpha(N) \rightarrow 0$ , the first part of the result implies that with probability tending to one, we have under-selection:

$$\mathbb{P} \{ \mathcal{E}_{\text{U}}(D) \text{ holds} \} \rightarrow 1.$$

By (S17) and Assumption 1, strict under-selection will not happen with high probability:

$$\mathbb{P} \{ \mathcal{E}_{\text{SU}}(D) \text{ holds} \} \rightarrow 0.$$

Therefore, we conclude the consistency of the screening procedure. □

### C.3 Proof of Theorem 1

We state and prove a more general version of Theorem 1:

**Theorem S3** (Bonferroni corrected marginal t test). *Let  $\tilde{\mathbb{M}}_d = \widehat{\mathbf{S}}(\mathbb{M}_{d,+}^*)$  where  $\mathbb{M}_{d,+}^* = \mathbf{P}(\mathbb{M}_{d-1}^*)$ . Assume Conditions 1, 2, 3 and 4. Then we have the following results for the screening procedure based on Bonferroni corrected marginal t-test:*

$$(i) \text{ (Validity) } \limsup_{N \rightarrow \infty} \sum_{d=1}^D \mathbb{P} \left\{ \tilde{\mathbb{M}}_d \cap \mathbb{M}_d^{*c} \neq \emptyset \right\} \leq \sum_{d=1}^D \alpha_d = \alpha.$$

$$(ii) \text{ (Consistency) } \limsup_{N \rightarrow \infty} D \sum_{d=1}^D \mathbb{P} \left\{ \tilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^* \neq \emptyset \right\} = 0.$$



(iii) (Type I error control) Overall the procedure achieves type I error rate control:

$$\limsup_{N \rightarrow \infty} \mathbb{P} \left( \widehat{\mathbb{M}} \cap (\cup_{d=1}^D \mathbb{M}_d^*)^c \neq \emptyset \right) \leq \alpha.$$

(iv) (Perfect screening) When  $\delta'$  is strictly positive, we have  $\max_{d \in [D]} \alpha_d \rightarrow 0$  and

$$\lim_{N \rightarrow \infty} \mathbb{P} \left( \widehat{\mathbb{M}} = \bigcup_{d=1}^D \mathbb{M}_d^* \right) = 1.$$

Part (i) and (ii) of Theorem 1 justified Assumption 1 and 2 respectively, which build up the basis for applying Theorem S2. Part (iii) guarantees type I error control under the significance level  $\alpha$ . When we let  $\alpha$  decay to zero, Part (iii) implies that we will not include redundant terms into the selected working model. Part (iv) further states a stronger result with vanishing  $\alpha$  - perfect selection can be achieved asymptotically.

*Proof of Theorem 1.* (i) First, we show validity:

$$\begin{aligned} \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d \cap \mathbb{M}_d^{*c} \neq \emptyset \right\} &= \mathbb{P} \left\{ \exists \mathcal{K} \in \mathbb{M}_{d,+}^* \setminus \mathbb{M}_d^*, \left| \frac{\widehat{\tau}_{\mathcal{K}}}{\widehat{v}_{\mathcal{K},R}} \right| \geq \Phi^{-1} \left( 1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^*|} \right) \right\} \\ &\leq \sum_{\mathcal{K} \in \mathbb{M}_{d,+}^* \setminus \mathbb{M}_d^*} \mathbb{P} \left\{ \left| \frac{\widehat{\tau}_{\mathcal{K}}}{\widehat{v}_{\mathcal{K},R}} \right| \geq \Phi^{-1} \left( 1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^*|} \right) \right\} \\ &\leq \sum_{\mathcal{K} \in \mathbb{M}_{d,+}^* \setminus \mathbb{M}_d^*} \left( \frac{\alpha_d}{|\mathbb{M}_{d,+}^*|} + \frac{\widetilde{C}}{(QN_0)^{1/3}} \right) \text{ (by Corollary S2)} \\ &\leq \left( \alpha_d + \frac{\widetilde{C}|\mathbb{M}_{d,+}^*|}{N^{1/3}} \right). \end{aligned}$$

Hence,

$$\sum_{d=1}^D \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d \cap \mathbb{M}_d^{*c} \neq \emptyset \right\} \leq \sum_{d=1}^D \left( \alpha_d + \frac{\widetilde{C}|\mathbb{M}_{d,+}^*|}{N^{1/3}} \right).$$

Due to the effect heredity condition 4, we have

$$|\mathbb{M}_{1,+}^*| = |\mathbb{M}_1^*|, \quad |\mathbb{M}_{d,+}^*| \leq K|\mathbb{M}_{d-1}^*|.$$

Hence

$$\limsup_{N \rightarrow \infty} \sum_{d=1}^D \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d \cap \mathbb{M}_d^{*c} \neq \emptyset \right\} \leq \alpha + \limsup_{N \rightarrow \infty} \frac{K\widetilde{C}|\mathbb{M}^*|}{N^{1/3}} = \alpha. \text{ (using Condition 2(iii))}$$

- (ii) Second, we show consistency. Assume the nonzero  $\tau_{\mathcal{K}}$ 's are positive. If some are negative one can simply modify the direction of some of the inequalities and still validate the proof.

$$\begin{aligned}
& \mathbb{P} \left\{ \tilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^* \neq \emptyset \right\} \\
&= \mathbb{P} \left\{ \exists \mathcal{K} \in \mathbb{M}_d^*, \left| \frac{\hat{\tau}_{\mathcal{K}}}{\hat{v}_{\mathcal{K},R}} \right| \leq \Phi^{-1} \left( 1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^*|} \right) \right\} \\
&\leq \sum_{\mathcal{K} \in \mathbb{M}_d^*} \mathbb{P} \left\{ \left| \frac{\hat{\tau}_{\mathcal{K}}}{\hat{v}_{\mathcal{K},R}} \right| \leq \Phi^{-1} \left( 1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^*|} \right) \right\} \\
&\leq \sum_{\mathcal{K} \in \mathbb{M}_d^*} \mathbb{P} \left\{ \left| \frac{\hat{\tau}_{\mathcal{K}}}{v_{\mathcal{K},R}} \right| \leq \frac{\hat{v}_{\mathcal{K},R}}{v_{\mathcal{K},R}} \Phi^{-1} \left( 1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^*|} \right) \right\} \\
&\leq \sum_{\mathcal{K} \in \mathbb{M}_d^*} \mathbb{P} \left\{ \left| \frac{\hat{\tau}_{\mathcal{K}}}{v_{\mathcal{K},R}} \right| \leq \left\{ 1 + \frac{\tilde{C}}{(QN_0)^{1/3}} \right\} \Phi^{-1} \left( 1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^*|} \right) \right\} + \mathbb{P} \left\{ \frac{\hat{v}_{\mathcal{K},R}}{v_{\mathcal{K},R}} > 1 + \frac{\tilde{C}}{(QN_0)^{1/3}} \right\}.
\end{aligned}$$

For simplicity, let

$$Z_d^* = \Phi^{-1} \left( 1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^*|} \right).$$

Then

$$\begin{aligned}
& \mathbb{P} \left\{ \tilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^* \neq \emptyset \right\} \\
&\leq \sum_{\mathcal{K} \in \mathbb{M}_d^*} \left( \mathbb{P} \left\{ -Z_d^* - \frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}} \leq \frac{\hat{\tau}_{\mathcal{K}}}{v_{\mathcal{K},R}} - \frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}} \leq Z_d^* - \frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}} \right\} + \frac{\tilde{C}}{(QN_0)^{1/3}} \right) \\
&= \sum_{\mathcal{K} \in \mathbb{M}_d^*} \Phi \left\{ r_{\mathcal{K}}^{-1} \left( Z_d^* - \frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}} \right) \right\} - \Phi \left\{ r_{\mathcal{K}}^{-1} \left( -Z_d^* - \frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}} \right) \right\} (\triangleq \otimes) + \frac{\tilde{C}|\mathbb{M}_d^*|}{(QN_0)^{1/3}}.
\end{aligned}$$

With Condition 2, we have

$$Z_d^* = \Theta \left( \sqrt{2 \ln \frac{2|\mathbb{M}_{d,+}^*|}{\alpha_d}} \right) = \Theta(\sqrt{(\delta' + \delta''/3) \ln N}),$$

$$\left| \frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}} \right| = \Theta(N^{1/2+\delta}) = \Theta(N^{\delta_0}) \text{ (by defining } \delta_0 = 1/2 + \delta > 0).$$

Because  $\delta > -1/2$  and  $\delta' \geq 0$ , we have  $|\frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}}| \rightarrow \infty$  and  $Z_d^*/(|\frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}}|) \rightarrow 0$ . Therefore,

$$\Phi \left\{ r_{\mathcal{K}}^{-1} \left( Z_d^* - \frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}} \right) \right\} - \Phi \left\{ r_{\mathcal{K}}^{-1} \left( -Z_d^* - \frac{\tau_{\mathcal{K}}}{v_{\mathcal{K},R}} \right) \right\} = \Theta(N^{-\delta_0} \exp\{-N^{2\delta_0}/2\}).$$

Now applying Condition 2 again, we have

$$D \sum_{d=1}^D \mathbb{P} \left\{ \tilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^* \neq \emptyset \right\} = \Theta \left( D|\mathbb{M}^*|N^{-\delta_0} \exp\{-N^{2\delta_0}/2\} + D|\mathbb{M}^*|/N^{1/3} \right) = o(1).$$

(iii) The Type I error rate control comes from Theorem S2.

(iv) The perfect selection result follows from Theorem S2.

□

## C.4 Proof of Theorem 2

Theorem 2 is a direct result of Theorem 1, Lemma S2 and the following Berry–Esseen bound:

**Lemma S5** (Berry–Esseen bound under perfect screening). *Assume (S12). Then*

$$\begin{aligned} & \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}(\widehat{\mathbb{M}}) - \gamma}{v(\mathbb{M}^*)} \leq t \right\} - \Phi(t) \right| \\ & \leq 2\mathbb{P} \left\{ \widehat{\mathbb{M}} \neq \mathbb{M}^* \right\} + 2C\sigma_w \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\sqrt{\underline{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})} \cdot \sqrt{N_0}} \cdot \frac{\|\mathbf{f}[\mathbb{M}^*]\|_\infty}{\|\mathbf{f}[\mathbb{M}^*]\|_2}. \end{aligned}$$

*Proof of Lemma S5.* This lemma is a direct application of Lemma S3. First we check that

$$\gamma(\mathbb{M}^*) = \gamma.$$

From the definition of  $\gamma$  (S7), we have

$$\begin{aligned} \gamma &= \mathbf{f}^\top \bar{Y} \\ &= \mathbf{f}^\top G\tau = \mathbf{f}^\top G(\cdot, \mathbb{M}^*)\tau(\mathbb{M}^*) \\ &= Q^{-1} \mathbf{f}^\top G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \bar{Y} = \gamma(\mathbb{M}^*). \end{aligned}$$

□

Now apply Lemma S3 with  $\mathbb{M} = \mathbb{M}^*$  to get the result of Theorem 2.

## C.5 Statement and proof of Lemma S6

The following lemma gives the closed form solution of the RLS estimator (4.7).

**Lemma S6.**  $\widehat{Y}_R$  from (4.7) can be expressed as:

$$\widehat{Y}_R = Q^{-1} G(\cdot, \widehat{\mathbb{M}}) G(\cdot, \widehat{\mathbb{M}})^\top \widehat{Y}.$$

If  $\widehat{\mathbb{M}} = \mathbb{M}^*$ ,  $\mathbb{E} \left\{ \widehat{Y}_R \right\} = \bar{Y}$ .

*Proof of Lemma S6.* Due to the orthogonality of  $G$ , we have the following decomposition:

$$\hat{Y} = Q^{-1}G(\cdot, \hat{\mathbb{M}})G(\cdot, \hat{\mathbb{M}})^\top \hat{Y} + Q^{-1}G(\cdot, \hat{\mathbb{M}}^c)G(\cdot, \hat{\mathbb{M}}^c)^\top \hat{Y}.$$

By the constraint in (4.7), we have

$$\|\hat{Y} - \mu\|^2 = \|Q^{-1}G(\cdot, \hat{\mathbb{M}}^c)G(\cdot, \hat{\mathbb{M}}^c)^\top \hat{Y}\|^2 + \|Q^{-1}G(\cdot, \hat{\mathbb{M}})G(\cdot, \hat{\mathbb{M}})^\top \hat{Y} - \mu\|^2,$$

which is minimized at

$$\hat{\mu} = \hat{Y}_R = Q^{-1}G(\cdot, \hat{\mathbb{M}})G(\cdot, \hat{\mathbb{M}})^\top \hat{Y}.$$

Besides,  $\hat{\mu}$  satisfies the constraint in (4.7).

Next we verify  $\mathbb{E}\{\hat{Y}_R\} = \bar{Y}$  if  $\hat{\mathbb{M}} = \mathbb{M}^\star$ . Utilizing the orthogonality of  $G$  again, we have

$$\bar{Y} = Q^{-1}G(\cdot, \mathbb{M}^\star)G(\cdot, \mathbb{M}^\star)^\top \bar{Y} + Q^{-1}G(\cdot, \mathbb{M}^{\star c})G(\cdot, \mathbb{M}^{\star c})^\top \bar{Y}$$

□

## C.6 Proof of Proposition 1

*Proof of Proposition 1.* (i) Based on the definition of  $v_R^2$  and  $v^2$ , we have

$$\frac{v_R^2}{v^2} = \frac{\mathbf{f}^{\star\top} V_{\hat{Y}} \mathbf{f}^\star}{\mathbf{f}^\top V_{\hat{Y}} \mathbf{f}} = \frac{\|\mathbf{f}^\star\|_2^2}{\|\mathbf{f}\|_2^2}$$

because  $\kappa(V_{\hat{Y}}) = 1$ . We further compute

$$\frac{\|\mathbf{f}^\star\|_2^2}{\|\mathbf{f}\|_2^2} = \frac{\mathbf{f}^\top \{Q^{-1}G(\cdot, \mathbb{M}^\star)G(\cdot, \mathbb{M}^\star)^\top\} \mathbf{f}}{\mathbf{f}^\top \mathbf{f}} \leq 1$$

where the inequality holds because of the following dominance relationship:

$$Q^{-1}G(\cdot, \mathbb{M}^\star)G(\cdot, \mathbb{M}^\star)^\top \preceq I_Q.$$

(ii) Because the order of the nonzero elements in  $\mathbf{f}$  is not crucial here, we assume the first  $s^\star$  coordinates of  $\mathbf{f}$  are nonzero while the rest are zero without loss of generality. We can compute

$$\frac{v_R^2}{v^2} = \frac{\mathbf{f}^{\star\top} V_{\hat{Y}} \mathbf{f}^\star}{\mathbf{f}^\top V_{\hat{Y}} \mathbf{f}} \leq \kappa(V_{\hat{Y}}) \cdot \frac{\|\mathbf{f}^\star\|_2^2}{\|\mathbf{f}\|_2^2}. \quad (\text{S22})$$

For  $\mathbf{f}^*$  we have

$$\begin{aligned}
\|\mathbf{f}^*\|_2 &= \|Q^{-1}G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \mathbf{f}\|_2 \\
&= \left\| Q^{-1}G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \left\{ \sum_{s=1}^{s^*} \mathbf{f}(s) \mathbf{e}_s \right\} \right\|_2 \\
&\leq \sum_{s=1}^{s^*} |\mathbf{f}(s)| \|Q^{-1}G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \mathbf{e}_s\|_2 \\
&= \left( \frac{|\mathbb{M}^*|}{Q} \right)^{1/2} \sum_{s=1}^{s^*} |\mathbf{f}(s)| = \left( \frac{|\mathbb{M}^*|}{Q} \right)^{1/2} \|\mathbf{f}\|_1.
\end{aligned}$$

Then we have

$$\frac{\|\mathbf{f}^*\|_2^2}{\|\mathbf{f}\|_2^2} \leq \frac{|\mathbb{M}^*|}{Q} \frac{\|\mathbf{f}\|_1^2}{\|\mathbf{f}\|_2^2} \leq \frac{s^* |\mathbb{M}^*|}{Q}. \quad (\text{S23})$$

Combining (S22) and (S23), we conclude the result.  $\square$

**Proposition S1** (Asymptotic length of confidence comparison). *Assume that both  $\hat{\gamma}$  and  $\hat{\gamma}_R$  converge to a normal distribution as the sample size tends to infinity. Assume the variance estimators are consistent:  $N(\hat{v}^2 - v_{\text{lim}}^2) = o_{\mathbb{P}}(1)$ ,  $N(\hat{v}_R^2 - v_{R,\text{lim}}^2) = o_{\mathbb{P}}(1)$ .*

(i) *If the condition number of  $D_{\hat{\gamma}}$  satisfies  $\kappa(D_{\hat{\gamma}}) = 1$ , we have*

$$\frac{v_{R,\text{lim}}^2}{v_{\text{lim}}^2} \leq 1.$$

(ii) *Let  $s^*$  denote the number of nonzero elements in  $\mathbf{f}$ , then we have*

$$\frac{v_{R,\text{lim}}^2}{v_{\text{lim}}^2} \leq \kappa(D_{\hat{\gamma}}) \cdot \frac{s^* |\mathbb{M}^*|}{Q}.$$

### C.7 Proof of Theorem 3

*Proof of Theorem 3.* According to Condition 5 and Theorem 1, with Strategy 1,

$$\mathbb{P} \left\{ \hat{\mathbb{M}} = \cup_{d=1}^{d^*} \mathbb{M}_d^* \right\} \rightarrow 1.$$

We will apply Lemma S5 with

$$\mathbb{M} = \underline{\mathbb{M}}^* = \cup_{d=1}^{d^*} \mathbb{M}_d^*.$$

We only need to verify  $\gamma = \gamma[\mathbb{M}]$  under the orthogonality condition (5.13).

$$\begin{aligned}
\gamma &= \mathbf{f}^\top \bar{\mathbf{Y}} \\
&= \mathbf{f}^\top G \tau = \mathbf{f}^\top G(\cdot, \underline{\mathbb{M}}^*) \tau(\underline{\mathbb{M}}^*) + \mathbf{f}^\top G(\cdot, \underline{\mathbb{M}}^{*c}) \tau(\underline{\mathbb{M}}^{*c}).
\end{aligned}$$

Now by (5.13),  $\mathbf{f}^\top G(\cdot, \mathbb{M}^c) = 0$ . Hence

$$\gamma = Q^{-1} \mathbf{f}^\top G(\cdot, \cup_{d=1}^{d^*} \mathbb{M}_d^*) G(\cdot, \cup_{d=1}^{d^*} \mathbb{M}_d^*)^\top \bar{Y} = \gamma.$$

□

## C.8 Proof of Theorem 4

*Proof of Theorem 4.* This proof can be finished by applying Lemma S3 and S4 with  $\mathbb{M} = \overline{\mathbb{M}}^*$  and checking  $\gamma[\overline{\mathbb{M}}^*] = \gamma$ , which is omitted here. □

## C.9 Proof of Proposition 1

*Proof of Proposition 1.* Because the order of the nonzero elements in  $\mathbf{f}^*$  is not crucial, we assume only the first  $s^*$  elements of  $\mathbf{f}$  are nonzero. That is,

$$\mathbf{f} = f_1 \mathbf{e}_1 + \cdots + f_{s^*} \mathbf{e}_{s^*}. \quad (\text{S24})$$

We can verify that

$$\|Q^{-1} G(\cdot, \mathbb{M}^*) G(\cdot, \mathbb{M}^*)^\top \mathbf{e}_k\|_2 = \frac{|\mathbb{M}^*|}{Q}, \quad \forall k \in [Q]. \quad (\text{S25})$$

Therefore,

$$\frac{v_{\text{R}}^2}{v^2} = \frac{\mathbf{f}^* V_{\hat{Y}} \mathbf{f}^*}{\mathbf{f} V_{\hat{Y}} \mathbf{f}} \leq \frac{\varrho_{\max}(V_{\hat{Y}}) \|\mathbf{f}^*\|_2^2}{\varrho_{\min}(V_{\hat{Y}}) \|\mathbf{f}\|_2^2} = \kappa(V_{\hat{Y}}) \cdot \frac{\|\mathbf{f}^*\|_2^2}{\|\mathbf{f}\|_2^2}.$$

On the one hand, using  $Q^{-1} G(\cdot, \mathbb{M}^*) G(\cdot, \mathbb{M}^*)^\top \preceq I_Q$ , we have

$$\frac{\|\mathbf{f}^*\|_2^2}{\|\mathbf{f}\|_2^2} \leq 1. \quad (\text{S26})$$

On the other hand, using (S24) and (S25), we have

$$\frac{\|\mathbf{f}^*\|_2^2}{\|\mathbf{f}\|_2^2} \leq \frac{\|\mathbf{f}\|_1^2}{\|\mathbf{f}\|_2^2} \cdot \frac{|\mathbb{M}^*|}{Q} \leq \frac{s^* |\mathbb{M}^*|}{Q}. \quad (\text{S27})$$

Combining (S26) and (S27) concludes the proof. □

## C.10 Proof of Theorem 5

For simplicity, we focus on the case given by (6.14). The general proof can be completed similarly. We begin by the following lemma:

**Lemma S7** (Consistency of the selected tie sets). *Assume Conditions 1, 3 and 6. There exists universal constants  $C, C' > 0$ , such that when  $N > n(\delta_1, \delta_2, \delta_3)$ , we have*

$$\mathbb{P}\left\{\widehat{\mathcal{T}}_1 = \mathcal{T}_1\right\} \geq 1 - \mathbb{P}\{\widehat{\mathbb{M}} \neq \mathbb{M}^*\} - C|\mathcal{T}'||\mathcal{T}_1| \left\{ \sqrt{\frac{\bar{c}\Delta|\mathbb{M}^*|}{N^{1+2\delta_2}}} \exp\left(-\frac{C'N^{1+2\delta_2}}{\bar{c}\Delta|\mathbb{M}^*|}\right) + \sigma \frac{\underline{c}^{-1/2} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N}} \right\}.$$

Lemma S7 establishes a finite sample bound to quantify the performance of the tie set selection step in Algorithm 2. The tail bound implies that the performance of tie selection depends on several elements:

- Quality of effect screening. Ideally we hope perfect screening can be achieved. In other words, the misspecification probability  $\mathbb{P}\{\widehat{\mathbb{M}} \neq \mathbb{M}^*\}$  is small in an asymptotic sense.
- Size of the tie  $|\mathcal{T}_1|$  and the number of factor combinations considered  $|\mathcal{T}'|$ . These two quantities play a natural role because one can expect the difficulty of selection will increase if there are too many combinations present in the first tie or involved in comparison.
- Size of between-group distance  $d_h^*$ . If the gap between  $\bar{Y}_{(1)}$  and the remaining order values are large,  $\eta_N = \Theta(N^{\delta_2})$  is allowed to take larger values and the term

$$\sqrt{\frac{\bar{c}\Delta|\mathbb{M}^*|}{N^{1+2\delta_2}}} \exp\left(-\frac{C'N^{1+2\delta_2}}{\bar{c}\Delta|\mathbb{M}^*|}\right)$$

can become smaller in magnitude.

- Population level property of potential outcomes. The scale of the centered potential outcomes  $|Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|$  should be controlled, and the population variance  $S(\mathbf{z}, \mathbf{z})$  should be non-degenerate.
- The relative scale between number of nonzero effects  $|\mathbb{M}^*|$  and the total number of units  $N$ . The larger  $N$  is compared to  $|\mathbb{M}^*|$ , the easier for us to draw valid asymptotic conclusions.

*Proof of Lemma S7.* The high level idea of the proof is: we first prove the non-asymptotic bounds over the random event  $\widehat{\mathbb{M}} = \mathbb{M}^*$ , then make up for the cost of  $\widehat{\mathbb{M}} \neq \mathbb{M}^*$ . Over  $\widehat{\mathbb{M}} = \mathbb{M}^*$ , we have

$$\widehat{Y}_R = \widehat{Y}_R^* = G(\cdot, \mathbb{M}^*)\widehat{\tau}(\mathbb{M}^*) = Q^{-1}G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \widehat{Y}.$$

We apply Lemma S3 to establish a Berry–Esseen bound for each  $\widehat{Y}_R^*(\mathbf{z})$ . Note that

$$\widehat{Y}_R^*(\mathbf{z}) = \mathbf{f}_z^\top \widehat{Y}, \quad \mathbf{f}_z^\top = Q^{-1}G(\mathbf{z}, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top.$$

By calculation we have

$$\|\mathbf{f}_z\|_\infty = Q^{-1}|\mathbb{M}^\star|, \quad \|\mathbf{f}_z\|_2 = \sqrt{Q^{-1}|\mathbb{M}^\star|}.$$

Also we can show that

$$\sum_{z'=1}^Q \mathbf{f}_z(z')^2 N_{z'}^{-1} S(z', z') \leq \sigma^2 v^2(\mathbb{M}).$$

and obtain

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{Y}_R^\star(z) - \bar{Y}(z)}{v_N} \leq t \right\} - \Phi(t) \right| \leq 2C\sigma \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \sqrt{\frac{|\mathbb{M}^\star|}{QN_0}}.$$

**A probabilistic bound on the ordered statistics.** We show a bound on

$$\mathbb{P} \left\{ \max_{z \in \mathcal{T}' \setminus \mathcal{T}_1} \hat{Y}_R^\star(z) < \min_{z \in \mathcal{T}_1} \hat{Y}_R^\star(z) \leq \max_{z \in \mathcal{T}_1} \hat{Y}_R^\star(z) \right\}.$$

It is known that (Wainwright, 2019, Exercise 2.2)

$$1 - \Phi(x) = \int_x^\infty \phi(t) dt \leq \frac{1}{x} \int_x^\infty t \phi(t) dt \leq \frac{1}{\sqrt{2\pi}x} \left\{ \exp \left( -\frac{x^2}{2} \right) \right\}.$$

Hence

$$\begin{aligned} & \mathbb{P} \left\{ \sqrt{N} \left| \hat{Y}_R^\star(z) - \bar{Y}(z) \right| \geq \sqrt{N} d_h^\star \right\} \\ & \leq \frac{v_N}{\sqrt{2\pi} d_h^\star} \cdot \exp \left( -\frac{d_h^{*2}}{2v_N^2} \right) + 2C\sigma \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^\star|}{N_0 Q}}. \end{aligned} \quad (\text{S28})$$

Therefore, for all  $z \in \mathcal{T}' \setminus \mathcal{T}$  and  $z' \in \mathcal{T}_1$ ,

$$\begin{aligned} & \mathbb{P} \left\{ \hat{Y}_R^\star(z') - \hat{Y}_R^\star(z) < 0 \right\} \\ & = \mathbb{P} \left\{ \sqrt{N}(\hat{Y}_R^\star(z') - \bar{Y}(z')) - \sqrt{N}(\hat{Y}_R^\star(z) - \bar{Y}(z)) < \sqrt{N}(\bar{Y}(z) - \bar{Y}(z')) \right\} \\ & \leq \mathbb{P} \left\{ \sqrt{N}(\hat{Y}_R^\star(z') - \bar{Y}(z')) - \sqrt{N}(\hat{Y}_R^\star(z) - \bar{Y}(z)) < -2\sqrt{N} d_h^\star \right\} \\ & = \mathbb{P} \left\{ \sqrt{N}(\hat{Y}_R^\star(z') - \bar{Y}(z')) - \sqrt{N}(\hat{Y}_R^\star(z) - \bar{Y}(z)) < -2\sqrt{N} d_h^\star, \right. \\ & \quad \left. \sqrt{N}(\hat{Y}_R^\star(z) - \bar{Y}(z)) < \sqrt{N} d_h^\star \right\} \\ & + \mathbb{P} \left\{ \sqrt{N}(\hat{Y}_R^\star(z') - \bar{Y}(z')) - \sqrt{N}(\hat{Y}_R^\star(z) - \bar{Y}(z)) < -2\sqrt{N} d_h^\star, \right. \\ & \quad \left. \sqrt{N}(\hat{Y}_R^\star(z) - \bar{Y}(z)) < \sqrt{N} d_h^\star \right\} \\ & \leq \mathbb{P} \left\{ \sqrt{N}(\hat{Y}_R^\star(z') - \bar{Y}(z')) < -\sqrt{N} d_h^\star \right\} + \mathbb{P} \left\{ \sqrt{N}(\hat{Y}_R^\star(z) - \bar{Y}(z)) \geq \sqrt{N} d_h^\star \right\}. \end{aligned}$$



Using (S28) we have

$$\begin{aligned} & \mathbb{P} \left\{ \hat{Y}_R^*(z') - \hat{Y}_R^*(z) < 0 \right\} \\ & \leq \frac{\sqrt{\bar{c}\Delta|\mathbb{M}^*|}}{\sqrt{2\pi N_0 Q d_h^*}} \cdot \exp \left( -\frac{N_0 Q d_h^{*2}}{2\bar{c}\bar{s}|\mathbb{M}^*|} \right) + 2C\sigma \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}}. \end{aligned}$$

Now a union bound gives

$$\begin{aligned} & \mathbb{P} \left\{ \hat{Y}_R^*(z') - \hat{Y}_R^*(z) < 0 \right\} \\ & \geq 1 - |\mathcal{T}_1| |\mathcal{T}'| \left\{ \frac{\sqrt{\bar{c}\bar{s}|\mathbb{M}^*|}}{\sqrt{2\pi N_0 Q d_h^*}} \cdot \exp \left( -\frac{N_0 Q d_h^{*2}}{2\bar{c}\bar{s}|\mathbb{M}^*|} \right) + 2C\sigma \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned}$$

Now using that  $d_h^* = \Theta(N^{\delta_1})$ ,  $N d_h^{*2} = \Theta(N^{1+2\delta_1})$  with  $1 + 2\delta_1 > 0$ . The first term in the bracket has the following order

$$\frac{\sqrt{\bar{c}\bar{s}|\mathbb{M}^*|}}{\sqrt{2\pi N_0 Q d_h^*}} \cdot \exp \left( -\frac{N_0 Q d_h^{*2}}{2\bar{c}\bar{s}|\mathbb{M}^*|} \right) = \Theta \left( \sqrt{\frac{\bar{c}\bar{s}|\mathbb{M}^*|}{N^{1+2\delta_1}}} \exp \left\{ -\frac{C' N^{1+2\delta_1}}{\bar{c}\bar{s}|\mathbb{M}^*|} \right\} \right)$$

where  $C' > 0$  is a universal constant due to Condition 2. Note that  $\delta_2 > \delta_1$ . Thus when  $N$  is large enough, we have

$$\begin{aligned} & \mathbb{P} \left\{ \hat{Y}_R^*(z') - \hat{Y}_R^*(z) < 0 \right\} \\ & \geq 1 - C |\mathcal{T}_1| |\mathcal{T}'| \left\{ \sqrt{\frac{\bar{c}\bar{s}|\mathbb{M}^*|}{N^{1+2\delta_1}}} \exp \left\{ -\frac{C' N^{1+2\delta_1}}{\bar{c}\bar{s}|\mathbb{M}^*|} \right\} + \sigma \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\} \quad (\text{S29}) \end{aligned}$$

**Nice separation.** Consider the following random index:

$$\tilde{z} \in \arg \max_{z \in \mathcal{T}'} \hat{Y}_R^*(z).$$

For any  $\bar{\epsilon} > 0$ ,

$$\begin{aligned} & \mathbb{P} \left\{ \min_{z \notin \mathcal{T}_1} |\hat{Y}_R^*(z) - \hat{Y}_R^*(\tilde{z})| / \eta_N \geq 2\bar{\epsilon} \right\} \\ & \geq \mathbb{P} \left\{ \min_{z \notin \mathcal{T}_1, z' \in \mathcal{T}_1} |\hat{Y}_R^*(z) - \hat{Y}_R^*(z')| / \eta_N \geq 2\bar{\epsilon}, \tilde{z} \in \mathcal{T}_1 \right\} \\ & \geq \mathbb{P} \left\{ \min_{z \notin \mathcal{T}_1, z' \in \mathcal{T}_1} |\hat{Y}_R^*(z) - \hat{Y}_R^*(z')| / \eta_N \geq 2\bar{\epsilon} \right\} + \mathbb{P} \{ \tilde{z} \in \mathcal{T}_1 \} - 1 \\ & \geq \mathbb{P} \{ \tilde{z} \in \mathcal{T}_1 \} - \sum_{z \notin \mathcal{T}_1, z' \in \mathcal{T}_1} \mathbb{P} \left\{ |\hat{Y}_R^*(z) - \hat{Y}_R^*(z')| / \eta_N \leq 2\bar{\epsilon} \right\}. \quad (\text{S30}) \end{aligned}$$

To proceed we have the following tail bound:

$$\begin{aligned}
& \mathbb{P} \left\{ |\hat{Y}_R^*(z) - \hat{Y}_R^*(z')|/\eta_N \leq 2\bar{\epsilon} \right\} \\
&= \mathbb{P} \left\{ |\{\hat{Y}_R^*(z) - \bar{Y}(z)\} - \{\hat{Y}_R^*(z') - \bar{Y}(z')\} - \{\bar{Y}(z) - \bar{Y}(z')\}| \leq 2\bar{\epsilon}\eta_N \right\} \\
&\leq \mathbb{P} \left\{ |\bar{Y}(z) - \bar{Y}(z')| - |\hat{Y}_R^*(z) - \bar{Y}(z)| - |\hat{Y}_R^*(z') - \bar{Y}(z')| \leq 2\bar{\epsilon}\eta_N \right\} \\
&\leq \mathbb{P} \left\{ |\hat{Y}_R^*(z) - \bar{Y}(z)| + |\hat{Y}_R^*(z') - \bar{Y}(z')| \geq 2d_h^* - 2\bar{\epsilon}\eta_N \right\} \\
&\leq \mathbb{P} \left\{ |\hat{Y}_R^*(z) - \bar{Y}(z)| \geq d_h^* - \bar{\epsilon}\eta_N \right\} + \mathbb{P} \left\{ |\hat{Y}_R^*(z') - \bar{Y}(z')| \geq d_h^* - \bar{\epsilon}\eta_N \right\} \\
&\quad (\text{because } z \notin \mathcal{T}_1 \text{ and } z' \in \mathcal{T}_1) \\
&\leq 4 \left\{ \frac{\sqrt{\bar{c}\Delta|\mathbb{M}^*|}}{\sqrt{2\pi N_0 Q}(d_h^* - \bar{\epsilon}\eta_N)} \cdot \exp \left( -\frac{N_0 Q(d_h^* - \bar{\epsilon}\eta_N)^2}{2\bar{c}\bar{s}|\mathbb{M}^*|} \right) \right. \\
&\quad \left. + 2C\sigma \frac{\bar{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \\
&\quad (\text{This is deduced analogously to the proof in the previous part})
\end{aligned}$$

By the conditions we imposed in the theorem, we know that when  $N$  is large enough,

$$d_h^* - \bar{\epsilon}\eta_N > d_h^*/2.$$

Hence, for  $N > N(\delta_1, \delta_2)$ , we have

$$\begin{aligned}
& \sum_{z \notin \mathcal{T}_1, z' \in \mathcal{T}_1} \mathbb{P} \left\{ |\hat{Y}_R^*(z) - \hat{Y}_R^*(z')|/\eta_N \leq 2\bar{\epsilon} \right\} \\
&\leq 4|\mathcal{T}_1||\mathcal{T}'| \left\{ \frac{\sqrt{2\bar{c}\bar{s}|\mathbb{M}^*|}}{\sqrt{\pi N_0 Q}d_h^*} \cdot \exp \left( -\frac{N_0 Q d_h^{*2}}{8\bar{c}\bar{s}|\mathbb{M}^*|} \right) + 2C\sigma \frac{\bar{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}.
\end{aligned}$$

Combined with (S30), we have:

$$\begin{aligned}
& \mathbb{P} \left\{ \min_{z \notin \mathcal{T}_1} |\hat{Y}_R^*(z) - \hat{Y}_R^*(\tilde{z})|/\eta_N \geq 2\bar{\epsilon} \right\} \\
&\geq \underbrace{\mathbb{P} \{ \tilde{m} \in \mathcal{T}_1 \} - 4|\mathcal{T}_1||\mathcal{T}'| \frac{\sqrt{2\bar{c}\bar{s}|\mathbb{M}^*|}}{\sqrt{\pi N_0 Q}d_h^*} \cdot \exp \left( -\frac{N_0 Q d_h^{*2}}{8\bar{c}\bar{s}|\mathbb{M}^*|} \right)}_{\text{Term I}} \\
&\quad - \underbrace{4|\mathcal{T}_1||\mathcal{T}'| 2C\sigma \frac{\bar{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}}}_{\text{Term II}}.
\end{aligned}$$

Analogous to the discussion in the previous part, when  $N$  is sufficiently large, we can show

$$\begin{aligned} & \mathbb{P} \left\{ \min_{\mathbf{z} \notin \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}})|/\eta_N \geq 2\bar{\epsilon} \right\} \\ & \geq \mathbb{P} \{ \tilde{m} \in \mathcal{T}_1 \} - C|\mathcal{T}_1||\mathcal{T}'| \left\{ \sqrt{\frac{\bar{c}\bar{s}|\mathbb{M}^*|}{N^{1+2\delta_2}}} \exp \left\{ -\frac{C'N^{1+2\delta_2}}{\bar{c}\bar{s}|\mathbb{M}^*|} \right\} + \sigma \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned}$$

Similarly we can show for any  $\mathbf{z} \in \mathcal{T}_1$  and  $\underline{\epsilon} > 0$ ,

$$\begin{aligned} & \mathbb{P} \left\{ \max_{\mathbf{z} \in \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}})|/\eta_N \leq 2\underline{\epsilon} \right\} \\ & \geq \mathbb{P} \{ \tilde{\mathbf{z}} \in \mathcal{T}_1 \} - \sum_{\mathbf{z} \neq \mathbf{z}' \in \mathcal{T}_1} \mathbb{P} \left\{ |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\mathbf{z}')|/\eta_N > 2\underline{\epsilon} \right\}. \end{aligned}$$

Then we have for  $\mathbf{z} \neq \mathbf{z}' \in \mathcal{T}_1$ ,

$$\begin{aligned} & \mathbb{P} \left\{ |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\mathbf{z}')|/\eta_N > 2\underline{\epsilon} \right\} \\ & \leq \mathbb{P} \left\{ |\hat{Y}_R^*(\mathbf{z}) - \bar{Y}(\mathbf{z})| \geq \underline{\epsilon}\eta_N - d_h \right\} + \mathbb{P} \left\{ |\hat{Y}_R^*(\mathbf{z}') - \bar{Y}(\mathbf{z}')| \geq \underline{\epsilon}\eta_N - d_h \right\} \\ & \leq 4 \left\{ \frac{\sqrt{\bar{c}\bar{s}|\mathbb{M}^*|}}{\sqrt{2\pi N_0 Q}(\underline{\epsilon}\eta_N - d_h)} \cdot \exp \left( -\frac{N_0 Q(\underline{\epsilon}\eta_N - d_h)^2}{2\bar{c}\bar{s}|\mathbb{M}^*|} \right) \right. \\ & \quad \left. + 2C\sigma \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned}$$

By the scaling of the parameters, when  $N_0$  is large enough  $N > N(\delta_2, \delta_3)$ ,  $\underline{\epsilon}\eta_N - d_h > \underline{\epsilon}\eta_N/2$ . That being said,

$$\begin{aligned} & \mathbb{P} \left\{ |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\mathbf{z}')|/\eta_N > 2\underline{\epsilon} \right\} \\ & \leq 4 \left\{ \frac{\sqrt{2\bar{c}\bar{s}|\mathbb{M}^*|}}{\sqrt{\pi N_0 Q}(\underline{\epsilon}\eta_N)} \cdot \exp \left( -\frac{N_0 Q(\underline{\epsilon}\eta_N)^2}{8\bar{c}\bar{s}|\mathbb{M}^*|} \right) + 2C\sigma \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned}$$

Hence we have:

$$\begin{aligned} & \mathbb{P} \left\{ \max_{\mathbf{z} \in \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}})|/\eta_N \leq 2\underline{\epsilon} \right\} \\ & \geq \mathbb{P} \{ \tilde{\mathbf{z}} \in \mathcal{T}_1 \} - \underbrace{4|\mathcal{T}_1||\mathcal{T}'| \frac{\sqrt{2\bar{c}\bar{s}|\mathbb{M}^*|}}{\sqrt{\pi N_0 Q}(\underline{\epsilon}\eta_N)} \cdot \exp \left( -\frac{N_0 Q(\underline{\epsilon}\eta_N)^2}{8\bar{c}\bar{s}|\mathbb{M}^*|} \right)}_{\text{Term I}} \\ & \quad - \underbrace{4|\mathcal{T}_1||\mathcal{T}'| 2C\sigma \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}}}_{\text{Term II}}. \end{aligned}$$

Again, by the conditions, we can show

$$\begin{aligned} & \mathbb{P} \left\{ \max_{\mathbf{z} \in \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}})| / \eta_N \leq 2\epsilon \right\} \\ & \geq \mathbb{P} \{ \tilde{\mathbf{z}} \in \mathcal{T}_1 \} - C|\mathcal{T}_1||\mathcal{T}'| \left\{ \sqrt{\frac{\bar{c}\bar{s}|\mathbb{M}^*|}{N^{1+2\delta_2}}} \exp \left\{ -\frac{C'N^{1+2\delta_2}}{\bar{c}\bar{s}|\mathbb{M}^*|} \right\} + \sigma \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned}$$

Applying (S29) we know that

$$\begin{aligned} & \mathbb{P} \{ \tilde{\mathbf{z}}_h \in \mathcal{T}_1 \} \\ & \geq 1 - C|\mathcal{T}'||\mathcal{T}_1| \left\{ \sqrt{\frac{\bar{c}\bar{s}|\mathbb{M}^*|}{N^{1+2\delta_2}}} \exp \left( -\frac{C'N^{1+2\delta_2}}{\bar{c}\bar{s}|\mathbb{M}^*|} \right) + \sigma \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned}$$

**Aggregating parts.** Aggregating all the results above, we can show that, when  $N$  is large enough, i.e.,  $N > n(\delta_1, \delta_2, \delta_3)$ ,

$$\begin{aligned} & \mathbb{P} \left\{ \max_{\mathbf{z} \in \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}})| \leq \epsilon \eta_N, \min_{\mathbf{z} \notin \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}})| \geq \bar{\epsilon} \eta_N \right\} \\ & \geq 1 - C|\mathcal{T}'||\mathcal{T}_1| \left\{ \sqrt{\frac{\bar{c}\bar{s}|\mathbb{M}^*|}{N^{1+2\delta_2}}} \exp \left( -\frac{C'N^{1+2\delta_2}}{\bar{c}\bar{s}|\mathbb{M}^*|} \right) + \sigma \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned}$$

**Bounding the factor level combination selection probability.** From the formulated procedure, we have

$$\begin{aligned} & \mathbb{P} \left\{ \hat{\mathcal{T}}_1 = \mathcal{T}_1 \right\} \\ & = \mathbb{P} \left\{ |\hat{Y}_R(\mathbf{z}) - \max_{\mathbf{z} \in \mathcal{T}'} \hat{Y}_R(\mathbf{z})| \leq \epsilon \eta_N, \text{ for } \mathbf{z} \in \mathcal{T}_1; \right. \\ & \quad \left. |\hat{Y}_R(\mathbf{z}) - \max_{\mathbf{z} \in \mathcal{T}'} \hat{Y}_R(\mathbf{z})| > \epsilon \eta_N, \text{ for } \mathbf{z} \notin \mathcal{T}_1 \right\} \\ & \geq \mathbb{P} \left\{ |\hat{Y}_R^*(\mathbf{z}) - \max_{\mathbf{z} \in \mathcal{T}'} \hat{Y}_R^*(\mathbf{z})| \leq \epsilon \eta_N, \text{ for } \mathbf{z} \in \mathcal{T}_1; \right. \\ & \quad \left. |\hat{Y}_R^*(\mathbf{z}) - \max_{\mathbf{z} \in \mathcal{T}'} \hat{Y}_R^*(\mathbf{z})| > \epsilon \eta_N, \text{ for } \mathbf{z} \notin \mathcal{T}_1 \right\} - \mathbb{P} \{ \hat{\mathbb{M}} \neq \mathbb{M}^* \} \\ & = \mathbb{P} \left\{ |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}}_h)| \leq \epsilon \eta_N, \text{ for } \mathbf{z} \in \mathcal{T}_1; \right. \\ & \quad \left. |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}}_h)| > \epsilon \eta_N, \text{ for } \mathbf{z} \notin \mathcal{T}_1 \right\} - \mathbb{P} \{ \hat{\mathbb{M}} \neq \mathbb{M}^* \} \end{aligned}$$

(where we introduce random index  $\tilde{\mathbf{z}}_h$  to record the position that achieves maximum)

$$\geq \mathbb{P} \left\{ |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}}_h)| \leq \epsilon \eta_N, \text{ for } \mathbf{z} \in \mathcal{T}_1; \right.$$

$$\begin{aligned}
& \left. |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}}_h)| > \bar{\epsilon}\eta_N, \text{ for } \mathbf{z} \notin \mathcal{T}_1 \right\} - \mathbb{P}\{\hat{\mathbf{M}} \neq \mathbf{M}^*\} \\
& \text{(simply using the fact that } \bar{\epsilon} > \underline{\epsilon}\text{)} \\
& = \mathbb{P}\left\{ \max_{\mathbf{z} \in \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}}_h)| \leq \underline{\epsilon}\eta_N; \min_{\mathbf{z} \notin \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}}_h)| > \bar{\epsilon}\eta_N \right\} \\
& \quad - \mathbb{P}\{\hat{\mathbf{M}} \neq \mathbf{M}^*\} \\
& \geq 1 - \sum_{h=1}^{H_0} \left( 1 - \mathbb{P}\left\{ \max_{\mathbf{z} \in \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}}_h)| \leq \underline{\epsilon}\eta_N; \min_{\mathbf{z} \notin \mathcal{T}_1} |\hat{Y}_R^*(\mathbf{z}) - \hat{Y}_R^*(\tilde{\mathbf{z}}_h)| > \bar{\epsilon}\eta_N \right\} \right) \\
& \quad - \mathbb{P}\{\hat{\mathbf{M}} \neq \mathbf{M}^*\} \\
& \geq 1 - \mathbb{P}\{\hat{\mathbf{M}} \neq \mathbf{M}^*\} \\
& \quad - C|\mathcal{T}'||\mathcal{T}_1| \left\{ \sqrt{\frac{\bar{c}\bar{s}|\mathbf{M}^*|}{N^{1+2\delta_2}}} \exp\left(-\frac{C'N^{1+2\delta_2}}{\bar{c}\bar{s}|\mathbf{M}^*|}\right) + \sigma \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \cdot \sqrt{\frac{|\mathbf{M}^*|}{N_0 Q}} \right\}.
\end{aligned}$$

□

Lemma S7 suggests that, under the conditions assumed in Theorem 5, we select the first tie set consistently as  $N \rightarrow \infty$ . Now Theorem 5 is a direct result of Lemma S5 and Lemma S7.