

This article was downloaded by: [University of Arizona]

On: 02 October 2014, At: 11:31

Publisher: Taylor & Francis

Informa Ltd Registered in England and Wales Registered Number: 1072954 Registered office: Mortimer House, 37-41 Mortimer Street, London W1T 3JH, UK

## Journal of the American Statistical Association

Publication details, including instructions for authors and subscription information:

<http://www.tandfonline.com/loi/uasa20>

### Interaction Screening for Ultrahigh-Dimensional Data

Ning Hao & Hao Helen Zhang

Accepted author version posted online: 18 Feb 2014. Published online: 02 Oct 2014.



**To cite this article:** Ning Hao & Hao Helen Zhang (2014) Interaction Screening for Ultrahigh-Dimensional Data, Journal of the American Statistical Association, 109:507, 1285-1301, DOI: [10.1080/01621459.2014.881741](https://doi.org/10.1080/01621459.2014.881741)

**To link to this article:** <http://dx.doi.org/10.1080/01621459.2014.881741>

PLEASE SCROLL DOWN FOR ARTICLE

Taylor & Francis makes every effort to ensure the accuracy of all the information (the "Content") contained in the publications on our platform. However, Taylor & Francis, our agents, and our licensors make no representations or warranties whatsoever as to the accuracy, completeness, or suitability for any purpose of the Content. Any opinions and views expressed in this publication are the opinions and views of the authors, and are not the views of or endorsed by Taylor & Francis. The accuracy of the Content should not be relied upon and should be independently verified with primary sources of information. Taylor and Francis shall not be liable for any losses, actions, claims, proceedings, demands, costs, expenses, damages, and other liabilities whatsoever or howsoever caused arising directly or indirectly in connection with, in relation to or arising out of the use of the Content.

This article may be used for research, teaching, and private study purposes. Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden. Terms & Conditions of access and use can be found at <http://www.tandfonline.com/page/terms-and-conditions>

# Interaction Screening for Ultrahigh-Dimensional Data

Ning HAO and Hao Helen ZHANG

In ultrahigh-dimensional data analysis, it is extremely challenging to identify important interaction effects, and a top concern in practice is computational feasibility. For a dataset with  $n$  observations and  $p$  predictors, the augmented design matrix including all linear and order-2 terms is of size  $n \times (p^2 + 3p)/2$ . When  $p$  is large, say more than tens of hundreds, the number of interactions is enormous and beyond the capacity of standard machines and software tools for storage and analysis. In theory, the interaction-selection consistency is hard to achieve in high-dimensional settings. Interaction effects have heavier tails and more complex covariance structures than main effects in a random design, making theoretical analysis difficult. In this article, we propose to tackle these issues by forward-selection-based procedures called iFOR, which identify interaction effects in a greedy forward fashion while maintaining the natural hierarchical model structure. Two algorithms, iFORT and iFORM, are studied. Computationally, the iFOR procedures are designed to be simple and fast to implement. No complex optimization tools are needed, since only OLS-type calculations are involved; the iFOR algorithms avoid storing and manipulating the whole augmented matrix, so the memory and CPU requirement is minimal; the computational complexity is linear in  $p$  for sparse models, hence feasible for  $p \gg n$ . Theoretically, we prove that they possess sure screening property for ultrahigh-dimensional settings. Numerical examples are used to demonstrate their finite sample performance. Supplementary materials for this article are available online.

KEY WORDS: Forward selection; GWAS; Heredity condition; Sure screening.

## 1. INTRODUCTION

Ultrahigh-dimensionality is a significant feature of data collected in contemporary scientific research, owing to rapid advances of technologies and computer power. Big data are abundant in many areas including biology, genetics, medicine, finance, social science, environmental science, and so on. One major challenge in dealing with big datasets is that, the number of predictors  $p$  is much larger than the sample size  $n$ . In this article, we allow  $p$  to be as large as  $O(\exp(n^\xi))$  for some  $\xi \in (0, 1/2)$ , which is described as *nonpolynomial* (NP) dimensionality in Fan and Song (2010). To extract useful information from such data and build an interpretable model with high-prediction power, variable selection or screening must be employed. A variety of variable selection methods have been developed and in common use, such as the LASSO (Tibshirani 1996), SCAD (Fan and Li 2001), Dantzig selector (Candes and Tao 2007), elastic net (Zou and Hastie 2005), minimax concave penalty (MCP) (Zhang 2010), and others (Zou 2006; Zou and Li 2008). Many methods possess favorable theoretical properties such as model selection consistency (Zhao and Yu 2006) and oracle properties (Fan and Lv 2011). When  $p$  is much larger than  $n$ , sure screening is a more realistic goal to achieve than oracle properties or selection consistency (Fan and Lv 2008; Wang 2009). Sure screening assures that all important variables are identified with a probability tending to 1, hence achieving effective dimension reduction without information loss and providing a reasonable starting point for low-dimensional methods to be applied.

Most existing methods for variable selection are designed for selecting main effects only. However, main effects may not be

sufficient to characterize the relationship between the response and predictors in complex situations, where predictors work together. Interaction models provide a better approximation to the response surface, improve prediction accuracy, and bring new insight on the interplay between predictors. They are useful in social, political, and economic problems to identify nontrivial interactions between covariates in modeling election results, product sales, social networks, stock market changes. One interesting application is to study the effects of combinations of various behaviors and exposures on disease rates, commonly needed in bioassay and epidemiology. In genome-wide association studies (GWAS), there is growing interest to identify the interaction (epistatic) effects of single-nucleotide polymorphisms (SNPs) (Evans et al. 2006; Manolio and Collins 2007; Kooperberg and LeBlanc 2008; Cordell 2009), since gene–gene interactions may provide critical insight on the complex biological pathways that underpin human diseases. A common class of linear models considering two-way interactions assume

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \beta_{11} X_1^2 + \beta_{12} X_1 X_2 + \cdots + \beta_{pp} X_p^2 + \varepsilon, \quad (1.1)$$

where  $Y$  is the response,  $X_1, \dots, X_p$  are covariates, and  $\varepsilon$  is the error. Marginality principle (Nelder 1977, 1994; McCullagh and Nelder 1989; McCullagh 2002) or heredity conditions (Hamada and Wu 1992; Chipman 1996; Chipman, Hamada, and Wu 1997) are generally employed to characterize the hierarchical structure between main and interaction effects. In particular, the strong heredity condition is

$$\beta_{k\ell} \neq 0 \Rightarrow \beta_k \beta_\ell \neq 0,$$

that is,  $X_k X_\ell$  is important only if its both parents  $X_k$  and  $X_\ell$  are important. The weak heredity is

$$\beta_{k\ell} \neq 0 \Rightarrow \beta_k^2 + \beta_\ell^2 \neq 0,$$

that is,  $X_k X_\ell$  is important only if at least one of  $X_k$  and  $X_\ell$  is important.

Ning Hao is Assistant Professor, Department of Mathematics, University of Arizona, Tucson, AZ 85721 (E-mail: [nhao@math.arizona.edu](mailto:nhao@math.arizona.edu)). Hao Helen Zhang is Associate Professor, Department of Mathematics, University of Arizona, Tucson, AZ 85721 (E-mail: [hzhang@math.arizona.edu](mailto:hzhang@math.arizona.edu)). The authors are partially supported by NSF Grants DMS-1309507 (Hao and Zhang), DMS-1347844 (Zhang), NIH Grants NIH/NCI R01 CA-085848 (Zhang) and P01 CA142538 (Zhang), AMS-Simons Travel Grant (Hao). The authors are grateful to Dr. Han Xiao and to the editors, associate editor, and four referees for their helpful comments and suggestions.

Color versions of one or more of the figures in the article can be found online at [www.tandfonline.com/r/jasa](http://www.tandfonline.com/r/jasa).

Interaction selection for (1.1) has lately drawn much attention in the literature. Recent works include Efron et al. (2004), Turlach (2004), Yuan, Joseph, and Lin (2007), Yuan, Joseph, and Zou (2009), Zhao, Rocha, and Yu (2009), and Choi, Li, and Zhu (2010), among others. In particular, Efron et al. (2004), Turlach (2004), and Yuan, Joseph, and Lin (2007) considered enforcing the strong heredity principle in the LARS; Yuan, Joseph, and Zou (2009) incorporated the structural relationship by imposing linear inequality constraints on coefficients; Zhao, Rocha, and Yu (2009) introduced the composite absolute penalties (CAP) to achieve hierarchy in variable selection. Choi, Li, and Zhu (2010) employed a special reparameterization of regression coefficients to enforce the heredity constraint. These procedures, except Efron et al. (2004), can be described as *joint analysis*, as they consider main and interaction effects in (1.1) altogether and make a global search over all candidate models. When  $p$  is small or moderate, joint analysis is effective in identifying important interaction effects. Some joint-analysis methods can produce consistency selection results under the strong heredity condition for a fixed  $p$  (Yuan, Joseph, and Zou 2009; Choi, Li, and Zhu 2010). However, joint-analysis methods become infeasible if  $p$  is very large. Two major limiting factors are memory requirement and computational cost. Joint analysis typically requires to store the entire augmented design matrix of size  $n \times (p^2 + 3p)/2$ . Take an example of  $n = 200$ ,  $p = 10,000$ , where the total number of entries is  $\approx 10^{10}$  and beyond the capacity of standard software such as R and MATLAB. Since sophisticated programming tools are needed to handle complex penalty structures (Zhao, Rocha, and Yu 2009; Choi, Li, and Zhu 2010) or multiple inequality constraints (Yuan, Joseph, and Zou 2009), joint analysis implementation can be extremely expensive. Furthermore, it is not clear whether selection consistency would still hold in ultrahigh-dimensional settings.

An alternative interaction selection tool is *two-stage* analysis: first select main effects only (by intentionally leaving interaction terms out), then select interactions of main effects which are previously identified. When the data dimension is very large, two-stage approaches are possibly only feasible choices for practitioners (Wu et al. 2009; Wu et al. 2010). Despite their computational advantages over joint analysis, two-stage procedures have been criticized for their theoretical validity, even for low-dimensional data with  $p < n$  (Turlach 2004).

Motivated by the above practical and theoretical concerns, we propose new greedy-type model selection procedures for high-dimensional interaction selection, study their numerical properties and performance, and provide rigorous theoretical justifications. In particular, we consider interaction-selection procedures featured with *FOR*ward selection, which are referred to as *iFOR*. Forward selection (FS) is a classical variable selection method in linear regression and it builds the model sequentially by adding one variable at a time. FS is easy to implement as it involves only simple OLS-type operations. Though the local search is suboptimal, it is a necessary compromise when dealing with high dimensionality for the sake of computation. In this article, we propose two algorithms: *iFORT* and *iFORM*. The *iFORT* is a two-stage procedure: at the first stage, it selects only main effects (all quadratic terms and interactions ignored) by FS; at the second stage, interaction terms generated under the heredity condition are considered. The *iFORM*, on the other hand,

selects main effects and interactions altogether in an iterative fashion. Compared to join-analysis procedures, the *iFOR* methods can incorporate the strong or weak heredity condition in a much simpler fashion. Their implementation does not require the storage of the entire augmented matrix, making them feasible for large problems. The memory and computational complexity are shown to be *linear* in  $p$ . In one simulation example with  $p = 10,000$  and  $n = 400$ , it takes *iFOR* fewer than 30 sec to complete the selection process. Numerical examples suggest promising performance of *iFOR* in terms of effective coverage. In addition to the new algorithms and numerical results, another major goal of this work is to investigate theoretical properties of *iFOR* estimators and understand their asymptotic behaviors. By rigorously analyzing the covariance structure between main effects and interaction terms, we prove that the *iFORT* has a sure screening property for ultrahigh-dimensional settings. This is the first theoretical justification of two-stage approaches.

The rest of this article is organized as follows. Section 2 introduces the basic model setup and the new procedures: *iFORT* and *iFORM*, under the strong heredity condition. Major theoretical results are presented in Section 3. Section 4 extends the *iFOR* to the context of the weak heredity condition. Numerical results are demonstrated in Sections 5 and 6. Final remarks are given in Section 7. All technical proofs are relegated to the Appendix.

## 2. METHODOLOGY

### 2.1 Model Setup and Notations

Given  $n$  IID observations  $(\mathbf{x}_1, Y_1), \dots, (\mathbf{x}_n, Y_n)$ , we consider a regression model with linear and second-order terms

$$Y_i = \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}^{(1)} + \mathbf{z}_i^\top \boldsymbol{\beta}^{(2)} + \varepsilon_i, \quad 1 \leq i \leq n, \quad (2.1)$$

where  $Y_i$  is a real-valued response,  $\mathbf{x}_i = (X_{i1}, \dots, X_{ip})$  is a  $p$ -dimensional vector, the vector  $\mathbf{z}_i = (X_{i1}^2, X_{i1}X_{i2}, \dots, X_{i1}X_{ip}, X_{i2}^2, X_{i2}X_{i3}, \dots, X_{ip}^2)^\top$  contains quadratic and two-way interaction terms,  $\beta_0$  is the intercept,  $\boldsymbol{\beta}^{(1)}$  and  $\boldsymbol{\beta}^{(2)}$  are, respectively, regression coefficients of linear effects and order-2 effects, and  $\varepsilon_i$  is the noise with mean zero and finite variance  $\sigma^2$ . The length of  $\mathbf{z}_i$  or  $\boldsymbol{\beta}^{(2)}$  is  $q = (p + p^2)/2$ . The entire parameter vector is  $\boldsymbol{\beta} = (\boldsymbol{\beta}^{(1)\top}, \boldsymbol{\beta}^{(2)\top})^\top$ . Throughout this article, we assume that  $E(X_{ij}) = 0$ ,  $\text{var}(X_{ij}) = 1$ ,  $E(Y_i) = 0$ ,  $\text{var}(Y_i) = 1$  in (2.1) for  $i = 1, \dots, n$  and  $j = 1, \dots, p$ . We also assume that all the quadratic effects and two-way interactions are centered, that is,  $\mathbf{z}_i = (\dots, X_{ik}X_{i\ell} - E(X_{ik}X_{i\ell}), \dots)^\top$ . This eliminates the need of the intercept term  $\beta_0$  in (2.1).

For convenience, denote  $([\mathbf{x}_i^\top]_{i=1}^n)$  as the design matrix containing only linear effects, and  $([\mathbf{x}_i^\top, \mathbf{z}_i^\top]_{i=1}^n)$  as the augmented design matrix. Define the index sets of linear and order-2 terms as

$$\mathcal{P}_1 = \{1, 2, \dots, p\}, \quad \mathcal{P}_2 = \{(k, \ell) : 1 \leq k \leq \ell \leq p\}.$$

In (2.1), any term  $\beta_j \neq 0$  or  $\beta_{jk} \neq 0$  is regarded as *relevant*; the corresponding predictor can be a linear, or quadratic, or interaction effect. We define the nonzero linear and order-2 effects as

$$\begin{aligned} \mathcal{T}_1 &= \{j : \beta_j \neq 0, j \in \mathcal{P}_1\}, \\ \mathcal{T}_2 &= \{(j, k) : \beta_{jk} \neq 0, (j, k) \in \mathcal{P}_2\}. \end{aligned}$$

The full model is  $\mathcal{F} = \mathcal{P}_1 \cup \mathcal{P}_2$  and the true model  $\mathcal{T} = \mathcal{T}_1 \cup \mathcal{T}_2$ . For any model  $\mathcal{M}$ , use  $|\mathcal{M}|$  to denote its model size, that is, the number of predictors contained in  $\mathcal{M}$ . We have  $|\mathcal{P}_1| = p$ ,  $|\mathcal{P}_2| = q$ , and  $|\mathcal{F}| = d = p + q$ . We assume  $|\mathcal{T}_1| = p_0$  and  $|\mathcal{T}_2| = q_0$ , and then the true model size  $|\mathcal{T}| = d_0 = p_0 + q_0$ . In the literature, variable selection for (2.1) has been studied by penalized least squares using the augmented matrix  $([\mathbf{x}_i^\top, \mathbf{z}_i^\top])_{i=1}^n$  as the covariates and conducting variable selection under heredity principles. They work quite well when  $p$  is moderate. But when  $p$  is big, their implementation becomes infeasible since the full model size  $d$  increases quadratically in  $p$ . For example,  $p = 50$  and  $d = 1325$ ,  $p = 500$  and  $d = 125,750$ , and  $p = 5000$  and  $d = 12,507,500$ .

Next we give a review of the FS solution path algorithm (Wang 2009), which is closely related to the interaction selection algorithms under consideration. For each  $1 \leq k \leq n$ , we use  $\mathcal{S}_k$  to denote the index of selected variables at the end of the  $k$ th step. Let  $\text{RSS}_{\mathcal{M}}$  be the residual sum of squares (RSS) using model  $\mathcal{M}$  to fit the data.

**2.1.1 Forward Selection (FS).** *Initial step:* Set  $k = 0$  and  $\mathcal{S}_0 = \emptyset$ .

*Iterative step:* set  $k = k + 1$ . If  $k > n$ , stop. Otherwise, given  $\mathcal{S}_{k-1}$ , for every  $j \in \mathcal{P}_1 \setminus \mathcal{S}_{k-1}$ , construct a candidate model  $\mathcal{M}_{j,k-1} = \mathcal{S}_{k-1} \cup \{j\}$ . Compute the  $\text{RSS}_{\mathcal{M}_{j,k-1}}$  for each  $j$ . Find  $a_k = \arg \min_{j \in \mathcal{P}_1 \setminus \mathcal{S}_{k-1}} \text{RSS}_{\mathcal{M}_{j,k-1}}$  and update  $\mathcal{S}_k = \mathcal{S}_{k-1} \cup \{a_k\}$ . Repeat this step until stop.

The FS algorithm produces a solution path consisting of  $n$ -nested models  $\mathcal{S}_1 \subset \dots \subset \mathcal{S}_n$ , where  $\mathcal{S}_k = \{a_1, \dots, a_k\}$  for  $1 \leq k \leq n$ . When  $p \gg n$ , the FS automatically terminates after  $n$  steps when RSS reduces to zero. Since the solution path of the FS depends only on the subspaces spanned by the predictor vectors (column vectors in the design matrix), centering and standardization does not change the solution path. Wang (2009) showed the screening consistency of the FS for main-effect selection under the ultrahigh-dimensional setup.

One straightforward way of extending the FS to the interaction selection is to apply FS directly to model (2.1), ignoring the hierarchical structure. We name this procedure FS2 to distinguish it from the usual FS for main effect selection. Based on our empirical experience, FS2 works well for small and moderate  $p$  in sparse settings. In Section 3, we prove that FS2 has a sure screening property for interaction selection under some regularity conditions. However, similar to joint-analysis methods, the implementation of FS2 requires to store the entire augmented design matrix or call the features repeatedly during computation, making it difficult for high-dimensional data analysis.

## 2.2 New Methods: iFOR

We propose two forward-selection-based algorithms for interaction selection. The new algorithms naturally incorporate the marginality or heredity principles (Yuan, Joseph, and Zou 2009; Zhao, Rocha, and Yu 2009; Choi, Li, and Zhu 2010), without invoking complex constraints or optimization tools as done in joint analysis. Throughout this section, we use  $\mathcal{C}$  to denote the candidate index set which consists of all the terms to be considered for selection in the immediately following step.

We first describe the two-stage approach (iFORT) algorithm. At Stage 1, only main effects are selected by FS, while all of the

order-2 terms left out of the model. Denote the selected main-effect set by  $\widehat{\mathcal{M}}$ . At Stage 2, we expand  $\widehat{\mathcal{M}}$  by adding all the two-way interactions within  $\widehat{\mathcal{M}}$  and then implement FS on the expanded set while forcing  $\widehat{\mathcal{M}}$  to stay in the final model.

### 2.2.1 Two-Stage iFOR (iFORT).

- Stage 1. Define  $\mathcal{C} = \mathcal{P}_1$ . Implement FS on  $\mathcal{C}$ . The resulting solution path is  $\{\mathcal{S}_t^{(1)}, t = 1, 2, \dots\}$ , and the selected main effects are  $\widehat{\mathcal{M}} = \{j_1, \dots, j_{t_1}\}$ .
- Stage 2. Update  $\mathcal{C} = \widehat{\mathcal{M}} \cup \{(k, l) : k \in \widehat{\mathcal{M}} \text{ and } l \in \widehat{\mathcal{M}}\}$ . Implement FS on  $\mathcal{C}$  by forcing-in  $\widehat{\mathcal{M}}$ . Denote the solution path by  $\{\mathcal{S}_{t_1+t}^{(2)}, t = 1, 2, \dots\}$ .

The iFORT is simple, fast, and feasible to implement for high-dimensional data analysis. It does not require complex optimization tools, and the strong heredity condition is automatically satisfied in the final true model by forcing-in  $\widehat{\mathcal{M}}$ . If the model is sparse, the number of important linear effects  $p_0$  would be small, so the number of terms considered at Stage 2 would be much smaller than  $(p^2 + p)/2$ . Theoretical properties of iFORT are studied in Section 3.

The iFORT separately selects main effects and order-2 terms at two stages. Alternatively, one may select them altogether under the marginality principle, and this leads to another new algorithm iFORM. The main idea of the iFORM is to apply FS to a submodel of model (2.1) indexed by a dynamic candidate set  $\mathcal{C}$ . At step  $t$ , we use  $\mathcal{S}_t$ ,  $\mathcal{M}_t$ , and  $\mathcal{C}_t$ , respectively, to represent the index set of all selected effects, selected main effects, and current candidate set. Initially,  $\mathcal{C} = \mathcal{P}_1$ , that is, all the main effects. Then the candidate set  $\mathcal{C}$  grows gradually by adding two-way interactions between the main effects already in the model. In other words, we update  $\mathcal{C}_t$  by defining  $\mathcal{C}_t = \mathcal{P}_1 \cup \{(k, \ell) : k, \ell \in \mathcal{M}_t\}$ .

### 2.2.2 iFOR Under Marginality Principle (iFORM).

- Step 1. (Initialization) Set  $\mathcal{S}_0 = \emptyset$ ,  $\mathcal{M}_0 = \emptyset$  and  $\mathcal{C}_0 = \mathcal{P}_1$ .
- Step 2. (Selection) In the  $t$ th step with given  $\mathcal{S}_{t-1}$ ,  $\mathcal{C}_{t-1}$ , and  $\mathcal{M}_{t-1}$ , forward regression is used to select one more predictor from  $\mathcal{C}_{t-1} \setminus \mathcal{S}_{t-1}$  into the model. We add the selected one into  $\mathcal{S}_{t-1}$  to get  $\mathcal{S}_t$ . We also update  $\mathcal{C}_t$  and  $\mathcal{M}_t$  if the newly selected predictor is a main effect. Otherwise,  $\mathcal{C}_t = \mathcal{C}_{t-1}$  and  $\mathcal{M}_t = \mathcal{M}_{t-1}$ .
- Step 3. (Solution path) Iterating Step 2, we get a solution path  $\{\mathcal{S}_t : t = 1, 2, \dots, D\}$ .

In the above algorithm,  $D$  is chosen as a reasonable upper bound of  $d_0$  (the total number of important effects), to terminate the procedure. A direct advantage of the iFORM is that it allows the interactions to enter the model early, making it easier to select weak relevant main effects. Moreover, when we decide the optimal model along the solution path, we only need to use model size selection criteria, say BIC, once, while for iFORT, we have to use BIC twice which may cause additional error in practice even if the solution path is correct. Our empirical experience also suggests the iFORM has better finite sample performance. The screening consistency of iFORM is shown in Section 3.3.

To select the optimal model from the FS path, we consider the use of BIC. There are two types of BIC proposed in the



literature, the standard BIC

$$\text{BIC}_1(\widehat{\mathcal{M}}) = \log \hat{\sigma}_{\widehat{\mathcal{M}}}^2 + n^{-1}|\widehat{\mathcal{M}}| \log(n),$$

and the BIC specially designed for high-dimensional data (Chen and Chen 2008)

$$\text{BIC}_2(\widehat{\mathcal{M}}) = \log \hat{\sigma}_{\widehat{\mathcal{M}}}^2 + n^{-1}|\widehat{\mathcal{M}}|(\log(n) + 2 \log d^*),$$

where  $d^*$  is the number of predictors in the full model. The only difference between two BICs is the extra term  $2 \log d^*$  in  $\text{BIC}_2$ . Chen and Chen (2008) derived  $\text{BIC}_2$  by controlling the false discovery rate (FDR) and showed that it is selection consistent if  $d^* = O(n^\xi)$  for some  $\xi > 0$ . Wang (2009) showed its selection consistency for FS under ultrahigh-dimensional setup  $d^* = O(\exp(n^\xi))$ . Since we deal with the ultrahigh-dimensional data, we use  $\text{BIC}_2$  for iFORM and the first stage of iFORT. At the second stage of iFORT, since the number of candidate predictors is already dramatically reduced after the first stage,  $\text{BIC}_1$  is more appropriate. Section 5 demonstrates their effective performance, in terms of coverage, false discovery control, and prediction accuracy.

### 2.3 Computational Complexity and Practical Issues

We show that the computational complexity of iFOR procedures is linear in  $p$ , which explains their feasibility for  $p \gg n$ . The FS algorithm described in Section 2.1 is equivalent to the following procedure. At each step, the response is regressed on the most correlated covariate, and the residual is calculated and used as the new response in next step. After the most correlated covariate (say,  $X_1$ ) is selected, all other covariates are regressed on  $X_1$ , and then the covariates are substituted by the corresponding normalized residuals, which are used as the new covariates in next step. Note that the computation complexity of each step is  $O(nm)$ , where  $n$  is the sample size and  $m$  is the number of predictors in the candidate set. First, the absolute correlations between the response and all covariates in the current candidate set are calculated at each step, so the complexity is  $O(nm)$ . Once the most correlated covariate is selected, the response and all other covariates are regressed on it, whose cost is also  $O(nm)$ . For the iFORT and iFORM algorithms, the number of steps to build the whole solution path is at most  $n$ , so the number of main effects selected is not larger than  $n$ . This implies that, at each step, there are at most  $p + n(n+1)/2$  predictors in the candidate set, that is,  $m \leq p + n(n+1)/2$  holds for any step. Therefore, the overall complexity is  $nO(n(p + n(n+1)/2)) = O(n^2p + n^4)$ , which is linear in  $p$ .

The parameter  $D$  controls the length of the solution path for the iFORM. Since the final model is chosen based on BIC by comparing all the models along the path, the final model select results is not sensitive to the exact value of  $D$  as long as it is reasonably large. In practice, though  $d_0$  is unknown, it is reasonable to assume that  $d_0$  is much smaller than  $n$  in high-dimensional sparse regression problems (Fan and Lv 2008). In our numerical study, we have tried  $D = n/2, n/3, n/4$  and obtained the same results since  $D > d_0$ . In general, we suggest  $D = n/2$ .

## 3. THEORETICAL RESULTS

We study theoretical properties of iFOR. In literature, a long-term concern about two-stage methods is their theoretical validity, as the main effect selection at Stage 1 is conducted under a misspecified working model. In Section 3.1, we first prove that the iFORT is able to capture all important main effects under ultrahigh-dimensional settings. This fundamental result provides rigorous justifications for two-stage methods. In Section 3.2, we prove that iFORT can identify all important interactions consistently with probability tending to one under heredity conditions. The screening consistency of iFORM is shown in Section 3.3.

### 3.1 Screening Consistency of iFORT for Main Effects

Recall the true model  $\mathcal{T} = \mathcal{T}_1 \cup \mathcal{T}_2$ , where  $\mathcal{T}_1 \subset \mathcal{P}_1$  and  $\mathcal{T}_2 \subset \mathcal{P}_2$ . For any square matrix  $A$ , denote its smallest and largest eigenvalues, respectively, by  $\lambda_{\min}(A)$  and  $\lambda_{\max}(A)$ . Denote the covariance matrices of main linear effects and interactions (i.e., all degree 2 monomials), respectively, by  $\Sigma^{(1)}$  and  $\Sigma^{(2)}$ . The total covariant matrix is  $\Sigma$ . The following regularity conditions are needed.

- (C1). Normality:  $X_{i1}, \dots, X_{ip}$  are jointly normal and marginally standard normal.  $\varepsilon_i \sim N(0, \sigma^2)$  is independent of  $X_{i1}, \dots, X_{ip}$ .
- (C2). Covariance matrix: We assume that there exist two constants  $0 < \tau_{\min} < \tau_{\max} < \infty$ , such that  $2\tau_{\min} < \lambda_{\min}(\Sigma^{(1)}) \leq \lambda_{\max}(\Sigma^{(1)}) < \tau_{\max}/2$ .
- (C3). Signal strength: We assume that  $\|\beta\| \leq C_\beta$  for some positive constant  $C_\beta$  and  $\beta_{\min} \geq \nu_\beta n^{-\xi_{\min}}$ , where  $\beta_{\min} = \min_{\kappa \in \mathcal{T}} |\beta_\kappa|$  and  $\xi_{\min} > 0$ .
- (C4). Dimensionality and sparsity: There exist positive constants  $\xi, \xi_0$ , and  $\nu$  such that  $\log p \leq \nu n^\xi$ ,  $d_0 \leq \nu n^{\xi_0}$ , and  $\xi + 6\xi_0 + 12\xi_{\min} < 1$ ,  $\xi < \frac{1}{2}$ .

*Remark 3.1.* Conditions (C1) to (C4) are standard in the literature of ultrahigh-dimensional inference (Fan and Lv 2008; Zhang and Huang 2008). The normality assumption (C1) is extensively used in the past literature to facilitate proof (Fan and Lv 2008; Zhang and Huang 2008; Wang 2009). (C2) requires the design matrix of main effects to be well-behaved. (C1) and (C2) together assure the Sparse Riesz Condition (Zhang and Huang 2008); see the proof in Appendix for more details. (C3) requires that the smallest signal should not decay too fast, otherwise they cannot be consistently identified; see (Fan and Peng 2004) for more discussions. (C4) allows the dimension  $p$  to diverge with  $n$  at an exponential rate, or the NP dimensionality (Fan and Lv 2008). Intuitively, one would expect that stronger conditions are needed to develop theory for interaction selection due to their heavier tails. However, to our satisfaction, conditions (C1) to (C4) are comparable to those used in the main-effect selection literature (Fan and Lv 2008; Wang 2009). The only difference is  $\xi < \frac{1}{2}$  in (C4) while  $\xi < 1$  is used in Wang (2009), due to heavier tails of interaction terms. Note if  $X_{1j}$ 's are sub-Gaussian with  $E(e^{aX_{1j}^2}) < b$  for positive constants  $a$  and  $b$ , typically, we can only bound a product term by  $E(e^{2aX_{1j}X_{1k}}) < b^2$ .

*Theorem 3.1. (sure screening of main effects)* Define  $K = 2\tau_{\max}\nu C_\beta^2 \tau_{\min}^{-2} \nu_\beta^{-4}$ . Under conditions (C1)-(C4), the first stage

of iFORT is screening consistent for the main effects. For  $t_1 \geq K \nu n^{2\xi_0 + 4\xi_{\min}}$ ,

$$\mathbf{P}\left(\mathcal{T}_1 \subset \mathcal{S}_t^{(1)}\right) \rightarrow 1 \quad \text{as } n \rightarrow \infty. \quad (3.1)$$

Next we give insight on why screening consistency (3.1) still holds for selection under a misspecified model. A key observation is Lemma 1 in Appendix, which says, under (C1),

$$\Sigma = \begin{pmatrix} \Sigma^{(1)} & 0 \\ 0 & \Sigma^{(2)} \end{pmatrix}.$$

The block structure of  $\Sigma$  guarantees that ignored important interaction terms have minimal affects to the procedure at Stage 1. Imagining if there are some nonzero terms on the right top corner of  $\Sigma$ , we have to put some strong and complicated conditions on  $\Sigma$  to guarantee screening consistency.

*Remark 3.2.* In general, as long as  $\Sigma$  has a block structure, Theorem 1 holds even without normality. Here (C1) is used as a convenient and sufficient condition to assure the covariance block structure. There are other weaker but sufficient conditions (C1)' or (C1)'', which can replace (C1):

- (C1)'.  $X_{i1}, \dots, X_{ip}$  are sub-Gaussian marginally, and their joint distribution is symmetric with respect to  $\mathbf{0}$ .  
 (C1)'.  $X_{i1}, \dots, X_{ip}$  are sub-Gaussian marginally, and their joint distribution has vanished third moments.

### 3.2 Screening Consistency of iFORT for Interaction Effects

After Stage 1, the iFORT essentially reduces the main effect dimensionality from  $p$  to  $t_1 = o(n^{\frac{1}{3}})$ , which is significant if  $p \gg n$ . Using (C4), it is straightforward to show  $2\xi_0 + 4\xi_{\min} < \frac{1}{3}$ . Next we study the asymptotic behaviors of iFORT for interaction selection under the strong heredity:

(H1). Strong heredity condition:  $\beta_{k\ell} \neq 0 \Rightarrow \beta_k \beta_\ell \neq 0$ .

Under (H1), the interaction selection of iFORT at Stage 2 does not need to deal with high-dimensional predictors any more, since the number of selected main effects is  $o(n^{\frac{1}{3}})$ . Even if we include all interactions within the selected model at Stage 1, the final model has cardinality  $o(n^{\frac{2}{3}})$ . Corollary 1 gives the fundamental result: the iFORT is screening consistent for interaction selection under the heredity condition for ultrahigh-dimensional settings.

*Corollary 3.1. (sure screening of interactions)* Conditional on (3.1) and (H1), for  $t_2 \geq K \nu n^{2\xi_0 + 4\xi_{\min}}$ ,

$$\mathbf{P}\left(\mathcal{T} \subset \mathcal{S}_{t_1+t_2}^{(2)}\right) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

*Remark 3.3.* The strong heredity is necessary to ensure the consistency of two-stage procedures for screening interaction terms. Otherwise, if  $X_1 X_2$  is important but neither  $X_1$  nor  $X_2$ , then the main effects are not guaranteed to be identified at Step 1, and consequently, their interaction  $X_1 X_2$  might not be considered at the second step. We also point out that the strong heredity condition is actually not that strong with a simple illustration. Consider the case  $p = 2$ , where the full model space (for simplicity, ignoring two quadratic terms) can be represented by the

parameter set  $(\beta_0, \beta_1, \beta_2, \beta_{12})^\top$  in  $\mathbb{R}^4$ . The strong heredity condition covers the entire  $\mathbb{R}^4$  except for a couple of subsets, such as  $\{\beta_0 = 0, \beta_1^2 + \beta_2^2 + \beta_{12}^2 > 0\}$  and  $\{\beta_1 \beta_2 = 0, \beta_{12}^2 > 0\}$ . The excluded subsets have zero mass in  $\mathbb{R}^4$ , so the strong heredity condition is met by most models. This implies that the iFORT methods work for a generic model.

### 3.3 Screening Consistency of FS2 and iFORM

Naively, we can use any one-stage variable selection tool to fit (1.1) directly (as long as computation is feasible), ignoring the hierarchical structure. Though the model consistency or screening consistency result (Zhao and Yu 2006; Wang 2009; Fan and Lv 2011) could be generalized to the context of interaction selection, the extension of earlier proofs is not straightforward due to the heavy tails of interaction effects. Actually, all the existing proof technique would require some regularity conditions on the eigenvalues of  $\Sigma^{(2)}$ . Next, we establish the screening consistency of FS2 under conditions that are related only to  $\Sigma^{(1)}$ .

- (C2a). Covariance matrix: Assume that there exist two constants  $0 < \tau_{\min} < \frac{1}{4} < 1 < \tau_{\max} < \infty$ , such that  $\sqrt{\tau_{\min}} < \lambda_{\min}(\Sigma^{(1)}) \leq \lambda_{\max}(\Sigma^{(1)}) < \sqrt{\tau_{\max}}/4$ .  
 (C4a). Dimensionality and sparsity: There exist positive constants  $\xi$ ,  $\xi_0$ , and  $\nu$ , such that  $\log p \leq \nu n^\xi$ ,  $d_0 \leq \nu n^{\xi_0}$ , and  $\xi + 6\xi_0 + 12\xi_{\min} < \frac{1}{2}$ .

There is no essential difference between (C2a) and (C2). (C2a) is used only for easy presentation. (C4a) is slightly stronger than (C4). Note that under (C1) and (C2a), the population and sample covariance matrices  $\Sigma$  and  $\hat{\Sigma}$  can be well controlled because  $\Sigma^{(2)}$  can be explicitly represented by  $\Sigma^{(1)}$ . See Lemma 3 in the appendix. On the other hand, the screening consistency result below strongly depends on the normality condition (C1) since there is no easy way to capture the structure of  $\Sigma^{(2)}$  by  $\Sigma^{(1)}$  without normality condition.

*Theorem 3.2.* Under conditions (C1), (C2a), (C3), and (C4a), FS2 is screening consistent. For  $t \geq K \nu n^{2\xi_0 + 4\xi_{\min}}$ ,

$$\mathbf{P}\left(\mathcal{T} \subset \mathcal{S}_t^{\text{FS2}}\right) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

The screening consistency of iFORM is implied in the proof of Theorem 2, as iFORM is similar to FS2 but with a restrictive candidate set each step.

*Corollary 3.2.* Under conditions (C1), (C2a), (C3), (C4a), and (H1), iFORM is screening consistent. For  $t \geq K \nu n^{2\xi_0 + 4\xi_{\min}}$ ,

$$\mathbf{P}(\mathcal{T} \subset \mathcal{S}_t) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

## 4. EXTENSIONS TO WEAK HEREDITY

In some real applications, the weak heredity provides a useful alternative for the underlying model structure. Under the weak heredity, for a two-way interaction effect to be active, at least one of the parent effects need to be effective. In this section, we generalize the iFOR algorithms described in Section 2 to satisfy the weak heredity condition. Similar to the strong heredity situation, both iFOR algorithms under the weak heredity are easy to implement.

(H2). Weak heredity condition:  $\beta_{k\ell} \neq 0 \Rightarrow \beta_k^2 + \beta_\ell^2 \neq 0$ .

#### 4.1 iFORT Under Weak Heredity (iFORT-w)

- Stage 1. Define  $\mathcal{C} = \mathcal{P}_1$ . Implement FS on  $\mathcal{C}$ . The resulting solution path is  $\{\mathcal{S}_t^{(1)}, t = 1, 2, \dots\}$ , and the selected main effects are  $\widehat{\mathcal{M}} = \{j_1, \dots, j_{t_1}\}$ .
- Stage 2. Update  $\mathcal{C} = \widehat{\mathcal{M}} \cup \{(k, l) : k \in \widehat{\mathcal{M}} \text{ or } l \in \widehat{\mathcal{M}}\}$ . Implement FS on  $\mathcal{C}$  by forcing-in  $\widehat{\mathcal{M}}$ . Denote the solution path by  $\{\mathcal{S}_{t_1+t}^{(2)}, t = 1, 2, \dots\}$ .

For the iFORM extension, after selecting any new linear term, we need to expand the candidate set by including all of its interactions with the other linear effects. Denote by  $\mathcal{M}_t$  the index set of selected linear effects at Step  $t$ . Under the weak heredity condition, we update  $\mathcal{C}_t$  as

$$\mathcal{C}_t = \mathcal{P}_1 \cup \{(k, \ell) : k \text{ or } \ell \in \mathcal{M}_t\}.$$

For each  $t$ , we use  $\mathcal{S}_t$ ,  $\mathcal{M}_t$ , and  $\mathcal{C}_t$  to represent the index set of selected model, selected main effects and candidates set at Step  $t$ , respectively.

#### 4.2 iFORM Under Weak Heredity (iFORM-w)

The iFORM-w algorithm is the same as the iFOR algorithm given in Section 2.2.2, except for the rule of updating  $\mathcal{C}_t$ .

*Remark 4.1.* The weak heredity condition is slightly more flexible than the strong heredity condition, and generally chooses a larger model. In practice, the weak heredity is more useful to identify important interactions with one weak parent effect (Yuan, Joseph, and Zou 2009). With regard to the computation speed, since the candidate set size at each step is larger than in the strong heredity case, the iFORT-w and iFORM-w are slower than the iFORT and iFORM.

### 5. NUMERICAL STUDIES

#### 5.1 Experiments and Setup

We demonstrate performance of the iFOR methods in various  $p \gg n$  scenarios, including the regression settings with independent predictors, predictors with autoregressive (AR) correlation structure, compound symmetry (CS) correlation, and more complex settings as considered in Fan and Song (2010). We consider forward-based joint analysis (FS2), and the proposed forward-based procedures iFORT, iFORM, iFORT-w, iFORM-w. In the literature, there are other two-step procedures which are not based on forward selection such as Mendel (Wu et al. 2009) and screen-and-clean (Wu et al. 2010). For comparison, we also include two such procedures, iMART1 and iMART2. The iMART1 screens main effects based on marginal correlation at Step 1, that is, those that exceed a threshold are retained as candidate predictors, and then conducts the LASSO penalized regression on the expanded dictionary consisting of all the candidate predictors and their pairwise interaction terms at Step 2. The iMART2 first screens main effects by marginal correlation, then screens the pairwise products of the main effect candidates by pairwise correlation, and then implements the LASSO to obtain the final model. The standard BIC is used to select the tuning parameter of LASSO. The oracle (ORACL) procedure is also presented as the gold standard, which is generally not available in practice.

Recall that the full model is  $\mathcal{F} = \mathcal{P}_1 \cup \mathcal{P}_2$ ,  $|\mathcal{P}_1| = p$ ,  $|\mathcal{P}_2| = q$ . The true model is  $\mathcal{T} = \mathcal{T}_1 \cup \mathcal{T}_2$ ,  $|\mathcal{T}_1| = p_0$ ,  $|\mathcal{T}_2| = q_0$ . We run  $M = 100$  Monte Carlo simulations and report their average performance in selecting linear effects and interactions, estimating coefficients, and making predictions. For the  $m$ th replication, let  $\widehat{\beta}^{(m)}$  denote the fitted regression coefficients,  $\widehat{\mathcal{T}}_1^{(m)}$  and  $\widehat{\mathcal{T}}_2^{(m)}$  respectively denote the selected linear effects and interactions. To evaluate linear effect selection, we report

1. Coverage probability (Cov)  $\sum_{m=1}^M I(\mathcal{T}_1 \subset \widehat{\mathcal{T}}_1^{(m)})/M$ ,
2. Percentage of correct zeros (Cor0)  $\sum_{m=1}^M \sum_{j=1}^p I(\widehat{\beta}_j^{(m)} = 0, \beta_j = 0)/[M(p - p_0)]$ ,
3. Percentage of incorrect zeros (Inc0):  $\sum_{m=1}^M \sum_{j=1}^p I(\widehat{\beta}_j^{(m)} = 0, \beta_j \neq 0)/[Mp_0]$ .
4. Exact selection probability (Ext)  $\sum_{m=1}^M I(\mathcal{T}_1 = \widehat{\mathcal{T}}_1^{(m)})/M$ .

For interaction selection, we report

1. Coverage probability (iCov)  $\sum_{m=1}^M I(\mathcal{T}_2 \subset \widehat{\mathcal{T}}_2^{(m)})/M$ ,
2. Percentage of correct zeros (iCor0)  $\sum_{m=1}^M \sum_{(j,k) \in \mathcal{P}_2} I(\widehat{\beta}_{jk}^{(m)} = 0, \beta_{jk} = 0)/[M(q - q_0)]$ ,
3. Percentage of incorrect zeros (iInc0)  $\sum_{m=1}^M \sum_{(j,k) \in \mathcal{P}_2} I(\widehat{\beta}_{jk}^{(m)} = 0, \beta_{jk} \neq 0)/[Mq_0]$ .
4. Exact selection probability (iExt)  $\sum_{m=1}^M I(\mathcal{T}_2 = \widehat{\mathcal{T}}_2^{(m)})/M$ .

The overall model selection is measured by the model size  $\sum_{m=1}^M |\widehat{\mathcal{T}}_1^{(m)} \cup \widehat{\mathcal{T}}_2^{(m)}|/M$ . For estimation, we report the squared root of mean-squared error (RMSE)  $\sum_{m=1}^M [\sum_{j=1}^p (\widehat{\beta}_j^{(m)} - \beta_j)^2 + \sum_{(j,k) \in \mathcal{P}_2} (\widehat{\beta}_{jk}^{(m)} - \beta_{jk})^2]^{1/2}/M$ . For the prediction error, we report the out-of-sample  $R^2$  (Rsqr):

$$100\% \times \left\{ 1 - \frac{\sum_{i=1}^n [Y_i^* - x_i^* \widehat{\beta}^{(1)} - z_i^* \widehat{\beta}^{(2)}]^2}{\sum_{i=1}^n (Y_i^* - Y^*)^2} \right\},$$

where the test data  $(\mathbf{X}_i^*, Y_i^*)$ ,  $i = 1, \dots, n$  are generated independently from the same distribution as the training set, and  $Y^* = \frac{1}{n} \sum_{i=1}^n Y_i^*$ . A larger Rsqr suggests a better prediction. The standard error of Rsqr is reported as well. We also report the average computation time.

#### 5.2 Simulation Results

In all the examples, we generate the response  $Y$  from model (2.1) with  $\sigma = 2, 3, 4$ .

*Example 5.1.* (Independent predictors) Let  $(n, p, p_0, q_0) = (100, 500, 4, 4)$ .  $\mathbf{X}$ 's are iid from  $MVN(\mathbf{0}, I_p)$ . The true  $\beta^{(1)} = (3, 0, 3, 0, 0, 3, 0, 0, 0, 3, 0_{490})$ , so  $\mathcal{T}_1 = \{1, 3, 6, 10\}$ . The important interaction set  $\mathcal{T}_2 = \{(1, 3), (1, 6), (3, 10), (6, 10)\}$  with coefficient 2.

*Example 5.2.* (Autoregressive correlation) Consider the same setup as Example 1, except that  $\mathbf{X}$  follows  $MVN$  with mean  $\mathbf{0}$  and  $\text{cov}(X_j, X_k) = 0.5^{|j-k|}$  for  $1 \leq j, k \leq p$ .

*Example 5.3.* (High dimensional: AR) Let  $(n, p, p_0, q_0) = (400, 5000, 10, 10)$ . We generate  $\mathbf{X}$  from  $MVN$  with mean  $\mathbf{0}$  and  $\text{cov}(X_j, X_k) = 0.5^{|j-k|}$ . The true  $\beta^{(1)} = (3, 3, 3,$

Table 1. Results of Example 1,  $(n, p, p_0, q_0) = (100, 500, 4, 4)$ , independent predictors

	Linear term selection				Interaction selection				Size and prediction			
	Cov	Cor0	Inc0	Ext	iCov	iCor0	iInc0	Iext	size	RMSE	Rsqr	sdR
$\sigma = 2$												
FS2	0.34	1.00	0.55	0.34	0.22	1.00	0.73	0.14	3.32	5.27	30.00	4.25
iMART1	0.71	1.00	0.07	0.42	0.71	1.00	0.17	0.29	10.02	2.39	77.21	2.07
iMART2	0.69	1.00	0.08	0.67	0.06	1.00	0.45	0.06	8.24	3.55	67.29	1.74
iFORT-w	0.84	1.00	0.06	0.82	0.77	1.00	0.17	0.60	7.48	1.63	81.99	2.12
iFORT	0.84	1.00	0.06	0.82	0.84	1.00	0.11	0.63	7.71	1.35	84.57	1.84
iFORM-w	0.98	1.00	0.01	0.98	0.90	1.00	0.08	0.60	8.13	1.05	88.24	1.32
iFORM	1.00	1.00	0.00	0.96	0.99	1.00	0.00	0.99	8.07	0.63	91.89	0.25
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	8.00	0.59	92.03	0.22
$\sigma = 3$												
FS2	0.03	1.00	0.83	0.03	0.01	1.00	0.95	0.01	1.06	6.91	4.93	2.18
iMART1	0.59	1.00	0.12	0.30	0.53	1.00	0.29	0.13	11.20	3.48	56.91	2.98
iMART2	0.58	1.00	0.12	0.56	0.06	1.00	0.53	0.04	8.46	4.08	53.72	2.16
iFORT-w	0.63	1.00	0.17	0.62	0.22	1.00	0.61	0.17	5.22	3.86	54.64	2.45
iFORT	0.63	1.00	0.17	0.60	0.63	1.00	0.29	0.48	6.45	2.62	65.52	2.73
iFORM-w	0.69	1.00	0.16	0.69	0.24	1.00	0.59	0.13	5.39	3.72	56.16	2.54
iFORM	0.98	1.00	0.01	0.95	0.74	1.00	0.14	0.74	7.52	1.48	79.31	1.08
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	8.00	0.89	83.61	0.43
$\sigma = 4$												
FS2	0.01	1.00	0.90	0.01	0.00	1.00	0.98	0.00	0.67	7.13	0.39	1.72
iMART1	0.49	1.00	0.15	0.28	0.34	1.00	0.44	0.05	11.63	4.49	39.12	2.73
iMART2	0.46	1.00	0.15	0.43	0.05	1.00	0.59	0.04	8.50	4.60	41.03	2.10
iFORT-w	0.41	1.00	0.33	0.41	0.02	1.00	0.90	0.01	3.36	5.17	33.99	2.15
iFORT	0.41	1.00	0.33	0.38	0.40	1.00	0.50	0.29	5.02	4.00	44.46	2.89
iFORM-w	0.22	1.00	0.50	0.22	0.03	1.00	0.90	0.03	2.53	5.67	25.92	2.28
iFORM	0.67	1.00	0.19	0.67	0.15	1.00	0.64	0.15	4.72	3.93	48.14	2.46
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	8.00	1.19	73.97	0.63

3, 3, 2, 2, 2, 2, 2,  $\mathbf{0}_{4990}$ ). The nonzero interaction set is  $\mathcal{T}_2 = \{(1, 2), (1, 3), (2, 3), (2, 5), (3, 4), (6, 8), (6, 10), (7, 8), (7, 9), (9, 10)\}$ , and their coefficients are (2, 2, 2, 2, 2, 1, 1, 1, 1, 1).

*Example 5.4.* (High dimensional: AR) We increase the dimension  $p = 10,000$  in Example 3.

*Example 5.5.* (High dimensional: FS2010) We use the same setup as in Example 4, except that  $\mathbf{X}$  has a more complex covariance structure as considered in Fan and Song (2010). First, we generate  $X_j, j = 1, \dots, 50$  independently from the standard normal distribution. Then we define

$$X_k = \sum_{j=1}^s X_j (-1)^{j+1} / 5 + \sqrt{25 - s} / 5 \epsilon_k, \quad k = p - 50, \dots, p,$$

with  $s = 10$  and  $\{\epsilon_k\}_{k=p-49}^{50}$  follow the standard normal distribution.

*Example 5.6.* (Weak heredity) We use the same setup as in Example 3, except the nonzero interaction set  $\mathcal{T}_2 = \{(1, 2), (1, 13), (2, 3), (2, 15), (3, 4), (6, 10), (6, 18), (7, 9), (7, 18), (10, 19)\}$  and the corresponding coefficients (2, 2, 2, 2, 2, 1, 1, 1, 1, 1). Note that the weak heredity condition holds here.

Three additional examples, Examples 7 to 9, are listed in the online available supplementary material due to the page limit. In particular, the compound symmetry (CS) correlation

is considered in Examples 7 and 9. The numerical results are summarized in the following Tables 1–6 and Tables S1–S3 in the supplementary material.

We first summarize the results for Examples 1–5, where the strong heredity condition holds. All the methods perform reasonably well in most of the settings, including the high-dimensional cases with  $p = 5000$  and  $p = 10,000$ , as long as the noise level is not too high. Overall speaking, the iFORM is the best among all the methods in terms of both model selection and prediction performance. The iFORM method has the smallest RMSE, the largest out-of-the-sample  $R^2$ , and the highest exact coverage probability for main effects and interactions. When  $\sigma = 2$ , the iFORM’s performance is quite close to the ORACL procedure. The performance of iFORT is sensitive to the dimensionality and noise level. In particular, when  $p$  is large and the noise level is high, it may miss some important main effects in Stage 1, although the result may be improved by using less aggressive selection criteria such as AIC and standard BIC. On the other hand, iFORM consistently gives higher coverage of important main effects and interactions than iFORT, which supports our motivation for the dynamic selection procedure. The FS2 has the worst performance, and it fails to run when  $p$  is 5000 or larger. Both iMART1 and iMART2 are reasonably fast and perform well, sometimes quite competitive in prediction. But when the covariance structure is complex, their performance is



Table 2. Results of Example 2,  $(n, p, p_0, q_0) = (100, 500, 4, 4)$ , AR(0.5) correlation

	Linear term selection				Interaction selection				Size and prediction			
	Cov	Cor0	Inc0	Ext	iCov	iCor0	iInc0	lext	size	RMSE	Rsqr	sdR
$\sigma = 2$												
FS2	0.55	1.00	0.31	0.55	0.43	1.00	0.52	0.24	5.14	3.76	58.43	3.77
iMART1	0.81	1.00	0.05	0.65	0.81	1.00	0.12	0.29	10.43	2.08	82.71	2.11
iMART2	0.80	1.00	0.05	0.80	0.11	1.00	0.41	0.10	7.96	3.23	75.54	1.54
iFORT-w	0.71	1.00	0.11	0.71	0.65	1.00	0.25	0.55	7.04	2.29	79.69	2.18
iFORT	0.71	1.00	0.11	0.71	0.71	1.00	0.20	0.58	7.28	1.96	81.80	2.08
iFORM-w	0.96	1.00	0.01	0.94	0.83	1.00	0.11	0.50	8.18	1.30	89.02	1.22
iFORM	0.98	1.00	0.01	0.90	0.98	1.00	0.02	0.94	8.05	0.73	92.45	0.70
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	8.00	0.61	93.36	0.18
$\sigma = 3$												
FS2	0.14	1.00	0.60	0.14	0.03	1.00	0.91	0.02	2.18	6.19	26.54	2.83
iMART1	0.75	1.00	0.07	0.46	0.59	1.00	0.21	0.11	11.95	3.15	70.70	1.95
iMART2	0.73	1.00	0.07	0.71	0.09	1.00	0.45	0.07	8.01	3.62	67.01	1.52
iFORT-w	0.58	1.00	0.18	0.58	0.21	1.00	0.60	0.19	5.21	3.93	59.93	2.27
iFORT	0.58	1.00	0.18	0.57	0.58	1.00	0.32	0.46	6.51	2.91	67.35	2.53
iFORM-w	0.68	1.00	0.15	0.68	0.29	1.00	0.54	0.20	5.77	3.59	62.39	2.32
iFORM	0.92	1.00	0.04	0.86	0.64	1.00	0.19	0.62	7.28	1.82	79.30	1.59
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	8.00	0.91	86.18	0.36
$\sigma = 4$												
FS2	0.03	1.00	0.76	0.03	0.00	1.00	0.97	0.00	1.24	6.91	11.72	2.24
iMART1	0.67	1.00	0.10	0.41	0.41	1.00	0.31	0.01	12.01	4.00	56.82	1.98
iMART2	0.67	1.00	0.09	0.67	0.05	1.00	0.50	0.03	7.94	4.03	57.53	1.53
iFORT-w	0.29	1.00	0.33	0.29	0.02	1.00	0.89	0.02	3.34	5.33	41.31	1.82
iFORT	0.29	1.00	0.33	0.29	0.27	1.00	0.56	0.22	4.97	4.41	48.54	2.40
iFORM-w	0.24	1.00	0.41	0.24	0.02	1.00	0.89	0.02	3.01	5.62	35.73	2.10
iFORM	0.61	1.00	0.18	0.60	0.15	1.00	0.59	0.15	4.99	3.88	55.68	2.08
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	8.00	1.21	77.73	0.54

Table 3. Results of Example 3,  $(n, p, p_0, q_0) = (400, 5000, 10, 10)$ , AR correlation

	Linear term selection				Interaction selection				Size and prediction			
	Cov	Cor0	Inc0	Ext	iCov	iCor0	iInc0	lext	size	RMSE	Rsqr	sdR
$\sigma = 2$												
iMART1	1.00	1.00	0.00	1.00	0.69	1.00	0.04	0.69	19.66	1.00	97.76	0.05
iMART2	0.99	1.00	0.00	0.99	0.02	1.00	0.34	0.02	16.89	2.29	95.17	0.20
iFORT-w	0.00	1.00	0.33	0.00	0.03	1.00	0.29	0.03	14.13	5.51	90.89	0.35
iFORT	0.00	1.00	0.33	0.00	0.00	1.00	0.57	0.00	14.80	6.45	86.86	0.55
iFORM-w	1.00	1.00	0.00	1.00	0.93	1.00	0.01	0.62	20.28	0.88	97.91	0.03
iFORM	1.00	1.00	0.00	1.00	0.98	1.00	0.00	0.37	20.74	0.82	97.93	0.03
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	0.79	97.94	0.03
$\sigma = 3$												
iMART1	1.00	1.00	0.00	1.00	0.31	1.00	0.10	0.25	19.40	1.49	95.30	0.08
iMART2	1.00	1.00	0.00	1.00	0.02	1.00	0.34	0.02	16.90	2.37	93.00	0.21
iFORT-w	0.00	1.00	0.36	0.00	0.00	1.00	0.37	0.00	13.14	5.97	87.35	0.42
iFORT	0.00	1.00	0.36	0.00	0.00	1.00	0.60	0.00	13.98	6.78	83.70	0.59
iFORM-w	1.00	1.00	0.00	1.00	0.22	1.00	0.14	0.11	19.12	1.70	95.34	0.07
iFORM	1.00	1.00	0.00	1.00	0.37	1.00	0.10	0.15	19.90	1.52	95.50	0.06
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	1.00	95.79	0.05
$\sigma = 4$												
iMART1	1.00	1.00	0.00	0.98	0.01	1.00	0.14	0.08	19.66	1.89	92.26	0.14
iMART2	0.97	1.00	0.00	0.97	0.01	1.00	0.36	0.01	16.84	2.56	89.96	0.23
iFORT-w	0.00	1.00	0.40	0.00	0.00	1.00	0.48	0.00	11.50	6.68	82.31	0.41
iFORT	0.00	1.00	0.40	0.00	0.00	1.00	0.66	0.00	12.72	7.38	79.12	0.59
iFORM-w	0.85	1.00	0.02	0.85	0.00	1.00	0.27	0.00	17.57	2.53	91.47	0.17
iFORM	0.94	1.00	0.01	0.94	0.02	1.00	0.22	0.01	18.81	2.24	91.97	0.13
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	1.24	92.91	0.09

Table 4. Results of Example 4 with ultrahigh-dimensional data,  $(n, p, p_0, q_0) = (400, 10,000, 10, 10)$

	Linear term selection				Interaction selection				Size and prediction			
	Cov	Cor0	Inc0	Ext	iCov	iCor0	iInc0	Iext	size	RMSE	Rsqr	sdR
$\sigma = 2$												
iMART1	0.98	1.00	0.00	0.98	0.54	1.00	0.06	0.54	19.44	1.13	97.63	0.88
iMART2	0.97	1.00	0.00	0.97	0.01	1.00	0.34	0.01	16.84	2.32	94.95	0.21
iFORT-w	0.00	1.00	0.35	0.00	0.00	1.00	0.29	0.00	13.98	5.67	90.50	0.34
iFORT	0.00	1.00	0.35	0.00	0.00	1.00	0.60	0.00	14.55	6.67	86.35	0.54
iFORM-w	1.00	1.00	0.00	0.99	0.95	1.00	0.01	0.68	20.29	0.85	97.91	0.04
iFORM	1.00	1.00	0.00	0.97	0.99	1.00	0.00	0.47	20.66	0.82	97.92	0.03
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	0.79	97.94	0.03
$\sigma = 3$												
iMART1	0.98	1.00	1.00	0.98	0.25	1.00	0.11	0.22	19.25	1.55	95.23	0.09
iMART2	0.97	1.00	0.00	0.97	0.02	1.00	0.36	0.02	16.86	2.46	92.78	0.22
iFORT-w	0.00	1.00	0.38	0.00	0.00	1.00	0.40	0.00	12.63	6.22	86.65	0.41
iFORT	0.00	1.00	0.38	0.00	0.00	1.00	0.65	0.00	13.57	7.09	82.92	0.59
iFORM-w	1.00	1.00	0.00	0.99	0.16	1.00	0.16	0.13	18.78	1.77	95.20	0.08
iFORM	1.00	1.00	0.00	0.98	0.35	1.00	0.11	0.18	19.63	1.58	95.39	0.07
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	1.01	95.78	0.06
$\sigma = 4$												
iMART1	0.97	1.00	1.00	0.97	0.12	1.00	0.16	0.07	19.66	2.01	92.05	0.13
iMART2	0.97	1.00	0.00	0.97	0.02	1.00	0.36	0.02	16.91	2.59	89.85	0.25
iFORT-w	0.00	1.00	0.42	0.00	0.00	1.00	0.52	0.00	11.03	6.91	81.24	0.48
iFORT	0.00	1.00	0.42	0.00	0.00	1.00	0.68	0.00	12.56	7.50	78.78	0.63
iFORM-w	0.83	1.00	0.02	0.83	0.00	1.00	0.29	0.00	17.44	2.63	91.35	0.14
iFORM	0.97	1.00	0.00	0.97	0.01	1.00	0.23	0.01	18.43	2.26	91.90	0.13
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	1.25	92.91	0.09

Table 5. Results of Example 5,  $(n, p, p_0, q_0) = (400, 10,000, 10, 10)$ , FS2010 correlation

	Linear term selection				Interaction selection				Size and prediction			
	Cov	Cor0	Inc0	Ext	iCov	iCor0	iInc0	Iext	size	RMSE	Rsqr	sdR
$\sigma = 2$												
iMART1	0.30	1.00	0.09	0.29	0.03	1.00	0.18	0.24	19.56	2.33	88.90	0.54
iMART2	0.30	1.00	0.09	0.30	0.00	1.00	0.38	0.00	16.69	2.95	85.97	0.55
iFORT-w	0.68	1.00	0.07	0.64	0.75	1.00	0.09	0.73	18.87	1.70	91.59	0.59
iFORT	0.68	1.00	0.07	0.64	0.68	1.00	0.11	0.30	19.48	1.78	91.17	0.67
iFORM-w	1.00	1.00	0.00	0.89	1.00	1.00	0.00	0.90	20.39	0.72	95.12	0.05
iFORM	0.98	1.00	0.01	0.90	0.98	1.00	0.02	0.98	19.98	0.86	94.46	0.47
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	0.70	95.13	0.05
$\sigma = 3$												
iMART1	0.28	1.00	0.10	0.21	0.03	1.00	0.20	0.10	20.78	2.66	82.91	0.76
iMART2	0.28	1.00	0.10	0.28	0.00	1.00	0.39	0.00	17.40	3.17	79.98	0.69
iFORT-w	0.60	1.00	0.08	0.60	0.21	1.00	0.24	0.20	17.20	2.40	84.39	0.62
iFORT	0.60	1.00	0.08	0.59	0.60	1.00	0.13	0.29	19.08	2.11	85.32	0.68
iFORM-w	0.96	1.00	0.01	0.93	0.32	1.00	0.15	0.31	18.89	1.60	87.99	0.33
iFORM	0.96	1.00	0.01	0.89	0.64	1.00	0.07	0.63	19.57	1.29	88.67	0.51
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	0.87	89.98	0.10
$\sigma = 4$												
iMART1	0.22	1.00	0.13	0.11	0.01	1.00	0.27	0.00	21.06	3.24	74.74	0.76
iMART2	0.22	1.00	0.13	0.22	0.00	1.00	0.43	0.00	17.77	3.61	71.77	0.75
iFORT-w	0.34	1.00	0.13	0.34	0.00	1.00	0.48	0.00	14.35	3.49	73.68	0.66
iFORT	0.34	1.00	0.13	0.32	0.33	1.00	0.21	0.13	17.77	2.85	76.51	0.68
iFORM-w	0.86	1.00	0.04	0.86	0.00	1.00	0.44	0.00	15.83	2.80	77.78	0.53
iFORM	0.91	1.00	0.02	0.91	0.02	1.00	0.34	0.02	16.85	2.41	79.46	0.47
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	1.07	83.61	0.15

Table 6. Results of Example 6 for the weak heredity case,  $(n, p, p_0, q_0) = (400, 5000, 10, 10)$

	Linear term selection				Interaction selection				Size and prediction			
	Cov	Cor0	Inc0	Ext	iCov	iCor0	iInc0	lext	size	RMSE	Rsqr	sdR
$\sigma = 2$												
iMART1	1.00	1.00	0.00	0.92	0.00	1.00	0.56	0.00	16.04	3.68	90.86	0.13
iMART2	1.00	1.00	0.00	0.93	0.00	1.00	0.67	0.00	13.79	3.82	89.99	0.17
iFORT-w	0.04	1.00	0.20	0.04	0.06	1.00	0.24	0.06	15.83	3.95	93.62	0.28
iFORT	0.04	1.00	0.20	0.04	0.00	1.00	0.65	0.00	13.53	5.36	87.11	0.33
iFORM-w	1.00	1.00	0.00	1.00	1.00	1.00	0.00	0.91	20.09	0.76	97.80	0.03
iFORM	0.96	1.00	0.00	0.96	0.00	1.00	0.60	0.00	14.61	3.71	90.73	0.12
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	0.75	97.81	0.03
$\sigma = 3$												
iMART1	0.99	1.00	0.00	0.87	0.00	1.00	0.57	0.00	16.58	3.84	88.33	0.17
iMART2	1.00	1.00	0.00	0.85	0.00	1.00	0.67	0.00	13.93	3.87	87.74	0.19
iFORT-w	0.00	1.00	0.24	0.00	0.00	1.00	0.34	0.00	14.51	4.49	89.96	0.33
iFORT	0.00	1.00	0.24	0.00	0.00	1.00	0.68	0.00	12.94	5.70	84.02	0.36
iFORM-w	1.00	1.00	0.00	1.00	0.21	1.00	0.15	0.19	18.66	1.60	94.85	0.07
iFORM	0.79	1.00	0.03	0.79	0.00	1.00	0.63	0.00	14.12	3.96	88.00	0.19
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	0.97	95.42	0.06
$\sigma = 4$												
iMART1	0.96	1.00	0.00	0.67	0.00	1.00	0.59	0.00	17.42	4.03	85.09	0.20
iMART2	1.00	1.00	0.00	0.67	0.00	1.00	0.67	0.00	14.30	3.96	84.62	0.22
iFORT-w	0.00	1.00	0.29	0.00	0.00	1.00	0.48	0.00	12.60	5.37	84.16	0.44
iFORT	0.00	1.00	0.29	0.00	0.00	1.00	0.72	0.00	11.91	6.30	79.57	0.43
iFORM-w	0.86	1.00	0.02	0.86	0.00	1.00	0.30	0.00	17.06	2.38	90.77	0.14
iFORM	0.48	1.00	0.08	0.48	0.00	1.00	0.68	0.00	13.05	4.47	83.94	0.29
ORACL	1.00	1.00	0.00	1.00	1.00	1.00	0.00	1.00	20.00	1.20	92.26	0.09

not very good. This can be seen in Example 5, and Examples 7 and 8 in the supplementary material.

In Example 6, the weak heredity condition holds, and therefore we expect that the iFOR under the weak heredity constraint should perform better than those under the strong heredity. The results in Table 6 confirm this pattern: iFORM-w (or iFORT-w) gives better performance than iFORM (or iFORT) in terms of both model selection and prediction accuracy. Since the strong heredity methods make an incorrect model structure assumption, they suffer by missing some important interactions. For example, if  $\sigma = 2$ , iFORM-w is the only method showing a high exact selection probability (91%) for important interactions.

Finally, we illustrate the quality of the solution path by the *hit-rate* plot. In each plot, the x-axis denotes the solution path steps  $\{1, 2, \dots, S\}$ , and the y-axis represents the “hit rate” which is defined as the percentage of important terms recovered up to step  $s$ . Denote the true model size by  $d_0$ . The ideal hit plot (given by ORACL) should show a linearly increasing trend with slope  $1/d_0$  within the first  $d_0$  steps and then stays at 1 afterward. For the graph clarity, we only draw the hit rates for the strong heredity methods. Figure 1 plots the hit-rates for Examples 1 and 2 with the moderate  $p = 500$ . Here  $d_0 = 8$ , so we choose  $S = 20$ . Based on Figure 1, the iFORM has the highest hit rate among all, very close to the oracle. For  $\sigma = 2$  and 3, its hit rate is more than 95% after 20 steps; for the more difficult case  $\sigma = 4$ , iFORM still achieves approximately 90% hit rate. The iFORT is slightly worse than iFORM, with rates 90%, 80%, 70%, respectively, for  $\sigma = 2, 3, 4$ . The FS2 has the lowest hit rate, only 20% when  $\sigma = 4$ . Figure 2 plots the hit rates for the large  $p$ . Since  $d_0 = 20$ ,

we choose  $S = 40$ . The FS2 is not shown in Figure 2, because it fails to run. Again, iFORM has the highest hit rate among all (except the oracle). The iFORT is slightly worse, about 80% hit rate in most cases.

Table 7 summarizes the average computation time (seconds per run) for each procedure. The machine we used equips Intel Core (TM) i7-2600 CPU @ 3.40GHZ with 4.00 GB ram. Since the time difference is small for varying  $\sigma$ , we only present the results for  $\sigma = 2$ . When  $p$  is moderately large, the FS2 is slowest, taking 16.40 sec in average for Example 1. The iFORT and iFORM are the fastest, taking 0.04 and 0.08 sec in Example 1, which is more than 100 times faster than FS2. The weak heredity methods are slower than their strong heredity counterparts. When  $p$  is large, the FS2 fails to run, while the iFOR procedures are still amazingly fast. When  $p = 5000$ , it takes 11.39 (and 16.06) sec for iFORT (and iFORM). When  $p = 10,000$ , it takes 22.13 (and 29.17) sec for iFORT (and iFORM). The weak heredity methods now take significantly more time. Overall, the iFORM appears the most promising in terms of both performance and speed.

6. REAL DATA ANALYSIS

We analyze two real datasets, the inbred mouse microarray gene expression dataset (Lan et al. 2006) and the supermarket data (Wang 2009). The inbred mouse microarray dataset contains 60 mouse arrays, with 31 from female mice and 29 from male mice, respectively. Each array measures the expression values of 22,690 genes. The response is a continuous phenotypic

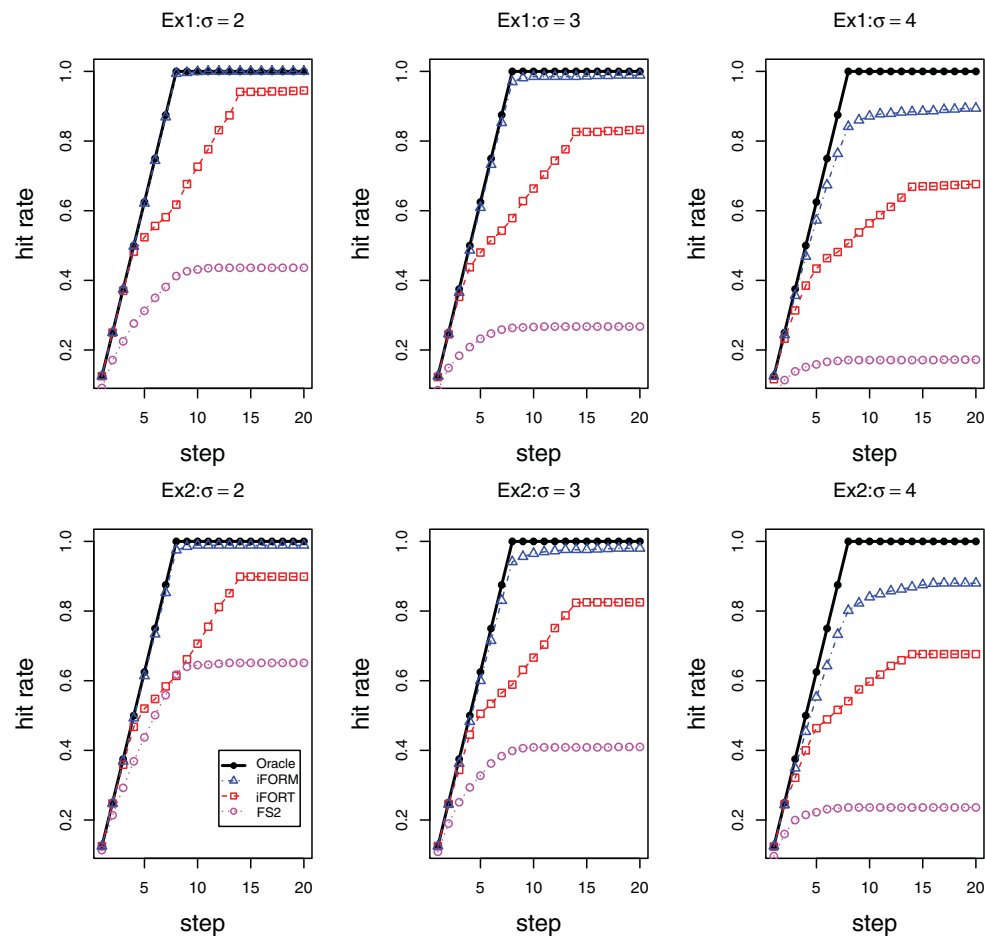


Figure 1. The hit-rate plots for the moderate  $p$  for ORACL, iFORM, iFORT, and FS2.

variable measured by real-time RT-PCR, stearyl-CoA desaturase 1 (SCD1). The supermarket dataset collects daily sale information of a major supermarket located in northern China, with  $n = 464$  and  $p = 6398$ . The response  $Y$  is the number of customers per day, and the predictors  $\mathbf{X}$  are sale volumes of various products. The supermarket manager is interested in the relationship between the number of customers and the sale volume of certain products. For convenience, the response and all predictors are centered to zero and standardized to have a unit variance prior to the analysis.

The proposed methods are applied to both datasets. To assess the prediction performance of the procedures, we randomly sample  $n_1$  observations to form the training set, and use the remaining  $n - n_1$  observations as the test data to compute the out-of-sample  $R^2$  for the final model. We use  $n_1 = 50$  in the inbred mouse data analysis and use  $n_1 = 400$  for the supermarket data analysis. The results are summarized in Table 8. It is observed that the iFOR methods give similar performance for both datasets.

7. DISCUSSION

In this article, we tackle the important problem of interaction selection for ultrahigh-dimensional data. The task is both computationally and theoretically challenging. We propose a new class of procedures, called iFOR, and study their numerical and theoretical properties. One major advantage of the proposed methods is their computation feasibility. The code is simple and fast. Theoretically we show that the iFOR can discover all relevant interactions consistently, even if the dimension increases exponentially fast with the sample size. Our numerical examples suggest that the new methods, especially iFORM, give promising performance for ultrahigh-dimensional data.

We use the extended BIC (Chen and Chen 2008) to select a final model from the solution path in this work. Since the motivation of the extended BIC is to control FDR, it tends to be conservative in real-data analysis. It would be interesting to study the performance of other selection criteria such as AIC and cross-validation for iFOR methods in the future. Other

Table 7. Average computation time (in seconds) for  $\sigma = 2$

Example	$n$	$p$	$(p_0, q_0)$	FS2	iFORT	iFORT-w	iFORM	iFORM-w
1	100	500	(4, 4)	16.40	0.04	0.36	0.09	0.51
2	100	500	(4, 4)	16.29	0.04	0.34	0.08	0.50
3	400	5000	(10, 10)	—	11.39	80.66	16.06	126.21
4	400	10,000	(10, 10)	—	22.13	144.29	29.17	209.65



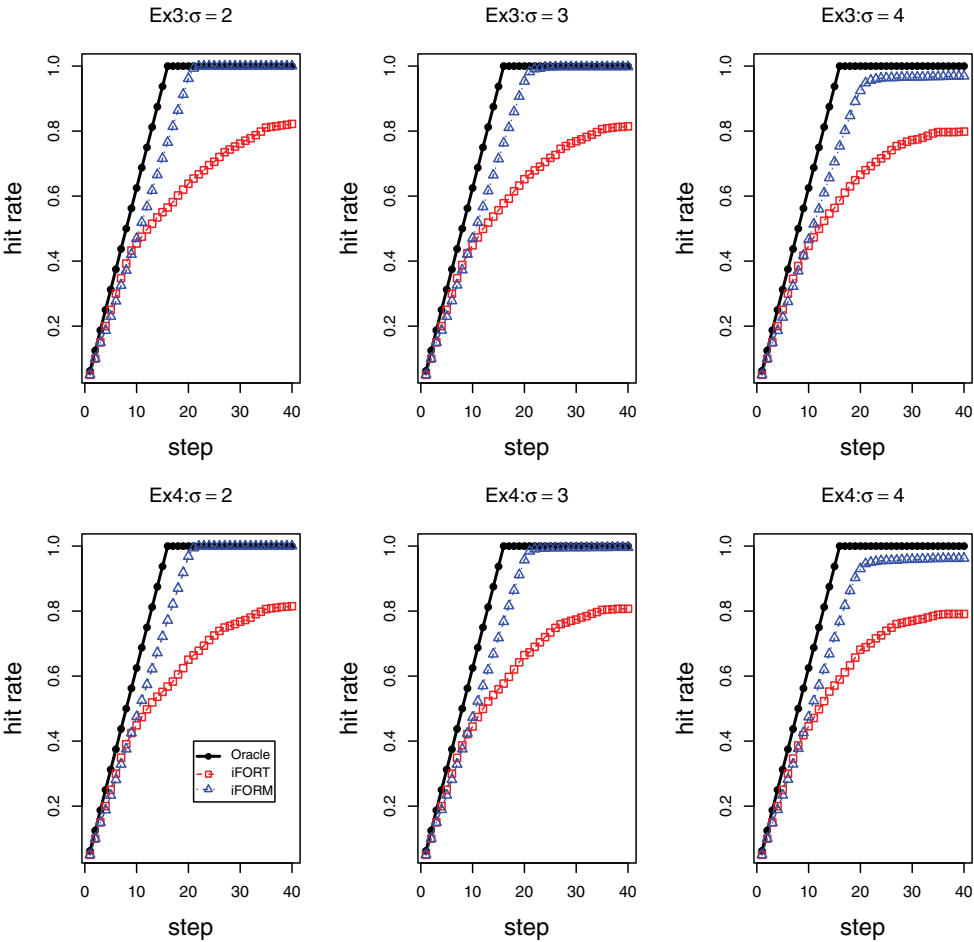


Figure 2. The hit-rate plots for large  $p$  for ORACL, iFORM, iFORT, and FS2.

works of interest include the generalization of the iFOR to other loss functions in GLM or nonparametric regression, and how to improve computational efficiency of penalized methods with the iFOR ideas.

In practice, higher-order interactions are useful to uncover multiway relationships among predictors for complex problems where two-way interactions are not sufficient. The proposed methods can be readily extended to selecting higher-order interactions, by including higher-order products of predictors in the candidate set. No essential change is needed in the computational algorithm, except that the enlarged candidate set will demand extra time. When considering higher order interaction models, one should tune the model properly to avoid the overfitting. The interpretation of higher order interactions should be cautious as well. The topic is worth a full investigation.

APPENDIX A: TOTAL COVARIANCE MATRIX

In this section, we work on the total covariance matrix  $\Sigma$  and show it is determined by the covariance matrix  $\Sigma^{(1)}$  of main effects under the Gaussian assumption (C1).

Let us temporarily ignore the index labeling the order of observations, and denote by  $X_j$  for  $1 \leq j \leq p$  the main effects and  $Z_{jk} = X_j X_k - E(X_j X_k)$  for  $(j, k) \in \mathcal{P}_2$  the interactions. Let  $\Sigma^{(1)} = (\sigma_{ij})$  denote covariance matrix of the main effects  $X_1, \dots, X_p$ . The first two lemmas help us to characterize the total covariance matrix  $\Sigma$ .

*Lemma A.1.* Under the normality condition (C1), for  $\forall j, k, \ell$ ,  $\text{cov}(X_j, Z_{k\ell}) = 0$ , which implies

$$\Sigma = \begin{pmatrix} \Sigma^{(1)} & 0 \\ 0 & \Sigma^{(2)} \end{pmatrix}.$$

*Proof.*  $\text{cov}(X_j, Z_{k\ell}) = \text{cov}(X_j, X_k X_\ell) = E(X_j X_k X_\ell) - E(X_j)E(X_k X_\ell) = 0$ . The conclusion still holds if the joint density of  $X_1, \dots, X_p$  is symmetric with respect to the original point  $\mathbf{0}$ .  $\square$

*Lemma A.2.* Under the normality condition (C1),

$$\text{cov}(Z_{ij}, Z_{k\ell}) = \text{cov}(X_i X_j, X_k X_\ell) = \sigma_{ik} \sigma_{j\ell} + \sigma_{i\ell} \sigma_{jk}. \tag{A.1}$$

Table 8. Prediction performance: The average out-of-sample  $R^2$  for iFOR methods

Dataset	iFORT	iFORT-w	iFORM	iFORM-w
Inbred mouse data	60.73 (1.15)	58.46 (1.37)	60.22 (1.15)	60.31 (1.28)
Supermarket data	88.91 (0.17)	88.42 (0.19)	88.66 (0.18)	86.61 (0.22)

*Proof.* This lemma follows directly from the following useful formula (Bar and Ditttrich 1971):

$$\begin{aligned} E(X_i X_j X_k X_\ell) &= E(X_i X_j)E(X_k X_\ell) + E(X_i X_k)E(X_j X_\ell) \\ &\quad + E(X_i X_\ell)E(X_j X_k) - 2E(X_i)E(X_j)E(X_k)E(X_\ell). \end{aligned}$$

□

Let  $A = (A_{ij})$  be an  $N \times N$  matrix. In linear algebra, a  $K \times K$  submatrix is called a principal submatrix if it is of the form  $A_{\mathcal{I}} = (A_{\ell_i \ell_j})$ , where  $\mathcal{I}$  is an index set  $\mathcal{I} = \{1 \leq \ell_1 < \dots < \ell_K \leq N\}$ . Here with slight abuse of this conception, we allow arbitrary order for the index set  $\mathcal{I}$ . For example, let  $\mathcal{I} = \{2, 1\}$  and

$$A_{\mathcal{I}} = \begin{pmatrix} A_{22} & A_{21} \\ A_{12} & A_{11} \end{pmatrix}$$

is still called a principal submatrix in this article.

Based on the formula (A.1), we can decompose  $\Sigma^{(2)}$  to a sum  $\Sigma_1^{(2)} + \Sigma_2^{(2)}$ . In fact, we have

**Lemma A.3.** Both  $\Sigma_1^{(2)}$  and  $\Sigma_2^{(2)}$  are principal submatrices of  $\Sigma^{(1)} \otimes \Sigma^{(1)}$ .

*Proof.* The Kronecker product (Laub 2005)  $\Sigma^{(1)} \otimes \Sigma^{(1)}$  is a  $p^2 \times p^2$  matrix whose rows and columns are both indexed by the set  $\mathcal{P}_1 \times \mathcal{P}_1$ . The entry corresponding to the index  $(ij, k\ell)$  is  $\sigma_{ij}\sigma_{k\ell}$ . By formula (A.1), both  $\Sigma_1^{(2)}$  and  $\Sigma_2^{(2)}$  are  $\frac{p(p+1)}{2} \times \frac{p(p+1)}{2}$  principal submatrices of  $\Sigma^{(1)} \otimes \Sigma^{(1)}$ . □

**Lemma A.4.** Under the conditions (C1) and (C2a), we have

$$2\tau_{\min} < \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) < \tau_{\max}/2. \quad (\text{A.2})$$

*Proof.* By Laub (2005) Theorem 13.12, the eigenvalues of  $\Sigma^{(1)} \otimes \Sigma^{(1)}$  are  $\lambda_i \lambda_j$ ,  $1 \leq i, j \leq p$ , if the eigenvalues of  $\Sigma^{(1)}$  are  $\lambda_1, \dots, \lambda_p$ . Therefore, under condition (C2a), we have

$$\tau_{\min} < \lambda_{\min}(\Sigma^{(1)} \otimes \Sigma^{(1)}) \leq \lambda_{\max}(\Sigma^{(1)} \otimes \Sigma^{(1)}) < \tau_{\max}/4.$$

By Lemma A.3, the eigenvalues of  $\Sigma_1^{(2)}$  and  $\Sigma_2^{(2)}$  are also bounded by  $\tau_{\min}$  and  $\tau_{\max}/4$ , so

$$2\tau_{\min} < \lambda_{\min}(\Sigma^{(2)}) \leq \lambda_{\max}(\Sigma^{(2)}) < \tau_{\max}/2.$$

It is straightforward to get (A.2). □

## APPENDIX B. A BERNSTEIN INEQUALITY AND ITS APPLICATION

In this section, we study a Bernstein-type inequality and its applications in bounding the eigenvalues of submatrices of sample covariance matrix  $\hat{\Sigma}$ , which is crucial in the proofs of theorems. For any index set  $\mathcal{M}$ ,  $\hat{\Sigma}_{\mathcal{M}}$  denotes the principal submatrix corresponding to  $\mathcal{M}$ .

**Lemma B.1.** Let  $W_1, \dots, W_n$  be independent random variables with mean zero and variances bounded by  $\sigma^2 \geq 1$ . Assume for some  $0 < \alpha < 1$ ,

$$E(|W_i|^{3(1-\alpha)} e^{t|W_i|^\alpha}) \leq A, \quad \text{for all } 1 \leq i \leq n, \quad 0 \leq t \leq T. \quad (\text{B.1})$$

Then for  $x > (\frac{2A}{T^2})^{\frac{1}{1-\alpha}}$ ,

$$\begin{aligned} \mathbf{P}\left(\left|\sum_{i=1}^n W_i\right| \geq x\right) &\leq 2 \exp\left\{-\frac{x^2}{2(n\sigma^2 + x^{2-\alpha}/T)}\right\} + \sum_{i=1}^n \mathbf{P}(|W_i| \geq x). \end{aligned} \quad (\text{B.2})$$

*Proof.* Let  $W_i^* = W_i \cdot I_{(-\infty, x]}(W_i)$ . Then

$$\mathbf{P}\left(\sum_{i=1}^n W_i \geq x\right) \leq \mathbf{P}\left(\sum_{i=1}^n W_i^* \geq x\right) + \sum_{i=1}^n \mathbf{P}(W_i \geq x). \quad (\text{B.3})$$

For  $W_i^* \geq 0$ , we have

$$e^{tW_i^*} \leq 1 + tW_i^* + \frac{t^2}{2}W_i^{*2} + \sum_{k=3}^{\infty} \frac{t^k}{k!}|W_i|^{k\alpha+3(1-\alpha)}x^{(k-3)(1-\alpha)}. \quad (\text{B.4})$$

Note that (B.4) is true also for  $W_i^* < 0$  because of the monotonicity of function  $f(u) = e^u - 1 - u - u^2/2$ . □

It is easy to get  $E|W_i|^{k\alpha+3(1-\alpha)} \leq \frac{k!A}{T^k}$  from (B.1). Moreover, we have  $E(W_i^*) \leq 0$ ,  $\text{var}(W_i^*) \leq \sigma^2$  from definition. Taking expectation of (B.4),

$$\begin{aligned} E(e^{tW_i^*}) &\leq 1 + \frac{t^2\sigma^2}{2} + \sum_{k=3}^{\infty} \frac{2A}{T^2x^{1-\alpha}} \frac{1}{2} \left(\frac{x^{1-\alpha}}{T}\right)^{k-2} t^k \\ &\leq 1 + \frac{t^2\sigma^2}{2} + \frac{t^2}{2} \sum_{k=3}^{\infty} \left(\frac{tx^{1-\alpha}}{T}\right)^{k-2} \\ &\leq 1 + \frac{t^2\sigma^2}{2(1 - tx^{1-\alpha}/T)}, \end{aligned} \quad (\text{B.5})$$

when  $|tx^{1-\alpha}/T| < 1$ .

Let  $t = \frac{x}{n\sigma^2 + x^{2-\alpha}/T}$ . By Markov inequality

$$\begin{aligned} \mathbf{P}\left(\sum_{i=1}^n W_i^* \geq x\right) &\leq e^{-tx} E(e^{t\sum_{i=1}^n W_i^*}) \\ &\leq e^{-tx} \prod_{i=1}^n E(e^{tW_i^*}) \\ &\leq e^{-tx} \left(1 + \frac{t^2\sigma^2}{2(1 - tx^{1-\alpha}/T)}\right)^n \\ &\leq \exp\left\{-\frac{x^2}{n\sigma^2 + x^{2-\alpha}/T}\right\} \\ &\quad \times \left(1 + \frac{1}{2n} \frac{x^2}{n\sigma^2 + x^{2-\alpha}/T}\right)^n \\ &\leq \exp\left\{-\frac{x^2}{2(n\sigma^2 + x^{2-\alpha}/T)}\right\}. \end{aligned}$$

Therefore,

$$\mathbf{P}\left(\sum_{i=1}^n W_i \geq x\right) \leq \exp\left\{-\frac{x^2}{2(n\sigma^2 + x^{2-\alpha}/T)}\right\} + \sum_{i=1}^n \mathbf{P}(W_i \geq x).$$

Apply the same technique to  $-W_i^*$  and combine the results, we can get (B.2).

The following is the Lemma 1 in Wang (2009), which is useful in the proof of Theorem 1.

**Lemma B.2.** Under condition (C1) and (C2), for  $m = o(n^{\frac{1}{3}-\frac{1}{3}\xi})$ ,  $\mathcal{M} \subset \mathcal{P}_1$ ,

$$\mathbf{P}\left(\tau_{\min} \leq \min_{|\mathcal{M}| \leq m} \lambda_{\min}(\hat{\Sigma}_{\mathcal{M}}) \leq \max_{|\mathcal{M}| \leq m} \lambda_{\max}(\hat{\Sigma}_{\mathcal{M}}) \leq \tau_{\max}\right) \rightarrow 1. \quad (\text{B.6})$$

Furthermore, under condition (C4), (B.6) holds for  $m = O(n^{2\xi_0+4\xi_{\min}}) = o(n^{\frac{1}{3}-\frac{1}{3}\xi})$ .

**Lemma B.3.** Let  $W_1, \dots, W_n$  be independent random variables with zero mean such that  $E(e^{T_0|W_i|^\alpha}) \leq A_0$  for constants  $T_0 > 0$ ,  $A_0 > 0$ , and  $0 < \alpha < 1$ . Then, for a sequence  $a_n \rightarrow \infty$  with  $a_n = o(n^{\frac{\alpha}{2(2-\alpha)}})$ , there exist constants  $c_1, c_2$ , such that

$$\mathbf{P}(|W_1 + \dots + W_n| > \sqrt{n}a_n) \leq c_1 \exp(-c_2 a_n^2). \quad (\text{B.7})$$

*Proof.* The condition  $E(e^{T_0|W_i|^\alpha}) \leq A_0$  implies  $\text{var}(W_i) \leq \sigma^2$ ,  $E(|W_i|^2 e^{T|W_i|^\alpha}) \leq A$  and  $E(|W_i|^{3(1-\alpha)} e^{T|W_i|^\alpha}) \leq A$  for some constants  $\sigma^2$ ,  $T$  and  $A$ . By Lemma B.1, we have

$$\mathbf{P}\left(\left|\sum_{i=1}^n W_i\right| \geq x\right) \leq 2 \exp\left\{-\frac{x^2}{2(n\sigma^2 + x^{2-\alpha}/T)}\right\} + \sum_{i=1}^n \mathbf{P}(|W_i| \geq x).$$

Let  $x = \sqrt{n}a_n$ . Then

$$\begin{aligned} \exp\left\{-\frac{x^2}{2(n\sigma^2 + x^{2-\alpha}/T)}\right\} &= \exp\left\{-\frac{na_n^2}{2(n\sigma^2 + n^{\frac{2-\alpha}{2}}a_n^{2-\alpha}/T)}\right\} \\ &= \exp\left\{-\frac{a_n^2}{2\sigma^2 + o(1)}\right\}. \end{aligned}$$

On the other hand, by Markov inequality

$$\begin{aligned} \mathbf{P}(|W_i| \geq x) &= \mathbf{P}(W_i^2 e^{T|W_i|^\alpha} \leq x^2 e^{Tx^\alpha}) \leq Ax^{-2} \exp\{-Tx^\alpha\} \\ &\leq \frac{A}{na_n^2} \exp\{-Ta_n^2/o(1)\}. \end{aligned}$$

Hence,  $\sum_{i=1}^n \mathbf{P}(|W_i| \geq x) \leq \frac{A}{a_n^2} \exp\{-Ta_n^2/o(1)\}$ . And (B.7) is easily obtained.  $\square$

*Remark B.1.* We are interested in the case that  $W_i = X_{ij}X_{ik}X_{i\ell}$ , where  $X_{ij}$ ,  $X_{ik}$ ,  $X_{i\ell}$  are joint normal and marginally standard normal. It is easy to see that  $W_i$  satisfies  $E(e^{\frac{1}{4}|W_i|^{\frac{2}{3}}}) \leq \sqrt{2}$  and  $\text{var}(W_i) \leq 30$ . Therefore, (B.7) holds for  $c_1 = 3$ ,  $c_2 = 1/61$  when  $n$  is sufficiently large.

To show Theorem 2, we have to obtain an analog of Lemma B.2 for arbitrary submodel  $\mathcal{M}$ . We start from a generalization of Lemma A3 in Bickel and Levina (2008).

*Lemma B.4.* Let  $W_1, \dots, W_n$  be independent random variables with zero mean such that  $E(e^{T_0|W_i|^\alpha}) \leq A_0$  for constants  $T_0 > 0$ ,  $A_0 > 0$ , and  $0 < \alpha \leq 1$ . Then there exist constants  $c_3, c_4$ , for  $0 < \epsilon \leq 1$

$$\mathbf{P}(|W_1 + \dots + W_n| > n\epsilon) \leq c_3 \exp(-c_4 n^\alpha \epsilon^2). \quad (\text{B.8})$$

*Proof.* The condition  $E(e^{T_0|W_i|^\alpha}) \leq A_0$  implies  $\text{var}(W_i) \leq \sigma^2$ ,  $E(|W_i|^2 e^{T|W_i|^\alpha}) \leq A$  and  $E(|W_i|^{3(1-\alpha)} e^{T|W_i|^\alpha}) \leq A$  for some constants  $\sigma^2$ ,  $T$  and  $A$ . When  $\alpha < 1$ , by Lemma B.1,

$$\mathbf{P}\left(\left|\sum_{i=1}^n W_i\right| \geq x\right) \leq 2 \exp\left\{-\frac{x^2}{2(n\sigma^2 + x^{2-\alpha}/T)}\right\} + \sum_{i=1}^n \mathbf{P}(|W_i| \geq x).$$

Let  $x = n\epsilon$ . Then

$$\begin{aligned} \exp\left\{-\frac{x^2}{2(n\sigma^2 + x^{2-\alpha}/T)}\right\} &= \exp\left\{-\frac{n^2\epsilon^2}{2(n\sigma^2 + n^{2-\alpha}\epsilon^{2-\alpha}/T)}\right\} \\ &= \exp\left\{-\frac{n^\alpha\epsilon^2}{2n^{\alpha-1}\sigma^2 + 2\epsilon^{2-\alpha}/T}\right\} \\ &\leq \exp\left\{-\frac{n^\alpha\epsilon^2}{o(1) + 2/T}\right\}. \end{aligned}$$

On the other hand, by Markov inequality

$$\begin{aligned} \mathbf{P}(|W_i| \geq x) &= \mathbf{P}(W_i^2 e^{T|W_i|^\alpha} \leq x^2 e^{Tx^\alpha}) \leq Ax^{-2} \exp\{-Tx^\alpha\} \\ &\leq \frac{A}{n^2\epsilon^2} \exp\{-Tn^\alpha\epsilon^\alpha\}. \end{aligned}$$

Hence,  $\sum_{i=1}^n \mathbf{P}(|W_i| \geq x) \leq \frac{A}{n^2\epsilon^2} \exp\{-\frac{1}{2}Tn^\alpha\epsilon^\alpha\} \exp\{-\frac{1}{2}Tn^\alpha\epsilon^\alpha\} \leq o(1) \exp\{-\frac{1}{2}Tn^\alpha\epsilon^2\}$ . And (B.8) is easily obtained.  $\square$

When  $\alpha = 1$ ,  $E(e^{T_0|W_i|}) \leq A_0$  implies  $\sum_{k=0}^\infty \frac{1}{k!} T_0^k E(|W_i|^k) \leq A_0$ . So  $E(|W_i|^k) \leq \frac{1}{2} k! (\frac{1}{T_0})^{k-2} \frac{2A_0}{T_0^2}$  for  $k \geq 2$ . By Bernstein's inequality, Lemma

2.2.11 in van der Vaart and Wellner (1996), we have

$$\begin{aligned} \mathbf{P}\left(\left|\sum_{i=1}^n W_i\right| \geq n\epsilon\right) &\leq 2 \exp\left\{-\frac{n^2\epsilon^2}{2(2nA_0/T_0^2 + n\epsilon/T_0)}\right\} \\ &\leq 2 \exp\left\{-\frac{n\epsilon^2}{4A_0/T_0^2 + 2/T_0}\right\}. \end{aligned}$$

*Lemma B.5.* Under condition (C1) and (C2), for  $0 < \epsilon < 1$ , we have

$$\begin{aligned} \mathbf{P}\left(\left|\sum_{s=1}^n (X_{si}X_{sj} - \sigma_{ij})\right| \geq n\epsilon\right) &\leq C_1 \exp(-C_2 n\epsilon^2) \\ \mathbf{P}\left(\left|\sum_{s=1}^n (X_{si}X_{sj}X_{sk} - 0)\right| \geq n\epsilon\right) &\leq C_3 \exp(-C_4 n^{\frac{2}{3}}\epsilon^2) \\ \mathbf{P}\left(\left|\sum_{s=1}^n (X_{si}X_{sj}X_{sk}X_{s\ell} - \sigma_{ij}\sigma_{k\ell} - \sigma_{ik}\sigma_{j\ell} - \sigma_{i\ell}\sigma_{jk})\right| \geq n\epsilon\right) \\ &\leq C_5 \exp(-C_6 n^{\frac{1}{2}}\epsilon^2), \end{aligned} \quad (\text{B.9})$$

where  $C_1, \dots, C_6$  are constants.

*Proof.* We show the last inequality here. The first two are similar. Let  $W_s = X_{si}X_{sj}X_{sk}X_{s\ell} - \sigma_{ij}\sigma_{k\ell} - \sigma_{ik}\sigma_{j\ell} - \sigma_{i\ell}\sigma_{jk}$ .

$$\begin{aligned} E\left(e^{\frac{1}{4}|W_s|^{\frac{1}{2}}}\right) &= E\left(e^{\frac{1}{4}|X_{si}X_{sj}X_{sk}X_{s\ell} - \sigma_{ij}\sigma_{k\ell} - \sigma_{ik}\sigma_{j\ell} - \sigma_{i\ell}\sigma_{jk}|^{\frac{1}{2}}}\right) \\ &\leq E\left(e^{\frac{1}{4}|X_{si}X_{sj}X_{sk}X_{s\ell}|^{\frac{1}{2}} + \frac{1}{4}|\sigma_{ij}\sigma_{k\ell} + \sigma_{ik}\sigma_{j\ell} + \sigma_{i\ell}\sigma_{jk}|^{\frac{1}{2}}}\right) \\ &\leq e^{\frac{\sqrt{3}}{4}} E\left(e^{\frac{1}{4}|X_{si}X_{sj}X_{sk}X_{s\ell}|^{\frac{1}{2}}}\right) \\ &\leq e^{\frac{\sqrt{3}}{4}} E\left(e^{\frac{1}{4}\frac{X_{si}^2 + X_{sj}^2 + X_{sk}^2 + X_{s\ell}^2}{4}}\right) \\ &\leq e^{\frac{\sqrt{3}}{4}} E\left[\left(e^{\frac{X_{si}^2}{4}} + e^{\frac{X_{sj}^2}{4}} + e^{\frac{X_{sk}^2}{4}} + e^{\frac{X_{s\ell}^2}{4}}\right)/4\right] \\ &= \sqrt{2}e^{\frac{\sqrt{3}}{4}}. \end{aligned}$$

The inequality follows directly from the last lemma.  $\square$

*Lemma B.6.* Under conditions (C1) and (C2a), for  $m = o(n^{\frac{1}{6} - \frac{1}{3}\xi})$ ,

$$\mathbf{P}\left(\tau_{\min} \leq \min_{|\mathcal{M}| \leq m} \lambda_{\min}(\hat{\Sigma}_{\mathcal{M}}) \leq \max_{|\mathcal{M}| \leq m} \lambda_{\max}(\hat{\Sigma}_{\mathcal{M}}) \leq \tau_{\max}\right) \rightarrow 1. \quad (\text{B.12})$$

Furthermore, under condition (C4), (B.12) holds for  $m = O(n^{2\xi_0 + 4\xi_{\min}}) = o(n^{\frac{1}{6} - \frac{1}{3}\xi})$ .

*Proof.* The proof is similar to Lemma 1 in Wang (2009), where the inequality (B.9) plays a crucial role. The inequality (B.9) implies  $\mathbf{P}(|\hat{\Sigma}_{ij}^{(1)} - \Sigma_{ij}^{(1)}| > \epsilon) \leq C_1 \exp(-C_2 n\epsilon^2)$  for  $\forall 1 \leq i, j \leq p$ . Since the distribution of interactions have heavier tails, we have

$$\mathbf{P}(|\hat{\Sigma}_{\kappa\gamma} - \Sigma_{\kappa\gamma}| > \epsilon) \leq C_7 \exp(-C_8 n^{\frac{1}{2}}\epsilon^2), \quad (\text{B.13})$$

for  $\forall \kappa, \gamma \in \mathcal{P}_1 \cup \mathcal{P}_2$ . For example, if  $\kappa = (i, j)$ ,  $\gamma = (k, \ell) \in \mathcal{P}_2$ ,

$$\begin{aligned}
& |\hat{\Sigma}_{\kappa\gamma} - \Sigma_{\kappa\gamma}| \\
&= \left| \frac{1}{n} \sum_{s=1}^n (X_{si} X_{sj} - \hat{\Sigma}_{ij})(X_{sk} X_{s\ell} - \hat{\Sigma}_{k\ell}) - (\sigma_{ik}\sigma_{j\ell} + \sigma_{i\ell}\sigma_{jk}) \right| \\
&= \left| \frac{1}{n} \sum_{s=1}^n X_{si} X_{sj} X_{sk} X_{s\ell} - \hat{\Sigma}_{ij}^{(1)} \hat{\Sigma}_{k\ell}^{(1)} - (\sigma_{ik}\sigma_{j\ell} + \sigma_{i\ell}\sigma_{jk}) \right| \\
&\leq \left| \frac{1}{n} \sum_{s=1}^n X_{si} X_{sj} X_{sk} X_{s\ell} - (\sigma_{ij}\sigma_{k\ell} + \sigma_{ik}\sigma_{j\ell} + \sigma_{i\ell}\sigma_{jk}) \right| \\
&\quad + |\hat{\Sigma}_{ij}^{(1)} \hat{\Sigma}_{k\ell}^{(1)} - \sigma_{ij}\sigma_{k\ell}| \\
&\leq \left| \frac{1}{n} \sum_{s=1}^n X_{si} X_{sj} X_{sk} X_{s\ell} - (\sigma_{ij}\sigma_{k\ell} + \sigma_{ik}\sigma_{j\ell} + \sigma_{i\ell}\sigma_{jk}) \right| \\
&\quad + |\hat{\Sigma}_{ij}^{(1)}(\hat{\Sigma}_{k\ell}^{(1)} - \sigma_{k\ell})| + |(\hat{\Sigma}_{ij}^{(1)} - \sigma_{ij})\sigma_{k\ell}| \\
&\leq \left| \frac{1}{n} \sum_{s=1}^n X_{si} X_{sj} X_{sk} X_{s\ell} - (\sigma_{ij}\sigma_{k\ell} + \sigma_{ik}\sigma_{j\ell} + \sigma_{i\ell}\sigma_{jk}) \right| + |\hat{\Sigma}_{k\ell}^{(1)} - \sigma_{k\ell}| \\
&\quad + |\hat{\Sigma}_{ij}^{(1)} - \sigma_{ij}|
\end{aligned}$$

Therefore,

$$\begin{aligned}
& \mathbf{P}(|\hat{\Sigma}_{\kappa\gamma} - \Sigma_{\kappa\gamma}| > \epsilon) \\
&\leq \mathbf{P}\left(\left|\frac{1}{n} \sum_{s=1}^n X_{si} X_{sj} X_{sk} X_{s\ell} - (\sigma_{ij}\sigma_{k\ell} + \sigma_{ik}\sigma_{j\ell} + \sigma_{i\ell}\sigma_{jk})\right| > \frac{\epsilon}{3}\right) \\
&\quad + \mathbf{P}\left(|\hat{\Sigma}_{k\ell}^{(1)} - \sigma_{k\ell}| > \frac{\epsilon}{3}\right) + \mathbf{P}\left(|\hat{\Sigma}_{ij}^{(1)} - \sigma_{ij}| > \frac{\epsilon}{3}\right) \\
&\leq C_5 \exp(-C_6 n^{\frac{1}{2}}(\epsilon/3)^2) + 2C_1 \exp(-C_2 n(\epsilon/3)^2) \\
&\leq C_7 \exp(-C_8 n^{\frac{1}{2}}\epsilon^2).
\end{aligned}$$

□

Let  $\mathbf{v} = (v_1, \dots, v_p, v_{p+1}, \dots, v_{p(p+1)/2})^\top$  be a  $p + p(p+1)/2$  dimensional vector and  $\mathbf{v}_{\mathcal{M}}$  be the subvector corresponding to index set  $\mathcal{M} \subset \mathcal{P}_1 \cup \mathcal{P}_2 = \mathcal{F}$ . Recall  $\Sigma_{\mathcal{M}}$  is the principle submatrix corresponding to  $\mathcal{M}$ . By Lemma 4, we have

$$2\tau_{\min} < \min_{\mathcal{M} \subset \mathcal{F}} \inf_{\|\mathbf{v}_{\mathcal{M}}\|=1} \mathbf{v}_{\mathcal{M}}^\top \Sigma_{\mathcal{M}} \mathbf{v}_{\mathcal{M}} \leq \max_{\mathcal{M} \subset \mathcal{F}} \sup_{\|\mathbf{v}_{\mathcal{M}}\|=1} \mathbf{v}_{\mathcal{M}}^\top \Sigma_{\mathcal{M}} \mathbf{v}_{\mathcal{M}} < \tau_{\max}/2.$$

To show (B.12), it suffices to show

$$\mathbf{P}\left(\max_{|\mathcal{M}| \leq m} \sup_{\|\mathbf{v}_{\mathcal{M}}\|=1} |\mathbf{v}_{\mathcal{M}}^\top (\hat{\Sigma}_{\mathcal{M}} - \Sigma_{\mathcal{M}}) \mathbf{v}_{\mathcal{M}}| > \epsilon\right) \rightarrow 0, \quad (\text{B.14})$$

for arbitrarily small positive number  $\epsilon$ . The left-hand side of (B.14) is bounded by

$$\sum_{|\mathcal{M}| \leq m} \sum_{\kappa, \gamma \in \mathcal{F}} \mathbf{P}\left(|\hat{\Sigma}_{\kappa\gamma} - \Sigma_{\kappa\gamma}| > \frac{\epsilon}{m}\right). \quad (\text{B.15})$$

Note that the number of possible models with sizes smaller than  $m$  is less than  $(p + p(p+1)/2)^m \leq p^{2m}$  when  $p \geq 3$ . Applying (B.13), we can bound (B.15) further

$$(B.15) \leq p^{2m} (p^2)^2 C_7 \exp(-C_8 n^{\frac{1}{2}} \epsilon^2 / m^2) \quad (\text{B.16})$$

$$= C_7 \exp((2m+4) \log p - C_8 n^{\frac{1}{2}} \epsilon^2 / m^2) \quad (\text{B.17})$$

$$\leq C_7 \exp\left(2m \nu n^\xi \left(1 - \frac{1}{2} C_8 \nu^{-1} \epsilon^2 n^{\frac{1}{2}-\xi} m^{-3}\right)\right), \quad (\text{B.18})$$

which converges to zero when  $n \rightarrow \infty$  and  $m = o(n^{\frac{1}{6}-\frac{1}{3}\xi})$ .

**Remark B.2.** Beyond normality. Lemmas B.2, B.3, B.6 play important roles in the proofs of Theorems 1 and 2. A key assumption is  $E(e^{T_0|W_i|^\alpha}) \leq A_0$  where  $W_i$  is (higher) product of predictors. It is easy to see that the condition still holds, using the argument of Lemma B.5, if the marginal distributions of  $\mathbf{X}$  is sub-Gaussian. In particular, Theorem 1 is still true if condition (C1') holds and Theorem 3.2 is still true if

(C1'') holds and the total covariance matrix  $\Sigma$  has bounded eigenvalues asymptotically.

## APPENDIX C: PROOFS OF THEOREM 1 AND 2

With slight abuse of notations, we denote by  $\mathbf{X}$  the total design matrix including main and interaction effects. For any index set  $\mathcal{M} \subset \mathcal{F}$ ,  $\mathbf{X}_{(\mathcal{M})}$  is the submatrix of  $\mathbf{X}$  whose columns correspond to  $\mathcal{M}$ ;  $\boldsymbol{\beta}_{(\mathcal{M})}$  is the subvector of  $\boldsymbol{\beta}$  corresponding to  $\mathcal{M}$ . If  $\mathcal{M} = \{j\}$ , we simply use  $\mathbf{X}_j$  and  $\boldsymbol{\beta}_j$ .

We first overview the general strategy (in the context of FS2) and then give proofs for theorems. The goal is to show that all important predictors in the candidate pool are selected within a number of steps, for FS2 and the first stage of iFORT. By the nature of FS, the predictors are selected sequentially, one at each step. Therefore, we divide the whole procedure into a sequence of stages, each of which consists of several steps, starting immediately after one important term is selected and ending when the next predictor is identified. If we can show that the length of each stage is less than some integer  $L$ , then after  $d_0 L$  steps, all important predictors would have been selected.

Assume that stage  $T$  is the earliest stage among all that lasts longer than  $L$  steps, and  $T < d_0$ . Working within stage  $T$ , we omit the stage-label  $T$ , and denote by  $\mathcal{S}_t$  the index set of all selected predictors up to step  $t$  of stage  $T$ . Define

$$\Omega(t) = \text{RSS}(\mathcal{S}_t) - \text{RSS}(\mathcal{S}_{t+1}),$$

where  $\text{RSS}(\mathcal{S}_t)$  is the residual sum of squares of  $\mathbf{Y}$  regressed on the predictor space spanned by  $\mathcal{S}_t$ . A key step is to show that

$$n^{-1} \Omega(t) \geq 2L^{-1}(1 - o(1)) \quad \text{for all } 1 \leq t \leq L. \quad (\text{C.1})$$

Therefore, we have  $n^{-1} \|\mathbf{Y}\|^2 \geq \sum_{t=1}^L n^{-1} \Omega(t) \geq 2(1 - o(1)) \rightarrow 2$ , which contradicts with the fact  $\text{var}(Y) = 1$ . Then we can conclude that every stage contains less than  $L$  steps.

The inequalities of type (C.1) are obtained in the following proofs, which lead to Theorems 1 and 2. We illustrate Theorem 3.2 first, because it is technically more straightforward.

*Proof of Theorem 3.2.* Given the regularity conditions and Lemma B.6, the proof of Theorem 2 is similar to that of Theorem 1 in Wang (2009). Let  $K = 2\tau_{\max} \nu C_\beta^2 \tau_{\min}^{-1} \nu_\beta^{-4}$  and  $L = K n^{\xi_0 + 4\xi_{\min}}$ . Note that  $|\mathcal{S}_t| < d_0 L \leq K \nu n^{2\xi_0 + 4\xi_{\min}}$ , so the eigenvalues of  $\Sigma_{\mathcal{M}}$  can be controlled by Lemma B.6. Following (B.1) and (B.2) in Wang (2009), we have

$$\Omega(t)^{\frac{1}{2}} \geq \max_{j \in \mathcal{T}} \|\mathbf{H}_j^{(t)} \mathbf{Q}_{(\mathcal{S}_t)} \mathbf{X}_{(T)} \boldsymbol{\beta}_{(T)}\| - \max_{j \in \mathcal{T}} \|\mathbf{H}_j^{(t)} \mathbf{Q}_{(\mathcal{S}_t)} \boldsymbol{\epsilon}\|, \quad (\text{C.2})$$

where  $\mathbf{Q}_{(\mathcal{S}_t)} = \mathbf{I}_n - \mathbf{H}_{(\mathcal{S}_t)} = \mathbf{I}_n - \mathbf{X}_{(\mathcal{S}_t)} (\mathbf{X}_{(\mathcal{S}_t)}^\top \mathbf{X}_{(\mathcal{S}_t)})^{-1} \mathbf{X}_{(\mathcal{S}_t)}^\top$ ,  $\mathbf{H}_j^{(t)} = \mathbf{X}_j^{(t)} \mathbf{X}_j^{(t)\top} \|\mathbf{X}_j^{(t)}\|^{-2}$  and  $\mathbf{X}_j^{(t)} = (\mathbf{I}_n - \mathbf{H}_{(\mathcal{S}_t)}) \mathbf{X}_j$ .

Following the procedure leading to (B.7) in Wang (2009), we have, with probability tending to 1,

$$\max_{j \in \mathcal{T}} \|\mathbf{H}_j^{(t)} \mathbf{Q}_{(\mathcal{S}_t)} \mathbf{X}_{(T)} \boldsymbol{\beta}_{(T)}\|^2 \geq \tau_{\max}^{-1} \nu^{-1} C_\beta^{-2} \tau_{\min}^2 \nu_\beta^4 n^{1-\xi_0-4\xi_{\min}}. \quad (\text{C.3})$$

Similar to (B.8) in Wang (2009),

$$\max_{j \in \mathcal{T}} \|\mathbf{H}_j^{(t)} \mathbf{Q}_{(\mathcal{S}_t)} \boldsymbol{\epsilon}\|^2 \leq \tau_{\min}^{-1} n^{-1} \max_{j \in \mathcal{T}} \max_{|\mathcal{M}| \leq m^*} (\mathbf{X}_j^\top \mathbf{Q}_{(\mathcal{M})} \boldsymbol{\epsilon})^2, \quad (\text{C.4})$$

where  $m^* \leq TL \leq d_0 L$ . Given  $\mathbf{X}, \mathbf{X}_{(\mathcal{M})}^\top \boldsymbol{\epsilon}$  is a normal random variable with mean 0 and variance  $\|\mathbf{Q}_{(\mathcal{M})} \mathbf{X}_j\|^2 \leq \|\mathbf{X}_j\|^2$ . So (C.4) is further bounded by

$$\leq \tau_{\min}^{-1} n^{-1} \max_{j \in \mathcal{T}} \|\mathbf{X}_j\|^2 \max_{j \in \mathcal{T}} \max_{|\mathcal{M}| \leq m^*} \chi_1^2,$$

where  $\chi_1^2$  represents a chi-square random variable with one degree of freedom. By Lemma B.6,  $n^{-1} \max_{j \in \mathcal{T}} \|\mathbf{X}_j\|^2 \leq \tau_{\max}$  with probability



tending to 1. Moreover, the total number of combinations for  $j \in \mathcal{T}$  and  $|\mathcal{M}| \leq m^*$  is no more than  $(p^2)^{m^*+2} = p^{2m^*+4}$ . Therefore,

$$\begin{aligned} \max_{j \in \mathcal{T}} \max_{|\mathcal{M}| \leq m^*} \chi_1^2 &\leq 2(2m^* + 4) \log p \\ &\leq 5d_0 L v n^\xi \\ &\leq 5K v^2 n^{\xi+2\xi_0+4\xi_{\min}} \end{aligned}$$

with probability tending to 1. Finally, we have

$$\begin{aligned} n^{-1} \Omega(t) &\geq n^{-1} \left( (\tau_{\max}^{-1} v^{-1} C_{\beta}^{-2} \tau_{\min}^2 v_{\beta}^4 n^{1-\xi_0-4\xi_{\min}})^{\frac{1}{2}} \right. \\ &\quad \left. - (\tau_{\min}^{-1} \tau_{\max} 5K v^2 n^{\xi+2\xi_0+4\xi_{\min}})^{\frac{1}{2}} \right)^2 \\ &\geq \tau_{\max}^{-1} v^{-1} C_{\beta}^{-2} \tau_{\min}^2 v_{\beta}^4 n^{-\xi_0-4\xi_{\min}} \\ &\quad \times \text{Big}(1 - 2(\tau_{\max}^2 v^3 C_{\beta}^2 \tau_{\min}^{-3} v_{\beta}^{-4} 5K n^{\xi+3\xi_0+8\xi_{\min}-1})^{\frac{1}{2}}) \\ &= 2L^{-1}(1 - o(1)). \end{aligned}$$

*Proof of Theorem 3.1.* Because we concentrate on only main effects in the first stage of iFORT, similar to (C.2), we have

$$\Omega(t)^{\frac{1}{2}} \geq \max_{j \in \mathcal{T}_1} \left| \left| \mathbf{H}_j^{(t)} \mathbf{Q}_{(S_r)} \mathbf{X}_{(T_1)} \boldsymbol{\beta}_{(T_1)} \right| \right| - \max_{j \in \mathcal{T}_1} \left| \left| \mathbf{H}_j^{(t)} \mathbf{Q}_{(S_r)} (\mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)} + \boldsymbol{\epsilon}) \right| \right|. \quad (\text{C.5})$$

The first term on the right-hand side can be bounded as

$$\max_{j \in \mathcal{T}_1} \left| \left| \mathbf{H}_j^{(t)} \mathbf{Q}_{(S_r)} \mathbf{X}_{(T_1)} \boldsymbol{\beta}_{(T_1)} \right| \right|^2 \geq \tau_{\max}^{-1} v^{-1} C_{\beta}^{-2} \tau_{\min}^2 v_{\beta}^4 n^{1-\xi_0-4\xi_{\min}}. \quad (\text{C.6})$$

Similar to (C.4),

$$\begin{aligned} \max_{j \in \mathcal{T}_1} \left| \left| \mathbf{H}_j^{(t)} \mathbf{Q}_{(S_r)} (\mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)} + \boldsymbol{\epsilon}) \right| \right|^2 &\leq \tau_{\min}^{-1} n^{-1} \max_{j \in \mathcal{T}_1} \max_{|\mathcal{M}| \leq m^*} (\mathbf{X}_j^{\top} \mathbf{Q}_{(\mathcal{M})} (\mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)} + \boldsymbol{\epsilon}))^2 \\ &\leq 3\tau_{\min}^{-1} n^{-1} \max_{j \in \mathcal{T}_1} \max_{|\mathcal{M}| \leq m^*} \left( (\mathbf{X}_j^{\top} \mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)})^2 + (\mathbf{X}_j^{\top} \mathbf{H}_{(\mathcal{M})} \mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)})^2 \right. \\ &\quad \left. + (\mathbf{X}_j^{\top} \mathbf{Q}_{(\mathcal{M})} \boldsymbol{\epsilon})^2 \right), \end{aligned} \quad (\text{C.7})$$

where  $m^* \leq TL \leq p_0 L$ .

For the first term in (C.7),

$$(\mathbf{X}_j^{\top} \mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)})^2 = \left( \sum_{\kappa \in \mathcal{T}_2} \mathbf{X}_j^{\top} \mathbf{X}_{\kappa} \boldsymbol{\beta}_{(\kappa)} \right)^2 \leq q_0 \left( \max_{\kappa \in \mathcal{T}_2} |\mathbf{X}_j^{\top} \mathbf{X}_{\kappa}| \right)^2 \|\boldsymbol{\beta}_{T_2}\|^2.$$

Therefore,

$$\begin{aligned} 3\tau_{\min}^{-1} n^{-1} \max_{j \in \mathcal{T}_1} \max_{|\mathcal{M}| \leq m^*} (\mathbf{X}_j^{\top} \mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)})^2 &\leq 3\tau_{\min}^{-1} n^{-1} q_0 C_{\beta} \max_{j \in \mathcal{T}_1} \max_{\kappa \in \mathcal{T}_2} (\mathbf{X}_j^{\top} \mathbf{X}_{\kappa})^2. \end{aligned} \quad (\text{C.8})$$

By Lemma B.3, Remark 1 and Bonferroni inequality,

$$\begin{aligned} \mathbf{P} \left( \max_{j \in \mathcal{T}_1} \max_{\kappa \in \mathcal{T}_2} (\mathbf{X}_j^{\top} \mathbf{X}_{\kappa}) > \sqrt{n} 20 \sqrt{\log n} \right) &\leq p_0 q_0 3 \exp(-400 \log n / 61) \\ &\leq \exp(2 \log v + 2\xi_0 \log n - 2 \log n) \rightarrow 0. \end{aligned}$$

Thus (C.8) can be bounded by  $1200 \tau_{\min}^{-1} C_{\beta} v n^{\xi_0} \log n$  with probability tending to 1.

For the second term,

$$\begin{aligned} (\mathbf{X}_j^{\top} \mathbf{H}_{(\mathcal{M})} \mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)})^2 &= \left( \sum_{\kappa \in \mathcal{T}_2} \mathbf{X}_j^{\top} \mathbf{X}_{(\mathcal{M})} (\mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{(\mathcal{M})})^{-1} \mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{\kappa} \boldsymbol{\beta}_{(\kappa)} \right)^2 \\ &\leq q_0 \left( \max_{\kappa \in \mathcal{T}_2} \mathbf{X}_j^{\top} \mathbf{X}_{(\mathcal{M})} (\mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{(\mathcal{M})})^{-1} \mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{\kappa} \right)^2 \|\boldsymbol{\beta}_{T_2}\|^2. \end{aligned}$$

Therefore,

$$\begin{aligned} 3\tau_{\min}^{-1} n^{-1} \max_{j \in \mathcal{T}_1} \max_{|\mathcal{M}| \leq m^*} (\mathbf{X}_j^{\top} \mathbf{H}_{(\mathcal{M})} \mathbf{X}_{(T_2)} \boldsymbol{\beta}_{(T_2)})^2 &\leq 3\tau_{\min}^{-1} n^{-1} q_0 C_{\beta} \max_{j \in \mathcal{T}_1} \max_{|\mathcal{M}| \leq m^*} \max_{\kappa \in \mathcal{T}_2} (\mathbf{X}_j^{\top} \mathbf{X}_{(\mathcal{M})} (\mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{(\mathcal{M})})^{-1} \mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{\kappa})^2 \\ &\leq 3\tau_{\min}^{-1} n^{-1} q_0 C_{\beta} \max_{j \in \mathcal{T}_1} \max_{|\mathcal{M}| \leq m^*} \max_{\kappa \in \mathcal{T}_2} (\|\mathbf{X}_j^{\top} \mathbf{X}_{(\mathcal{M})} (\mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{(\mathcal{M})})^{-1}\|_{\infty} m^* \\ &\quad \times \max_{\ell \in \mathcal{M}} |\mathbf{X}_{\ell}^{\top} \mathbf{X}_{\kappa}|)^2 \\ &\leq 3\tau_{\min}^{-1} n^{-1} q_0 C_{\beta} m^{*2} \max_{j \in \mathcal{T}_1} \max_{|\mathcal{M}| \leq m^*} \|\mathbf{X}_j^{\top} \mathbf{X}_{(\mathcal{M})} (\mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{(\mathcal{M})})^{-1}\|_{\infty}^2 \max_{\kappa \in \mathcal{T}_2} \\ &\quad \times \max_{\ell \in \mathcal{P}_1} (\mathbf{X}_{\ell}^{\top} \mathbf{X}_{\kappa})^2, \end{aligned} \quad (\text{C.9})$$

where  $\|\cdot\|_{\infty}$  denote the vectorized infinity norm. By Lemma B.2,

$$\begin{aligned} \|\mathbf{X}_j^{\top} \mathbf{X}_{(\mathcal{M})} (\mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{(\mathcal{M})})^{-1}\|_{\infty} &\leq \left\| \frac{\mathbf{X}_j^{\top} \mathbf{X}_{(\mathcal{M})}}{n} \right\|_2 \left\| \left( \frac{\mathbf{X}_{(\mathcal{M})}^{\top} \mathbf{X}_{(\mathcal{M})}}{n} \right)^{-1} \right\|_2 \\ &\leq \tau_{\max} \tau_{\min}^{-1}, \end{aligned}$$

with probability tending to one. By Lemma B.3,

$$\begin{aligned} \mathbf{P} \left( \max_{\kappa \in \mathcal{T}_2} \max_{\ell \in \mathcal{P}_1} (\mathbf{X}_{\ell}^{\top} \mathbf{X}_{\kappa}) > \sqrt{n} \sqrt{100 v n^{\xi}} \right) &\leq p q_0 3 \exp \left( -\frac{200 v n^{\xi}}{61} \right) \\ &\leq 3 \exp \left( v n^{\xi} + \log v + \xi_0 \log n - \frac{100}{61} v n^{\xi} \right) \rightarrow 0. \end{aligned}$$

Thus, with probability tending to 1, (C.9) is further bounded by

$$300 \tau_{\max}^2 \tau_{\min}^{-3} C_{\beta} m^{*2} v^2 n^{\xi_0+\xi} \leq 300 \tau_{\max}^2 \tau_{\min}^{-3} C_{\beta} v^4 K^2 n^{5\xi_0+8\xi_{\min}+\xi}. \quad (\text{C.10})$$

Following the same steps after (C.4), the third term in (C.7) can be controlled by,

$$15 \tau_{\min}^{-1} \tau_{\max} K v^2 n^{\xi+2\xi_0+4\xi_{\min}}. \quad (\text{C.11})$$

Finally, combining all results, we have

$$\begin{aligned} \Omega(t)^{\frac{1}{2}} &\geq (\tau_{\max}^{-1} v^{-1} C_{\beta}^{-2} \tau_{\min}^2 v_{\beta}^4 n^{1-\xi_0-4\xi_{\min}})^{\frac{1}{2}} - (1200 \tau_{\min}^{-1} C_{\beta} v n^{\xi_0} \log n \\ &\quad + 300 \tau_{\max}^2 \tau_{\min}^{-3} C_{\beta} v^4 K^2 n^{5\xi_0+8\xi_{\min}+\xi} \\ &\quad + 15 \tau_{\min}^{-1} \tau_{\max} K v^2 n^{\xi+2\xi_0+4\xi_{\min}})^{\frac{1}{2}} \\ &\geq n^{\frac{1}{2}} (\tau_{\max}^{-1} v^{-1} C_{\beta}^{-2} \tau_{\min}^2 v_{\beta}^4 n^{-\xi_0-4\xi_{\min}})^{\frac{1}{2}} \times (1 - A_1 - A_2 - A_3)^{\frac{1}{2}}, \end{aligned}$$

where  $A_1 = 1200 \tau_{\min}^{-3} \tau_{\max} C_{\beta}^3 v^2 v_{\beta}^{-4} n^{2\xi_0+4\xi_{\min}-1} \log n$ ,  $A_2 = 300 \tau_{\max}^3 \tau_{\min}^{-5} C_{\beta}^3 v_{\beta}^{-4} v^5 K^2 n^{6\xi_0+12\xi_{\min}+\xi-1}$ ,  $A_3 = 15 \tau_{\min}^{-3} \tau_{\max}^2 K v^3 C_{\beta}^2 v_{\beta}^{-4} n^{\xi+3\xi_0+8\xi_{\min}-1}$ . Therefore,

$$n^{-1} \Omega(t) \geq 2L^{-1}(1 - o(1)).$$

## SUPPLEMENTARY MATERIALS

Additional numerical experiments (Examples 7–9) to illustrate performance of the iFOR methods for interaction selection under various settings.

[Received June 2012. Revised October 2013.]

## REFERENCES

- Bar, W., and Dittrich, F. (1971), "Useful Formula for Moment Computation of Normal Random Variables With Nonzero Means," *Automatic Control, IEEE Transactions on*, 16, 263–265. [1297]  
Bickel, P. J., and Levina, E. (2008), "Regularized Estimation of Large Covariance Matrices," *The Annals of Statistics*, 36, 199–227. [1298]

- Candes, E., and Tao, T. (2007), "The Dantzig Selector: Statistical Estimation When  $p$  is Much Larger Than  $n$ ," *The Annals of Statistics*, 35, 2313–2351. [1285]
- Chen, J., and Chen, Z. (2008), "Extended Bayesian Information Criteria for Model Selection With Large Model Spaces," *Biometrika*, 95, 759–771. [1288,1295]
- Chipman, H. (1996), "Bayesian Variable Selection With Related Predictors," *The Canadian Journal of Statistics / La Revue Canadienne de Statistique*, 24, 17–36. [1285]
- Chipman, H., Hamada, M., and Wu, C. F. J. (1997), "A Bayesian Variable-Selection Approach for Analyzing Designed Experiments With Complex Aliasing," *Technometrics*, 39, 372–381. [1285]
- Choi, N. H., Li, W., and Zhu, J. (2010), "Variable Selection With the Strong Heredity Constraint and its Oracle Property," *Journal of the American Statistical Association*, 105, 354–364. [1286,1287]
- Cordell, H. J. (2009), "Detecting Gene-Gene Interactions that Underlie Human Diseases," *Nature Reviews Genetics*, 10, 392–404. [1285]
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004), "Least Angle Regression," *The Annals of Statistics*, 32, 407–451. [1286]
- Evans, D. M., Marchini, J., Morris, A. P., and Cardon, L. R. (2006), "Two-Stage Two-Locus Models in Genome-Wide Association," *PLoS Genetics*, 2, e157. [1285]
- Fan, J., and Li, R. (2001), "Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties," *Journal of the American Statistical Association*, 96, 1348–1360. [1285]
- Fan, J., and Lv, J. (2008), "Sure Independence Screening for Ultrahigh Dimensional Feature Space," *Journal of the Royal Statistical Society, Series B*, 70, 849–911. [1285,1288]
- (2011), "Nonconcave Penalized Likelihood With np-Dimensionality," *Information Theory, IEEE Transactions on*, 57, 5467–5484. [1285,1289]
- Fan, J., and Peng, H. (2004), "On Non-Concave Penalized Likelihood With Diverging Number of Parameters," *The Annals of Statistics*, 32, 928–961. [1288]
- Fan, J., and Song, R. (2010), "Sure Independence Screening in Generalized Linear Models With np-Dimensionality," *The Annals of Statistics*, 38, 3567–3604. [1285,1290,1291]
- Hamada, M., and Wu, C. F. J. (1992), "Analysis of Designed Experiments With Complex Aliasing," *Journal of Quality Technology*, 24, 130–137. [1285]
- Kooperberg, C., and LeBlanc, M. (2008), "Increasing the Power of Identifying GxG Interactions in Genome-Wide Association Studies," *Genetic Epidemiology*, 32, 255–263. [1285]
- Lan, H., Chen, M., Flowers, J. B., Yandell, B. S., Stapleton, D. S., Mata, C. M., Mui, E. T., Flowers, M. T., Schueler, K. L., Manly, K. F., et al. (2006), "Combined expression trait correlations and expression quantitative trait locus mapping," *PLoS Genet*, 2, e6. [1294]
- Laub, A. (2005), *Matrix Analysis for Scientists and Engineers*, Philadelphia, PA: Society for Industrial and Applied Mathematics. [1297]
- Manolio, T. A., and Collins, F. S. (2007), "Genes, Environment, Health, and Disease: Facing up to Complexity," *Human Heredity*, 63, 63–66. [1285]
- McCullagh, P. (2002), "What is a Statistical Model?" *The Annals of Statistics*, 30, 1225–1267. [1285]
- McCullagh, P., and Nelder, J. (1989), *Generalized Linear Models*, Monographs on Statistics and Applied Probability, London, UK: Chapman and Hall. [1285]
- Nelder, J. A. (1977), "A Reformulation of Linear Models," *Journal of the Royal Statistical Society, Series A*, 140, 48–77. [1285]
- (1994), "The Statistics of Linear Models: Back to Basics," *Statistics and Computing*, 5, i–i. doi:10.1007/BF00143933. [1285]
- Tibshirani, R. (1996), "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society, Series B*, 58, 267–288. [1285]
- Turlach, B. (2004), Discussion of "Least Angle Regression," *The Annals of Statistics*, 32, 481–490. [1286]
- van der Vaart, A. and Wellner, J. (1996), *Weak Convergence and Empirical Processes*, Springer Series in Statistics, New York: Springer. [1298]
- Wang, H. (2009), "Forward Regression for Ultra-High Dimensional Variable Screening," *Journal of the American Statistical Association*, 104, 1512–1524. [1285,1287,1288,1289,1297,1298,1299]
- Wu, J., Devlin, B., Ringquist, S., Trucco, M., and Roeder, K. (2010), "Screen and Clean: A Tool for Identifying Interactions in Genome-Wide Association Studies," *Genetic Epidemiology*, 34, 275–285. [1286,1290]
- Wu, T. T., Chen, Y. F., Hastie, T., Sobel, E., and Lange, K. (2009), "Genome-Wide Association Analysis by Lasso Penalized Logistic Regression," *Bioinformatics*, 25, 714–721. [1286,1290]
- Yuan, M., Joseph, V., and Lin, Y. (2007), "An Efficient Variable Selection Approach for Analyzing Designed Experiments," *Technometrics*, 49, 430–439. [1286]
- Yuan, M., Joseph, V. R., and Zou, H. (2009), "Structured Variable Selection and Estimation," *Annals of Applied Statistics*, 3, 1738. [1286,1287,1290]
- Zhang, C.-H. (2010), "Nearly Unbiased Variable Selection Under Minimax Concave Penalty," *The Annals of Statistics*, 38, 894–942. [1285]
- Zhang, C.-H., and Huang, J. (2008), "The Sparsity and Bias of the Lasso Selection in High-Dimensional Linear Regression," *The Annals of Statistics*, 36, 1567–1594. [1288]
- Zhao, P., Rocha, G., and Yu, B. (2009), "The Composite Absolute Penalties Family for Grouped and Hierarchical Variable Selection," *The Annals of Statistics*, 37, 3468–3497. [1286,1287]
- Zhao, P., and Yu, B. (2006), "On Model Selection Consistency of Lasso," *Journal of Machine Learning Research*, 7, 2541–2563. [1285,1289]
- Zou, H. (2006), "The Adaptive Lasso and its Oracle Properties," *Journal of the American Statistical Association*, 101, 1418–1429. [1285]
- Zou, H., and Hastie, T. (2005), "Regularization and Variable Selection via the Elastic Net," *Journal of the Royal Statistical Society, Series B*, 67, 301–320. [1285]
- Zou, H., and Li, R. (2008), "One-step Sparse Estimates in Nonconcave Penalized Likelihood Models," *The Annals of Statistics*, 36, 1509–1533. [1285]