

Forward factorial screening with statistical inference

Abstract

Modern factorial designs may involve a large number of factors and limited number of replications in each treatment level, which leads to a large volume of causal parameters and invalidates classical inference results. We propose a forward factorial screening procedure, which identify salient factorial effects in a forward fashion following the natural heredity structure of factorial experiments. Computationally, the proposed procedure is simple to implement and is compatible with many existing model selection methods. Theoretically, we prove the forward screening procedure selects the non-zero effects consistently in asymptotic regimes. Furthermore, by combining the proposed forward screening framework with factor-based regressions, we construct point estimators for general causal effects in factorial experiments and propose variance estimation schemes, which can incorporate information from multiple treatment levels and improve asymptotic efficiency for certain choices of target effects. Our numerical study illustrates the finite sample performance of the proposed screening and inferential frameworks.

Keywords: causal inference; potential outcomes; design-based inference; heredity; interaction; variable selection

1 Introduction

1.1 Motivation

Factorial designs have been widely utilized in many fields, ranging from agricultural, industrial, to biomedical sciences (Wu and Hamada, 2011; Zhao and Ding, 2021b; Egami and Imai, 2018). The popularity of factorial designs partially ascribed to their ability to simultaneously accommodate multiple factors and provide informative assessments for the magnitudes of main causal effects and their interactions.

Statistical inference in classical factorial experiments has been well understood when the number of factors, denoted as K , is a fixed number. However, in modern factorial designs involving a large number of factors, it is debatable if the fixed K assumption accurately reflects the data collection

mechanism, and classical statistical inference in factorial experiments faces two key challenges. On the one hand, in an ocean of causal parameters (e.g., when $K = 10$, we have $2^K - 1 = 1,023$ factorial effects), practitioners may hope to target on a few salient ones. The choice of these target parameters need to be carefully made not only to guarantee the validity of statistical inference but also to ensure the selected parameters respect the design principles of factorial experiments. On the other hand, it is unknown how the selection step affects estimation and inference for general causal effects.

To address these challenges, we aim to leverage three unique structural principles first capitalized by Wu and Hamada (2011), and propose a forward factor screening procedure that screens out noisy factors in a principled fashion. These three principles are:

- *Effect Hierarchy Principle.* (i) Lower-order effects are more likely to be important than higher-order effects. (ii) Effects of the same order are equally likely to be important.
- *Effect Sparsity Principle.* The number of relatively important effects in a factorial experiment is small.
- *Effect Heredity Principle.* In order for an interaction to be significant, at least one of its parent main effects should be significant.

Forward screening can naturally incorporate these principles and produce a few interpretable parameters that respect the structure of factorial experiments. This procedure will allow people to work on an inherent parameter space with lower dimension, combine information across multiple treatment levels and mitigate the impact of limited replications due to the large Q setup. Moreover, we can take advantage of such benefits to make inference on many general interesting causal quantities.

1.2 Our contribution

This article makes several contributions. First, we propose a forward factorial screening procedure, which identifies salient factorial effects in a forward fashion following the general principles of factorial experiments. The procedure produces a working model with high interpretability because it fully respects the hierarchical structure of factorial experiments. Besides, combined with factor-based regressions, the proposed forward screening framework is easy to implement and demonstrates high computational efficiency. Theoretically, we also justify the screening consistency property of the forward screening procedure in asymptotic regimes. Second, we did thorough efficiency

comparison for inference using factor based regression with/without effect screening, both from theoretical justification and numerical experimentation. We show that effect screening can bring a great efficiency gain for causal effects defined by sparse contrasts while maintaining desirable performance for general contrast. Moreover, we extend the analysis to comparison of multiple causal parameters. Third, we propose strategies that are practically appealing in marginal cases where the effect sizes are small and perfect screening is hard to achieve. We introduce an under-screening strategy that excludes high order interactions and an over-selection strategy that includes high order interactions following the heredity principle. These two strategies can circumvent the potential challenges of under-powered effect size and deliver valid inference.

1.3 Literature review

In the realm of factorial experiments, the factor-based regression typically serves as a dominant strategy for delivering point estimators and confidence regions, due to its simplicity and flexibility in real-life applications. For example, Dasgupta et al. (2015) extended the classical notion of factorial effects to causal counterparts by introducing potential outcome framework and contrast designs. Zhao and Ding (2021b) studied the use of both saturated and unsaturated linear models for estimating the factorial causal effects. They discussed the parameter specifications of the regression models and justified the commonly used ordinary least squares (OLS) practice from a theoretical perspective. Pashley and Bind (2019) highlighted the desirable property of regression schemes combined with fractional factorial designs when full designs are possible due to constraints on resources such as units or cost. Zhao and Ding (2021a) explores the possibility of incorporating covariate information and applying restricted least squares (RLS) for multiple treatment experimental designs, including factorial studies as a special instance.

A closely related thread of research in factorial designs focuses on variable screening and screening. Powerful variable selection procedures can significantly reduce the complexity of the working model and lead to additional benefits in statistical estimation and inference. In practice pre-screening serves as an appealing scheme for optimizing allocation and utilization of resources. To this end, Wang (2009) introduced forward regression for main effects screening and proves its screening consistency property. Hao and Zhang (2014) further included second-order interactions into the linear model and proposes a two-step procedure for ultra-high dimensional variable screening. Meanwhile, to save resources and build an interpretable model with high prediction power, variable selection or screening must be employed. Haris et al. (2016) considered convex modelling of the factorial effects estimation and introduces strong heredity condition to achieve adaptive selection. Hao et al.

(2018) utilized a regularization scheme to tackle the curse of high dimensionality and perform valid variable screening with quadratic regression. Other works including Lim and Hastie (2015); Bien et al. (2013), proposed procedures for learning interactions based on ℓ_1 regularized least squares based on a purely algorithmic perspective without statistical guarantee.

1.4 Notations

We adopt the following notations throughout the manuscript. 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 a matrix V , define $\varrho_{\max}(V)$, $\varrho_{\min}(V)$ as the largest and smallest eigenvalue, respectively, and use $\kappa(V) = \varrho_{\max}(V)/\varrho_{\min}(V)$ to denote its condition number. For two positive semi-definite matrices V_1 and V_2 , we use $V_1 \preceq V_2$ to indicate V_1 is dominated by V_2 .

For analyzing factorial effects, we work with different levels of sets. For an integer K , let $[K] = \{1, \dots, K\}$. We use \mathcal{K} in calligraphic to denote a subset of $[K]$. For subsets of the power set of $[K]$, we use blackboard bold font for presentation,. For example, we denote $\mathbb{M} \subset \{\mathcal{K} \mid \mathcal{K} \subset [K]\}$ and denote the power set of $[K]$ as \mathbb{K} .

2 Factorial experiment: A general setup

In this section, we introduce four key components and their corresponding notations in factorial experiments, including the general framework of factorial experiments, potential outcomes and treatment assignment mechanisms, the definition of various causal effects, and factor based regressions. To demonstrate how our setup can be conveniently integrated in factorial experiments, we end the section with a concrete example (Example 1) with three factors.

First, we introduce the general framework of a 2^K factorial experimental design, for some $K \geq 2$. In this design, we have K factors each with binary levels, denoted by $z_k \in \{0, 1\}$, $k = 1, \dots, K$. Let the vector $\mathbf{z}_{\mathcal{K}} = (z_k)_{k \in \mathcal{K}}$ denotes the treatment combining the factors in the set $\mathcal{K} \subset [K] = \{1, \dots, K\}$; let \mathbf{z} denotes the treatment combining all K factors (we drop the subscript $[K]$ in $\mathbf{z}_{[K]}$ for simplicity). The K factors in total define a collection of $Q = 2^K$ treatment combinations, and we denote the set of of these treatment combinations as

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

Second, we introduce potential outcomes and treatment assignment mechanisms in a factorial

experiment. Suppose N units are enrolled in a factorial experiment, with $N(\mathbf{z})$ units assigned the treatment $\mathbf{z} \in \mathcal{T}$. Unit i has potential outcome $Y_i(\mathbf{z})$ if assigned to the treatment \mathbf{z} . As each unit has Q potential treatment combinations, we aggregate these Q potential outcomes for unit i into a vector $\mathbf{Y}_i = \{Y_i(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}$ using lexicographic order. The potential outcomes have mean $\bar{Y} = (\bar{Y}(\mathbf{z}))_{\mathbf{z} \in \mathcal{T}}$ and covariance $S = (S(\mathbf{z}, \mathbf{z}'))_{\mathbf{z}, \mathbf{z}' \in \mathcal{T}}$ with elements defined as follows:

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

Let Z_i encode the observed treatment for unit i , and the probability of observing this treatment combination is

$$\mathbb{P}\{Z_i = \mathbf{z}, i \in [N], \mathbf{z} \in \mathcal{T}\} = \frac{\prod_{\mathbf{z} \in \mathcal{T}} (N(\mathbf{z}))!}{N!}.$$

The observed outcome is $Y_i = Y_i(Z_i) = \sum_{\mathbf{z} \in \mathcal{T}} Y_i(\mathbf{z}) \mathbf{1}\{Z_i = \mathbf{z}\}$. Furthermore, we use N_i to denote the number of units for the treatment group in which the i -th individual is assigned.

Third, we formally define treatment effects frequently used in factorial experiments, and our definition is in line with the existing literature Dasgupta et al. (2015); Zhao and Ding (2021b); Wu and Hamada (2011). For a subset $\mathcal{K} \subset [K] = \{1, \dots, K\}$ of the K factors, we use the following “contrast vector” notation for better defining the causal effects later. To start with, we define the main causal effect for factor k . To do so, 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 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 interaction causal effect of multiple factors in the case with $|\mathcal{K}| \geq 2$, the contrast vector $g_{\mathcal{K}} \in \mathbb{R}^{Q \times 1}$ can be similarly defined 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 $g_{\emptyset} = \mathbf{1}_Q$.

Equipped with the contrast vector notation, we are ready to introduce the main effects and the k -way interaction causal effects in factorial experiments. More concretely, the main causal effect of a single factor and the k -way interaction causal effect of multiple factors ($k \geq 2$) are defined by the inner product of the contrast vector $g_{\mathcal{K}}$ and the averaged potential outcome \bar{Y} , that is

$$\tau_{\mathcal{K}} = Q^{-1} \cdot g_{\mathcal{K}}^\top \bar{Y}, \quad \mathcal{K} \subset [K]. \quad (2.3)$$

We call the effect $\tau_{\mathcal{K}}$ a *parent* of $\tau_{\mathcal{K}'}$ if $\mathcal{K} \subset \mathcal{K}'$ and $|\mathcal{K}| = |\mathcal{K}'| - 1$.

To more conveniently summarize these causal parameters of interest using synthesized notation, we stack the contrast vectors into an orthogonal matrix $G \in \mathbb{R}^{Q \times Q}$ (i.e., $G^\top G = Q \cdot I_Q$), in which the first column is each column is g_\emptyset and the rest of the columns are contrast vectors $g_{\mathcal{K}}$'s. We refer to matrix G as the contrast matrix. The entire collection of causal parameters in factorial experiments can be more handily written as

$$\tau = (\tau_{\mathcal{K}})_{\mathcal{K} \subset [K]} = Q^{-1} \cdot G^\top \bar{Y}, \text{ where } G = (g_{\mathcal{K}})_{\mathcal{K} \subset [K]}. \quad (2.4)$$

As an additional notation which can be useful for our future derivation, we use $\tau_\emptyset = Q^{-1} g_\emptyset^\top \bar{Y}$ to capture the total average of potential outcomes. We realize that the above notations can be rather abstract, and we have provided an illustrative example in Example 1 with three factors $K = 3$.

Fourth, to estimate the above causal parameters, we introduced two types of factor-based regression commonly adopted in the literature Wu and Hamada (2011); Dasgupta et al. (2015); Zhao and Ding (2021a). The first type is the *saturated regression* in which we regress the observed outcome Y_i on the regressor t_i . Here, t_i is one row vector in the contrast matrix G indexed by the unit i 's observed treatment status $Z_i = (z_{i,1}, \dots, z_{i,K})^\top$; or equivalently, we can write $t_i = \{g_{\mathcal{K}}(Z_i)\}_{\mathcal{K} \subset [K]} \in \mathbb{R}^{Q \times 1}$, where $g_{\mathcal{K}}(Z_i) = \prod_{k \in \mathcal{K}} g_{\{k\}}(z_k)$ is defined in (2.2). The second type of the *unsaturated regression* in which we regress the observed outcome Y_i on a subvector of t_i . We denote the subvector of t_i as $t_{i,\mathbb{M}}$, and \mathbb{M} is the subset of the power set of all factor combinations (i.e., $\mathbb{M} \subset \mathbb{K} = \{\mathcal{K} \mid \mathcal{K} \subset [K]\}$). Throughout this manuscript, we refer to any subset of \mathbb{K} , \mathbb{M} , as a *working model*. The estimated coefficients of the unsaturated regression is denoted as $\hat{\tau}(\mathbb{M})$, and the collection of true causal effects in the set \mathbb{M} is denoted as $\tau(\mathbb{M})$, i.e., $\tau(\mathbb{M}) = \{\tau_{\mathcal{K}}\}_{\mathcal{K} \in \mathbb{M}}$. Moreover, we use $G(\cdot, \mathbb{M})$ to denote the columns in G indexed by \mathbb{M} .

In what follows, we provide a concrete example of our notations with three factors in the uniform factorial experiment. By uniform, we mean that each treatment level contains the same number of units.

Example 1 (An explanation in the uniform 2^3 factorial design). *Suppose we have three binary factors z_1, z_2 and z_3 . This three factors amount to 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*

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

There are $N = \sum_{z_1, z_2, z_3} N(z_1 z_2 z_3) = 2^3 N_0$ units from this design, where $N(z_1 z_2 z_3) = N_0$ denotes the group size under treatment $(z_1 z_2 z_3)$ under the uniform design.

Each unit i corresponds to 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 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

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

We observe the pair (Y_i, Z_i) for the 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 covariate t_i can be constructed from Z_i :

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

For instance, when $Z_i = (101)$, t_i corresponds to the row (101) of the contrast matrix G . Then, a saturated regression is to regress Y_i on t_i :

$$Y_i \sim t_i.$$

For the unsaturated regression, if we only include indices \emptyset (the intercept), $\{1\}, \{1, 2\}, \{1, 3\}, \{1, 2, 3\}$ in our regression, we can form the working model $\mathbb{M} = \{\emptyset, \{1\}, \{1, 2\}, \{1, 3\}, \{1, 2, 3\}\}$ and perform

$$Y_i \sim t_{i,\mathbb{M}}, \text{ where } t_{i,\mathbb{M}} = \left[1, g_{\{1\}}(Z_i), g_{\{1,3\}}(Z_i), g_{\{1,2\}}(Z_i), g_{\{1,2,3\}}(Z_i) \right].$$

3 Forward screening in factorial experiments

In factorial experiments, when the number of factors K is small (hence small Q), classical statistical inference for factorial effects has been well studied Zhao and Ding (2021b). While modern factorial designs often involve a rather large number of factors, we propose a principled procedure to carefully decide an unsaturated working model \mathbb{M} and scientifically screen out potential noisy factors, followed by an investigation of its statistical property.

3.1 Proposed procedure

In this section, we introduce a principle 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. The forward screening procedure is summarized in Algorithm 1 below:

Algorithm 1: Forward factorial screening

Input: Factorial data $\{(Y_i, Z_i)\}_{i=1}^N$; predetermined integer D ; initial working model $\widehat{\mathbb{M}} = \{\emptyset\}$; significance level $\{\alpha_d\}_{d=1}^D$.

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 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 defined in (3.5) and (3.6), and denote the screened working model as $\widehat{\mathbb{M}}'$.
- 5 Run the unsaturated factor-based regression on the working model $\widehat{\mathbb{M}}'$:

$$Y_i \sim t_{i, \widehat{\mathbb{M}}'}, \text{ with weights } w_i = N/N_i.$$
- 6 Obtain coefficients $\widehat{\tau}(\widehat{\mathbb{M}}')$ and robust covariance estimation:

$$\widehat{\Sigma}(\widehat{\mathbb{M}}') = \frac{1}{Q^2} G(\cdot, \widehat{\mathbb{M}}')^\top \text{Diag} \left\{ N(\mathbf{z})^{-1} \widehat{S}(\mathbf{z}, \mathbf{z}) \right\} G(\cdot, \widehat{\mathbb{M}}').$$
- 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-test using the above $\widehat{\tau}_{\mathcal{K}}(\widehat{\mathbb{M}}')$ and $\widehat{\sigma}_{\mathcal{K}}(\widehat{\mathbb{M}}')$ under 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** $\widehat{\mathbb{M}}$

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

First, our algorithm adheres to the “Effect 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.

Second, our algorithm operates under the “Effect Sparsity Principle” as it removes potentially unimportant effects using marginal t-test with Bonferroni correction (Step 8). This step induces a sparse working model and helps us to identify essential factorial effects, echoing the “Effect Sparsity Principle.”

Third, our algorithm incorporates the “Effect Heredity Principle” as it screens out the interaction effects, 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). Here, by strong and weak heredity (Hao and Zhang, 2014; Lim and Hastie, 2015), mathematically speaking, we mean that:

- Weak heredity: remove all the d -way interaction term indexed by \mathcal{K} from $\widehat{\mathbb{M}}'$ if

$$\mathcal{K}' \notin \widehat{\mathbb{M}}' \text{ for all } \mathcal{K}' \subset \mathcal{K}, |\mathcal{K}'| = |\mathcal{K}| - 1. \quad (3.5)$$

- Strong heredity: remove all the d -way interaction term indexed by \mathcal{K} from $\widehat{\mathbb{M}}'$ if

$$\mathcal{K}' \notin \widehat{\mathbb{M}}' \text{ for some } \mathcal{K}' \subset \mathcal{K}, |\mathcal{K}'| = |\mathcal{K}| - 1. \quad (3.6)$$

Lastly, we note that our forward screening procedure enhances the interpretability of the selected working model by iterating between the “Sparsity-screening” step (shorted as 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 (shorted as 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,

$$\widehat{\mathbb{M}}_1 \xrightarrow{\mathbf{H}} \cdots \xrightarrow{\widehat{\mathbf{S}}} \widehat{\mathbb{M}}_{d-1} \xrightarrow{\mathbf{H}} \widehat{\mathbb{M}}_{d,+} \xrightarrow{\widehat{\mathbf{S}}} \widehat{\mathbb{M}}_d \rightarrow \cdots \xrightarrow{\widehat{\mathbf{S}}} \widehat{\mathbb{M}}_D. \quad (3.7)$$

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

3.2 Forward screening consistency

We are now ready to analyze the screening consistency property of Algorithm 1. To avoid redundancy, we focus on the algorithm that carries out the S-step via the marginal t-test. We shall show

that the proposed algorithm selects the “true working model” up to level D with probability tending to one as the sample size goes to infinity. Here, the true working model at level $k \in [K]$, denoted as \mathbb{M}_k^* , is the collection of \mathcal{K} ’s where $|\mathcal{K}| = k$ and $\tau_{\mathcal{K}} \neq 0$. While the key technique in our proof is to utilize Berry-Esseen bounds to analyze the statistical property of moment estimators for factorial effects, we defer interested readers for Supplementary Materials for details.

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

Condition 1 (Nearly uniform design). *There exists an positive integer $N_0 > 0$ and absolute constants $\underline{c} \leq \bar{c}$, such that*

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

This condition generalizes the classical assumption in the fixed Q regime (meaning that Q is fixed and N tends to infinity)

$$Q \text{ is fixed, and } N(\mathbf{z})/N \rightarrow e(\mathbf{z}) \in (0, 1), \text{ as } N \rightarrow \infty, \quad (3.8)$$

where $e(\mathbf{z})$ is a positive constant.

Next, we quantify the order of the size of the true effects $\tau_{\mathcal{K}}$ ’s and the tuning parameters α_d ’s adopted in Bonferroni correction, and we allow them to change with the sample size N :

Condition 2 (Order of parameters). *The true parameters and tuning parameters have the following order:*

- (i) *True parameter: $|\tau_{\mathcal{K}}| = \Theta(N^\delta)$ for some $-1/2 < \delta \leq 0$ and all $\mathcal{K} \in \mathbb{M}^*$.*
- (ii) *Tuning parameter in Bonferroni correction: $\alpha_d = \Theta(N^{-\delta'})$ for all $d \in [D]$ with some $\delta' > 0$.*

Condition 2 specifies the allowable range of the true factorial effects. The tuning parameter α_d converges to zero, which ensures no Type I error asymptotically. Wasserman and Roeder (2009) (Theorem 4.1 and 4.2) assumed similar conditions in high-dimensional model selection settings for linear models.

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

Condition 3 (Regularity conditions on $Y_i(\mathbf{z})$ ’s). *The potential outcomes satisfy*

- (i) *Nondegenerate correlation matrix. Let V^* be the correlation matrix of \hat{Y} . There exists a $\sigma > 0$, such that*

$$\varrho_{\min}(V^*)/\varrho_{\max}(V^*) \geq \sigma^2. \quad (3.9)$$

(ii) Bounded fourth central moments. *There exists a universal constant $\Delta > 0$ such that*

$$\max_{z \in [Q]} \frac{1}{N} \sum_{i=1}^N \{Y_i(z) - \bar{Y}(z)\}^4 \leq \Delta^4. \quad (3.10)$$

(iii) Bounded standardization scales. *Assume M_N is bounded by a constant $M > 0$, where M_N is:*

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

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.

Lastly, we assume the following structural condition 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 $K' \subset K$, $|K'| = |K| - 1$ such that $\tau_{K'} \neq 0$.*
- *Strong heredity: $\tau_K \neq 0$ only if for all $K' \subset K$, $|K'| = |K| - 1$, $\tau_{K'} \neq 0$.*

We then present the screening consistency property of our algorithm:

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

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

4 Inference under perfect screening

In this section, we start with presenting statistical inference framework (including causal parameters, estimators, and confidence interval constructions) when the perfect screening property established in (1) is satisfied. We then study the theoretical properties of the proposed framework in the second subsection.

4.1 Statistical inference

In factorial experiments, a causal parameter of interest can be expressed as a weighted combination of average potential outcomes, that is

$$\gamma = \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z}) \bar{Y}(\mathbf{z}) \triangleq \mathbf{f}^\top \bar{Y}, \quad (4.11)$$

where $\mathbf{f} = \{\mathbf{f}(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}$ 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). When we have multiple causal parameters defined through different weighting vectors, we later present a method to handle the multiple comparison issues in this context.

Without factor screening, a well-studied plug-in estimator of γ in the existing literature is to replace with its sample analogue (Li and Ding, 2017; Zhao and Ding, 2021b; Shi and Ding, 2022):

$$\hat{\gamma} = \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z}) \hat{Y}(\mathbf{z}) \triangleq \mathbf{f}^\top \hat{Y}, \quad (4.12)$$

where $\hat{Y}(\mathbf{z}) = \{N(\mathbf{z})\}^{-1} \sum_{Z_i=\mathbf{z}} Y_i$. When $N(\mathbf{z}) \geq 2$, its variance can be estimated by:

$$\hat{v}^2 = \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z})^2 N(\mathbf{z})^{-1} \hat{S}(\mathbf{z}, \mathbf{z}) \triangleq \mathbf{f}^\top \hat{V}_{\hat{Y}} \mathbf{f}, \quad (4.13)$$

where $\hat{S}(\mathbf{z}, \mathbf{z}) = \{N(\mathbf{z}) - 1\}^{-1} \sum_{Z_i=\mathbf{z}} (Y_i - \hat{Y}(\mathbf{z}))^2$, and $\hat{V}_{\hat{Y}} = \text{Diag} \left\{ N(\mathbf{z})^{-1} \hat{S}(\mathbf{z}, \mathbf{z}) \right\}_{\mathbf{z} \in \mathcal{T}}$.

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 restricted least squares

$$\hat{Y}_R = \arg \min \left\{ \|\hat{Y} - \mu\|_2^2 : G(\cdot, \hat{\mathbb{M}}^c)^\top \mu = 0, \mu \in \mathbb{R}^Q \right\} = Q^{-1} G(\cdot, \hat{\mathbb{M}}) G(\cdot, \hat{\mathbb{M}})^\top \hat{Y}. \quad (4.14)$$

The closed form solution is obtained by exploiting the orthogonality property of G . We provide a justification of this result in Lemma 6. Under perfect screening, \hat{Y}_R is also an unbiased estimator for \bar{Y} . This allows us to construct a restricted least squares (RLS) based estimator of γ with variance estimation:

$$\hat{\gamma}_R = \mathbf{f}^\top \hat{Y}_R \triangleq \mathbf{f}[\hat{\mathbb{M}}]^\top \hat{Y}, \quad \text{and} \quad \hat{v}_R^2 = \mathbf{f}[\hat{\mathbb{M}}]^\top \hat{V}_{\hat{Y}} \mathbf{f}[\hat{\mathbb{M}}], \quad (4.15)$$

where $\mathbf{f}[\hat{\mathbb{M}}] = Q^{-1} G(\cdot, \hat{\mathbb{M}}) G(\cdot, \hat{\mathbb{M}})^\top \mathbf{f}$. Leverage these two estimators and Theorem 2 to be presented, we are able to construct a level- $(1 - \alpha)$ confidence interval for γ :

$$\left[\hat{\gamma}_R \pm z_{1-\alpha/2} \times \hat{v}_R \right], \quad (4.16)$$

where $z_{1-\alpha/2}$ is $(1 - \alpha/2)$ th quantile of a standard normal distribution.

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 $\mathbf{f}^* = \mathbf{f}[\mathbb{M}^*] = Q^{-1}G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \mathbf{f}$. Given \mathbb{M}^* is the true sparse model, we have $(\mathbf{f}^*)^\top \bar{Y} = \mathbf{f}^\top \bar{Y}$, for all $\mathbf{f} \in \mathbb{R}^Q$. We make the following condition for \mathbf{f}^* :

Condition 5 (Condition on \mathbf{f}^*). *The weighting vector $\mathbf{f}^* = Q^{-1}G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \mathbf{f}$ satisfies*

$$N_0^{-1/2} \cdot \frac{\|\mathbf{f}^*\|_\infty}{\|\mathbf{f}^*\|_2} \longrightarrow 0. \quad (4.17)$$

Because $N_0Q = \Theta(N)$ under Condition 1 and the ratio $\|\mathbf{f}^*\|_\infty/\|\mathbf{f}^*\|_2$ is upper bounded by (see Lemma 1 for a rigorous statement):

$$N_0^{-1/2} \cdot \frac{\|\mathbf{f}^*\|_\infty}{\|\mathbf{f}^*\|_2} \leq \left(\frac{|\mathbb{M}^*|}{N_0Q} \right)^{1/2}, \quad (4.18)$$

Condition 5 says that the cardinality of the true working model $|\mathbb{M}^*|$ is asymptotically much smaller compared to the sample size N .

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). *Under Conditions 1-5, we have, as the sample size N tends to infinity,*

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

where $v_R^2 = \mathbf{f}^{*\top} V_{\hat{Y}} \mathbf{f}^*$ and $V_{\hat{Y}} = \text{Diag}\{N(\mathbf{z})^{-1}S(q, q)\} - N^{-1}S$. Furthermore, without loss of generality, assume $\|\mathbf{f}^*\|_\infty = O(Q^{-1})$, the variance estimator \hat{v}_R^2 is consistent and robust:

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

where $v_{R,\text{lim}}^2 = \mathbf{f}^{*\top} \text{Diag}\{N(\mathbf{z})^{-1}S(\mathbf{z}, \mathbf{z})\} \mathbf{f}^*$ is the robust asymptotic variance of $\hat{\gamma}_R$.

The above theorem guarantees that the proposed confidence interval in (4.16) for γ attains the nominal coverage probability asymptotically. We note that the above result applies to any model selection procedure with the oracle property (Fan and Li, 2001) and does not restrict to the marginal t-test with Bonferroni correction adopted in Algorithm 1. 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). *To establish the central limit theorem of the classical plug-un estimator $\hat{\gamma}$ (4.12), Shi and Ding (2022) provided the following condition similar to our Condition 5*

$$N_0^{-1/2} \cdot \frac{\|\mathbf{f}\|_\infty}{\|\mathbf{f}\|_2} \rightarrow 0. \quad (4.19)$$

In a case where N_0 is small (see Condition 1 for the definition of N_0) and both \mathbf{f} and \mathbb{M}^ are sparse, Condition (4.19) fails while Condition (4.17) continues to hold, suggesting that our proposed estimator $\hat{\gamma}_R$ can be more robust than $\hat{\gamma}$ in practice.*

To elaborate the benefits of conducting forward factorial screening in terms of asymptotic efficiency, we make a simple comparison in Proposition 1 based on a concrete choice of \mathbf{f} . Based on the result established in Shi and Ding (2022), the plug-in estimator $\hat{\gamma}$ satisfies

$$\frac{\hat{\gamma} - \gamma}{v} \rightsquigarrow \mathcal{N}(0, 1), \quad \text{where } v^2 = \mathbf{f}^\top V_{\hat{\gamma}} \mathbf{f},$$

under appropriate conditions.

Proposition 1 (Asymptotic relative efficiency comparison between $\hat{\gamma}$ and $\hat{\gamma}_R$ with sparse \mathbf{f}). *Define the condition number of $V_{\hat{\gamma}}$ as $\kappa(V_{\hat{\gamma}})$, and let s^* denote the number of nonzero elements in \mathbf{f} . Under the condition that both $\hat{\gamma}$ and $\hat{\gamma}_R$ converge to a normal distribution as the sample size tends to infinity, the asymptotic relative efficiency between $\hat{\gamma}$ and $\hat{\gamma}_R$ is upper bounded by*

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

provided that $s^ |\mathbb{M}^*| < Q$.*

Now we add some interpretation for Proposition 1. The condition number $\kappa(V_{\hat{\gamma}})$ 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} . The higher the variability is, the larger the condition number becomes. When the variance of $\hat{Y}(\mathbf{z})$ does not change with its group membership \mathbf{z} , the condition number $\kappa(V_{\hat{\gamma}}) = 1$. This suggests that the proposed RLS-based estimator $\hat{\gamma}_R$ is more efficient than the plug-in estimator $\hat{\gamma}$. When the variance of $\hat{Y}(\mathbf{z})$ does change with its group membership \mathbf{z} , but the variability of such changes is limited in the sense that $\kappa(V_{\hat{\gamma}}) < Q/s^* |\mathbb{M}^*|$, the proposed estimator is also more efficient than $\hat{\gamma}$. The above result can be extend to compare the length of the confidence intervals as well; see Proposition 2 in the Supplementary Material for details.

5 Inference under imperfect screening

In factorial experiments, the perfect screening property can be too much to hope for, this is partially because the factorial experimental data are driven by the Effect Hierarchy Principle (Wu and Hamada, 2011). Implied by this principle, the main factorial effects and lower-order factorial effects are more likely to be non-negligible than the higher-order factorial effects. Once the higher factorial effects are marginal, the perfect screening property may no longer hold. Take our selection procedure in Algorithm 1 (marginal t-test with Bonferroni correction) for example, whenever Condition 2(i) is violated, we may lose the perfect screening property established in Theorem 6 with high chance.

In response to this challenge, we propose two alternative strategies along with proposed causal effects estimators (summarized in Figure 1) in the first subsection. We then present their theoretical guarantees in the subsequent section.

5.1 Statistical inference with two alternative strategies

The two proposed strategies are built on the belief that perfect screening is more plausible for selecting the main factorial effects and lower-order factorial effects up to level d^* .

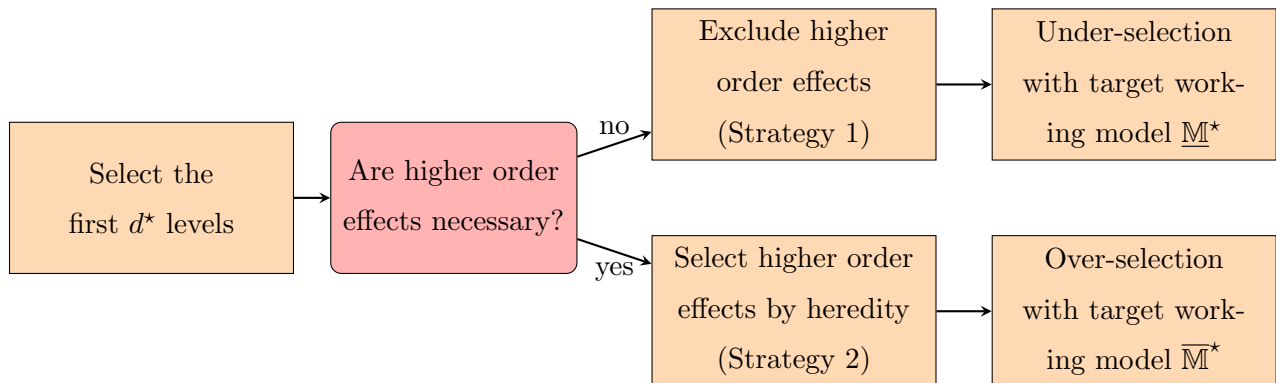


Figure 1: General strategy for factorial screening

For the first strategy, when the higher order factorial effects beyond level d^* are believed to be negligible, we may stop our forward screening procedure in Algorithm 1 at $d = d^*$ (instead of $d = D$). Similar strategies have been applied in the existing hierarchical model selection literature. For example, Hao and Zhang (2014); Lim and Hastie (2015); Yuan et al. (2007), among others, propose screening procedures that only include main effect selection and two-way interaction screening, assuming that the higher order interactions (beyond level two) are negligible. As this strategy

targets a true working model $\underline{\mathbb{M}}^*$ up to level d^* , that is,

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

this strategy leads to an under-selected parsimonious working model. We summarize this strategy in as follows:

Strategy 1 (Under selection by excluding higher 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 $\widehat{\underline{\mathbb{M}}} = \cup_{d=1}^{d^*} \widehat{\mathbb{M}}_d$.*

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

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

For the second strategy, rather than excluding all higher order interactions with negligible effects, we may further leverage the Effect 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 true working model $\overline{\mathbb{M}}^*$ that includes the true model \mathbb{M}^* , that is,

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

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 Effect 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*

$$\widehat{\overline{\mathbb{M}}} = \cup_{d=1}^D \widehat{\overline{\mathbb{M}}}_d; \quad \widehat{\overline{\mathbb{M}}}_d = \text{H}^{(d-d^*)}(\widehat{\mathbb{M}}_{d^*}), d \geq d^* + 1.$$

Here $\text{H}^{(d-d^*)}$ 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 $\gamma = \mathbf{f}^\top \bar{\mathbf{Y}}$ based on the restricted least squares:

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

In the following subsection, we study the theoretical properties of $\widehat{\gamma}_{\text{RO}}$ and $\widehat{\gamma}_{\text{RU}}$ and demonstrate their trade-offs for statistical inference.

5.2 Theoretical properties under imperfect screening

Throughout this section, as we are in a scenario where the higher order factorial effects beyond level d^* are negligible, we work under a relaxed condition of Condition 2 as follows:

Condition 6 (Order of parameters up to level d^*). *The true parameters and tuning parameters have the following order:*

- (i) *True parameter: $|\tau_{\mathcal{K}}| = \Theta(N^\delta)$ for some $-1/2 < \delta \leq 0$ and all $\mathcal{K} \in \underline{\mathbb{M}}^*$.*
- (ii) *Tuning parameter in Bonferroni correction: $\alpha_d = \Theta(N^{-\delta'})$ for all $d \leq d^*$ with some $\delta' > 0$.*

Condition 6 no longer impose any restriction on the order of the parameters beyond level d^* . By Theorem 6, Condition 6 guarantees that Algorithm 1 perfectly screens the first d^* levels of factorial effects in the sense that

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

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 unbiased estimator of γ . 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). *Define $\underline{\mathbf{f}}^* = Q^{-1}G(\cdot, \underline{\mathbb{M}}^*)G(\cdot, \underline{\mathbb{M}}^*)^\top \mathbf{f}$. Assume Conditions 1, 3, 4 and 6. Also assume \mathbf{f} satisfies the following orthogonality condition:*

$$G(\cdot, \mathbb{M}_d^*)^\top \mathbf{f} = 0, \ d^* + 1 \leq d \leq D^*. \quad (5.20)$$

If as N tends to infinity,

$$N_0^{-1/2} \cdot \frac{\|\underline{\mathbf{f}}^*\|_\infty}{\|\underline{\mathbf{f}}^*\|_2} \longrightarrow 0,$$

then

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

where $v_{\text{RU}}^2 = \underline{\mathbf{f}}^{\star\top} V_{\hat{\mathbf{Y}}} \underline{\mathbf{f}}^{\star}$. Furthermore, without loss of generality, assume $\|\underline{\mathbf{f}}^{\star}\|_{\infty} = O(Q^{-1})$, the variance estimator \hat{v}_{RU}^2 is consistent and robust:

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

where $v_{\text{RU,lim}}^2 = \underline{\mathbf{f}}^{\star\top} \text{Diag}\{N(\mathbf{z})^{-1}S(\mathbf{z}, \mathbf{z})\} \underline{\mathbf{f}}^{\star}$ is the robust asymptotic variance of $\hat{\gamma}_{\text{RU}}$.

The orthogonality condition presented in (5.20) restricts the weighting vector \mathbf{f} to be orthogonal to the higher order contrasts, because the higher order interactions are excluded from the model and making inference on a weighted combination of those excluded interactions is infeasible. One example satisfying (5.20) is an arbitrary nonzero linear combination of lower-order contrasts, given by $\mathbf{f} = G(\cup_{d=1}^{d^*} \mathbb{M}_d^{\star}) \mathbf{b}$, for any real \mathbf{b} with dimension $\sum_{d=1}^{d^*} |\mathbb{M}_d^{\star}|$. When \mathbf{b} is one of canonical bases, γ is the corresponding non-zero factorial effects.

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

Theorem 4 (Guarantee for Strategy 2). *Assume Conditions 1, 3, 4 and 6. Define the weighting vector $\overline{\mathbf{f}}^{\star} = Q^{-1}G(\cdot, \overline{\mathbb{M}}^{\star})G(\cdot, \overline{\mathbb{M}}^{\star})^{\top} \mathbf{f}$. If as N tends to infinity,*

$$N_0^{-1/2} \cdot \frac{\|\overline{\mathbf{f}}^{\star}\|_{\infty}}{\|\overline{\mathbf{f}}^{\star}\|_2} \longrightarrow 0,$$

then

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

where $v_{\text{RO}}^2 = \overline{\mathbf{f}}^{\star\top} V_{\hat{\mathbf{Y}}} \overline{\mathbf{f}}^{\star}$. Furthermore, without loss of generality, assume $\|\overline{\mathbf{f}}^{\star}\|_{\infty} = O(Q^{-1})$, the variance estimator \hat{v}_{RO}^2 is consistent and robust:

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

where $v_{\text{RO,lim}}^2 = \overline{\mathbf{f}}^{\star\top} \text{Diag}\{N(\mathbf{z})^{-1}S(\mathbf{z}, \mathbf{z})\} \overline{\mathbf{f}}^{\star}$ is the robust asymptotic variance of $\hat{\gamma}_{\text{RO}}$.

When analyzing Strategy 1 and 2, Algorithm 1 recovers a target 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 lose efficiency because

one might include too 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 homoscedasticity, 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_{RU}^2 and v_{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 derive*

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

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

6 An extension to comparing multiple causal effects

In the previous sections, we consider the problem on making inference on a single factorial causal effect $\gamma = \mathbf{f}^\top \bar{Y}$. We now provide an extension of our framework that delivers valid inference on multiple causal parameters.

In this section, we study the problem of “inference on best causal effects” as an application of forward factorial screening and inference based on restricted least squares. The goal of the problem is to compare a set of candidate causal effects in the form of (4.11) and identify the maximal effect. More concretely, suppose we have a set of effects defined by pre-specified weighting vectors:

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

We hope to perform statistical inference on the maximal effects in Γ :

$$\gamma_{(1)} = \max_{l \in [L]} \gamma_l.$$

As one specific example, if we choose $\{\mathbf{f}_l\}_{l \in [L]} = \{\mathbf{e}(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}$ to be a subset of the canonical bases $\{\mathbf{e}(\mathbf{z})\}_{\mathbf{z} \in \mathcal{T}}$, then our inferential target is

$$\bar{Y}_{(1)} = \max_{\mathbf{z} \in \mathcal{T}'} \bar{Y}(\mathbf{z}). \tag{6.21}$$

In this case, the goal is to identify the treatment combination that demonstrates the “best” performance measured by level of average potential outcome and do inference on the average potential outcome for the selected arms.

The challenge in identifying $\gamma_{(1)}$ lies in that similar effects might not be statistically distinguishable with limited sample size. We introduce the following notion \mathcal{L}_1 to include all effects that achieve or are statistically “close” to $\gamma_{(1)}$:

$$\mathcal{L}_1 = \left\{ l \in [L] \mid |\gamma_l - \gamma_{(1)}| = O(N^{-\delta_3}) \right\}, \text{ for some } \delta_3 > 0. \quad (6.22)$$

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

With the above basic setup, we introduce the following Algorithm 2 for selecting the most effective effects. The algorithm consists of three major components:

- (i) Construct $\hat{\gamma}_l = \mathbf{f}_l^\top \hat{Y}_R$ with feature screening (Step 1-2). These RLS based estimators enjoy great benefits for large Q and small N_0 regimes.
- (ii) Construct $\hat{\mathcal{L}}_1$ as an estimator for \mathcal{L}_1 (Step 3). $\hat{\mathcal{L}}_1$ aggregates the arms whose estimated averages are close to the empirical maximum. We will show that with proper tuning, this procedure recovers \mathcal{L}_1 with high probability.
- (iii) Construct estimators by averaging over $\hat{\mathcal{L}}_1$ (Step 4). The difference in values across arms in \mathcal{L}_1 are mainly due to randomness of the study. Therefore, we average the estimates over the selected $\hat{\mathcal{L}}_1$ to alleviate the impact of randomness and obtain accurate estimates.

We provide some theoretical analysis of Algorithm 2. Define

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

which we refer to as within-group diameter and between-group distance respectively. Condition 7 below characterizes the order of the involved quantities:

Condition 7 (Order of d_h , d_h^* and η_N). *Assume the following scaling of parameters:*

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

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

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

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

Algorithm 2: Select the most effective arms

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 η_N .

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

- 1 Perform effects screening with Algorithm 1 and obtain working model $\widehat{\mathbb{M}}$.
- 2 Obtain RLS based estimates: let \widehat{Y}_R be defined as (4.14), and compute

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

- 3 Select the best factor level combinations:

$$\widehat{\mathcal{L}}_1 = \left\{ l \in [L] \mid |\widehat{\gamma}_l - \max_{l \in [L]} \widehat{\gamma}_l| \leq \eta_N \right\}.$$

- 4 Define

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

Then generate point estimates and variance estimator for the effect size over the selected tie:

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

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

Theorem 5 (Asymptotic results on the estimated effects with screening). *Assume Condition 1, 3 and 7. Further assume that the perfect screening holds with probability tending to one (i.e., $\mathbb{P}\{\widehat{\mathbb{M}} = \mathbb{M}^*\} \rightarrow 1$). Assume $|\mathbb{M}^*| = \Theta(N^{\delta_4})$ for some $\delta_4 \geq 0$ and*

$$N^{-(1+2\delta_2-\delta_4)} \rightarrow 0, \quad L \cdot |\mathcal{L}_1| \cdot N^{-\frac{1-\delta_4}{2}} \rightarrow 0.$$

Then the point estimates are asymptotically jointly normal:

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

where $v_{(1)}^2 = \mathbf{f}_{(1)}^{\star\top} V_Y \mathbf{f}_{(1)}^{\star}$. Moreover, $\hat{v}_{(1)}^2$ is consistent and robust:

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

where $v_{(1),\text{lim}}^2 = \mathbf{f}_{(1)}^{\star\top} \text{Diag} \{N(\mathbf{z})^{-1} S(\mathbf{z}, \mathbf{z})\} \mathbf{f}_{(1)}^{\star}$.

From Theorem 5, with factorial screening, when the size of the true working model is small (say $\delta_4 = 0$) and screening is consistent, $N^{-(1+2\delta_2)}$ always converge to 0. Thus one only needs

$$L|\mathcal{L}_1| (|\mathbb{M}^*|/N)^{1/2} \rightarrow 0.$$

The condition relies on the scaling of N instead of a particular set of $N(\mathbf{z})$'s. In other words, we can incorporate information from other treatment arms to facilitate inference. On the contrary, without screening, one requires Q to be small compared to N or $\{\mathbf{f}_l\}_{l \in [L]}$ are dense, which might not be true in large Q setups and practical scenarios such as (6.21).

7 Simulation

In this section, we use numerical experiments to study the finite sample performance of the proposed forward screening framework and the inferential properties of the RLS based estimator.

7.1 Varying sample size

In the first setup, we evaluate the performance of different screening procedure under varying sample sizes. We set up a 2^8 uniform 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 potential outcomes independently 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} . $\mu(\mathbf{z})$ are generated 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 two-way interactions are zero.

- Higher order interactions: all the higher order interactions $\tau_{\mathcal{K}}$ where $|\mathcal{K}| \geq 3$ are zero.

The above rules guarantee that the factorial effects follow the strong heredity principle and demonstrate a sparse pattern. More details can be found in the R code attached in the support materials.

To perform effect screening, we compare four methods in our study: (i) *Forward Bonferroni*. screening in a forward fashion based on Bonferroni corrected margin t tests; (ii) *Forward Lasso*. screening in a forward fashion based on Lasso; (iii) *Naive Bonferroni*. screening with the full working model based on Bonferroni corrected margin t tests; (iv) *Naive Lasso*. screening with the full working model based on Lasso. For each of the screening methods, we use several criteria to compare the performance of the methods. (i) *Perfect selection probability*. We report the perfect selection probability $\mathbb{P}\{\hat{\mathbb{M}} = \mathbb{M}^*\}$. (ii) *Power of $\hat{\gamma}_{\text{R}}$* . We evaluate the power of the RLS based estimators $\hat{\gamma}_{\text{R}}$ and \hat{v}_{R} for testing a causal effect γ_{target} specified by a sparse vector:

$$\gamma_{\text{target}} = \mathbf{f}^{\top} \bar{\mathbf{Y}}, \quad \mathbf{f}(\mathbf{z}) = \begin{cases} -1, & z_1 = \cdots = z_8 = 0; \\ 1, & z_1 = 0, z_2 = \cdots = z_8 = 1; \\ 0, & \text{for other } \mathbf{z}. \end{cases}$$

Intuitively speaking, γ_{target} reflects the causal effects of setting all the factors as 1 except for the first factor compared to the baseline treatment level. (iii) *Coverage probability of the RLS based confidence interval*. We report the coverage probability for γ_{target} defined above with the RLS based confidence interval (4.16) with confidence level 0.95. The results are average over 500 Monte Carlo runs and are presented in Figure 2.

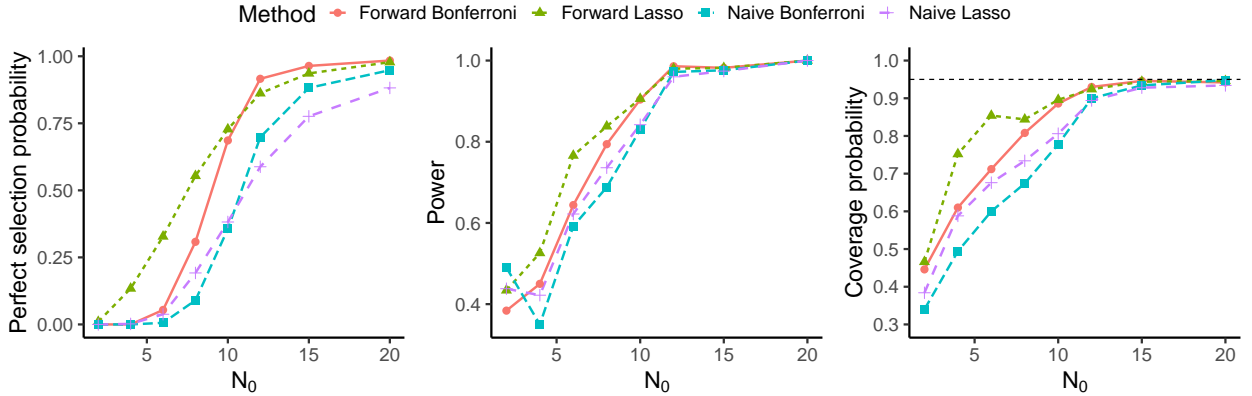


Figure 2: Performance evaluation of effect screening under varying sample size. The left panel presents the probability of perfect model selection. The middle panel shows the power curve for testing $H_0 : \gamma_{\text{target}} = 0$ with the RLS based estimator $\hat{\gamma}_{\text{R}}$. The right panel presents the coverage probability of $\hat{\gamma}_{\text{R}}$ for the target causal effect γ_{target} .

From the left panel of Figure 2, all four effect screening methods lead to perfect selection with probability tending to 1 as the number of replications N_0 increases. Nevertheless, with the forward procedure, the probability of perfect screening are higher than naive screening procedure for all tested N_0 . 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 $H_0 : \gamma_{\text{target}} = 0$, the middle panel of Figure 2 shows that all four methods have asymptotic rejection probability 1, while forward screening possesses higher power when the number of replications are limited. It is worthy of mentioning that in our numerical study we have also compared the behavior of the naive moment estimator $\hat{\gamma}$ and \hat{v} , whose power is much lower than the the RLS based estimators and thus omitted in the figure. This echoes with our conclusion that effect screening can incorporate information across treatment levels and improve the power for testing general causal effects. Analogously, from the right panel of Figure 2, we can witness an improvement in coverage probability of the RLS based confidence intervals with the forward screening procedure.

7.2 Varying effect size

In this section, we study the finite sample performance of effect screening under varying effect sizes. The data generating mechanism is similar to the previous section, except that we fix N_0 to be 2 but vary the size of the factorial effects. Analogously, we still use the four effect screening procedure (Forward Bonferroni, Forward Lasso, Naive Bonferroni and Naive Lasso) and the three criteria (perfect selection probability, power of $\hat{\gamma}_R$ and coverage probability of the RLS based confidence interval) for comparison. The results are reported in Figure 3.

From Figure 3, we can see that perfect effect screening and valid inference is more plausible when the effect size is large enough. Small effect size is a challenging setting, especially when coupled with limited number of replications (recall $N_0 = 2$ in this simulation setting). Nevertheless, forward screening can increase the perfect selection probability and improve the performance of the RLS based inference even for small effect size.

8 Discussion

In this paper we focus on randomized factorial experiments with binary-level factors. We will leave possible extensions to multi-level factors as future endeavor. Besides, it is of interest to discuss effect screening in observational factorial experiments where the data generating mechanism

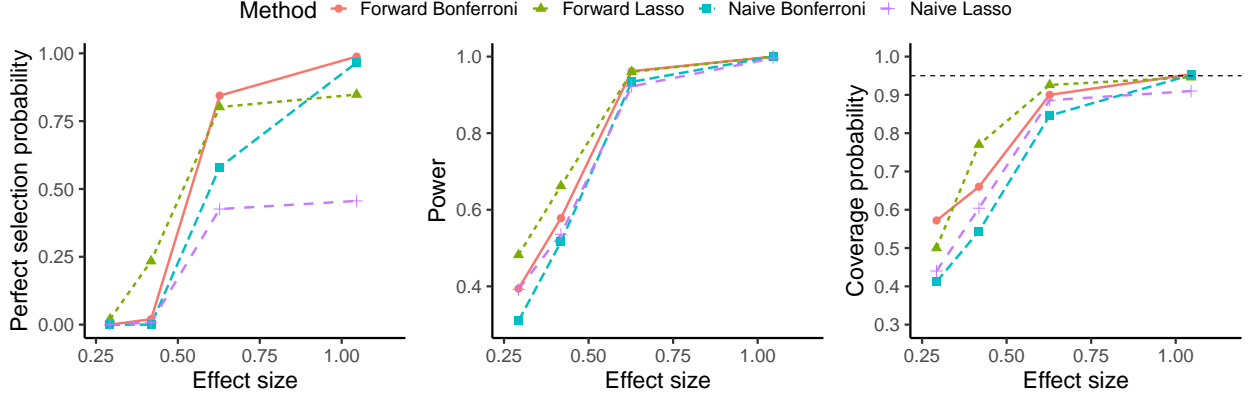


Figure 3: Performance evaluation of effect screening under varying effect sizes. The left panel presents the probability of perfect model selection. The middle panel shows the power curve for testing $H_0 : \gamma_{\text{target}} = 0$ with the RLS based estimator $\hat{\gamma}_R$. The right panel presents the coverage probability of $\hat{\gamma}_R$ for the target causal effect γ_{target} .

is typically more complicated. Moreover, many classical super population literature have propose various post-selection inference strategies, such as sample splitting and selective inference. It is an open research question whether these strategies work under complete randomization. We leave it to future research.

References

- Andrews, I., Kitagawa, T., and McCloskey, A. (2019). Inference on winners. Technical report, National Bureau of Economic Research.
- Bien, J., Taylor, J., and Tibshirani, R. (2013). A lasso for hierarchical interactions. *Annals of statistics*, 41(3):1111.
- 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*, pages 727–753.
- Egami, N. and Imai, K. (2018). Causal interaction in factorial experiments: Application to conjoint analysis. *Journal of the American Statistical Association*.
- Fan, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360.

- 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(522):615–625.
- Hao, N. and Zhang, H. H. (2014). Interaction screening for ultrahigh-dimensional data. *Journal of the American Statistical Association*, 109(507):1285–1301.
- Haris, A., Witten, D., and Simon, N. (2016). Convex modeling of interactions with strong heredity. *Journal of Computational and Graphical Statistics*, 25(4):981–1004.
- 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(520):1759–1769.
- Lim, M. and Hastie, T. (2015). Learning interactions via hierarchical group-lasso regularization. *Journal of Computational and Graphical Statistics*, 24(3):627–654.
- Pashley, N. E. and Bind, M.-A. C. (2019). Causal inference for multiple non-randomized treatments using fractional factorial designs. *arXiv e-prints*, pages arXiv–1905.
- 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*.
- Wang, H. (2009). Forward regression for ultra-high dimensional variable screening. *Journal of the American Statistical Association*, 104(488):1512–1524.
- Wasserman, L. and Roeder, K. (2009). High dimensional variable selection. *Annals of statistics*, 37(5A):2178.
- Wei, W., Zhou, Y., Zheng, Z., and Wang, J. (2022). Inference on the best policies with many covariates. *arXiv preprint arXiv:2206.11868*.
- Wu, C. J. and Hamada, M. S. (2011). *Experiments: planning, analysis, and optimization*, volume 552. John Wiley & Sons.
- Yuan, M., Joseph, V. R., and Lin, Y. (2007). An efficient variable selection approach for analyzing designed experiments. *Technometrics*, 49(4):430–439.
- Zhao, A. and Ding, P. (2021a). Covariate-adjusted fisher randomization tests for the average treatment effect. *Journal of Econometrics*, 225(2):278–294.
- Zhao, A. and Ding, P. (2021b). Regression-based causal inference with factorial experiments: estimands, model specifications, and design-based properties. *arXiv preprint arXiv:2101.02400*.

A General results on consistency of forward screening

We are now ready to show the guarantee for our procedure. Our presentation starts from a theorem (Theorem 6) quantifying the performance of the forward procedure (3.7) under general assumptions on the S-step and P-step. Then we derive a corollary (Theorem 1) by specifying the S-step as Bonferroni corrected marginal t tests and the P-step as screening based on heredity. The high level of the proof proceeds through mathematical induction:

1. (Base case) Show that forward screening selects the correct main effects with probability tending to one.
2. (Induction step) Show that if we correctly screen the effects up to k -way interactions (main effects if $k = 1$), then we are able to correctly detect the non-nulls among all $(k + 1)$ -way interactions.

In order to achieve satisfactory screening results, some regularization conditions need to be imposed to characterize a “good” layer-wise S-step, and the P-step should ensure that the procedure progress in a way that is compatible with the structure of the true factorial effects. In light of this, we use $\mathbb{M}_{d,+}^*$ to denote the pruned set of effects on the d -th layer based on the true model \mathbb{M}_{d-1}^* on the previous layer; that is,

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

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

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

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

where $\mathbb{M}_{d,+}^* = \mathbf{H}(\mathbb{M}_{d-1}^*)$ 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 $\mathbf{S}_N(\cdot)$: for $d = 1, \dots, D$, we have

$$\text{Validity: } \limsup_{N \rightarrow \infty} \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d \cap \mathbb{M}_d^{*c} \neq \emptyset \right\} \leq \alpha_d, \quad (1.23)$$

$$\text{Consistency: } \limsup_{N \rightarrow \infty} \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^* \neq \emptyset \right\} = 0. \quad (1.24)$$

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 ℓ_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 corresponding heredity principle such as (3.5) or (3.6) can improve screening accuracy as well as interpretability of the screening results.

Theorem 6 (Screening consistency). *Assume $\mathbb{M}^* \neq \emptyset$. Assume Assumption 1 and 2. Then the forward screening procedure (3.7) has the following properties:*

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

$$\limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}}_d \cap \mathbb{M}_d^{*c} \neq \emptyset \text{ for some } d \in [D] \right) \leq \alpha = \sum_{d=1}^D \alpha_d. \quad (1.25)$$

(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:*

$$\limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^* \text{ for all } d \in [D] \right) = 1. \quad (1.26)$$

Theorem 6 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 6 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 α_N plays an important role in theoretical discussion. To achieve perfect selection, we hope α_N decays as fast as possible; ideally if α_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 α can only decay at certain rates, because a fast decaying α means

higher possibility of rejection, thus can lead to strict under-selection. Therefore, in the tuning process, α_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.

B More results on inference under perfect screening

Lemma 1. *For $\mathbf{f}^* \neq 0$, we have*

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

Proof of Lemma 1. WLOG assume $\|\mathbf{f}\|_2 = 1$. Due to the orthogonality of G , we can decompose \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

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

Hence

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

□

Proposition 2 (Asymptotic length of confidence comparison with sparse \mathbf{f}). *Define the condition number of $V_{\hat{Y}}$ as $\kappa(V_{\hat{Y}})$. Let s^* be the number of nonzero elements in \mathbf{f} . We have the following upper bound:*

$$\frac{v_{\text{R},\text{lim}}^2}{v_{\text{lim}}^2} \leq \kappa(V_{\hat{Y}}) \cdot \min\left\{ \frac{s^* |\mathbb{M}^*|}{Q}, 1 \right\}.$$

C Technical proofs

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

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 estimation

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

It is known that (Li and Ding, 2017)

$$\mathbb{E}\{\hat{Y}\} = \bar{Y}, \quad V_{\hat{Y}} = \text{Var}\{\hat{Y}\} = \text{Diag}\{N(\mathbf{z})^{-1}S(\mathbf{z}, \mathbf{z})\} - N^{-1}S, \quad (3.27)$$

where $S \in \mathbb{R}^{Q \times Q}$ is the covariance matrix for potential outcomes. Then (3.27) further leads to the following facts:

$$\mathbb{E}\{\hat{\gamma}\} = \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z}) \bar{Y}(\mathbf{z}) = \gamma, \quad (3.28)$$

$$\text{Var}\{\hat{\gamma}\} = \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z})^2 N(\mathbf{z})^{-1} S(\mathbf{z}, \mathbf{z}) - N^{-1} \mathbf{f}^\top S \mathbf{f}, \quad (3.29)$$

$$\mathbb{E}\{\hat{v}^2\} = \sum_{\mathbf{z} \in \mathcal{T}} \mathbf{f}(\mathbf{z})^2 N(\mathbf{z})^{-1} S(\mathbf{z}, \mathbf{z}). \quad (3.30)$$

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

Lemma 2 (Variance concentration and Berry-Esseen bounds in finite population). *Denote $\gamma = \mathbb{E}\{\hat{\gamma}\}$, $v^2 = \text{Var}(\hat{\gamma})$ and $v_R^2 = \mathbb{E}\{\hat{v}_R^2\}$. Suppose the following conditions hold:*

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

$$Q^{-2} \sum_{\mathbf{z}=1}^Q w(\mathbf{z})^2 N_{\mathbf{z}}^{-1} S(\mathbf{z}, \mathbf{z}) \leq \sigma_w^2 v^2. \quad (3.31)$$

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

$$\max_{\mathbf{z} \in [Q]} \frac{1}{N} \sum_{i=1}^N \{Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})\}^4 \leq \Delta^4. \quad (3.32)$$

1. The variance estimator is robust for the true variance: $v_R \geq v$. Besides, the following tail bound holds:

$$\mathbb{P}\{N|\hat{v}_R - v_R| > t\} \leq \frac{C\bar{c}^3\bar{c}^{-4}\|w\|_\infty^2\Delta^4}{QN_0} \cdot \frac{1}{t^2}. \quad (3.33)$$

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

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P}\left\{\frac{\hat{\gamma} - \gamma}{v} \leq t\right\} - \Phi(t) \right| \leq 2C\sigma_w \frac{\bar{c}^{-1}\|w\|_\infty \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\|w\|_2 \sqrt{\bar{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})} \cdot \sqrt{N_0}}. \quad (3.34)$$

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

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

Proof of Lemma 2. 1. See Lemma S13 of Shi and Ding (2022).

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

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

$$\sup_{t \in \mathbb{R}} |\Phi\{(1+a)t + b\} - \Phi\{t\}| \leq |a| + |b|. \quad (3.35)$$

This can be proved by a simple step of intermediate value theorem: for any $t \in \mathbb{R}$,

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

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

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

Then we can show that

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

For the first term, we have

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

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

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

Besides, by (3.35), when $\epsilon_N \leq 1/2$, we also have

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

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

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

On the other hand, we can show that

$$\begin{aligned}\mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{\hat{v}_R} \leq t\right\} &\geq \mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{v} \leq \frac{\hat{v}_R}{v}t, \left|\frac{\hat{v}_R-v_R}{v}\right| \leq \epsilon_N\right\} \\ &\geq \mathbb{P}\left\{\frac{\hat{\gamma}-\gamma}{v} \leq \left(\frac{v_R}{v} - \epsilon_N\right)t\right\} - \mathbb{P}\left\{\left|\frac{\hat{v}_R-v_R}{v}\right| \geq \epsilon_N\right\}.\end{aligned} \quad (3.36)$$

By (3.35), when $\epsilon_N \leq 1/2$, we also have

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

So we can derive a lower bound analogous to (3.36). Note that the results can be analogously generalized to $t \leq 0$. Putting pieces together, we can show that for any $t \geq 0$,

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

□

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

Corollary 1. *Assume Condition (3.31) and (3.32) hold. Let $w = g_K$ for some $K \in \mathbb{K}$. Then we have a Berry-Esseen bound with the estimated variance:*

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

Proof of Corollary 1. Lower bound for Nv^2 . Note that $\|w\|_2^2 = Q$, $\|w\|_\infty = 1$. Using Condition (3.31), we have

$$\begin{aligned} Nv^2 &\geq N\sigma_w^{-2} Q^{-2} \sum_{\mathbf{z}=1}^Q w(\mathbf{z})^2 N_{\mathbf{z}}^{-1} S(\mathbf{z}, \mathbf{z}) \\ &\geq (\underline{c}QN_0) \cdot \sigma_w^{-2} \bar{c}^{-1} Q^{-1} N_0^{-1} \min_{\mathbf{z} \in \mathcal{T}} S(\mathbf{z}, \mathbf{z}) \cdot (Q^{-1} \|w\|_2^2) \\ &= \sigma_w^{-2} \underline{c} \bar{c}^{-1} \min_{\mathbf{z} \in \mathcal{T}} S(\mathbf{z}, \mathbf{z}). \end{aligned}$$

Therefore, the Berry-Esseen bound becomes

$$\begin{aligned} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{\tau}_K - \tau_K}{\hat{v}_R} \leq t \right\} - \Phi \left(\frac{v_R}{v} t \right) \right| &\leq \epsilon_N + \frac{C\sigma_w^4 \bar{c}^5 \underline{c}^{-6} \Delta^4}{(QN_0) \{\min_{\mathbf{z} \in \mathcal{T}} S(\mathbf{z}, \mathbf{z})\}^2} \cdot \frac{1}{\epsilon_N^2} \\ &\quad + 2C\sigma_w \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\sqrt{\bar{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})} \cdot \sqrt{QN_0}}. \end{aligned}$$

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

$$\epsilon_N = \left(\frac{2C\sigma_w^4 \bar{c}^5 \underline{c}^{-6} \Delta^4}{(QN_0) \{\min_{\mathbf{z} \in \mathcal{T}} S(\mathbf{z}, \mathbf{z})\}^2} \right)^{1/3},$$

the upper bound is minimized and

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

□

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

Lemma 3 (Berry Esseen bound with screening). *Let*

$$\mathbf{f}[\mathbb{M}] = Q^{-1}G(\cdot, \mathbb{M})G(\cdot, \mathbb{M})^\top \mathbf{f}. \quad (3.37)$$

Assume there exists $\sigma_w > 0$ such that

$$\sum_{\mathbf{z}=1}^Q \mathbf{f}[\mathbb{M}](\mathbf{z})^2 N_{\mathbf{z}}^{-1} S(\mathbf{z}, \mathbf{z}) \leq \sigma_w^2 v^2(\mathbb{M}). \quad (3.38)$$

Then it holds

$$\begin{aligned} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}[\widehat{\mathbb{M}}] - \gamma[\mathbb{M}]}{v(\mathbb{M})} \leq t \right\} - \Phi(t) \right| \\ \leq 2\mathbb{P} \left\{ \widehat{\mathbb{M}} \neq \mathbb{M} \right\} + 2C\sigma_w \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\sqrt{\underline{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})} \cdot \sqrt{N_0}} \cdot \frac{\|\mathbf{f}[\mathbb{M}]\|_\infty}{\|\mathbf{f}[\mathbb{M}]\|_2}. \end{aligned} \quad (3.39)$$

Proof of Lemma 3. 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:

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

□

C.2 Proof of Theorem 6

Proof of Theorem 6. Induction proof of a basic fact. 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. By induction we hope to prove the following fact under the given assumptions:

For all $D_0 \leq D$, we have

$$\left| \mathbb{P} \left(\widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*, d = 1, \dots, D_0 \right) - \mathbb{P} \left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d = 1, \dots, D_0 \right) \right| \rightarrow 0. \quad (3.40)$$

Because for any $D_0 \in [D]$, we always have:

$$\left\{ \widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d = 1, \dots, D_0 \right\} \subset \left\{ \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*, d = 1, \dots, D_0 \right\},$$

(3.40) is equivalent to: for all $D_0 \leq D$,

$$\mathbb{P} \left(\text{for all } d \in [D_0], \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*; \text{ there exists } d \in [D_0], \widehat{\mathbb{M}}_d \subsetneq \mathbb{M}_d^* \right) \rightarrow 0. \quad (3.41)$$

1. **Main effects.** First, because we assume the tests are consistent (Assumption 1), meaning asymptotically no false negatives:

$$\mathbb{P} \left(\widehat{\mathbb{M}}_1^c \cap \mathbb{M}_1^* \neq \emptyset \right) \rightarrow 0 \Leftrightarrow \mathbb{P} \left(\mathbb{M}_1^* \subset \widehat{\mathbb{M}}_1 \right) \rightarrow 1.$$

Therefore,

$$\mathbb{P} \left(\widehat{\mathbb{M}}_1 \subsetneq \mathbb{M}_1^* \right) \rightarrow 0. \text{ (no under selection for main effects)} \quad (3.42)$$

2. **Induction validity.** Generally speaking, the induction proceeds based on the following idea:

The case for $D_0 = 1$ has been shown in the previous part. Now assume (3.40) or (3.41) for some $D_0 \geq 1$. For $D_0 + 1$, the following holds:

$$\begin{aligned}
0 &\leq \mathbb{P} \left(\left\{ \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*, d \leq D_0 + 1 \right\} \right) - \mathbb{P} \left(\left\{ \widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subset \mathbb{M}_{D_0+1}^* \right\} \right) \\
&= \mathbb{P} \left(\left\{ \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*, d \leq D_0 + 1 \right\} - \left\{ \widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subset \mathbb{M}_{D_0+1}^* \right\} \right) \\
&\leq \mathbb{P} \left(\forall d \in [D_0 + 1], \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*; \exists d \in [D_0], \widehat{\mathbb{M}}_d \subsetneq \mathbb{M}_d^* \right) \\
&\leq \mathbb{P} \left(\forall d \in [D_0], \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*; \exists d \in [D_0], \widehat{\mathbb{M}}_d \subsetneq \mathbb{M}_d^* \right) \rightarrow 0. \text{ (by (3.41))}
\end{aligned}$$

Hence

$$\left| \mathbb{P} \left(\left\{ \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^*, d \leq D_0 + 1 \right\} \right) - \mathbb{P} \left(\left\{ \widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subset \mathbb{M}_{D_0+1}^* \right\} \right) \right| \rightarrow 0. \quad (3.43)$$

Now $\widehat{\mathbb{M}}_{D_0+1}$ 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

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

Hence we can compute

$$\begin{aligned}
0 &\leq \mathbb{P} \left(\left\{ \widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subset \mathbb{M}_{D_0+1}^* \right\} \right) - \mathbb{P} \left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0 + 1 \right) \\
&= \mathbb{P} \left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widehat{\mathbb{M}}_{D_0+1} \subsetneq \mathbb{M}_{D_0+1}^* \right) \\
&= \mathbb{P} \left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d \leq D_0; \widetilde{\mathbb{M}}_{D_0+1} \subsetneq \mathbb{M}_{D_0+1}^* \right) \\
&\leq \mathbb{P} \left(\widetilde{\mathbb{M}}_{D_0+1}^c \cap \mathbb{M}_{D_0+1}^* \neq \emptyset \right) \rightarrow 0.
\end{aligned}$$

The last convergence holds because of the consistency of the test.

The induction can be proceeded.

Proof of the first result. Now it follows

$$\begin{aligned}
&\limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}}_d \cap (\mathbb{M}_d^*)^c \neq \emptyset \text{ for some } d \in [D] \right) \\
&= \limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}}_1 \cap \mathbb{M}_1^{*c} \neq \emptyset \right) + \sum_{D_0=2}^D \mathbb{P} \left(\widehat{\mathbb{M}}_d \cap \mathbb{M}_d^{*c} = \emptyset, d = 1, \dots, D_0 - 1; \widehat{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{*c} \neq \emptyset \right) \\
&= \limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}}_1 \cap \mathbb{M}_1^{*c} \neq \emptyset \right) + \sum_{D_0=2}^D \mathbb{P} \left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^*, d = 1, \dots, D_0 - 1; \widehat{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{*c} \neq \emptyset \right) \\
&\text{(using (3.40) and the fact that } D \text{ is a fixed integer)}
\end{aligned}$$

$$\begin{aligned}
&\leq \limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}}_1 \cap \mathbb{M}_1^{\star c} \neq \emptyset \right) + \sum_{D_0=2}^D \mathbb{P} \left(\widehat{\mathbb{M}}_d = \mathbb{M}_d^{\star}, d = 1, \dots, D_0 - 1; \widetilde{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{\star c} \neq \emptyset \right) \\
&\quad (\text{on the given event, } \widehat{\mathbb{M}}_{D_0,+} = \mathcal{P}(\widehat{\mathbb{M}}_{D_0-1}) = \mathcal{P}(\mathbb{M}_{D_0-1}^{\star}) = \mathbb{M}_{D_0,+}^{\star} \text{ and } \widehat{\mathbb{M}}_{D_0} = \mathcal{S}_N(\widehat{\mathbb{M}}_{D_0,+}) = \widetilde{\mathbb{M}}_{D_0}) \\
&\leq \limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}}_1 \cap \mathbb{M}_1^{\star c} \neq \emptyset \right) + \sum_{D_0=2}^D \mathbb{P} \left(\widetilde{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{\star c} \neq \emptyset \right) \leq \sum_{D_0=1}^D \alpha_{D_0} = \alpha. \tag{3.44}
\end{aligned}$$

Therefore the target probability gets controlled under α .

Proof of the second result. Under $\alpha = \alpha_N \rightarrow 0$, (3.44) implies that with probability tending to one,

$$\widehat{\mathbb{M}}_d \cap (\mathbb{M}_d^{\star})^c = \emptyset, \text{ for } d = 1, \dots, D \Leftrightarrow \widehat{\mathbb{M}}_d \subset \mathbb{M}_d^{\star} = \emptyset, \text{ for } d = 1, \dots, D.$$

Now apply (3.40), we obtain

$$\widehat{\mathbb{M}}_d = \mathbb{M}_d^{\star}, \text{ for } d = 1, \dots, D,$$

with probability tending to one, which concludes the proof. □

C.3 Proof of Theorem 1

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

Theorem 7 (Bonferroni corrected marginal t test). *Let $\widetilde{\mathbb{M}}_d = \widehat{\mathbb{S}}(\mathbb{M}_{d,+}^{\star})$ where $\mathbb{M}_{d,+}^{\star} = \mathbb{P}(\mathbb{M}_{d-1}^{\star})$. 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) (Validity) $\limsup_{N \rightarrow \infty} \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d \cap \mathbb{M}_d^{\star c} \neq \emptyset \right\} \leq \alpha_d$ for all $d = 1, \dots, D$.
- (ii) (Consistency) $\limsup_{N \rightarrow \infty} \mathbb{P} \left\{ \widetilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^{\star} \neq \emptyset \right\} = 0$ for all $d = 1, \dots, D$.
- (iii) (Type I error control) Overall the procedure achieves type I error rate control:

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

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

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

Part (i) and ii of Theorem 1 justified Assumption 1 and 2 respectively, which build up the basis for applying Theorem 6. Part (iii) guarantees type I error control under the significance level α . When we let α 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 α - perfect selection can be achieved asymptotically.

Proof of Theorem 1. 1. First, we show validity:

$$\begin{aligned} \limsup_{N \rightarrow \infty} \mathbb{P} \left\{ \tilde{\mathbb{M}}_d \cap \mathbb{M}_d^{\star c} \neq \emptyset \right\} &= \limsup_{N \rightarrow \infty} \mathbb{P} \left\{ \exists \mathcal{K} \in \mathbb{M}_{d,+}^{\star} \setminus \mathbb{M}_d^{\star}, \left| \frac{\hat{\tau}_{\mathcal{K}}}{\hat{\sigma}_{\mathcal{K}}} \right| \geq \Phi^{-1} \left(1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^{\star}|} \right) \right\} \\ &\leq \limsup_{N \rightarrow \infty} \sum_{\mathcal{K} \in \mathbb{M}_{d,+}^{\star} \setminus \mathbb{M}_d^{\star}} \mathbb{P} \left\{ \left| \frac{\hat{\tau}_{\mathcal{K}}}{\hat{\sigma}_{\mathcal{K}}} \right| \geq \Phi^{-1} \left(1 - \frac{\alpha_d}{2|\mathbb{M}_{d,+}^{\star}|} \right) \right\} \\ &\leq \limsup_{N \rightarrow \infty} \sum_{\mathcal{K} \in \mathbb{M}_{d,+}^{\star} \setminus \mathbb{M}_d^{\star}} \left(\frac{\alpha_d}{|\mathbb{M}_{d,+}^{\star}|} + \frac{\tilde{C}}{(QN_0)^{1/3}} \right) \leq \alpha_d. \end{aligned}$$

2. Second, we show consistency.

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

For simplicity, let

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

Then

$$\begin{aligned} \limsup_{N \rightarrow \infty} \mathbb{P} \left\{ \tilde{\mathbb{M}}_d^c \cap \mathbb{M}_d^{\star} \neq \emptyset \right\} &\leq \limsup_{N \rightarrow \infty} \sum_{\mathcal{K} \in \mathbb{M}_d^{\star}} \left(\mathbb{P} \left\{ -Z_d^{\star} - \frac{\tau_{\mathcal{K}}}{\sigma_{\mathcal{K}}} \leq \frac{\hat{\tau}_{\mathcal{K}}}{\hat{\sigma}_{\mathcal{K}}} - \frac{\tau_{\mathcal{K}}}{\sigma_{\mathcal{K}}} \leq Z_d^{\star} - \frac{\tau_{\mathcal{K}}}{\sigma_{\mathcal{K}}} \right\} + \frac{\tilde{C}}{(QN_0)^{1/3}} \right) \\ &= \limsup_{N \rightarrow \infty} \sum_{\mathcal{K} \in \mathbb{M}_d^{\star}} \Phi \left\{ r_{\mathcal{K}}^{-1} \left(Z_d^{\star} - \frac{\tau_{\mathcal{K}}}{\sigma_{\mathcal{K}}} \right) \right\} - \Phi \left\{ r_{\mathcal{K}}^{-1} \left(-Z_d^{\star} - \frac{\tau_{\mathcal{K}}}{\sigma_{\mathcal{K}}} \right) \right\}. \end{aligned} \quad (3.45)$$

With Condition 2, we have

$$Z_d^* = \Theta \left(\sqrt{2 \ln \frac{2|\mathbb{M}_{d,+}^*|}{\alpha_d}} \right) = \Theta(\max\{\sqrt{\delta' \ln N}, \sqrt{\ln(2|\mathbb{M}_{d,+}^*|)}\}, \left| \frac{\tau_{\mathcal{K}}}{\sigma_{\mathcal{K}}} \right| = \Theta(N^{1/2+\delta}).$$

Because $\delta > -1/2$ and $\delta' \geq 0$, we have $|\frac{\tau_{\mathcal{K}}}{\sigma_{\mathcal{K}}}| \rightarrow \infty$ and $Z_d^*/(|\frac{\tau_{\mathcal{K}}}{\sigma_{\mathcal{K}}}|) \rightarrow 0$. Hence the above limit (3.45) converges to zero. This concludes the proof.

3. Based on the above two parts and Theorem 6, it suffices to conclude the Type I error rate control. A more delicate analysis in this particular setup can actually lead to sharper bound. Based on (3.44), we directly compute

$$\begin{aligned} & \limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}} \cap \mathbb{M}^{*c} \neq \emptyset \right) \\ & \leq \limsup_{N \rightarrow \infty} \mathbb{P} \left(\widehat{\mathbb{M}}_1 \cap \mathbb{M}_1^{*c} \neq \emptyset \right) + \sum_{D_0=2}^D \mathbb{P} \left(\widetilde{\mathbb{M}}_{D_0} \cap \mathbb{M}_{D_0}^{*c} \neq \emptyset \right) \\ & \leq \frac{\alpha_1}{K} \cdot |\mathbb{M}_1^*| + \sum_{D_0=2}^D \frac{\alpha_{D_0}}{|\mathbb{M}_{D_0,+}^*|} \cdot |\mathbb{M}_{D_0}^*| \leq \alpha. \end{aligned}$$

4. The perfect selection result follows from Part 1,2 and Theorem 6.

□

C.4 Proof of Theorem 2

Theorem 2 is a direct result of Lemma 2 and the following Berry Esseen result:

Lemma 4 (Berry-Esseen bound under perfect screening). *Let $\mathbf{f}[\mathbb{M}]$ be given by (3.37). Assume (3.38). Then it holds*

$$\begin{aligned} & \sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\widehat{\gamma}(\widehat{\mathbb{M}}) - \gamma}{v(\mathbb{M}^*)} \leq t \right\} - \Phi(t) \right| \\ & \leq 2\mathbb{P} \left\{ \widehat{\mathbb{M}} \neq \mathbb{M}^* \right\} + 2C\sigma_w \frac{\underline{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\sqrt{\bar{c}^{-1} \min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})} \cdot \sqrt{N_0}} \cdot \frac{\|\mathbf{f}[\mathbb{M}^*]\|_{\infty}}{\|\mathbf{f}[\mathbb{M}^*]\|_2}. \end{aligned} \quad (3.46)$$

Proof of Lemma 4. This lemma is a direct application of Lemma 3. First we check that

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

From the definition of γ (3.28), we have

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

Now apply Lemma 3 with $\mathbb{M} = \mathbb{M}^\star$ gives the result. \square

C.5 Proof of Theorem 3

Proof of Theorem 3. According to Condition 6 and Theorem 1, with Strategy 1,

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

We will apply Lemma 3 with

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

We only need to verify $\gamma = \gamma[\underline{\mathbb{M}}]$ in this case.

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

Now by (5.20), $\mathbf{f}^\top G(\cdot, \mathbb{M}^c) = 0$. Hence

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

\square

C.6 Proof of Theorem 4

Proof of Theorem 4. This proof can be finished by applying Lemma 3 with $\mathbb{M} = \overline{\mathbb{M}}^\star$ and checking $\gamma[\overline{\mathbb{M}}^\star] = \gamma$, which is omitted here. \square

C.7 Proof of Proposition 1

Proof of Proposition 1. WLOG we assume only the first s^\star elements of \mathbf{f} are nonzero. That is,

$$\mathbf{f} = f_1 \mathbf{e}_1 + \cdots + f_{s^\star} \mathbf{e}_{s^\star}. \quad (3.47)$$

We can verify that

$$\|Q^{-1}G(\cdot, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top \mathbf{e}_k\|_2 = \frac{|\mathbb{M}^*|}{Q}, \quad \forall k \in [Q]. \quad (3.48)$$

Therefore,

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

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

$$\frac{\|\mathbf{f}^*\|_2^2}{\|\mathbf{f}\|_2^2} \leq 1. \quad (3.49)$$

On the other hand, using (3.47) and (3.48), we have

$$\frac{\|\mathbf{f}^*\|_2^2}{\|\mathbf{f}\|_2^2} \leq \frac{\|\mathbf{f}\|_1^2}{\|\mathbf{f}\|_2^2} \cdot \frac{|\mathbb{M}^*|}{Q} \leq \frac{s^* |\mathbb{M}^*|}{Q}. \quad (3.50)$$

Combining (3.49) and (3.50) concludes the proof. \square

C.8 Proof of Theorem 5

For simplicity, we focus on the case given by (6.21). The general proof can be completed similarly.

We begin by the following lemma:

Lemma 5 (Consistency of the selected tie sets). *Assume Conditions 1, 3 and 7. There exists universal constants $C, C' > 0$, such that when $N > n(\delta_1, \delta_2, \delta_3)$,*

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

Lemma 5 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 5. The high level idea of the proof is: we first prove the non-asymptotic bounds over the random event $\hat{\mathbb{M}} = \mathbb{M}^*$, then make up for the cost of $\hat{\mathbb{M}} \neq \mathbb{M}^*$. Over $\hat{\mathbb{M}} = \mathbb{M}^*$, we have

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

We apply Lemma 3 to establish a Berry-Esseen bound for each $\hat{Y}_r^*(\mathbf{z})$. Note that

$$\hat{Y}_r^*(\mathbf{z}) = \mathbf{f}_z^\top \hat{Y}, \quad \mathbf{f}_z^\top = Q^{-1}G(\mathbf{z}, \mathbb{M}^*)G(\cdot, \mathbb{M}^*)^\top. \quad (3.51)$$

By calculation we have

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

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

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \frac{\hat{Y}_r^*(\mathbf{z}) - \bar{Y}(\mathbf{z})}{v_N} \leq t \right\} - \Phi(t) \right| \leq 2C\sigma \frac{\bar{c}^{-1} \max_{i \in [N], \mathbf{z} \in [Q]} |Y_i(\mathbf{z}) - \bar{Y}(\mathbf{z})|}{\bar{c}^{-1/2} \{\min_{\mathbf{z} \in [Q]} S(\mathbf{z}, \mathbf{z})\}^{1/2}} \sqrt{\frac{|\mathbb{M}^*|}{QN_0}}. \quad (3.52)$$

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

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

We can show that

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

Hence we know that

$$\begin{aligned} & \mathbb{P} \left\{ \sqrt{N} \left| \hat{Y}_r^*(z) - \bar{Y}(z) \right| \geq \sqrt{N} d_h^* \right\} \\ & \leq \frac{v_N}{\sqrt{2\pi} d_h^*} \cdot \exp\left(-\frac{d_h^{*2}}{2v_N^2}\right) + 2C\sigma \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}}. \end{aligned} \quad (3.53)$$

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

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

Using (3.53) 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

$$\begin{aligned} & \mathbb{P} \left\{ \hat{Y}_r^*(z') - \hat{Y}_r^*(z) < 0 \right\} \\ & \geq 1 - C|\mathcal{T}_1||\mathcal{T}'| \left\{ \sqrt{\frac{\bar{c}\bar{s}|\mathbb{M}^*|}{N^{1+2\delta_1}}} \exp \left\{ -\frac{C'N^{1+2\delta_1}}{\bar{c}\bar{s}|\mathbb{M}^*|} \right\} + \sigma \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned} \quad (3.54)$$

Nice separation. Suppose we are working on a random coordinate \tilde{z} . For $z \notin \mathcal{T}_1$ and any $\bar{\epsilon} > 0$,

$$\begin{aligned} & \mathbb{P} \left\{ \min_{z \notin \mathcal{T}_1} |\hat{Y}_r^*(z) - \hat{Y}_r^*(\tilde{z})|/\eta_N \geq 2\bar{\epsilon} \right\} \\ & \geq \mathbb{P} \left\{ \min_{z \notin \mathcal{T}_1, z' \in \mathcal{T}_1} |\hat{Y}_r^*(z) - \hat{Y}_r^*(z')|/\eta_N \geq 2\bar{\epsilon}, \tilde{m} \in \mathcal{T}_1 \right\} \\ & \geq \mathbb{P} \left\{ \min_{z \notin \mathcal{T}_1, z' \in \mathcal{T}_1} |\hat{Y}_r^*(z) - \hat{Y}_r^*(z')|/\eta_N \geq 2\bar{\epsilon} \right\} + \mathbb{P} \{ \tilde{m} \in \mathcal{T}_1 \} - 1 \\ & \geq \mathbb{P} \{ \tilde{m} \in \mathcal{T}_1 \} - \sum_{z \notin \mathcal{T}_1, z' \in \mathcal{T}_1} \mathbb{P} \left\{ |\hat{Y}_r^*(z) - \hat{Y}_r^*(z')|/\eta_N \leq 2\bar{\epsilon} \right\}. \end{aligned} \quad (3.55)$$

To proceed we have the following tail bound:

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

(This is deduced analogously to the proof in the previous part)

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 (3.55), we have:

$$\begin{aligned} & \mathbb{P} \left\{ \min_{z \notin \mathcal{T}_1} |\hat{Y}_r^*(z) - \hat{Y}_r^*(\tilde{z})|/\eta_N \geq 2\bar{\epsilon} \right\} \\ & \geq \mathbb{P} \{ \tilde{m} \in \mathcal{T}_1 \} - \underbrace{4|\mathcal{T}_1||\mathcal{T}'| \frac{\sqrt{2\bar{c}\bar{s}|\mathbb{M}^*|}}{\sqrt{\pi N_0 Q d_h^*}} \cdot \exp \left(-\frac{N_0 Q d_h^{*2}}{8\bar{c}\bar{s}|\mathbb{M}^*|} \right)}_{\text{Term I}} \\ & \quad - \underbrace{4|\mathcal{T}_1||\mathcal{T}'| 2C\sigma \frac{\bar{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}}}_{\text{Term II}}. \end{aligned} \tag{3.56}$$

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

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

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

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

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

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

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,

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

Hence we have:

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

Again, by the conditions, we can show

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

Specifying the random indices. Introduce the following random indices:

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

Applying (3.54) we know that

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

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

$$\begin{aligned} & \mathbb{P} \left\{ \max_{z \in \mathcal{T}_1} |\hat{Y}_r^*(z) - \hat{Y}_r^*(\tilde{z})| \leq \underline{\epsilon}\eta_N, \min_{z \notin \mathcal{T}_1} |\hat{Y}_r^*(z) - \hat{Y}_r^*(\tilde{z})| \geq \bar{\epsilon}\eta_N \right\} \\ & \geq 1 - C|\mathcal{T}'||\mathcal{T}_1| \left\{ \sqrt{\frac{\bar{c}\bar{s}|\mathbb{M}^*|}{N^{1+2\delta_2}}} \exp \left(-\frac{C'N^{1+2\delta_2}}{\bar{c}\bar{s}|\mathbb{M}^*|} \right) + \sigma \frac{\underline{c}^{-1} \max_{i \in [N], z \in [Q]} |Y_i(z) - \bar{Y}(z)|}{\bar{c}^{-1/2} \{\min_{z \in [Q]} S(z, z)\}^{1/2}} \cdot \sqrt{\frac{|\mathbb{M}^*|}{N_0 Q}} \right\}. \end{aligned} \tag{3.57}$$

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

$$\begin{aligned}
& \mathbb{P}\left\{\widehat{\mathcal{T}}_1 = \mathcal{T}_1\right\} \\
&= \mathbb{P}\left\{|\widehat{Y}_R(\mathbf{z}) - \max_{\mathbf{z} \in \mathcal{T}'} \widehat{Y}_R(\mathbf{z})| \leq \underline{\epsilon} \eta_N, \text{ for } \mathbf{z} \in \mathcal{T}_1; \right. \\
&\quad \left. |\widehat{Y}_R(\mathbf{z}) - \max_{\mathbf{z} \in \mathcal{T}'} \widehat{Y}_R(\mathbf{z})| > \underline{\epsilon} \eta_N, \text{ for } \mathbf{z} \notin \mathcal{T}_1\right\} \\
&\geq \mathbb{P}\left\{|\widehat{Y}_R^*(\mathbf{z}) - \max_{\mathbf{z} \in \mathcal{T}'} \widehat{Y}_R^*(\mathbf{z})| \leq \underline{\epsilon} \eta_N, \text{ for } \mathbf{z} \in \mathcal{T}_1; \right. \\
&\quad \left. |\widehat{Y}_R^*(\mathbf{z}) - \max_{\mathbf{z} \in \mathcal{T}'} \widehat{Y}_R^*(\mathbf{z})| > \underline{\epsilon} \eta_N, \text{ for } \mathbf{z} \notin \mathcal{T}_1\right\} - \mathbb{P}\{\widehat{\mathbf{M}} \neq \mathbf{M}^*\} \\
&= \mathbb{P}\left\{|\widehat{Y}_R^*(\mathbf{z}) - \widehat{Y}_R^*(\widetilde{\mathbf{z}}_h)| \leq \underline{\epsilon} \eta_N, \text{ for } \mathbf{z} \in \mathcal{T}_1; \right. \\
&\quad \left. |\widehat{Y}_R^*(\mathbf{z}) - \widehat{Y}_R^*(\widetilde{\mathbf{z}}_h)| > \underline{\epsilon} \eta_N, \text{ for } \mathbf{z} \notin \mathcal{T}_1\right\} - \mathbb{P}\{\widehat{\mathbf{M}} \neq \mathbf{M}^*\} \\
&\quad (\text{where we introduce random index } \widetilde{\mathbf{z}}_h \text{ to record the position that achieves maximum}) \\
&\geq \mathbb{P}\left\{|\widehat{Y}_R^*(\mathbf{z}) - \widehat{Y}_R^*(\widetilde{\mathbf{z}}_h)| \leq \underline{\epsilon} \eta_N, \text{ for } \mathbf{z} \in \mathcal{T}_1; \right. \\
&\quad \left. |\widehat{Y}_R^*(\mathbf{z}) - \widehat{Y}_R^*(\widetilde{\mathbf{z}}_h)| > \bar{\epsilon} \eta_N, \text{ for } \mathbf{z} \notin \mathcal{T}_1\right\} - \mathbb{P}\{\widehat{\mathbf{M}} \neq \mathbf{M}^*\} \\
&\quad (\text{simply using the fact that } \bar{\epsilon} > \underline{\epsilon}) \\
&= \mathbb{P}\left\{\max_{\mathbf{z} \in \mathcal{T}_1} |\widehat{Y}_R^*(\mathbf{z}) - \widehat{Y}_R^*(\widetilde{\mathbf{z}}_h)| \leq \underline{\epsilon} \eta_N; \min_{\mathbf{z} \notin \mathcal{T}_1} |\widehat{Y}_R^*(\mathbf{z}) - \widehat{Y}_R^*(\widetilde{\mathbf{z}}_h)| > \bar{\epsilon} \eta_N\right\} \\
&\quad - \mathbb{P}\{\widehat{\mathbf{M}} \neq \mathbf{M}^*\} \\
&\geq 1 - \sum_{h=1}^{H_0} \left(1 - \mathbb{P}\left\{\max_{\mathbf{z} \in \mathcal{T}_1} |\widehat{Y}_R^*(\mathbf{z}) - \widehat{Y}_R^*(\widetilde{\mathbf{z}}_h)| \leq \underline{\epsilon} \eta_N; \min_{\mathbf{z} \notin \mathcal{T}_1} |\widehat{Y}_R^*(\mathbf{z}) - \widehat{Y}_R^*(\widetilde{\mathbf{z}}_h)| > \bar{\epsilon} \eta_N\right\}\right) \\
&\quad - \mathbb{P}\{\widehat{\mathbf{M}} \neq \mathbf{M}^*\} \\
&\geq 1 - \mathbb{P}\{\widehat{\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{\bar{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 5 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 2 and Lemma 3.

D More discussions

D.1 Statement and proof of Lemma 6

Lemma 6. \hat{Y}_R from (4.14) can be expressed as:

$$\hat{Y}_R = Q^{-1}G(\cdot, \hat{\mathbb{M}})G(\cdot, \hat{\mathbb{M}})^\top \hat{Y}. \quad (4.58)$$

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

Proof of Lemma 6. 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 (4.14), 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 (4.14).

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

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

□

D.2 Discussion of several screening methods in finite population factorial experiments

- Bonferroni correction:
- LASSO:

$$\min_{\tau \in \mathbb{R}^H} \sum_{i=1}^N (y_i - f_i^\top \tau) + \lambda \|\tau\|_1,$$

- AIC:

$$\min_{\tau \in \mathbb{R}^H, \mathbb{M} \subset [K]} \sum_{i=1}^N (y_i - f_i^\top \tau_{\mathbb{M}}) + \lambda |\mathbb{M}|,$$

- BIC:

$$\min_{\tau \in \mathbb{R}^H, \mathbb{M} \subset [K]} \sum_{i=1}^N (y_i - f_{i,\mathbb{M}}^\top \tau_{\mathbb{M}}) + \lambda |\mathbb{M}| \log(N),$$

D.3 Discussion: more general centering values

D.3.1 Unsaturated weighted least square: a closed form expression

In this section we first derive the closed form expression for unsaturated WLS estimation, then verify the nice targeting property we mentioned in the previous section.

First we need to introduce a transformation matrix $\mathbf{P}_{\Delta\delta_{[K]}}$, with columns and rows indexed by subsets $\{\mathcal{K} \subset [K]\}$ of the K factors. Generally it is used to reveal the relationship between designs with different configurations of centering factors $\delta_{[K]}$ and $\delta'_{[K]} = \delta_{[K]} + \Delta\delta_{[K]}$. The transformation is actually linear:

$$\left(f_{\delta'_{[K]}}(z_{\mathcal{K}}^*) \right)_{\mathcal{K} \subset [K]} = \left(f_{\delta_{[K]}}(z_{\mathcal{K}}^*) \right)_{\mathcal{K} \subset [K]} \mathbf{P}_{\Delta\delta_{[K]}}. \quad (4.59)$$

The closed form of $\mathbf{P}_{\Delta\delta_{[K]}}$ is easy to derive. Note that for all $\mathcal{K}' \subset [K]$, we have

$$f_{\delta'_{[K]}}(z_{\mathcal{K}'}^*) = \sum_{\mathcal{K} \subset \mathcal{K}'} f_{\delta_{[K]}}(z_{\mathcal{K}}^*) \prod_{k \in \mathcal{K}' \setminus \mathcal{K}} (\Delta\delta)_k,$$

which implies the element of $\mathbf{P}_{\Delta\delta_{[K]}}$ indexed by $(\mathcal{K}, \mathcal{K}')$ is given by

$$\mathbf{P}_{\Delta\delta_{[K]}}(\mathcal{K}, \mathcal{K}') = \begin{cases} \prod_{k \in \mathcal{K}' \setminus \mathcal{K}} (\Delta\delta)_k & , \mathcal{K} \subset \mathcal{K}', \\ 0 & , \mathcal{K} \not\subset \mathcal{K}'. \end{cases} \quad (4.60)$$

Define $\mathbf{Q}_{\Delta\delta_{[K]}} = \mathbf{P}_{\Delta\delta_{[K]}}^{-1}$ to be the inverse. Note that $\mathbf{Q}_{\Delta\delta_{[K]}}$ is simply taking out a $\Delta\delta_{[K]}$ vector from a group of centering factors, so by symmetry we have

$$\mathbf{Q}_{\Delta\delta_{[K]}}(\mathcal{K}, \mathcal{K}') = \begin{cases} (-1)^{|\mathcal{K}'| - |\mathcal{K}|} \prod_{k \in \mathcal{K}' \setminus \mathcal{K}} (\Delta\delta)_k & , \mathcal{K} \subset \mathcal{K}', \\ 0 & , \mathcal{K} \not\subset \mathcal{K}'. \end{cases} \quad (4.61)$$

We shall give an example of the above matrix in the three-factor case, which appears (incompletely) in the appendix of Zhao and Ding (2021b). Let $A' = A - \delta_A$, $B' = B - \delta_B$, $C' = C - \delta_C$.

$$\begin{pmatrix} 1 \\ A \\ B \\ C \\ AB \\ AC \\ BC \\ ABC \end{pmatrix} = \mathbf{P}_{\Delta\delta_{[K]}}^\top \begin{pmatrix} 1 \\ A' \\ B' \\ C' \\ A'B' \\ A'C' \\ B'C' \\ A'B'C' \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \delta_A & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ \delta_B & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ \delta_C & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ \delta_A\delta_B & \delta_B & \delta_A & 0 & 1 & 0 & 0 & 0 \\ \delta_A\delta_C & \delta_C & 0 & \delta_A & 0 & 1 & 0 & 0 \\ \delta_B\delta_C & 0 & \delta_C & \delta_B & 0 & 0 & 1 & 0 \\ \delta_A\delta_B\delta_C & \delta_B\delta_C & \delta_A\delta_C & \delta_A\delta_B & \delta_C & \delta_B & \delta_A & 1 \end{pmatrix} \begin{pmatrix} 1 \\ A' \\ B' \\ C' \\ A'B' \\ A'C' \\ B'C' \\ A'B'C' \end{pmatrix}.$$

The following theorem shows that $\mathbf{P}_{\Delta\delta_{[K]}}$ and $\mathbf{Q}_{\Delta\delta_{[K]}}$ totally determines the structure of \mathbf{D}_h .

Theorem 8. *Consider weighted least squares with centering factors $\delta_{[K]}$ and weights proportional to size of each stratum. Let $\Delta\delta_{[K]} = \delta_{[K]} - (1/2)_{k=1}^K$. The unsaturated regression on up to all m -level main/interactions terms has coefficient vector:*

$$(\tilde{\tau}_{\mathcal{K}})_{\{|\mathcal{K}| \leq m\}} = (\tau_{\mathcal{K}})_{\{|\mathcal{K}| \leq m\}} + \mathbf{D}_h \cdot (\tau_{\mathcal{K}})_{\{|\mathcal{K}| > m\}}, \quad (4.62)$$

where \mathbf{D}_h is given by

$$\mathbf{D}_h = \mathbf{P}_{\Delta\delta_{[K]}} (\{\mathcal{K} \subset [m]\}, \{\mathcal{K} \subset [m]\}) \cdot \mathbf{Q}_{\Delta\delta_{[K]}} (\{\mathcal{K} \subset [m]\}, \{\mathcal{K} \subset [K] \setminus [m]\}).$$

Corollary 2. *The matrix \mathbf{D} has a closed form expression:*

1. For $\mathcal{K} \subsetneq \mathcal{K}'$,

$$\mathbf{D}_h(\mathcal{K}, \mathcal{K}') = 0. \quad (4.63)$$

2. For $\mathcal{K} \subset \mathcal{K}'$, let $|\mathcal{K}| = k$, $|\mathcal{K}'| = k'$, with $k \leq m < k'$,

$$\mathbf{D}_h(\mathcal{K}, \mathcal{K}') = \sum_{l=0}^{m-k} (-1)^{k'-k+1-l} \binom{k'-k+1}{l} \prod_{t \in \mathcal{K}' \setminus \mathcal{K}} \left(\delta_t - \frac{1}{2} \right). \quad (4.64)$$

Proof. This result can be derived through careful calculation based on the definition of \mathbf{P} and \mathbf{Q} from (4.60) and (4.61) along with Theorem 8 thus omitted here. \square

D.3.2 A sufficient condition for sign consistency in population WLS regression

Definition 1 (Active interaction number). *For every z_k of the K factors, there are s_k factors that have nonzero interaction with z_k , where $s_k \in [K - 1]$ is a nonnegative integer associated with K . We call s_k the active interaction number of factor z_k . The maximal active interaction number is subsequently defined as $s_K = \max_{k \in [K]} s_k$.*

This definition is mainly devoted to finer technical purposes in Theorem 9.

Theorem 9. *Assume we run weighted least square under the setting depicted in Theorem 8. Define the maximal decaying rate $c_K = \max_{l \in [K]} c_l$. Recall the predefined maximal active interaction number s_K from Definition 1. If we have*

$$s_K c_K \max_{k=1, \dots, K} |\delta_k - 1/2| < \ln 2, \quad (4.65)$$

then the unsaturated regression coefficients $(\tilde{\tau}_{\mathcal{K}})_{\{|\mathcal{K}| \leq m\}}$ and the corresponding saturated regression coefficients $(\tau_{\mathcal{K}})_{\{|\mathcal{K}| \leq m\}}$ from (4.62) have same signs on every term.

Condition (4.65) unifies the property of factorial effects and the information of the design pattern (the centering factors $\delta_{[K]}$). The product of s_K and c_K demonstrates a trade-off between the active interaction number and the hierarchy structure. Sparser interactions require slower decaying rate and vice versa. Besides, the product of $s_K c_K$ and $\max_{k=1, \dots, K} |\delta_k - 1/2|$ shows that if δ_k lies more close to $1/2$, less restriction are needed on the effect structure. This aligns with the result in Zhao and Ding (2021b): when $\delta_k = 1/2$ holds for all $k = 1, \dots, K$, $\mathbf{D}_h = \mathbf{0}$, so that forward selection always works.