## False Discovery Rate Control Under General Dependence By Symmetrized Data Aggregation

Lilun Du<sup>1</sup>, Xu Guo<sup>2</sup>, Wenguang Sun<sup>3</sup> and Changliang Zou<sup>4</sup>

#### **Abstract**

We develop a new class of distribution—free multiple testing rules for false discovery rate (FDR) control under general dependence. A key element in our proposal is a symmetrized data aggregation (SDA) approach to incorporating the dependence structure via sample splitting, data screening and information pooling. The proposed SDA filter first constructs a sequence of ranking statistics that fulfill global symmetry properties, and then chooses a data—driven threshold along the ranking to control the FDR. The SDA filter substantially outperforms the knockoff method in power under moderate to strong dependence, and is more robust than existing methods based on asymptotic p-values. We first develop finite—sample theories to provide an upper bound for the actual FDR under general dependence, and then establish the asymptotic validity of SDA for both the FDR and false discovery proportion (FDP) control under mild regularity conditions. The procedure is implemented in the R package sdafilter. Numerical results confirm the effectiveness and robustness of SDA in FDR control and show that it achieves substantial power gain over existing methods in many settings.

Keywords: Empirical distribution; Integrative multiple testing; Moderate deviation theory; Sample-splitting; Uniform convergence.

<sup>&</sup>lt;sup>1</sup>Hong Kong University of Science and Technology, Hong Kong

<sup>&</sup>lt;sup>2</sup>Beijing Normal University, Beijing, China

<sup>&</sup>lt;sup>3</sup>University of Southern California. Corresponding Email: wenguans@marshall.usc.edu.

<sup>&</sup>lt;sup>4</sup>Nankai University, Tianjin, China

## 1 Introduction

Multiple testing provides a useful approach to identifying sparse signals from massive data. Recent developments on false discovery rate (FDR; Benjamini and Hochberg, 1995) methodologies have greatly influenced a wide range of scientific disciplines including genomics (Tusher et al., 2001; Roeder and Wasserman, 2009), neuroimaging (Pacifico et al., 2004) Schwartzman et al., 2008), geography (Caldas de Castro and Singer, 2006; Sun et al., 2015) and finance (Barras et al., 2010). Conventional FDR procedures, such as the Benjamini—Hochberg (BH) procedure, adaptive p-value procedure (Benjamini and Hochberg, 1997) and adaptive z-value procedure based on local FDR (Efron et al., 2001; Sun and Cai, 2007), are developed under the assumption that the test statistics are independent. However, data arising from large—scale testing problems are often dependent. FDR control under dependence is a critical problem that requires much research. Two key issues include (a) how the dependence may affect existing FDR methods, and (b) how to properly incorporate the dependence structure into inference.

## 1.1 FDR control under dependence

The impact of dependence on FDR analysis was first investigated by Benjamini and Yekutieli (2001), who showed that the BH procedure, when adjusted at level  $\alpha/(\frac{p}{j-1} 1/j)$  with p being the number of tests, controls the FDR at level  $\alpha$  under arbitrary dependence among the p-values. However, this adjustment is often too conservative in practice. Benjamini and Yekutieli (2001) further proved that applying BH without any adjustment is valid for FDR control for correlated tests satisfying the PRDS property. This result was strengthened by Sarkar (2002), who showed that the FDR control theory under positive dependence holds for a generalized class of step-wise methods. Storey et al. (2004), Wu (2008) and Clarke and Hall (2009) respectively showed that, in the asymptotic sense, BH is valid under weak dependence, Markovian dependence and linear process models. Although controlling the FDR does not always require independence, some key quantities in FDR analysis, such as the expectation and variance of the number of false positives, may possess substantially different properties under dependence (Owen) 2005; Finner et al. 2007). This implies that conventional FDR methods such as BH can suffer from low power and high variability under strong

dependence. Efron (2007) and Schwartzman and Lin (2011) showed that strong correlations degrade the accuracy in both estimation and testing. In particular, positive/negative correlations can make the empirical null distributions of z-values narrower/wider, which has substantial impact on subsequent FDR analyses. These insightful findings suggest that it is crucial to develop new FDR methods tailored to capture the structural information among dependent tests.

Intuitively high correlations can be exploited to aggregate weak signals from individuals to increase the signal to noise ratio (SNR). Hence informative dependence structures can become a bless for FDR analysis. For example, the works of Benjamini and Heller (2007), Sun and Cai (2009) and Sun and Wei (2011) showed that incorporating functional, spatial, and temporal correlations into inference can improve the power and interpretability of existing methods. However, these methods are not applicable to general dependence structures. Efron (2007), Efron (2010) and Fan et al. (2012) discussed how to obtain more accurate FDR estimates by taking into account arbitrary dependence. For a general class of dependence models, Leek and Storey (2008), Friguet et al. (2009), Fan et al. (2012) and Fan and Han (2017) showed that the overall dependence can be much weakened by subtracting the common factors out, and factor-adjusted p-values can be employed to construct more powerful FDR procedures. The works by Hall and Jin (2010), Jin (2012) and Li and Zhong (2017) showed that, under both the global testing and multiple testing contexts, the covariance structures can be utilized, via transformation, to construct test statistics with increased SNR, revealing the beneficial effects of dependence. However, the above methods, for example by Fan and Han (2017) and Li and Zhong (2017), rely heavily on the accuracy of estimated models and the asymptotic normality of the test statistics. Under the finite-sample setting, poor estimates of model parameters or violations of normality assumption may lead to less powerful and even invalid FDR procedures. This article aims to develop a robust and assumption-lean method that effectively controls the FDR under general dependence with much improved power.

#### 1.2 Model and problem formulation

We consider a setup where p-dimensional vectors  $\xi_i = (\xi_{i1}, \dots, \xi_{ip})^>$ ,  $i = 1, \dots, n$ , follow a multivariate distribution with mean  $\mu = (\mu_1, \dots, \mu_p)^>$  and covariance matrix  $\Sigma$ . The problem of interest

is to test p hypotheses simultaneously:

$$H_i^0: \mu_j = 0$$
 versus  $H_i^1: \mu_j = 0$ , for  $j = 1, ..., p$ .

The summary statistic  $\bar{\xi} = n^{-1} \sum_{i=1}^{p} \xi_i$  obeys a multivariate normal (MVN) model asymptotically

$$\bar{\xi} \stackrel{d}{\approx} MVN(\mu, n^{-1}\Sigma).$$
 (1)

Denote  $\Omega = \Sigma^{-1}$  the precision matrix. We first assume that  $\Omega$  is known. For the case with unknown precision matrix, a data-driven methodology and its theoretical properties are discussed in Section  $\square$  The problem of multiple testing under dependence can be recast as a variable selection problem in linear regression. Specifically, by taking a "whitening" transformation, Model ( $\square$  is equivalent to the following model:

$$Y = X \mu + , \approx^{d} MVN(0, n^{-1}I_{p}),$$
 (2)

where Y =  $\Omega^{1/2}\bar{\xi}$   $\mathbb{R}^p$  is the pseudo response, X =  $\Omega^{1/2}$   $\mathbb{R}^{p \times p}$  is the design matrix,  $I_p$  is a p-dimensional identity matrix and =  $(1, \ldots, p)^p$  are noise terms that are approximately inde-pendent and normally distributed. The connection between model selection and FDR was discussed in Abramovich et al. (2006) and Bogdan et al. (2015), respectively under the normal means model and regression model with orthogonal designs.

$$FDP = \frac{P_{j=1}^{p} (1 - \theta_{j}) \delta_{j}}{(P_{j=1}^{p} \delta_{j}) ? 1}, \quad TDP = \frac{P_{j=1}^{p} \theta_{j} \delta_{j}}{(P_{j=1}^{p} \theta_{j}) ? 1},$$
(3)

where a  $\square$  b = max(a, b). The FDR is the expectation of the FDP: FDR = E(FDP). The average power is defined as AP = E(TDP).

## 1.3 FDR control by symmetrized data aggregation

This article introduces a new information pooling strategy, the symmetrized data aggregation (SDA), for handling the dependence issue in multiple testing. The SDA involves splitting and re-

assembling data to construct a sequence of statistics fulfilling symmetry properties. Our proposed SDA filter for FDR control consists of three steps:

- The first step splits the sample into two parts, both of which are utilized to construct statistics to assess the evidence against the null.
- The second step aggregates the two statistics to form a new ranking statistic fulfilling symmetry properties.
- The third step chooses a threshold along the ranking by exploiting the symmetry property between positive and negative null statistics to control the FDR.

To get intuitions on how the idea works, we start with the independent case [Zou et al. (2020)]. The more interesting but complicated dependent case will be described shortly, with detailed discussions, refinements and justifications deferred to later sections. Suppose the vectors  $\xi_i = (\xi_{i1}, \ldots, \xi_{ip})^{>}$  are i.i.d. obeying MVN( $\mu$ , I<sub>p</sub>). The proposed SDA method first splits the full sample into two disjoint subsets D<sub>1</sub> and D<sub>2</sub>, with sizes n<sub>1</sub> and n<sub>2</sub> and n = n<sub>1</sub> + n<sub>2</sub>. A pair of statistics, both of which follow N(0, 1) under the null, are then calculated to test H<sup>0</sup><sub>j</sub>:

$$(T_{1j},T_{2j}) = \begin{array}{c} P \\ \hline \begin{array}{c} \frac{i \mathbb{P} D_1}{\sqrt{n_1}} \xi_{j} \\ \hline \end{array}, \begin{array}{c} P \\ \frac{i \mathbb{P} D_2}{\sqrt{n_2}} \xi_{j} \\ \hline \end{array}.$$

The product  $W_j = T_{1j}T_{2j}$  is used to aggregate the evidence across the two groups. If  $|\mu_j|$  is large, then both  $T_{1j}$  and  $T_{2j}$  tend to have large absolute values with the same sign, thereby leading to a positive and large  $W_j$ . By contrast,  $W_j$  fulfills the symmetry property under  $H_i^0$ , i.e.

$$Pr(W_j \ge t \mid H_i^0) = Pr(W_j \le -t \mid H_i^0), \text{ for any } t \ ? R.$$
 (4)

This motivates one to consider the following selection procedure  $A^b = \{j : W_j \ge L\}$ , where L is the threshold chosen to control the FDR at level  $\alpha$ :

$$L = \inf t > 0: \frac{\#\{j : W_j \le -t\}}{\#\{j : W_j \ge t\} \supseteq 1} \le \alpha .$$
 (5)

According to the symmetry property (4), the count of negative  $W_j$ 's below –t strongly resembles the count of false positives in the selected subset (i.e. the null  $W_j$ 's above t). It follows that the fraction in Equation (5) provides a good estimate of the FDP.

The dependent case involves a more carefully designed SDA filter. After sample splitting, we apply variable selection techniques such as LASSO to  $D_1$  to construct  $T_{1j}$ .  $T_{1j}$ , which is calculated based on linear model (2), can effectively capture the dependence structure. Before using  $D_2$  to construct  $T_{2j}$ , we carry out a data screening step to narrow down the focus. We show that the screening step can significantly increase the SNR of  $T_{2j}$  under strong dependence, hence the correlations are exploited again to increase the power. The ranking statistic  $W_j$  is constructed by combining  $T_{1j}$  and  $T_{2j}$  with proven asymptotic symmetry properties. The theory of the proposed SDA filter is divided into two parts: the finite sample theory provides an upper bound for the FDR under general dependence, while the asymptotic theory shows that both the FDR and FDP can be controlled at  $\alpha$  + o(1) under mild regularity conditions.

## 1.4 Connections to existing work and our contributions

The SDA is closely related to existing ideas of sample–splitting (Wasserman and Roeder, 2009; Meinshausen et al., 2009) and data carving (Eithian et al., 2014) Lei et al., 2021), both of which firstly divide the data into two independent parts, secondly use one part to narrow down the focus (or rank the hypotheses) and finally use the remainder to perform inference tasks such as variable selection, estimation or multiple testing. These ideas have a common theme with covariate–assisted multiple testing (Lei and Fithian, 2018; Cai et al., 2019) Li and Barber, 2019, where the primary statistic plays the key role to assess the significance while the side information plays an auxiliary role to assist inference [see also the discussion by Ramdas (2019)]. SDA provides a novel way of data aggregation where both parts of data, which are combined under the symmetry principle, play essential roles in both ranking and selection. This substantially reduces the information loss in conventional sample–splitting methods, while the symmetry principle, which is fulfilled by construction, enables the development of an effective and assumption-lean FDR filter.

The SDA is inspired by the elegant knockoff filter for FDR control (Barber and Candès, 2015), which creates knockoff features that emulate the correlation structure in original features, to form symmetrized ranking statistics for selecting important variables via the same mechanism (5). The knockoff method, which is originally developed under regression models, can be applied for FDR

control in Model (1) via the equivalent Model (2). The knockoff filter employs local pairwise contrasts: the ranking variable is constructed to capture the differential evidences against the null exhibited by the pair (i.e. the original feature vs. its knockoff copy). While it is desirable to make the pair as "independent" as possible, high correlations will greatly restrict the geometric space in which the knockoff can be constructed; see Appendix B.1 for detailed discussions and illustrations. This would significantly increase the difficulty for distinguishing the variable and its knockoff and hence lower the power. By contrast, the SDA filter, which does not rely on pairwise contrasts, will not suffer from high correlations.

To visualize the correlation effects, we consider a setup similar to Figure 5 in Barber and Candès (2015), where correlated normal, t, and exponential data are generated based on an autoregressive model  $\Sigma = (\rho^{|j-i|})$  (see Section 5.2) for more details about the setup). We vary  $\rho$  from -0.9 to 0.9 and apply BH, knockoff and SDA at FDR level  $\alpha = 0.2$ . The actual FDRs and APs based on 500 replications are summarized in Figure 1. Our first column (normal data) shows that knockoff outperforms BH in some situations, but both the FDR and AP of the knockoff method decrease when correlations grow higher. By contrast, SDA controls the FDR near the nominal level consistently, and the power of SDA increases sharply with growing correlations. This pattern corroborates the insights by Benjamini and Heller (2007), Sun and Cai (2009) and Hall and Jin (2010) that high correlations, which can be exploited to increase the SNR, may become a bless in large—scale inference.

The proposed research improves the previous work by Zou et al. (2020) in several ways. First, Zou et al. (2020) has mainly focused on the independent and weak dependent case, with the major goal of deriving convergence rate of false discovery proportions when simultaneously performing thousands of t-tests. The methodology in Zou et al. (2020), which does not utilize LASSO and does not include the data screening step, becomes highly inefficient under strong dependence. See Appendix B.2 for an illustration. Second, our new theories for FDR and FDP control under dependence and the robustness of the SDA filter under model misspecification substantially depart from the theory in Zou et al. (2020).

The SDA filter provides a model–free framework that overcomes the limitations of many selective inference procedures, for example, the methods in Lockhart et al. (2014) and Javanmard and Javadi

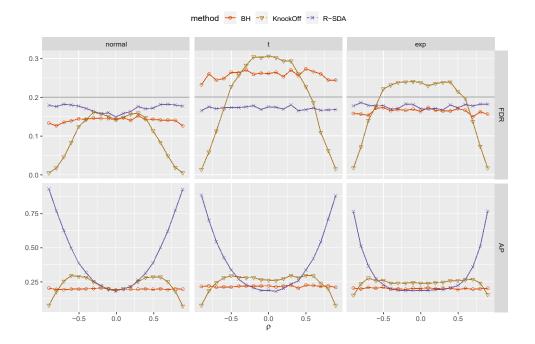


Figure 1: Impacts of correlation on different FDR procedures: "t" denotes the t distribution with 3 df and "exp" denotes the exponential distribution with scale parameter 2. In both cases the models have been misspecified as normal when computing the p-values.

(2019), which require strong assumptions about the conditional distribution to construct asymptotic p-values. Our numerical results show that the methods in Fan and Ham (2017) and Li and Zhong (2017), which require correctly specified models, accurate estimates of parameters and normality assumptions, are in general not robust for FDR control. The SDA filter, which employs empirical distributions instead of asymptotic distributions, only requires the global symmetry of the ranking statistics. It is more robust than its competitors for a wide range of scenarios since the asymptotic symmetry property is much easier to achieve in practice compared to asymptotic normality. As illustrated by the second column (multivariate t data) of Figure 1. BH fails to control the FDR under heavy—tailed models. The failure in accounting for the deviations from normality may result in misleading empirical null and severe bias in FDR analysis (Efron, 2004) belaigle et al., 2011; Liu and Shao, 2014). Finally, our Theorem 1. which develops a finite—sample upper bound of FDR

<sup>&</sup>lt;sup>1</sup>For example, the average of several t-variables fulfills the symmetry property perfectly but violates the normality assumption. For asymmetric distributions such as exponential, we usually need a smaller sample size to achieve asymptotic symmetry compared to asymptotic normality – the latter is stronger than the former since it requires an additional accurate approximation in the tail areas.

under dependence, is closely connected to robust knockoffs theory and is established utilizing key arguments from Barber et al. (2020). More specifically, we employ the leave-one-out technique suggested in Barber et al. (2020) to analyze the effect on the SDA filter of possible deviations from normality and the sure screening property, similarly to the analysis of the effect on the Model-X knockoff filter of errors in estimating the true covariance structure. This important connection sheds lights on how the model uncertainty can affect the actual FDR level and how the error bound in FDR can be explicitly quantified using appropriate deviation measures; a detailed discussion is provided in Section B.3 of the Supplementary Material.

## 1.5 Organization

The remainder of our paper is structured as follows. In Section we introduce the SDA filter for FDR control and discuss the effects of dependence on multiple testing. We develop finite sample and asymptotic theories for FDR control in Section Methodology and theory for the unknown dependence case are discussed in Section Simulation and real data analysis are presented in Sections and respectively. The extensions, proofs of theories and additional comparisons are provided in the Supplementary Material.

## 2 The SDA Filter for FDR Control

We start with the assumption that the covariance matrix  $\Sigma$  is known and then move to the case with unknown  $\Sigma$  in Section  $\Delta$  Our discussion is mainly based on regression model ( $\Sigma$ ); an equivalent

description of the methodology via model (1) follows similarly. We first outline in Section 2.1 the steps for constructing the ranking statistics, then provide intuitive explanations on how the SDA filter works in Sections 2.2 and 2.3 The detailed SDA algorithm is provided in Section A.4.

## 2.1 Construction of ranking statistics and the symmetry property

SDA first splits the data into two independent parts  $D_1$  and  $D_2$ , which are respectively used to construct statistics  $T_{1j}$  and  $T_{2j}$ . The information in the two parts is then combined to form the ranking statistic  $W_j = T_{1j}T_{2j}$ . A wide class of pairs may be constructed from the sample. This section presents a specific pair  $(T_{1j}, T_{2j})$ , which is used in all numerical studies. Examples of other possible pairs are presented in Section A.2 in the Supplementary Material.

We propose to use LASSO (Tibshirani, 1996) to extract information from  $D_1$  as it simultaneously takes into account the sparsity and dependency structures. Let  $\bar{\xi}_1 = n_1^{-1} P_{i \otimes D_1} \xi_i$  and  $y_1 = X\bar{\xi}_1$ . The LASSO estimator is given by  $\mu b_1 = (\mu_{11}, \dots, \mu_{1p})^{>} = \arg\min L(\mu)$ , where

$$L(\mu) = (y_1 - X\mu)^{>} (y_1 - X\mu) + \lambda k \mu k_1.$$
 (6)

Let  $S = \{j : \mathbf{h}_{1j} = 0\}$  denote the subset of coordinates selected by LASSO and  $S^c = \{1, \dots, p\} \setminus S$  its complement.

Remark 1 Following Wasserman and Roeder (2009), we suggest using  $n_1 = d2n/3e$ , which provides stable performance across a wide range of settings. To obtain asymptotically unbiased estimator in the next step, it is required that S contains all the signals with high probability. In practice, this can be achieved by deliberately choosing an overfitted model that includes most true signals and many false positives; see also Barber and Candès (2019) and Remark 2 in Section 3.2.

Next we use  $D_2$  to obtain the least–squares estimates (LSEs). Let  $\xi_2 = n_2^{-1} P_{i \boxtimes D_2} \xi_i$ ,  $y_2 = X \xi_2$ ,  $X_S = (X_j : j \boxtimes S) \text{ and } e_j = (0, \dots, 0, 1, 0, \dots, 0)^{>2}.$  The LSEs are only calculated for

<sup>&</sup>lt;sup>2</sup>Specifically,  $e_j$  is an |S|-vector with 1 in the jth coordinate and 0 elsewhere.

coordinates on the narrowed subset S. Let  $\mu_2 = (\mu_{21}, \ldots, \mu_{2p})^>$ , where

$$\mathbf{p}_{2j} = \begin{cases} \frac{?}{?} e^{>}_{j} (X^{>}X_{S})^{-1}X^{>}Y_{S}, & j ? S; \\ ? & 0, & j ? S^{c}. \end{cases}$$
(7)

Section 2.3 provides insights on why this data screening step can lead to increased SNR.

To aggregate information across both  $D_1$  and  $D_2$ , let  $W_j = T_{1j}T_{2j}$ , where

$$(\mathsf{T}_{1j},\mathsf{T}_{2j}) = \frac{\sqrt{n_1 | \mathbf{m}_{1j}|} \sqrt{n_2 | \mathbf{m}_{2j}|}}{\sigma_{S,j}}$$
(8)

and  $\sigma_{S,j}^2$ 's are the diagonal elements of  $(X_S^> X_S)^{-1}$ . A multiple testing procedure consists of two steps: ranking and thresholding. Next we show that  $W_j$ 's play key roles in both steps. Intuitively, the positive  $W_j$ 's can be used for ranking because a large and positive  $W_j$  indicates strong evidence against the null. Meanwhile, the negative  $W_j$ 's, which usually correspond to null cases, can be used for thresholding. The key idea is to exploit the following asymptotic symmetry property:

$$\sup_{0 \le t \le c \log p} \frac{P^{j \ge S \cap A^c} I(W_j \ge t)}{P^{j \ge S \cap A^c} I(W_j \le -t)} - 1 = o(1) \quad \text{for some } c > 0,$$
 (9)

#### 2.2 FDR thresholding

The asymptotic symmetry property ( $\bigcirc$ ) motivates us to choose the following data—driven threshold to control the FDR at level  $\alpha$ :

$$L = \inf t > 0: \frac{\#\{j : W_j \le -t\}}{\#\{j : W_j \ge t\} \ ! \ ! \ !} \le \alpha . \tag{10}$$

Our decision rule is given by  $\delta = (\delta_j : 1 \le j \le p)^> = \{I(W_j \ge L) : 1 \le j \le p\}^>$ . Denote  $A = \{j : \delta_j = 1\}$  the discovery set. To see why (10) makes sense, note that  $\{j : W_j \le -t\}$  is an overestimation of  $\{j : W_j \le -t, j \ge A^c\}$ , which is asymptotically equal to  $\{j : W_j \ge t, j \ge A^c\}$ , the number of false positives, due to the asymptotic symmetry property (9). It follows that the

We shall see that S contains all signals, then the LSEs of the null coordinates are symmetrically distributed around 0. Hence  $W_j$ 's satisfy (4). It is easy to see that (9) is an asymptotic version of the symmetry property given by (4); see Lemmas (5.1) (5.2) in Section (4) of the Supplementary Material for a rigorous discussion.

fraction in (10) provides an overestimate of the FDP, which (desirably) leads to a conservative FDR control. Moreover, the empirical FDR level is typically very close to  $\alpha$  because the gap between the fraction in (10) and the actual FDP is usually small in practice, where, for a suitably chosen L, most cases in  $\{j: W_j \leq -L\}$  should come from the null.

The operation of the SDA filter can be visualized in Figure [2]. We generate  $\{\xi_i: i=1,\ldots,90\}$  from an MVN distribution with  $\mu$   $\mathbb{R}$   $R^{p=1000}$  and  $\Sigma=(0.8^{\lfloor i-j \rfloor})_{1 \le i,j \le p}$ . We randomly set 10% of the coordinates in  $\mu$  to be 0.2 and 0 elsewhere. Panel (a) presents the scatter plot of 288 nonzero  $W_j$ 's with red triangles and black dots respectively denoting true signals and nulls. Panel (d) plots the normalized knockoff statistics that are constructed according to (1.7) in Barber and Candès (2015) 1. We can see that both SDA and knockoff fulfill the symmetry property approximately for the null  $W_j$ 's (black dots). However, SDA achieves a more clearcut separation of signals and noise. As explained in Section B.1 of the Supplement, the symmetrized knockoff statistics suffers from high correlations. By contrast, the construction of SDA statistic, which does not depend pairwise contrasts, eliminates the needs for creating fake variables. We can see from Panel (a) that the SDA ranking places most true signals above 0, and many true signals stay well above the majority of the null cases. However, in Panel (d) that illustrates the knockoff ranking, the true signals are not well separated from the nulls, and many true signals even fall below 0. Since the threshold must be positive, signals with negative  $W_j$ 's will be missed, which leads to substantial power loss.

The impacts on the FDP processes are shown in the second column in Figure  $\mathbb{Z}$ . We can see that the estimated FDP process [ $\mathbb{E}DP(t)$ ] of SDA approximates the true FDP process [FDP(t)] fairly accurately. However, the knockoff method yields overly conservative estimates of the true FDPs, which leads to overly conservative thresholds (marked by blue vertical lines). The last column in Figure  $\mathbb{Z}$  compares the TDP processes of SDA and knockoff. At the FDR level 0.2, the TDP of SDA is 0.87 (threshold L = 0.62), which is much higher than that of knockoff (TDP=0.03 with threshold L = 6.80). The low TDP of knockoff is due to the decreased power in distinguishing the signal from noise [Panel (d)] and an overly conservative threshold [Panel (e)].

The normalization, which makes the plot easier to read, does not affect the results of the knockoff method. This is because only the relative magnitudes of W<sub>i</sub> matter in the thresholding step of the knockoff method.

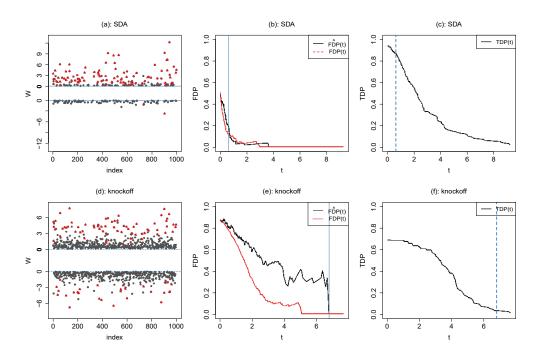


Figure 2: (a): Scatter plot of the 288 nonzero  $W_j$ s from the SDA filter with red triangles and black dots denoting true signals and nulls respectively. A vertical space is added to the middle of the plot to better contrast positive and negative  $W_j$ 's. (b): the corresponding estimate of FDP curve (against t) along with the true FDP for the SDA filter; (c): the true power curve (against t) for the SDA filter. (d)-(f): the scatter plot of  $p = 1000 W_j$ s, the corresponding FDP estimate, and the true power for the knockoff method.

## 2.3 Power and effects of dependence

The impact of dependence on FDR analysis has been extensively studied but most discussions have focused on the validity issue. This section first discusses the impact of dependence on power, and then provides insights on the information loss of conventional data splitting methods.

Under the SDA framework, many possible pairs of  $(T_{1j}, T_{2j})$  may be constructed. It is easy to show that  $W_i$  constructed via the pairs of sample averages

$$(\mathsf{T}_{1j}^{\,0},\mathsf{T}_{2j}^{\,0}) = (\sqrt[4]{\mathsf{n}_{1}\bar{\xi}_{1}},\sqrt[4]{\mathsf{n}_{2}\bar{\xi}_{2}}) \tag{11}$$

also fulfill the asymptotic symmetry property. However, the pair in (11), which falls into the class of marginal testing techniques, can be highly inefficient since it completely ignores the dependence structure. Next we provide intuitions on how the dependence structure is incorporated into the SDA filter to improve the efficiency of existing methods.

First,  $T_{1j}$  is superior to  $T_{1j}$  by everaging joint modeling techniques. The merit of joint modeling has been carefully illustrated by Barber and Candès (2015) through extensive simulations. Candès et al. (2018) further argued that the conditional testing techniques are in general more powerful in recovering sparse signals than marginal testing methods.  $T_{1j}$  is constructed based on LASSO (a conditional inference technique) and serves as a more suitable building block than  $T_{1j}^0$  for constructing  $W_{1j}$ . Second,  $T_{2j}$  enjoys a higher SNR than  $T_{2j}^0$  by exploiting the dependence between  $\xi_S$  and  $\xi_{S^c}$ . Clearly, the expectations of both  $\mu_{2S}$  and  $\bar{\xi}_{2S}$  are  $\mu_{2S}$ . The covariance of  $\mu_{2S}$  is  $n_2^{-1}Q$ , where  $Q=(X_S^{\times}X_S)^{-1}$ . By the inversion formula of a block matrix, we have  $X_S^{\times}X_S = \Omega_{S,S} = \Sigma_{S,S} - \Sigma_{S,S} \varepsilon \Sigma_{S^c} \Sigma_{S^c} \varepsilon_{S^c} \varepsilon_{S^c} S_{S^c} S_{S^c}$ 

Finally, both knockoff and SDA achieve the symmetry property at the expense of possibly reduced SNR: the former increases the dimension of the design matrix by adding noise variables while the latter involves sample splitting. In contrast with the sample splitting method in Wasserman and Roeder (2009), where  $D_1$  is thrown away after model selection, SDA provides a new aggregation strategy:  $T_{1j}$  is kept and combined with  $T_{2j}$  to form the ranking statistic  $W_j$ . This substantially reduces the information loss in conventional sample splitting methods.

## 2.4 Effects of data screening

The data screening step is always beneficial as long as the tests are correlated. Intuitively, the smaller the set S, the larger amount of uncertainty can be explained by the variables in  $S^c$ . Hence a more effective dimension reduction implies increased SNR and higher power. Meanwhile, our theory on FDR control requires that  $P(A \ B \ S)$  holds with high probability, indicating that an overly aggressive data screening step can hurt the FDR procedure. In practice, we recommend

deliberately choosing an overfitted model to ensure the validity in FDR control; this would slightly compromise the power. To illustrate the tradeoff, Figure presents a numerical study to investigate how the size of S may affect both the FDR and power. We can see that the actual FDRs of SDA may deviate from the nominal level when S is too small. By contrast, a large S (overfitted model) has little impact on the FDR levels, but affects the power negatively.

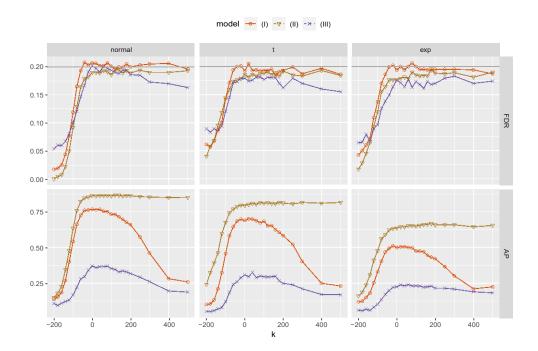


Figure 3: The effects of data screening. We choose n = 90, p = 500, and  $\mu = \pm 0.2$ . The proportion of non-nulls is 10% and  $\alpha = 0.2$ . We investigate the performance of SDA over 3 distributions and 3 covariance structures described in Section 5. Here k denotes the excess counts of |S| with  $\lambda$  selected by the ATC criterion (k can be negative).

## 3 Theoretical Properties of the SDA Filter

This section first establishes finite sample theory for FDR bounds (Section 8.1), and then develops asymptotic theories for FDR and FDP control.

## 3.1 Finite-sample theory on FDR control

Our finite—sample theory, which requires no model assumptions, establishes an upper bound for the FDR under general dependence. We emphasize that the upper bound holds for both known and estimated covariance matrices.

Our theory is developed for a modified SDA filter (SDA+) which chooses the threshold

$$1 + \#\{j : W_j \le -t\}$$
   
 L = inf t > 0:  $\#\{j : W_j \ge t\} \boxtimes 1 \le \alpha$ .

SDA+ is slightly more conservative than SDA but their difference is negligible when the number of rejections is large. Recall  $S = \{j : | \mathbf{b}_{1j} = 0\}$ . Denote  $W_S = (W_j : j \otimes S)^{>}$  and  $W_{-j} = W_S \setminus W_j$ . The key quantity that controls the upper bound is

$$\Delta_{i} = |Pr(W_{i} > 0 | |W_{i}|, W_{-i}) - 1/2|, \qquad (12)$$

which can be interpreted as a measure of the extent to which the "flip-sign" property of  $W_j$  is violated. Our finite sample theory for FDR control is given by Theorem  $\square$ 

Theorem 1 For any  $\alpha \ 2 \ (0,1)$ , the FDR of the SDA+ method satisfies

$$FDR \leq \min_{j \geq 0} \alpha(1+5) + Pr \max_{j \in A} \sum_{i \leq j} \Delta_{ij} > .$$
 (13)

Our theorem is closely connected to Theorem 1 in Barber et al. (2020). Both theorems involve assessing how the deviations from the "idealized situation" would affect the actual FDR level. However, the interpretations are very different. In model-X knockoff the deviation (from the assumption of a known X matrix) comes from the estimation errors of the X matrix whereas in SDA the deviation (from the perfect symmetry property) comes from the possible violations of the normality assumption and sure screening property. Our theorem shows that a tight control of  $\Delta_j$ 's leads to effective FDR control. Next we carefully interpret the bound and present several important settings in which the upper bound in (13) exactly achieves or is very close to the nominal level  $\alpha$ .

For a null variable (i.e.  $j \boxtimes A^c$ ), the flip-sign property means that  $W_j$  is equally likely to be positive or negative conditioning on its magnitude and other  $W_k$ 's in S.

Consider the ideal case where (a) the error distribution is symmetric, (b) S contains all signals and (c)  $W_j$ 's are independent of each other for  $j \ \mathbb{Z} \ S$ . We can show that  $\Delta_j = 0$  for all  $j \ \mathbb{Z} \ A^c \cap S$ . The upper bound achieves the nominal level  $\alpha$  exactly since  $Pr(W_j > 0 \mid |W_j|, W_{-j}) = Pr(W_j > 0 \mid |W_j|) = 1/2$  and hence we can set = 0. Even when the error distribution is asymmetric, we expect that  $\Delta_j$ 's would become vanishingly small for moderate sample size n due to the convergence of  $\mathbf{p}_{2j}$  to a symmetric distribution (Lemma  $\mathbf{S}.\mathbf{1}$ ). Hence the FDR bound would be close to  $\alpha$ .

Next we turn to the dependent case. For simplicity, assume that  $\xi_i$ 's come from a multivariate normal distribution. Let  $Q=(X_S^*X_S)^{-1}:=(Q_{jk})_{q_n\times q_n}$  with  $q_n=|S|$ . The matrix  $Q=\Sigma_{S,S}-\Sigma_{S^c,S^c}\Sigma_{S^c,S^c}\Sigma_{S^c,S^c}$  is the conditional covariance matrix of  $\xi_S$  given  $\xi_{S^c}$ . The following lemma shows that the magnitude of  $\Delta_j$  is controlled by the matrix Q.

Lemma 1 (Flip-sign property under Gaussian dependence). Assume that  $\xi_i$ 's obey a multivariate normal distribution. Denote  $Q_{-j,j}$  the jth column of Q excluding  $Q_{jj}$ . If  $Q_{-j,j} = 0$ , then  $\Delta_j = 0$ .

To provide some intuitions on how close the bound is to  $\alpha$  in practice, consider the autoregressive (AR) structure  $\Sigma = (\sigma_{j,l}) = (\rho^{\lfloor j-l \rfloor})$ . Since the precision matrix of AR structure is tridiagonal, only consecutive coordinates are correlated with each other conditional on remaining variables. Suppose sparse signals are randomly distributed on the p coordinates and the dimension reduction via S is performed effectively, e.g.  $q_n$  p. Let E be an event such that for any null variable  $j \ \mathbb{P} \ S \cap A^c$ , remaining variables in S are conditionally uncorrelated with it. We expect E to occur with high probability since for large tridiagonal precision matrices, there is a small chance that two consecutive coordinates are selected into a small set S simultaneously. On event E, we have  $Q_{-j,j} = 0$  and it follows from Lemma  $\mathbb{T}$  that  $\Delta_j = 0$ . Consequently the FDR bound would converge to  $\alpha$  when  $\Pr(E) \to 1$ . In the same vein, we expect that the bound would be close to  $\alpha$  for the class of power decay covariance matrices and the class of sparse precision matrices.

#### 3.2 Asymptotic theory on FDP control

Under the asymptotic paradigm we can prove that the FDR can be controlled at  $\alpha$  + o(1) under suitable conditions (asymptotic validity). Denote  $\epsilon_i$  = X( $\xi_i$  -  $\mu$ ). Let d<sub>n</sub> = |A|, q<sub>n</sub> = |S|, q<sub>0n</sub> =

 $|S \cap A^c|$ , and  $A(S) := (X_S^> X_S)^{-1} X_S^> = (a_{jk})_{q_n \times p}$ . Assume that  $q_n$  is uniformly bounded above by some non-random sequence  $\bar{q}_n$  that will be specified later. We start with some regularity conditions.

Condition 1 (Sure screening property) As  $n \to \infty$ ,  $Pr(A \square S) \to 1$ .

Remark 2 Condition 1 ensures that  $\mathbf{p}_{2j}$  is unbiased for j 2 S. This pre–selection property, which has been commonly used (Wasserman and Roeder) 2009; Meinshausen et al., 2009; Barber and Candès, 2019), can be fulfilled with suitably chosen  $\lambda$  under the "zonal" assumption (Bühlmann and Mandozzi, 2014). In practice, we recommend applying AIC to deliberately choose an overfitted model. The sure screening property may not hold exactly but missing small  $\mu_j$ 's is inconsequential. For example, if we ignore "unimportant" signals, then Condition 1 is fulfilled by LASSO for large signals exceeding the rate of  $d_n$   $\log p/n$ . Asymptotically unbiased estimators are usually sufficient for effective FDR control. This has been corroborated by our empirical results in Section 5.

Remark 3 Condition 2 assumes that  $\mu_1$  is a reasonable estimator of  $\mu$ ; this condition typically holds with  $c_{np} = d_n \frac{p}{\log p/n}$  for the LASSO solution (Van de Geer and Bühlmann) 2009).

The next two conditions are standard: Condition  $\blacksquare$  imposes constraints on the diverging rates of  $\bar{q}_n$  and p, both of which depend on the existence of certain moments; Condition  $\blacksquare$  requires that the eigenvalues of the design matrix are doubly bounded by two constants.

Condition 4 (Covariance) There exist positive constants  $\bar{\kappa}$  and  $\underline{\kappa}$  such that with probability one,

$$\underline{\kappa} \leq \lim\inf_{n \to \infty} \lambda_{\min}(X_S^> X_S) < \lim\sup_{n \to \infty} \lambda_{\max}(X_S^> X_S) \leq \bar{\kappa}.$$

Condition 5 (Signals) As  $n, p \rightarrow \infty$ ,  $\eta_n = |C_{\mu}| \rightarrow \infty$ , where

$$C_{\mu} = \{j \ \mathbb{Z} \ A : \mu_i^2 / \{\max(c_{np}^2, \log \bar{q}_n / n)\} \rightarrow \infty \}.$$

Remark 4 Condition Simplies that the number of identifiable effect sizes should not be too small as  $p \to \infty$ . This seems to be a necessary condition for FDP control. For example, Liu and Shao (2014) showed that if a multiple testing method controls the FDP with high probability, then its number of true alternatives must diverge when the number of tests goes to infinity.

Condition 6 (Dependence) Let  $\rho_{j\,k} = Q_{j\,k}/p^{p} \overline{Q_{j\,j}\,Q_{k\,k}}$ . Assume that for each j, Card $\{1 \le k \le q_n : |\rho_{j\,k}| \ge C(\log n)^{-2-\nu}\} \le r_p$ , where C > 0,  $\nu > 0$  is any small constant, and  $r_p/\eta_n \to 0$  as  $n, p \to \infty$ .

Remark 5 Condition allows  $\xi_j$  to be correlated with all others but requires that the number of large correlations cannot diverge too fast. The condition appears to be similar to the regularity conditions in Fan et al. (2012) and Xia et al. (2020) but in fact our condition is much weaker. For instance, the correlation between  $\mathbf{p}_{2j_1}$  and  $\mathbf{p}_{2j_2}$  is just the partial correlation of  $\xi_{j_1}$  and  $\xi_{j_2}$  given the rest variables. In particular, large correlations would be highly unlikely after data screening for a wide range of popular models, such as the class of power decay covariance matrices and the class of moderately sparse precision matrices. This reveals the advantage of SDA, which effectively de–correlates the strong dependence via data screening and conditioning.

Our main theoretical result on the asymptotic validity of the SDA method for both FDP and FDR control is given by the next theorem.

Theorem 2 Suppose Conditions 16 hold. For any  $\alpha \ 2 \ (0,1)$ , the FDP of the SDA method satisfies

$$FDP_{W}(L) := \frac{\#\{j : W_{j} \ge L, j ? A^{c}\}}{\#\{j : W_{j} \ge L\} ? 1} \le \alpha + o_{p}(1).$$
 (14)

It follows that  $\limsup_{(n,p)\to\infty} FDR \leq \alpha$ .

## 4 Unknown dependence

Now we turn to the case where the covariance structure is unknown. When  $\Omega$  is unknown, the SDA filter operates in the same way except that we substitute the estimate  $\Phi$  in place of  $\Omega$ .

We propose to estimate  $\Omega$  using only the first part of the sample  $D_1$ . Denote  $\Omega$  the corresponding estimator. Then the SDA filter can be readily constructed via the steps in Sections 2.1-2.2 with  $X = \Omega^{1/2}$ . Various high-dimensional precision matrix estimation methods, such as the graphical LASSO (Friedman et al., 2008) and CLIME (Cai et al., 2011), can be used to obtain  $\Omega$ . An attractive feature of the SDA filter under unknown dependence is its robustness for FDR control. We next show that the SDA filter is robust for FDR control if  $\Omega$  is constructed based only on  $D_1$ . We first state a modified version of Condition 6, which uses  $\Omega$ 0 in place of  $\Omega$ .

Condition 6' Let  $Q^0 = (X_S^> X_S)^{-1} X_S^> X_S^{-1} X_S^> X_S^{-1} X_S^> X_S^{-1} := (Q_{jk}^0)_{q_n \times q_n}$  and  $\rho_{jk}^0 = Q_{jk}^0 / Q_{jj}^0 Q_{kk}^C$ . Assume that for each j,  $Card\{1 \le k \le q_n : |\rho_{jk}^0| \ge C(\log n)^{-2-v}\} \le r_p$ , where C > 0, v > 0 is any small constant, and  $r_p / q_n \to 0$  as  $n, p \to \infty$ .

The following theorem, which is in parallel with Theorem 2 establishes the asymptotic validity of the SDA filter for estimated covariance.

Theorem 3 Let  $\Omega$  denote an estimator based on  $D_1$ . Suppose Conditions  $\Omega$  and 6' hold. Then the FDP of the SDA method utilizing  $X = \Omega^{1/2}$  satisfies FDP  $\leq \alpha + o_p(1)$ . It follows that  $\limsup_{(n,p)\to\infty} \mathsf{FDR} \leq \alpha$ .

Remark 6 Our FDR theory does not require an accurate estimator for  $\Omega$ . The accuracy of the estimator only affects the power but not the validity. Consider a working covariance structure that "estimates"  $\Omega$  as the identity matrix. Then it can be shown that the FDP can still be controlled. This is more attractive than the FDR theories in, for example, Fan and Han (2017) and Li and Zhong (2017) that critically depend on the accuracy of the covariance estimators.

The key step in the proof is to verify the validity of  $(\underline{\mathbb{P}})$ . This amounts to addressing two major issues: the asymptotic symmetry of  $W_i$  under the null and the uniform convergence of

 $q_{0n}^{-1} \stackrel{P}{\underset{j \boxtimes S \cap A^c}{}} I(W_j \geq t). \text{ Because } \Phi \text{ is obtained from } D_1, \text{ then } \blacktriangle_{2j} \text{ is unbiased conditional on } D_1 \text{ and } P \\ \text{thus } \stackrel{P}{\underset{j \boxtimes S \cap A^c}{}} P(W_j > t) \text{ is approximately equal to } \stackrel{P}{\underset{j \boxtimes S \cap A^c}{}} P(W_j < -t), \text{ establishing the symmetry property. The dependence assumption on } Q^0 \text{ ensures the convergence of } q_{0n}^{-1} \stackrel{P}{\underset{i \boxtimes S \cap A^c}{}} I(W_j \geq t).$ 

While sample–splitting ensures the independence between  $\not b_1$  and  $\not b_2$  and hence the robustness of the SDA filter, as one would expect, a more accurate estimate of  $\Omega$  yields better power. Previously we have proposed to estimate  $\Omega$  using  $D_1$  and construct the LSE ( $\square$  using  $D_2$ . In practice one may consider using  $D_1$  to construct  $T_{1j}$ , and then obtaining the LSE via the full sample estimator, denoted  $\not D_F$ , that is estimated using  $\{D_1,D_2\}$ . The caveat is that, although  $X=\not D_F^{1/2}$  can potentially increase the power, stronger conditions will be needed to guarantee the asymptotic validity of the "full–sample" SDA method. As pointed out by an insightful referee, the asymptotic theory requires that  $\not D_F$  must converge to  $\Omega$  at a very fast rate, which can be impractical in applications. We recommend the robust SDA filter that estimates  $\Omega$  using only  $D_1$ . Next we specify the requirements on the estimation accuracy of  $\not D_F$ .

Condition 7 The estimated precision matrix  $\Phi_F$  satisfies  $k\Phi_F - \Omega k_\infty = O_p(a_{np})$  with  $a_{np} \to 0$ .

The following theorem shows that the FDR and FDP can be controlled asymptotically when  $\Phi_F$  is sufficiently close to  $\Omega$ . Let  $s_n = k\Omega k_\infty$ .

Theorem 4 Consider a modified SDA procedure where we use  $D_1$  to construct  $T_{1j}$  and the full sample estimator  $\mathfrak{Q}_F$  to construct the LSE (7). Suppose Conditions 1-6 hold and  $\mathfrak{Q}_F$  satisfies Condition 7. Then, if

$$c_{np}a_{np}s_{n}\bar{q}_{n} \frac{p}{n \log p(\log \bar{q}_{n})^{1+\gamma}} \rightarrow 0$$
(15)

for a small  $\gamma$  > 0, the results in Theorem 2 hold for the procedure with  $\Phi_{\text{F}}$  .

This theorem, which is a complementary result to Theorem 3, provides conditions that warrant the implementation of a more efficient version of SDA. It is worth further investigating the condition (15), which seems to be unavoidable because  $T_{1j}$  and  $T_{2j}$  are no longer independent when the whole sample is used to estimate  $\Omega$ . To fix ideas, suppose that  $\Omega = (\omega_{ij})_{p \times p}$  is  $k_n$ -sparse, i.e.

 $\max_{1\leq i\leq p} \bigcap_{j=1}^p I(\omega_{ij}=0) \leq k_n$ , and that all its elements  $\omega_{ij}s$  are bounded. First, standard arguments in, for example, Yuan (2010) and Liu et al. (2012) indicate that  $a_{np}=O_p(k_n^p \overline{\log p/n})$ . Accordingly, with  $c_{np}=d_n^p \overline{\log p/n}$ , Equation (15) is equivalent to the condition  $d_n k_n s_n \overline{q}_n / n^{1/2} \to 0$  if p is of a polynomial rate of n. The condition above imposes restrictions on the diverging rates of  $d_n$ ,  $k_n$ ,  $s_n$  and  $\overline{q}_n$  Assume that  $d_n$ ,  $k_n$  and  $s_n$  are all bounded. Then we must require that  $\overline{q}=o(n^{1/2})$ . Alternatively, if we only assume that  $k_n$  and  $s_n$  are bounded, then a sufficient condition for (15) is  $\overline{q}\equiv o(n^{1/4})$  (since  $d_n\leq \overline{q}$ ). These rates are consistent with those in the literature; see, for example, Portnoy et al. (1984) and Fan and Peng (2004).

## 5 Simulation

This section first introduces the R package sdafilter (Section 5.1), followed by simulation designs (Section 5.2) and comparison results (Section 5.3). Additional results for comparisons with unknown covariance matrix and other correlation structures are provided in the Supplementary Material.

## 5.1 Implementation details

We describe the <u>implementation details of the R package sdafilter</u>. For sample–splitting, we follow the strategy in Wasserman and Roeder (2009), which uses  $n_1 = [2/3n]$  for selecting variables, and the rest  $n_2 = n - n_1$  for obtaining the LSEs. The AIC is used to select the tuning parameter in LASSO. If the number of the variables selected by AIC exceeds [p/3], then only the first [p/3] variables will be retained. For the case with unknown  $\Omega$ , our default option is to apply the R package glasso to  $D_1$ , where the tuning parameter is set by the R package huge. If prior knowledge suggests a nonsparse  $\Omega$ , the "nonsparse" option in our package can be used. This option first estimates the covariance matrix using the R package POET and then takes its inverse as the input. The stable option implements the R-SDA method described in Section A.1 of the Supplementary Material. The kwd option enables the usage of different estimators to summarizes the information in the first part of data, including the de-biased LASSO, innovated transformation of the sample means (Hall and Jin) 2010), and factor-adjusted sample means (Fan and Han) 2017).

## 5.2 Simulation settings

We consider three types of covariance structures: (I) Autoregressive (AR) structure:  $\Sigma = (\rho^{\lfloor j-i \rfloor})$ . (II) Compound symmetry structure: all off-diagonal elements of the  $\Sigma$  are  $\rho$ , which can be regarded as a factor model with one principal component. (III) Sparse covariance structure:  $\Sigma = \Gamma \Gamma^{>} + I_{p}$ , where  $\Gamma$  is a p × p matrix and each row of  $\Gamma$  has only one position with nonzero value sampled from uniform distribution [1, 2].

The diagonal elements are normalized as unity for all three settings. To investigate the robustness of different methods, we consider three error distributions: (i) multivariate normal; (ii) t-distribution with df = 3 and (iii) exponential distribution with scale parameter 2. The observations are then standardized to have mean zero and standard deviation one. The correlation structure remains nearly unchanged after transformation. The following six methods will be compared:

- (a) The Benjamini-Hochberg (BH) procedure with the p-values transformed from the t statistics.
- (b) The principal factor approximation (PFA) procedure proposed by Fan et al. (2012) for known covariance and Fan and Han (2017) for estimated covariance. Two versions of the PFA procedure using the unadjusted p-values and adjusted p-values are implemented using the R package pfa, denoted as PFA<sub>U</sub> and PFA<sub>A</sub> respectively. We only report the results for PFA<sub>A</sub> as it generally outperforms PFA<sub>U</sub>.
- (c) The sample-splitting method (SS; Wasserman and Roeder, 2009), which conducts data screening using LASSO and then applies BH to the p-values calculated based on  $\mu_2$ .
- (d) The knockoff method (Knockoff; Barber and Candès, 2015), which is implemented using function "create.fixed" in the R package knockoff.
- (e) The DATE method (DATE; Li and Zhong, 2017), which we implemented by ourselves.
- (f) The stability-refined SDA filter (R-SDA) implemented using our package sdafilter with the "stable" option. We only presented R-SDA, which we recommend to use in practice, to make the plots easier to read. SDA has similar performance to R-SDA.

Let n be the sample size, p the number tests, and  $\pi_1$  the proportion of signals. For each combination  $(n, p, \pi_1)$ , we generate data and apply the six methods at FDR level  $\alpha$ . The FDR and AP are calculated by averaging the proportions from 500 replications.

#### 5.3 Comparison results for known covariance structures

We fix  $(n, p, \pi_1, \alpha) = (90, 500, 0.1, 0.2)$  and generate  $\mu_j$  from the following random mixture model:

$$\mu_j \stackrel{i.j.d}{\supseteq} (1 - \pi_1)\delta_0 + \pi_1 g(\cdot), \quad j = 1, \dots, p,$$

where  $\delta_0$  is the dirac delta function (denoting a point mass at 0), and  $g(\cdot)$  is the density of the non-null distribution, specified as a uniform distribution  $[\mu_0 - 0.1, \mu_0 + 0.1]$ . The signals  $\mu_j$ 's are then randomly multiplied by a flip-sign. To assess the effect of signal strength, we vary  $\mu_0$  from 0.1 to 0.3 and apply the six methods to simulated data. The results for Structures (I) and (III) are summarized in Figure 4 where in the top row we fix  $\rho = 0.8$ . The results for Structure (II) with  $\rho = 0.8$  are shown in Figure 5 of the Supplementary Material. The following observations can be made.

- (a) For the Gaussian error case, BH, knockoff, R-SDA and SS control the FDR at the nominal level. The FDR levels of PFA $_{\rm A}$  and DATE are inflated when signals are weak.
- (b) For the non-Gaussian error case, BH, DATE, SS and PFA<sub>A</sub> fail to control the FDR under various settings and the FDR levels can be much higher than the nominal level. Knockoff controls the FDR in all settings but can be very conservative. R-SDA has the most accurate and stable FDR levels among all methods.
- (c) R-SDA vs SS and BH. As expected, SS and BH control the FDR under the Gaussian case but are not robust for non-Gaussian errors. R-SDA has much higher power than both methods (even when the FDR levels of R-SDA are much lower). It is interesting to note that although SS only uses the second part of the data, its power can be much higher than BH when the correlation structure is highly informative [Normal case under Structure (I) on top left]. This is because the data screening step can significantly increase the SNR (Section 2.3).

- (d) R-SDA vs Knockoff. R-SDA and knockoff, both of which are distribution—free, are the only methods that can control the FDR at the nominal level across all scenarios. The knockoff method is overly conservative in Setting (I) due to the high correlation. The conservativeness become less severe under Setting (III). By contrast, R-SDA controls the FDR more accurately near the target level and has significantly higher power than knockoff.
- (e) R-SDA vs DATE and PFA<sub>A</sub>. In some scenarios, DATE and PFA<sub>A</sub> can outperform SDA in power. However, the higher power may be attributed to the severely inflated FDRs. The numerical results reveal the promise of extending the SDA framework by employing other methods, such as factor–adjusted z-scores or innovated transformations, as alternatives to the LASSO estimates, to construct T<sub>1i</sub>.

Next we turn to investigate how the six methods are affected by the strength of correlation. For covariance structures (I) and (II), we fix  $\mu = 0.2$  under alternative and vary the magnitude of correlation  $\rho$  from independence ( $\rho = 0$ ) to strong dependence ( $\rho = 0.9$ ). The results are summarized in Figure 5. In addition to the observations that we have made based on the previous graph, the following additional patterns are worthy of mentioning.

- (a) The knockoff method becomes more conservative when correlations become higher. Note that the average correlations in Structure (II) is much higher than that in Structure (I), the power of the knockoff method deteriorates faster for Structure (II) as ρ increases. For Structure (II), the FDR of BH also decreases as ρ increases.
- (b) In contrast with BH and knockoff, both of which suffer from high correlations, the FDR of R-SDA remains at the nominal level consistently, and the power increases with the correlation.

  The power grows faster for Structure (II). This corroborates the insights that high correlations can be useful in FDR analysis (Benjamini and Heller, 2007; Sun and Cai, 2009).
- (c) In Column 2 of Figure 5 knockoff fails to control the FDR for heavy tailed distributions when correlation is low. By contrast, SDA controls the FDR accurately under non-Gaussian errors.

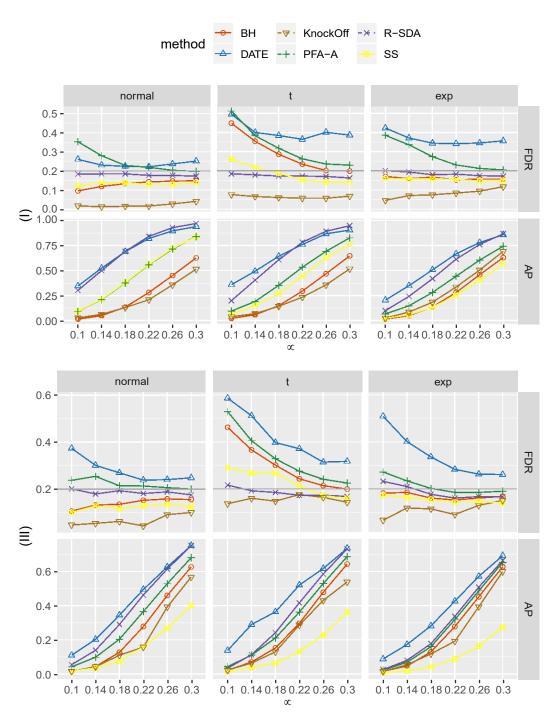


Figure 4: FDR and AP comparison for varying  $\mu$  in Settings (I) and (III) with known variance.

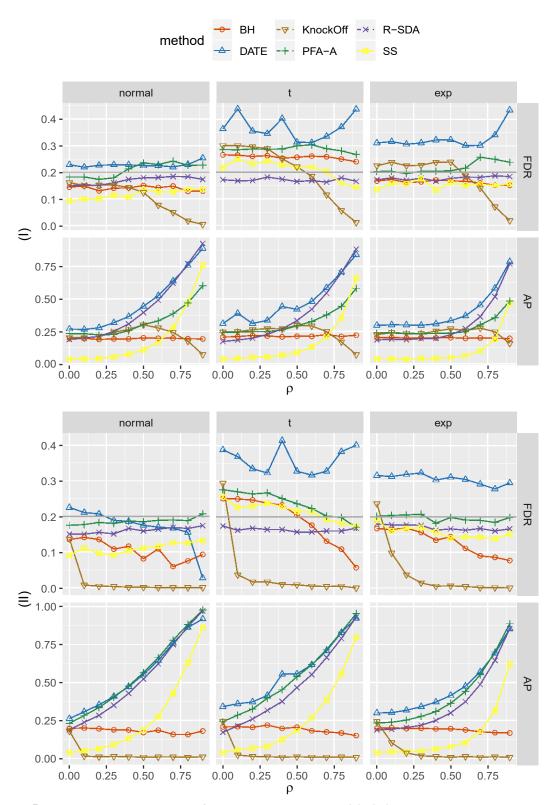


Figure 5: FDR and AP comparison for varying  $\rho$  in Settings (I)–(II) with known covariance matrix.

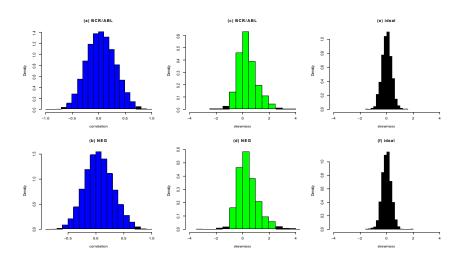


Figure 6: (a)-(b): Histograms of the off-diagonal elements of the sample correlation matrix for BCR/ABL and NEG; (c)-(d): Histogram of the skewness of the p = 1263 genes for BCR/ABL and NEG; (e)-(f): the ideal patterns of (c)-(d) when the data are normal.

## 6 A real-data example

This section illustrates the SDA filter for analysis of high-density oligonucleotide microarrays. The data set, which contains 12,625 probe sets from 128 adult patients enrolled in the Italian GIMEMA multi-center clinical trial, has been used in Chiaretti et al. (2005) and Bourgon et al. (2010) for identifying genetic factors that are associated with acute lymphoblastic leukemia (ALL). The ALL dataset is available at <a href="http://www.bioconductor.org">http://www.bioconductor.org</a>.

We focus on a subset of 79 patients with B-cell differentiation because existing research reveals that malignant cells in B-lineage ALL are often associated with genetic abnormalities that have significant impacts on the clinical course of the disease. The patients are divided into two groups based on the molecular heterogeneity of the B-lineage ALL: 37 with the BCR/ABL mutation and 42 with NEG. We further narrow down the focus to 10% of the genes (i.e., p = 1,263) before carrying out the FDR analysis. Specifically, the uncorrelated screening method (Bourgon et al., 2010) has been used to remove probe sets with small overall sample variances since they are unlikely to be differentially expressed.

We apply a two-sample version of R-SDA (see Section A.3 for details), BH, SS, PFAA, Knockoff

and DATE at several significance levels for identifying differentially expressed genes across the two groups. Table I summarizes the number of significant probe sets for each method. In Figure (a)-(b), we plot the pairwise correlations of the genes. We can see that a significant proportion of the correlations exceed 0.4. These correlations can jointly exhibit non-negligible dependence effect. This explains why the knockoff method is overly conservative. R-SDA is more powerful than SS by exploiting additional information from the second part of data. BH, PFA<sub>A</sub> and DATE claims more significant genes than R-SDA. However, some caveats need to be given regarding the reliability of BH, PFA<sub>A</sub> and DATE, which all require normality assumptions (and the latter two require accurate estimates of the unknown covariance matrices).

Next we conduct a preliminary analysis to investigate the normality assumption, which seems to have been severely violated in this data set. From Column 2 of Figure 1 we can see that the skewness scores of many genes exceed the conventional cutoff ±1. As a comparison, we display in Column 3 of Figure 1 the "ideal" pattern where the normality assumption holds. The histograms in Column 2 are much wider than the histograms in Column 3, indicating a possibly highly skewed error distribution. One possible explanation for the difference in power is that BH, PFA-A and DATE may have inflated FDR levels under violation of normality. This has been observed in our simulation studies (e.g. last column in Figure 3. By contrast, SDA and knockoff are distribution—free methods, which tend to produce more reliable and replicable findings. The lists of 19 highest ranked probe sets by the six methods are presented in Table 1 of Appendix 1.

Table 1: The number of rejections for six multiple testing procedures and various significance levels.

	R-SDA	SS	вн	PFA-A	Knockoff	DATE
$\alpha = 0.01$	19	7	29	98	2	364
$\alpha = 0.05$	33	15	146	182	2	452
$\alpha = 0.10$	56	37	229	252	2	501
α = 0.20	139	68	350	339	7	546

## Acknowledgments

The authors thank the Editor, Associate Editor and two anonymous referees for their many helpful comments that have resulted in significant improvements of the article.

## References

- Abramovich, F., Benjamini, Y., Donoho, D. L., and Johnstone, I. M. (2006), "Adapting to unknown sparsity by controlling the false discovery rate," The Annals of Statistics, 34, 584–653.
- Barber, R. F. and Candès, E. J. (2015), "Controlling the false discovery rate via knockoffs," The Annals of Statistics, 43, 2055–2085.
- (2019), "A knockoff filter for high-dimensional selective inference," The Annals of Statistics, 47, 2504–2537.
- Barber, R. F., Candès, E. J., and Samworth, R. J. (2020), "Robust inference with knockoffs," The Annals of Statistics, 48, 1409–1431.
- Barras, L., Scaillet, O., and Wermers, R. (2010), "False discoveries in mutual fund performance: Measuring luck in estimated alphas," The journal of finance, 65, 179–216.
- Benjamini, Y. and Heller, R. (2007), "False discovery rates for spatial signals," Journal of the American Statistical Association, 102, 1272–1281.
- Benjamini, Y. and Hochberg, Y. (1995), "Controlling the false discovery rate: a practical and powerful approach to multiple testing," Journal of the Royal Statistical Society: Series B (Methodological), 57, 289–300.
- (1997), "Multiple Hypotheses Testing with Weights," Scandinavian Journal of Statistics, 24, 407–418.
- Benjamini, Y. and Yekutieli, D. (2001), "The control of the false discovery rate in multiple testing under dependency," The Annals of Statistics, 29, 1165–1188.
- Bickel, P. J. and Levina, E. (2008), "Regularized estimation of large covariance matrices," The Annals of Statistics, 36, 199–227.
- Bogdan, M., Van Den Berg, E., Sabatti, C., Su, W., and Candès, E. J. (2015), "SLOPE-daptive variable selection via convex optimization," The Annals of Applied Statistics, 9, 1103.
- Bourgon, R., Gentleman, R., and Huber, W. (2010), "Independent filtering increases detection power for high-throughput experiments," Proceedings of the National Academy of Sciences, 107, 9546–9551.
- Bühlmann, P. and Mandozzi, J. (2014), "High-dimensional variable screening and bias in subsequent inference, with an empirical comparison," Computational Statistics, 29, 407–430.
- Cai, T. and Liu, W. (2016), "Large-scale multiple testing of correlations," Journal of the American Statistical Association, 111, 229–240.
- Cai, T., Liu, W., and Luo, X. (2011), "A constrained I<sub>1</sub> minimization approach to sparse precision matrix estimation," Journal of the American Statistical Association, 106, 594–607.
- Cai, T. T., Sun, W., and Wang, W. (2019), "CARS: Covariate assisted ranking and screening for large-scale two-sample inference (with discussion)," Journal of the Royal Statistical Society: Series B (Methodological), 81, 187–234.
- Caldas de Castro, M. and Singer, B. H. (2006), "Controlling the false discovery rate: a new application to account for multiple and dependent tests in local statistics of spatial association," Geographical Analysis, 38, 180–208.

- Candès, E., Fan, Y., Janson, L., and Lv, J. (2018), "Panning for gold: model-X knockoffs for high dimensional controlled variable selection," Journal of the Royal Statistical Society: Series B (Methodological), 80, 551–577.
- Chiaretti, S., Li, X., Gentleman, R., Vitale, A., Wang, K. S., Mandelli, F., Foa, R., and Ritz, J. (2005), "Gene expression profiles of B-lineage adult acute lymphocytic leukemia reveal genetic patterns that identify lineage derivation and distinct mechanisms of transformation," Clinical cancer research, 11, 7209–7219.
- Clarke, S. and Hall, P. (2009), "Robustness of multiple testing procedures against dependence," The Annals of Statistics, 37, 332–358.
- Delaigle, A., Hall, P., and Jin, J. (2011), "Robustness and accuracy of methods for high dimensional data analysis based on Student's t-statistic," Journal of the Royal Statistical Society: Series B (Methodological), 73, 283–301.
- Efron, B. (2004), "Large-scale simultaneous hypothesis testing: the choice of a null hypothesis," Journal of the American Statistical Association, 99, 96–104.
- (2007), "Correlation and large-scale simultaneous significance testing," Journal of the American Statistical Association, 102, 93–103.
- (2010), "Correlated z-values and the accuracy of large-scale statistical estimates," Journal of the American Statistical Association, 105, 1042–1055.
- Efron, B., Tibshirani, R., Storey, J. D., and Tusher, V. (2001), "Empirical Bayes analysis of a microarray experiment," Journal of the American Statistical Association, 96, 1151–1160.
- Fan, J. and Han, X. (2017), "Estimation of the false discovery proportion with unknown dependence," Journal of the Royal Statistical Society: Series B (Methodological), 79, 1143–1164.
- Fan, J., Han, X., and Gu, W. (2012), "Estimating false discovery proportion under arbitrary covariance dependence," Journal of the American Statistical Association, 107, 1019–1035.
- Fan, J., Liao, Y., and Mincheva, M. (2013), "Large covariance estimation by thresholding principal orthogonal complements," Journal of the Royal Statistical Society: Series B (Methodological), 75, 603–680.
- Fan, J. and Peng, H. (2004), "Nonconcave penalized likelihood with a diverging number of parameters," The Annals of Statistics, 32, 928–961.
- Finner, H., Dickhaus, T., and Roters, M. (2007), "Dependency and false discovery rate: asymptotics," The Annals of Statistics, 35, 1432–1455.
- Fithian, W., Sun, D., and Taylor, J. (2014), "Optimal inference after model selection," arXiv preprint arXiv:1410.2597.
- Friedman, J., Hastie, T., and Tibshirani, R. (2008), "Sparse inverse covariance estimation with the graphical lasso," Biostatistics, 9, 432–441.
- Friguet, C., Kloareg, M., and Causeur, D. (2009), "A factor model approach to multiple testing under dependence," Journal of the American Statistical Association, 104, 1406–1415.
- Hall, P. and Jin, J. (2010), "Innovated higher criticism for detecting sparse signals in correlated noise," The Annals of Statistics, 38, 1686–1732.
- Javanmard, A. and Javadi, H. (2019), "False discovery rate control via debiased lasso," Electronic Journal of Statistics, 13, 1212–1253.
- Jin, J. (2012), "Comment," Journal of the American Statistical Association, 107, 1042-1045.
- Leek, J. T. and Storey, J. D. (2008), "A general framework for multiple testing dependence," Proceedings of the National Academy of Sciences, 105, 18718–18723.
- Lei, L. and Fithian, W. (2018), "AdaPT: an interactive procedure for multiple testing with side information," Journal of the Royal Statistical Society: Series B (Methodological), 80, 649–679.

- Lei, L., Ramdas, A., and Fithian, W. (2021), "A general interactive framework for false discovery rate control under structural constraints," Biometrika, 108, 253–267.
- Li, A. and Barber, R. F. (2019), "Multiple testing with the structure-adaptive Benjamini–Hochberg algorithm," Journal of the Royal Statistical Society: Series B (Methodological), 81, 45–74.
- Li, J. and Zhong, P.-S. (2017), "A rate optimal procedure for recovering sparse differences between high-dimensional means under dependence," The Annals of Statistics, 45, 557–590.
- Liu, H., Han, F., Yuan, M., Lafferty, J., and Wasserman, L. (2012), "High-dimensional semiparametric Gaussian copula graphical models," The Annals of Statistics, 40, 2293–2326.
- Liu, W. and Shao, Q.-M. (2014), "Phase transition and regularized bootstrap in large-scale t-tests with false discovery rate control," The Annals of Statistics, 42, 2003–2025.
- Lockhart, R., Taylor, J., Tibshirani, R. J., and Tibshirani, R. (2014), "A significance test for the lasso," The Annals of Statistics, 42, 413–468.
- Meinshausen, N., Meier, L., and Bühlmann, P. (2009), "P-values for high-dimensional regression," Journal of the American Statistical Association, 104, 1671–1681.
- Owen, A. B. (2005), "Variance of the number of false discoveries." Journal of the Royal Statistical Society: Series B (Methodological), 67, 411–426.
- Pacifico, M. P., Genovese, C., Verdinelli, I., and Wasserman, L. (2004), "False Discovery Control for Random Fields," Journal of the American Statistical Association, 99, 1002–1014.
- Petrov, V. (2002), "On probabilities of moderate deviations," Journal of Mathematical Sciences, 109, 2189–2191.
- Portnoy, S. et al. (1984), "Asymptotic behavior of M-estimators of p regression parameters when p<sup>2</sup>/n is large. I. Consistency," The Annals of Statistics, 12, 1298–1309.
- Ramdas, A. (2019), "Discussion of CARS: Covariate assisted ranking and screening for large-scale two-sample inference," Journal of the Royal Statistical Society: Series B (Methodological), 81, 228.
- Roeder, K. and Wasserman, L. (2009), "Genome-wide significance levels and weighted hypothesis testing," Statistical science: a review journal of the Institute of Mathematical Statistics, 24, 398–413.
- Sarkar, S. K. (2002), "Some results on false discovery rate in stepwise multiple testing procedures," The Annals of Statistics, 30, 239–257.
- Schwartzman, A., Dougherty, R. F., and Taylor, J. E. (2008), "False discovery rate analysis of brain diffusion direction maps," The Annals of Applied Statistics, 2, 153–175.
- Schwartzman, A. and Lin, X. (2011), "The effect of correlation in false discovery rate estimation," Biometrika, 98, 199–214.
- Storey, J. D., Taylor, J. E., and Siegmund, D. (2004), "Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach," Journal of the Royal Statistical Society: Series B (Methodological), 66, 187–205.
- Sun, W. and Cai, T. (2009), "Large-scale multiple testing under dependence," Journal of the Royal Statistical Society: Series B (Methodological), 71, 393–424.
- Sun, W. and Cai, T. T. (2007), "Oracle and adaptive compound decision rules for false discovery rate control," Journal of the American Statistical Association, 102, 901–912.
- Sun, W., Reich, B. J., Cai, T. T., Guindani, M., and Schwartzman, A. (2015), "False discovery control in large-scale spatial multiple testing," Journal of the Royal Statistical Society: Series B (Methodological), 77, 59–83.
- Sun, W. and Wei, Z. (2011), "Large-Scale multiple testing for pattern identification, with applications to time-course microarray experiments," Journal of the American Statistical Association, 106, 73–88.

- Tibshirani, R. (1996), "Regression shrinkage and selection via the lasso," Journal of the Royal Statistical Society: Series B (Methodological), 58, 267–288.
- Tusher, V. G., Tibshirani, R., and Chu, G. (2001), "Significance analysis of microarrays applied to the ionizing radiation response," Proceedings of the National Academy of Sciences of the United States of America, 98, 5116–5121.
- Van de Geer, S. A. and Bühlmann, P. (2009), "On the conditions used to prove oracle results for the Lasso," Electronic Journal of Statistics, 3, 1360–1392.
- Wasserman, L. and Roeder, K. (2009), "High dimensional variable selection," The Annals of Statistics, 37, 2178–2201.
- Wu, W. B. (2008), "On false discovery control under dependence," The Annals of Statistics, 36, 364-380.
- Xia, Y., Cai, T. T., and Sun, W. (2020), "Gap: A general framework for information pooling in two-sample sparse inference," Journal of the American Statistical Association, 115, 1236–1250.
- Yuan, M. (2010), "High dimensional inverse covariance matrix estimation via linear programming," The Journal of Machine Learning Research, 11, 2261–2286.
- Zou, C., Ren, H., Guo, X., and Li, R. (2020), "A New Procedure for Controlling False Discovery Rate in Large-Scale t-tests," arXiv preprint arXiv:2002.12548.

# Supplementary Material for "False Discovery Rate Control Under General Dependence By Symmetrized Data Aggregation"

This supplement contains some refinements and extensions of the SDA filter (Appendix A), comparisons of the SDA filter with related ideas in the literature (Appendix B), the proofs of main theorems (Appendix C), other theoretical results (Appendix D), and additional numerical results (Appendix E).

## A Refinements and Extensions

SDA provides a general framework for constructing symmetrized statistics to aggregate structural information from dependent data. In this section, we discuss some extensions to illustrate how this framework can be implemented in different scenarios.

## A.1 A stability refinement

To improve the stability in selection and avoid "p-value lottery" occurred in a single sample splitting (Meinshausen et al., 2009), we propose a modified SDA algorithm that employs the "bagging" technique to aggregate results from multiple sample—splitting procedures.

Denote  $A_k^0$ ,  $k=1,\ldots,B$ , the discovery sets from repeatedly applying B times the SDA filter at level  $\alpha$  via random sample splittings. The decisions are aggregated by  $A_k^0 = \#\{j: P_{k=1}^B | (j) \mathbb{Z}\}$   $A_k^0 > dB/2e\}$ , the set of variables that are consistently selected in at least 50% of the replications. The stability refinement picks  $A_k^0$  having the biggest overlap with  $A_k^0$ :

$$k^{2} = \underset{1 \leq k \leq B}{\operatorname{arg \, max}} \quad \underset{j=1}{\mathsf{N}} (j \ 2 \ \mathsf{A} \ \mathsf{A} \ \mathsf{A} \ \mathsf{A} ) + \mathsf{I} (j \ 2 \ \mathsf{A} \ \mathsf{B}_{k} \cap \mathsf{A} \ \mathsf{B}_{k} ) \ . \tag{S.1}$$

The new method with stability refinement is denoted R-SDA. The asymptotic theory for the R-SDA filter is presented and proven in Section  $\square$ . Our theory implies that the FDPs of  $A^{\triangleright}_{R}$  can be controlled uniformly for all k. Hence the discovery set  $A^{\triangleright}_{R}$  produces more stable results with guaranteed FDR

control. Our numerical studies show that compared to SDA, R-SDA generally yields similar FDR and power but smaller variations in the FDP.

## A.2 Other types of ranking statistics

The SDA filter utilizes  $W_j = T_{1j}T_{2j}$  to rank the hypotheses. The asymptotic symmetry property (9) is fulfilled as long as  $T_{2j}$  are constructed as the LSEs on a subset S that includes all signals with high probability. This leaves much flexibility for constructing  $T_{1j}$ . We provide a few examples.

- 1)  $T_{1j} = \mathbf{p}_{1j}$ , where  $\mathbf{p}_{1j}$  is the LASSO estimate. In contrast with the scaled version  $\mathbf{p}_{1j}/\sigma_{S,j}$ , using  $\mathbf{p}_{1j}$  directly reflects the preference of selecting large effect sizes over significant ones. In our numerical studies the two methods seem to perform similarly.
- 2) If there is prior knowledge that the covariance structure can be well described by a factor model, then we can substitute the factor-adjusted statistics (Fan and Han, 2017) in place of  $T_{1j}$ .
- 3)  $T_{1j}$  is the de-biased estimate of  $\mu_j$  (or its scaled version) based on inverse regression method (Xia et al., 2020).
- 4)  $T_{1j}$  is the innovated transformation of the sample means (Hall and Jin, 2010; Jin, 2012).

In our simulation studies, we found LASSO works well and stably in a wide range of settings but can be outperformed by other choices of  $T_{1j}$  in special situations. How to develop more powerful ranking statistics is an interesting and challenging problem that requires further research. The main message of this section is that in applications practitioners may develop new types of ranking statistics tailored to problem contexts and prior knowledge about the data structure.

Finally we stress that our theory requires that  $T_{2j}$  must be chosen so that the asymptotic symmetry property is fulfilled. For example, it is not allowed to use the LASSO estimate again to construct  $T_{2j}$  because this improper choice would lead to a violation of the symmetry property, which no longer guarantees that the FDR can be controlled at the nominal level.

## A.3 Two-sample inference

Suppose we are interested in identifying features that exhibit differential levels across two conditions. Let  $\xi^{(k)} = (\xi_1^{(k)}, \dots, \xi_p^{(k)})^>$ , k = 1, 2, be two p-dimensional random vectors. The population mean vectors and covariance matrices are  $\mu^{(k)}$  and  $\Sigma^{(k)}$ , k = 1, 2, respectively. Consider the following two-sample multiple testing problem:

$$H_i^0: \mu_i^{(1)} = \mu_i^{(2)}$$
 versus  $H_j^1: \mu_i^{(1)} = \mu_i^{(2)}$ , for  $j = 1, ..., p$ .

The SDA filter can be easily generalized to handle the two-sample situation. Denote  $D^{(k)}=\{\xi_i^{(k)}=(\xi_{i1}^{(k)},\ldots,\xi_p^{(k)})^>, i=1,\cdots,n^{(k)}\}$ . First, we split  $D^{(k)}$  into two disjoint groups  $D_1^{(k)}=(\xi_1^{(k)})$  and  $D_2^{(k)}=(\xi_2^{(k)})$ , with sizes  $n_1^{(k)}$  and  $n_2^{(k)}$ , respectively. Denote  $n_l=\eta^{(1)}+\eta_l^{(2)}$ ,  $D_l=D_l^{(1)} \square D_l^{(2)}$ , I=1,2. Based on  $D_1$ , the LASSO estimator can be obtained via minimizing  $(y_1-X\omega)^>(y_1-X\omega)+\lambda k\omega k_1$ , where  $y_1=X(\xi_1^{(1)}-\xi_1^{(2)})$ ,  $X=\Omega^{1/2}$ , and  $\Omega=(n_1/n_1^{(1)}\Sigma^{(1)}+n_1/n_1^{(2)}\Sigma^{(2)})^{-1}$ . Denote S the selected subset by LASSO. Next we calculate the LSEs, using data  $D_2$ , for coordinates in S. The formula is identical to  $\square$  except that now we take  $y_2=X(\xi_2^{(1)}-\xi_2^{(2)})$  and  $X=\Omega^{1/2}$ . Finally, we can calculate  $W_j$  and determine the threshold L using ( $\square$ ). This procedure is implemented in Section  $\square$  in the main text to identify differentially expressed genes in microarray studies. Asymptotic theories for the two–sample SDA method, which are presented in Appendix  $\square$ , can be established similarly as done for the standard SDA method.

## A.4 The SDA algorithm: detailed steps

We summarize the operation of the SDA algorithm in this subsection.

- Step 1: Split the data set into two parts  $D_1$  and  $D_2$ . If the precision matrix  $\Omega$  is unknown, use  $D_1$  to obtain its estimate  $\Omega$ .
- Step 2: Let  $X = \Omega^{1/2}$ . Compute  $\mathfrak{b}_1$  by (a) and find the narrowed subset S. Record the estimated coefficients  $\mathfrak{b}_{1i}$ .
- Step 3: Compute  $p_2$  by  $\square$  by restricting on the coordinates in the subset S.

- Step 4: Compute the ranking statistic W<sub>i</sub> by (8).
- Step 5: Find the threshold L using ( $\boxed{10}$  and output  $A^b = \{j : W_j \ge L\}$  as the selected features.

# B Comparisons with Existing Literature

This section presents comparisons of SDA with existing literature. The goal is to provide insights on the limitations of existing works and highlight some key features of SDA.

#### B.1 SDA vs. Knockoff

We present some theoretical insights on why the knockoff method suffers from power loss under dependence. The whitening transformation from Model ( $\tilde{L}$ ) to Model ( $\tilde{L}$ ) implies that the fixed-design knockoff filter in Barber and Candès (2015) is directly applicable to our problem with the Gram matrix  $X > X = \Omega$ , where  $\Omega$  is the precision matrix. The augmented design matrix can accordingly be constructed as  $(\Omega^{1/2}, 0)^>$  (c.f. Section 2.1.2 of Barber and Candès, 2015). The knockoffs  $\tilde{X}$  must fulfill  $\tilde{X} > \tilde{X} = \Omega$  and  $X > \tilde{X} = \Omega$  – diag{s}, where  $s = (s_1, \ldots, s_p)^>$  is a p-dimensional nonnegative vector. Denote  $X_j$  the jth column of the design matrix and  $\tilde{X}_j$  its knockoff copy. In a setting where the features are normalized, i.e.  $\Omega_{jj} = 1$  for all j, the correlation between  $X_j$  and  $X_j$  is  $1 - s_j$ , where  $0 \le s_j \le 1$ . Intuitively, it is desirable to make the entries of s as large as possible; this ensures that  $X_j$  would deviate from its knockoff copy as much as possible (hence we will hopefully have sufficient power to distinguish the true signals from faked ones).

Consider two settings where the correlation structures are respectively AR(1) [Corr( $X_j$ ,  $X_k$ ) =  $\rho^{|j-k|}$ , j=k] and compound symmetric [Corr( $X_j$ ,  $X_k$ ) =  $\rho$ , j=k]. We consider two approaches, namely equi-correlated and SDP knockoffs, both of which were considered in Barber and Candès (2015) for optimizing  $s_j$ 's. Figure S1 depicts the "average similarity score" 1 – s as a function of different correlation levels  $\rho$ , where  $s=\rho^{-1} \frac{P}{j=1} s_j$  is calculated using both the equi-correlated (left column) and SDP (right column) optimizers. The plots for AR(1) and compound symmetric structures are shown in the top and bottom rows, respectively. We can see that the similarity score 1 – s increases rapidly in  $\rho$ . For example, 1 – s has already exceeded 95% when  $\rho$  is only 0.25 under

the compound symmetric structure. Consequently, it becomes extremely difficult to distinguish the original variables and their faked copies. This leads to substantial power loss of the knockoff filter. The relationship between the similarity scores and the correlation levels are consistent with the patterns in the power loss of the knockoff method as noted in Fig.5 of Barber and Candès (2015) and Figure 1 in the main text of this article.

In contrast with the knockoff filter, the operation of SDA does not rely on pairwise contrasts. It only utilizes the global symmetry property among all  $W_j$ 's. The sample-splitting approach eliminates the needs for constructing fake variables under a possibly highly restricted geometric space. This explains why the SDA does not suffer from high correlations.

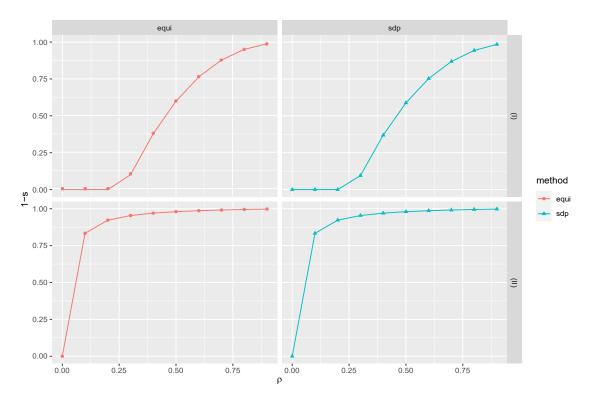


Figure S1: The knockoff filter suffers from power loss under moderate to strong dependence. The average similarity score (i.e., 1-s) between the original variable and its knockoff as a function of  $\rho$ . Top row: AR(1) structure; bottom row: compound symmetric structure. Both equi-correlated knockoff (left) and SDP knockoff (right) have been considered. The number of tests is p = 100.

### B.2 SDA vs. RESS

The reflection via sample-splitting (RESS) method in Zou et al. (2020) was developed for independent two-sample t-tests. It can be substantially improved by SDA that effectively exploits the informative dependence structure. For illustration, Figure 52 compares the FDR levels and average powers (AP) for SDA vs. BH and RESS in Zou et al. (2020) at different correlation levels. The simulation settings are the same as those in Figure 1 in the main text. We can see that the average powers of RESS and BH remain roughly the same across all correlation levels since the dependence structure has been ignored. In contrast, the power of SDA increases sharply with growing correlation levels. Section 2.3 in the main text provides high-level ideas on how the dependence is incorporated into the SDA filter to improve the power.

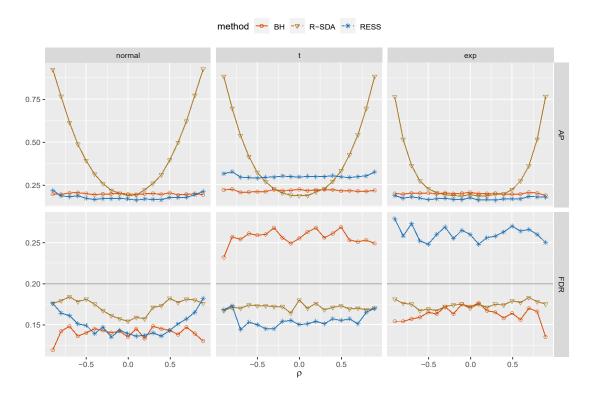


Figure S2: Impacts of correlation on different FDR procedures. Here RESS refers to the Refection via Sample Splitting procedure in Zou et al. (2020).

### B.3 Model uncertainty and error bound for FDR analysis

This section highlights the important connection of our theory to the robust knockoff theory in Barber et al. (2020), as pointed out by an insightful referee.

The model-X knockoff assumes that the distribution of the feature vector X is known exactly. However, in practical situations the X distribution must be estimated. In Theorem 1 of Barber et al. (2020), the KL divergence between the true distribution and its estimate is employed to quantify the effect of estimation errors on FDR control. The KL divergence can be interpreted as a measure of the extent to which the pairwise exchangeability property of the model-X knockoff is violated.

$$\Delta_i = |Pr(W_i > 0 | |W_i|, W_{-i}) - 1/2|,$$

which can be interpreted as a measure of the extent to which the flip-sign property is violated. We subsequently use  $\Delta_j$ 's to quantify the effect of asymmetry (i.e. deviation from the perfect symmetry assumption) on FDR control.

Barber et al. (2020) introduced an elegant leave-one-out argument to establish the upper bound for the actual FDR level of the model-X knockoff where the X matrix must be estimated from data. The analysis of SDA in Section 3.1 reveals that the technique can be readily extended to other important settings where the issue on model uncertainty must be addressed. In summary, the work of Barber et al. (2020) provides a set of useful technical tools for developing finite sample theory on (a) how the FDR control can be affected by the model uncertainty and (b) how the

In model-X knockoff the model uncertainty comes from the estimation errors whereas in SDA the model uncertainty corresponds to the possible deviation from normality and sure screening property.

error bound can be explicitly quantified using appropriate deviation measures. The connection of our theory to the robust knockoff theory also provides insights on the impact of deviation from symmetry on the performance of the SDA filter.

## C Proofs of Main Theorems

### C.1 Finite Sample Theory

This section proves Theorem 1. The proof of this theorem has extensively used the techniques developed by Barber et al. (2020), which shows that the Model-X knockoff (Candès et al., 2018) incurs an inflation of the FDR that is proportional to the errors in estimating the distribution of each feature conditional on the remaining features.

Fix > 0 and for any t > 0, define

$$R (t) = \frac{P}{P} \left[ (W_j \ge t, \Delta_j \le ) \right]$$

$$+ \frac{P}{P} \left[ (W_j \ge t, \Delta_j \le ) \right]$$

Consider the event that  $A = \{\Delta := \max_{j \supseteq A^c} \Delta_j \le \}$ . Furthermore, consider a thresholding rule L = T(W) that maps statistics W to a threshold  $L \ge 0$ . For each index j = 1, ..., p, by adopting the leave-one-out argument in Barber et al. (2020), define

$$L_{j} \; = \; T \; (W_{1}, \ldots, W_{j-1}, \, | \, W_{j} \, | \, , W_{j+1}, \ldots, W_{p}) \geq \; 0.$$

For the SDA filter with threshold L, we can write

$$\frac{P_{j \boxtimes A^{c}} | (W_{j} \ge L, \Delta_{j} \le)}{\mathbb{P}_{j} | (W_{j} \ge L)} = \frac{1 + \frac{P_{j}}{p^{j}} | (W_{j} \le -L)}{1 \mathbb{P}_{j} | (W_{j} \ge L)} \cdot \frac{P_{j \boxtimes A^{c}} | (W_{j} \ge L, \Delta_{j} \le)}{1 + \frac{P_{j}}{p^{j}} | (W_{j} \le -L)}$$

$$\le \alpha R(L).$$

Next we derive an upper bound for  $E\{R(L)\}$ . Note that

$$E\{R (L)\} = X = \begin{cases} & I(W_{j} \ge L, \Delta_{j} \le ) \\ & I + \int_{j} I(W_{j} \le -L) \\ & I + \int_{j} I(W_{j} \le -L) \\ & I + \int_{j} I(W_{j} \le L_{j}, \Delta_{j} \le ) \\ & I + \int_{j} I(W_{j} \ge L_{j}, \Delta_{j} \le ) \\ & I + \int_{j} I(W_{k} \le -L_{j}) \\ & I + \int_{k \ge A^{c}, k = j} I(W_{k} \le -L_{j})$$

The last step (S.2) holds since, after conditioning on  $(|W_j|, W_{-j})$ , the only unknown quantity is the sign of  $W_j$ . By the definition of  $\Delta_j$ , we have  $\Pr(W_j > 0 \mid |W_j|, W_{-j}) \le 1/2 + \Delta_j$ . Hence,

$$\begin{split} & E\{R(L)\} \\ & \times X \quad \begin{pmatrix} 1 & + & \Delta_{\frac{1}{2}} \end{pmatrix} I \left( \left| W_{j} \right| \geq L_{j}, \Delta_{j} \leq \right)_{c} \\ & \leq & 1^{\frac{1}{2}} + & \frac{P}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right) \\ & \geq & \left( \frac{1}{2} + \right) & \nearrow X \quad E \quad \frac{I \left( W_{j} \geq L_{j}, \Delta_{j} \leq \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad + \quad X \quad E \quad \frac{I \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} + \right) ? E\{R(L)\} + \quad X \quad E \quad \frac{P \left( W_{j} \leq -L_{j} \right)}{1 + \frac{1}{k \geq A^{c}, k = j} I \left( W_{k} \leq -L_{j} \right)} \quad ? \\ & = & \left( \frac{1}{2} +$$

The sum in the last expression can be simplified. If for all null j,  $W_j > -L_j$ , then the sum is equal to zero. Otherwise

where the first equality holds because for any j, k, if  $W_j \le -\min(L_j, L_k)$  and  $W_k \le -\min(L_j, L_k)$ , then  $L_j = L_k$ . Accordingly, we have

$$E\{R(L)\} \le \frac{1/2 + 1}{1/2 - 1} \le 1 + 5,$$

which proves the theorem.

### C.2 Asymptotic Theory with Known $\Omega$

We present the proofs of Theorem 2 here along with two key lemmas. The lemmas play key roles in our technical arguments and may be of independent interest in their own rights. Other technical lemmas and proofs are provided in Appendix 2.

For notational convenience, throughout this section, we consider variables that are included in the set S, and suppress "j  $\mathbb{Z}$  S" in all the summations with respect to j. Let  $\mathfrak{G}(x) = 1 - \Phi(x)$ ,  $G(t) = \bar{q}_{0n}^{1} \mathop{\stackrel{P}{\underset{j \otimes A^{c}}{\cap}}} \Pr(W_{j} \geq t \mid D_{1}), \ G_{-}(t) = \bar{q}_{0n}^{1} \mathop{\stackrel{P}{\underset{j \otimes A^{c}}{\cap}}} \Pr(W_{j} \leq -t \mid D_{1}) \ \text{and} \ G^{-1}(y) = \inf\{t \geq 0 : G(t) \leq y\} \ \text{for} \ 0 \leq y \leq 1.$ 

The first lemma characterizes the closeness between G(t) and  $G_{-}(t)$ .

Lemma S.1 Suppose Conditions 1 3 and 4 hold. We have

$$G(t) - 1 \rightarrow 0.$$

uniformly for all  $0 \le t \le G_{-}^{-1}(\alpha \eta_n/q_{0n})$ .

Proof. Define  $b_n = \sigma^{\sqrt[4]{C} \log \bar{q}_n}$  where C > 4. Denote  $T_{kj} = \sqrt[4]{n_k \mu_{kj}}/\sigma_j$  for  $j = 1, \ldots, q_n$  and  $\sigma^2 = Q_{jj}/\sigma_i^2$ . Observe that

$$\begin{split} \frac{G(t)}{G_{-}(t)} - 1 &= \frac{\Pr\left( \Pr(T_{1j}T_{2j} \geq t, |T_{2j}| \leq b_n \mid D_1) - \Pr(T_{1j}T_{2j} \leq -t, |T_{2j}| \leq b_n \mid D_1) \right)}{q_{0n}G_{-}(t)} \\ &+ \frac{p_{\mathbb{C}} \left\{ \Pr(T_{1j}T_{2j} \geq t, |T_{2j}| > b_n \mid D_1) - \Pr(T_{1j}T_{2j} \leq -t, |T_{2j}| > b_n \mid D_1) \right\}}{q_{0n}G_{-}(t)} \end{split}$$

 $:=\Delta_1 + \Delta_2.$ 

Firstly, for the term  $\Delta_2$ , by Lemma 5.8 we obtain that

$$\frac{\frac{P}{j \, \mathbb{P}_{A^c} \, Pr(T_{1j}T_{2j} \geq \, t, \, |T_{2j}| > \, b_n \, |D_1)}{q_{0n}G_-(t)} \leq \frac{\frac{P}{j \, \mathbb{P}_{A^c} \, Pr(|T_{2j}| > \, b_n \, |D_1)}{\alpha \eta_n} \, . \quad \frac{\bar{q_n} \times \, o(1/\bar{q_n})}{\eta_n}.$$

It follows that  $\Delta_2 = o(1)$ .

By Lemma S.7 it can be verified that

$$\frac{\Pr(T_{1j}T_{2j} \geq t, |T_{2j}| \leq b_n |D_1)}{\Pr(T_{1j}Z \geq t, |Z| \leq b_n |D_1)} \to 1,$$

where Z  $\square$  N(0,  $\sigma^2$ ) which is independent of T<sub>1i</sub>. Recall that

$$\mathbf{p}_{2j}/\sigma_{j} = n_{\bar{2}}^{-1} \sum_{i=1}^{X^{n_{2}}} e_{j}^{2} \quad X \geq X_{2S}^{-1} X_{2S} \epsilon_{i}/\sigma_{j} := n_{\bar{2}}^{-1} \sum_{i=1}^{X^{n_{2}}} (\sigma_{j} - \sigma_{i})^{2}$$

Note that  $B_n = n_2 \sigma^2$  and  $L_n = B_n^{-\theta/2} P_{n_2 \atop i=1} E(|_{ij}|^{\theta}) \le C n_2^{1-\theta/2} K_{n_2}^{\theta}$ . We have

$${2 \log(1/L_n)}^{1/2} \ge [2 \log{n_2^{\theta/2-1}/(K_{n^2}^{\theta})}]^{1/2} \ge {\rho 4 \overline{\log \bar{q}_n}},$$

according to Condition 4. The result follows by applying Lemma 5.7.

Similarly we get

$$\frac{\Pr(\mathsf{T}_{1j}\mathsf{T}_{2j} \le -\mathsf{t}, \, |\mathsf{T}_{2j}| \le \, b_n \, |\mathsf{D}_1)}{\Pr(\mathsf{T}_{1j}\mathsf{Z} \le -\mathsf{t}, \, |\mathsf{Z}| \le \, b_n \, |\mathsf{D}_1)} \to 1.$$

Note that

$$Pr(T_{1j}Z \le -t, |Z| \le b_n |D_1) = Pr(T_{1j}Z \ge t, |Z| \le b_n |D_1).$$

This implies that  $\Delta_1 = o(1)$ , which completes the proof.

The next lemma establishes the uniform convergence of  $\prod_{i \in A^c}^{P} I(W_i \ge t)/(q_{0n}G(t))$ .

Lemma S.2 Suppose Conditions 3 4 and 6 hold. Then, conditional on  $D_1$ , we have

$$\sup_{0 \le t \le G^{-1}(\alpha \eta_n / q_{0n})} \frac{\frac{\int_{j \ge A^c} I(W_j \ge t)}{q_{0n}G(t)} - 1 = o(f_j),$$

$$\sup_{0 \le t \le G^{-1}(\alpha \eta_n / q_{0n})} \frac{\int_{j \ge A^c} I(W_j \le -t)}{q_{0n}G(t)} - 1 = o(f_j).$$
(S.3)

$$\sup_{0 \le t \le G^{-1}(\alpha \eta_n/q_{0n})} \frac{\int_{j \mathbb{R}^{d}} I(W_j \le -t)}{q_{0n}G(t)} - 1 = o(1).$$
 (S.4)

Proof. We only prove the first formula; the second can be proven similarly. In the proof of Lemma S.1 we show that

$$G(t) = \ q_{0n}^{-1} \ \underset{j \, \text{\tiny{$\mathbb{Z}$A$}}^c}{X} \ \text{Pr}(T_{1j} T_{2j} \geq \ t, |T_{2j}| \leq \ b_n \ | \, D_1) \{1 + \ o(1)\} := \ \mathfrak{F}(t) \{1 + \ o(1)\}.$$

Similarly we can show that

$$q_{0n}^{-1} \underset{j \text{ } \mathbb{D} \text{ } A^{c}}{X} \text{ } I(W_{j} \geq t) = q_{0n}^{-1} \underset{j \text{ } \mathbb{D} \text{ } A^{c}}{X} \text{ } I(W_{j} \geq t, |T_{2j}| \leq b_{n})\{1 + o_{p}(1)\}.$$

Hence, it suffices to show that

$$\sup_{0 \le t \le G^{-1}(\alpha \eta_n/q_{0n})} \frac{P_{j \boxtimes A^c} I(W_j \ge t, |T_{2j}| \le b_n)}{g^{0n} \mathfrak{G}(t)} - 1 = o_p(1).$$

Note that the G(t) is a decreasing and continuous function. Let  $a_p = \alpha \eta_n$ ,  $z_0 < z_1 < \cdots < z_{h_n} \le 1$  and  $t_i = G^{-1}(\mathcal{E}_i)$ , where  $z_0 = a_p/q_{0n}$ ,  $z_i = a_p/q_{0n} + b_p \exp(i^{\zeta})/q_{0n}$ ,  $h_n = \{\log((q_{0n} - a_p)/b_p)\}^{1/\zeta}$  with  $b_p/a_p \to 0$  and  $0 < \zeta < 1$ . Note that  $G(t_i)/G(t_{i+1}) = 1 + o(1)$  uniformly in i. It is therefore enough to derive the convergence rate of

Note that

$$D(t) \leq r_p q_{0n} G(t) + \sum_{j \boxtimes A^c \ k \boxtimes M_j^c}^{X \ X} \left\{ Pr(W_k > t, W_j > t \mid D_1, B) - Pr(W_k > t \mid D_1, B) Pr(W_j > t \mid D_1, B) \right\}.$$

However, for each j  $\ 2$  A  $^c$  and k  $\ 2$  M  $^c$ , conditional on D<sub>1</sub>, the Pearson correlation coefficient between W<sub>j</sub> and W<sub>k</sub> is  $\rho_{jk}$ . By Lemma 1 in Cai and Liu (2016),

$$\frac{\Pr(W_k > t, W_j > t \mid D_1, B) - \Pr(W_k > t \mid D_1, B) \Pr(W_j > t \mid D_1, B)}{\Pr(W_k > t \mid D_1, B) \Pr(W_j > t \mid D_1, B)} \leq A_n,$$

uniformly holds, where  $A_n = (\log n)^{-1-v_1}$  for  $v_1 = \min(v, 1/2)$ 

From the above results, we can get

$$\begin{split} \Pr(D_{n} \geq \eta \mid D_{1}) \leq & \underset{0}{\overset{X^{h_{n}}}{\nearrow}} \Pr \overset{P}{\underset{j \geq A^{c}}{\nearrow}} [I(W_{j} > t_{i}, |T_{2j}| \leq b_{n}) - \Pr(W_{j} > t_{i}, |T_{2j}| \leq b_{n} \mid D_{1})] \geq & |D_{1}| \\ \leq & \underset{1}{\overset{X^{h_{n}}}{\nearrow}} \frac{1}{\overset{2}{\nearrow} e} D(t_{i}) \\ \leq & \underset{i=0}{\overset{X^{h_{n}}}{\nearrow}} \frac{1}{q_{0n} G^{2}(t_{i})} + h_{n} A_{n}) \; . \end{split}$$

Moreover, observe that

$$\frac{X^{h_n}}{\sum_{i=0}^{n} q_{0n} GP(t_i)} = \frac{1}{a_p} + \frac{X^{h_n}}{\sum_{i=1}^{n} a_p + b_p e^{i\zeta}} \cdot b_p^{-1}.$$

Finally, note that (a)  $\zeta$  can be arbitrarily close to 1 such that  $h_n A_n \to 0$ , and (b)  $b_p$  can be made arbitrarily large as long as  $b_p/a_p \to 0$ , we conclude that  $D_n = o_p(1)$  when  $r_p/\eta_n \to 0$ . This completes the proof.

In Lemma S.1 and Lemma S.2 we have established the symmetry property and uniform consistency for  $W_i$ 's. Now we are ready to present the proof of Theorem 2.

Proof of Theorem 2 By definition, SDA selects the jth variable if  $W_i \ge L$ , where

We need to establish an asymptotic bound for L so that Lemmas 5.15.2 can be applied.

Let  $t^2 = G_-^{-1}(\alpha \eta_n/q_{0n})$ . It follows from Lemma 5.2 that

$$\alpha \eta_n / q_{0n} = G_-(t^2) = \frac{1}{q_{0n}} \frac{X}{|f| \Delta^c} I(W_j < -t^2) \{1 + o(1)\}.$$

$$\begin{split} &\text{Pr}\left(W_{j} < \, t^{\mathbb{P}}, \, \text{for some } j \, \, \mathbb{P} \, C_{\mu}\right) \\ &\leq \, \eta_{n} \, \text{Pr} \, \, T_{1j} T_{2j} - \sqrt[V]{n_{1} n_{2} \mu_{j}^{2} / \sigma_{j}^{2}} < \, t^{\mathbb{P}} - \sqrt[V]{n_{1} n_{2} \mu_{j}^{2} / \sigma_{j}^{2}} \\ &\leq \, \eta_{n} \, \text{Pr} \, \, |\mu_{j}| \, (|\mu_{1j} - \mu_{j}| + |\mu_{2j} - \mu_{j}|) + \, |\mu_{1j} - \mu_{j}| \, |\mu_{2j} - \mu_{j}| > \, \mu_{j}^{2} - \, t^{\mathbb{P}} \sigma_{j}^{2} / \sqrt[V]{n_{1} n_{2}} \rightarrow 0. \end{split}$$

To see the last equation, denote  $d_j = \mu_j^2 - t^{\mathbb{Z}} \sigma_j^2 / \sqrt[4]{n_1 n_2}$ . Under Condition  $D_j$  it follows that  $d_j = \mu_j^2 \{1 + o(1)\}$ . We then get

$$\begin{split} & \Pr \left( |\mu_j| \left( |\pmb{p}_{1j} - \mu_j| + |\pmb{p}_{2j} - \mu_j| \right) + |\pmb{p}_{1j} - \mu_j| |\pmb{p}_{2j} - \mu_j| > d_j \right) \\ & \leq \Pr \left( |\mu_j| \left( |\pmb{p}_{1j} - \mu_j| + |\pmb{p}_{2j} - \mu_j| \right) > d_j/2 \right) + \Pr \left( |\pmb{p}_{1j} - \mu_j| |\pmb{p}_{2j} - \mu_j| > d_j/2 \right) =: H_1 + H_2. \end{split}$$

Note that  $d_j/|\mu_j|$  =  $|\mu_j|\{1 + o(1)\}$ . We observe that

$$\begin{split} &H_1 \leq \, \text{Pr} \, ( \, | \, \pmb{b}_{1j} - \, \mu_j \, | \, > \, d_j / (4 \, | \, \mu_j \, | \, ) ) + \, \text{Pr} \, ( \, | \, \pmb{b}_{2j} - \, \mu_j \, | \, > \, d_j / (4 \, | \, \mu_j \, | \, ) ) \, , \\ &H_2 \leq \, \text{Pr} \, ( \, | \, \pmb{b}_{1j} - \, \mu_j \, | \, > \, c_{np} ) + \, \text{Pr} \, | \, | \pmb{b}_{2j} - \, \mu_j \, | \, > \, C \, \frac{p}{\log \bar{q}_n / n} \, \, . \end{split}$$

Then the result follows from Lemmas 5.8 and Condition 2.

Consequently, we have  $Pr({P\atop j}|I(W_j>t^{\mathbb{Z}})\geq\eta_n)\rightarrow 1$ . We conclude that  ${P\atop j}|I(W_j<-t^{\mathbb{Z}})$ .  $\alpha\eta_n\leq\alpha^{P\atop j}|I(W_j>t^{\mathbb{Z}})$ , and hence L .  $t^{\mathbb{Z}}$ . By Lemmas S.1-S.2, we get

$$\frac{P}{P^{j \geq A^{c}} I(W_{j} \geq L)} - 1 \Rightarrow 0.$$

$$(S.5)$$

Next write

$$FDP = \frac{P}{1 \cdot P \cdot A^{c} \cdot I(W_{j} \geq L)} = \frac{P}{1 \cdot P \cdot I(W_{j} \leq -L)} \times \frac{P}{P \cdot P \cdot A^{c} \cdot I(W_{j} \geq L)} \times \frac{P}{P \cdot P$$

Note that  $R(L) \le \frac{P}{\int \mathbb{D} A^c} I(W_j \ge L) / \frac{P}{\int \mathbb{D} A^c} I(W_j \le -L)$ , and thus  $\limsup_{n \to \infty} FDP \le \alpha$  by (5.5) Then, for any > 0,

$$FDR \leq (1 + )\alpha R(L) + Pr(FDP \geq (1 + )\alpha R(L))$$
,

which proves the second part of this theorem.

# C.3 Asymptotic Theory with unknown $\Omega$ : Proof of Theorems 3 and 4

Proof of Theorem 3 The proof follows similar lines as those of Theorem 2 except that we now establish Lemmas 5.1 and 5.2 under Conditions 1-5 and 6'. Note that Lemma 5.8 still holds under Conditions 1 3 and 4 With unknown  $\Omega$ , conditional on  $D_1$ , the Pearson correlation coefficient between  $W_j$  and  $W_k$  is changed to  $\rho_{jk}^c$ . The rest of the proof is essentially the same as that of Theorem 2 and thus omitted.

Proof of Theorem 4 To establish this theorem, we consider another SDA procedure with the statistics  ${}^f\!W_j = \sqrt[4]{n_1n_2}|_{{\bf b}_{1j}|_{{\bf E}_{2j}}/\sigma_j^2}$ , where  ${\bf e}_2$  is the least–squares estimate that uses  ${\bf x}^2={\bf x}^2$  and  ${\bf e}_2={\bf x}^2$ . We choose a threshold  ${\bf E}>0$  by setting

$$e = \inf t > 0: \frac{\#\{j : w_j \le -t\}}{\#\{j : w_j \ge t\} \ \ 1} \le \alpha .$$

The proof of this theorem involves a careful investigation of the difference between  $W_j$  and  $\hat{W}_j$ . The main results are summarized by Lemmas 5.3-5.5. Define  $G = \{j : \mu_j = o(c_{np})\}$ .

From Lemma S.3 we have, for any j,

$$W_{j} - v_{j} = \sqrt[4]{n_{1}n_{2}} \mathbf{b}_{1j} (\mathbf{b}_{2j} - \mathbf{p}_{2j}) / \sigma_{j}^{2} = O_{p}(n \times s_{n}\bar{q}_{n}a_{np}) \overline{\log p/n} \times \{\mu_{j} + O_{p}(c_{np})\}.$$

$$\hat{W}_{j} = W_{j} \quad 1 + \frac{\mu_{2j} - \mu_{2j}}{\mu_{2j}} = W_{j} \quad 1 + \frac{O_{p}(s_{n}\bar{q}_{n}a_{np} - \frac{p}{\log p/n})}{\mu_{j} + O_{p}(\frac{p}{\log \bar{q}_{n}/n})} = W_{j}\{1 + o_{p}(1)\}.$$

In fact, under conditions  $c_{np}a_{np}s_n\bar{q}_n \sqrt{n\log p}(\log \bar{q}_n)^{1+\gamma} \rightarrow 0$  and  $1/(\sqrt[4]{nc_{np}}) = O(1)$ , we have:

$$\frac{s_n q_n a_{np} \frac{p}{\log p/n}}{c_{np}} = o(1), \quad \frac{s_n \bar{q}_n a_{np} \frac{p}{\log \bar{q}_n/n}}{p \frac{1}{\log \bar{q}_n/n}} = o(1).$$

From Lemma S.4 and Lemma S.5 given below, we conclude that

$$FDP_{\hat{W}}(E) := \frac{\#\{j : W_j \ge E, j ? A^c\}}{\#\{j : W_j \ge E\} ? 1} = FDP_W(L)\{1 + o_p(1)\}.$$

Under Conditions 1-6, similar to the proof of Theorem 2, we can show that  ${\rm FDP}_{\hat{\mathbb{W}}}$  (£) is controlled at the nominal level asymptotically. Thus the claimed result follows.

Lemma S.3 If Conditions 1 3 4 and 7 hold, then we have  $\mathbf{p}_j = \mathbf{p}_j + O_p(a_{np}s_n\bar{q}_n^p \log p/n)$  uniformly in j 2 S.

Proof. Note that

$$\begin{split} | \, b_j - \, \mu_j | &= \, e^{>}_j \, \left( \, X e^{>}_S X e^{>}_S \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X e^{>}_S X e^{-} \, \left( \, X^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X_{\,S} \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X_{\,S} \, \left( \, \xi^{-} - \, \mu \right) \\ &\leq \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X_{\,S} \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X_{\,S} \, \left( \, X e^{>}_S X_{\,S} \, \right)^{-1} \, X^{>}_S X_{\,S} \, \left( \, X e^{>}_S X_{\,S} \, \right)$$

Similar to Lemma S.8, we get  $k\xi - \mu k_{\infty} = O_p(\log p/n)$ . For the analysis of  $\Delta$ , we note the following fact

$$\begin{split} & k \, \Re_S^{>} \, \mathscr{R}_S \, - \, X_S^{>} \, X_S \, k_\infty \, \leq \, k \, \mathop{\raisebox{.4ex}{$\underline{\circ}$}} \, F \, - \, \Omega \, k_\infty \, = \, O_p(a_{np}), \\ & k \, \Re_S^{>} \, \mathscr{R}_S \, c \, - \, X_S^{>} \, X_S \, c \, k_\infty \, \leq \, k \, \mathop{\raisebox{.4ex}{$\underline{\circ}$}} \, P_F \, - \, \Omega \, k_\infty \, = \, O_p(a_{np}), \\ & k \, (\Re_S^{>} \, \mathscr{R}_S)^{-1} \, - \, (X_S^{>} \, X_S)^{-1} \, k_\infty \, \leq \, k \, (\Re_S^{>} \, \mathscr{R}_S)^{-1} \, k_\infty \, k \, (X_S^{>} \, X_S)^{-1} \, k_\infty \, k \, (X_S^{>} \, X_S^{>} \, X_S)^{-1} \, k_\infty \, k \, (X_S^{>} \, X_S^{>} \, X_S^{>} \, X_S^{>} \, X_S^{>} \, X_S^{>} \, ($$

Thus, by triangle inequality, we can conclude that

$$\begin{split} k\Delta k_{\infty} &= k (\Re_{S}^{>} \Re_{S})^{-1} \Re_{S}^{>} \Re_{S^{c}} - (X_{S}^{>} X_{S})^{-1} X_{S}^{>} X_{S^{c}} k_{\infty} \\ &\leq k (\Re_{S}^{>} \Re_{S})^{-1} - (X_{S}^{>} X_{S})^{-1} k_{\infty} k X_{S}^{>} X_{S^{c}} k_{\infty} + k \Re_{S}^{>} \Re_{S^{c}} - X_{S}^{>} X_{S^{c}} k_{\infty} k (\Re_{S}^{>} X_{S^{c}}^{>})^{-1} k_{\infty} \\ &= O_{p} (\bar{q}_{n} s_{n} a_{np}) \end{split}$$

and accordingly  $\max_{j} |\mathbf{p}_{j} - \mathbf{p}_{j}| = O_{p}(\bar{q}_{n}s_{n}a_{np}^{p}|\overline{\log p/n}).$ 

The next lemma establishes the approximation result of  $W_i$  to  $\hat{W}_j$  for those  $j \ 2 \ G$ .

Lemma S.4 Suppose Conditions 1 2 3 4 and 7 hold and

 $c_{np}a_{np}s_n\bar{q}_n \sqrt[V]{n\log p}(\log \bar{q}_n)^{1+\gamma} \rightarrow 0$  for a small  $\gamma > 0$ . Then, for any M > 0,

$$\sup_{\substack{M \leq t \leq G^{-1}(\alpha\eta_n/q_{0n}) \\ M \leq t \leq G^{-1}(\alpha\eta_n/q_{0n})}} \frac{P^{\frac{1}{p} \boxtimes G} I(W^f_j \geq t)}{\text{j$\mathbb{P}_G$ } I(W_j \geq t)} - 1 = o_p(1),$$
 
$$\sup_{\substack{M \leq t \leq G^{-1}(\alpha\eta_n/q_{0n}) \\ M \leq t \leq G^{-1}(\alpha\eta_n/q_{0n})}} P^{\frac{1}{p} \boxtimes G} I(W^f_j \leq -t)}_{\frac{1}{p} \boxtimes G} \frac{1}{p} = o_p(1),$$

Proof. By Lemma S.3 with probability tending to one,

$$\begin{array}{l} X \\ I \left( W_{j} \geq t \right) - X \\ j \mathbb{P}G \end{array}$$

$$\begin{array}{l} X \\ j \mathbb{P}G \end{array}$$

$$\begin{array}{l} X \\ j \mathbb{P}G \end{array}$$

$$\begin{array}{l} X \\ \{I \left( W_{j} \geq t + I_{n} \right) - I \left( W_{j} \geq t \right) \} + X \\ j \mathbb{P}G \end{array}$$

$$\begin{array}{l} \{I \left( W_{j} \geq t - I_{n} \right) - I \left( W_{j} \geq t \right) \} + X \\ j \mathbb{P}G \end{array}$$

$$\begin{array}{l} \{I \left( W_{j} \geq t - I_{n} \right) - I \left( W_{j} \geq t \right) \} + X \\ j \mathbb{P}G \end{array}$$

$$\begin{array}{l} \{I \left( W_{j} \geq t - I_{n} \right) - I \left( W_{j} \geq t \right) \} + X \\ j \mathbb{P}G \end{array}$$

where  $I_n/(c_{np}a_{np}s_n\bar{q}_n^{\ \ \ } \frac{1}{n\ log\ p}) \rightarrow \infty$  as  $n,p \rightarrow \infty.$  We will deal with  $\Delta_1$  only and the part of  $\Delta_2$  is

similar. Define the events  $C_t = \{|T_{1j}| > t/(C^{\sqrt{\log q_n}}), |T_{2j}| > t/(\sqrt[q]{nc_{np}}), j \ \mathbb{Z} \ G\}.$ 

$$\begin{split} E(\Delta_{1}) &= E \begin{bmatrix} ? & & ? \\ ? & X \\ ? & \end{bmatrix} \\ &= X \\ &\leq Pr(t \leq W_{j} \leq t + I_{n} | C_{t}) + Pr(t \leq W_{j} \leq t + I_{n}, C_{t}^{c}) \\ &= X \\ &\leq Pr(t \leq W_{j} \leq t + I_{n} | C_{t}) + Pr(t \leq W_{j} \leq t + I_{n}, C_{t}^{c}) \\ &= X \\ &\leq Pr(t \leq W_{j} \leq t + I_{n} | C_{t}) + o(1), \\ &= Pr(t \leq W_{j} \leq t + I_{n} | C_{t}) + o(1), \end{split}$$

where we use Lemmas S.8 and Condition 2 to get  $_{j \boxtimes G}^{P}$  Pr(t  $\leq W_{j} \leq t + I_{n}, C_{t}$ ) = o(1). Further note that under the event, {t  $\leq W_{j} \leq t + I_{n}, C_{t}$ }, we have

$$\left|T_{2j}\right| \leq \frac{t|+|I_n|}{T_{1j}} \leq \frac{C(t+|I_n|)^{\sqrt{\log \bar{q}_n}}}{t} = C^p \overline{\log \bar{q}_n} + \frac{I_n^{\sqrt{\log \bar{q}_n}}}{M} \leq C^p \overline{\log \bar{q}_n} = b_n,$$

under condition that  $I_n \to 0$ . Let  $T_{2j}^{\mathbb{P}} = \sqrt[4]{n_2(\mathbf{p}_{2j} - \mu_j)/\sigma_j}$  and  $U_j = \sqrt{n_2\mu_j/\sigma_j}$ . Thus from Lemma S.1 we conclude that

$$\begin{array}{l} X \\ P \, r \, (\,t \, - \, T_{1j} \, U_j \, \leq \, T_{1j} \, Z \, \leq \, t \, + \, I_n \, - \, T_{1j} \, U_j \, \, \big| \, C_t \big) \\ j \, \mathbb{P}G \\ X \\ = & E \, \big\{ \, \Phi \, (\,(\,t \, + \, I_n) / \big| \, T_{1j} \big| \, - \, U_j \, \big) \, - \, \Phi(t / \big| \, T_{1j} \big| \, - \, U_j \, \big) \, \, \big| \, C_t \big\} \\ j \, \mathbb{P}G \\ X \\ \leq & I_n \, E \, \big| \, T_{1j} \, \big|^{-1} \, \Phi(t / \big| \, T_{1j} \big| \, - \, U_j \, \big) \, \, \big| \, C_t \\ \leq & I_n \, X \quad E \quad (t / T_{1j}^2 \, - \, U_j / \big| \, T_{1j} \big| \, \big) \, + \, \frac{1}{t \, - \, U_j \, \big| \, T_{1j} \big|} \, \Phi(t / \big| \, T_{1j} \big| \, - \, U_j \, \big) \, \, \big| \, C_t \\ \vdots \\ I_n \, M \, & \qquad X \quad n \quad \Phi(t / \big| \, T_{1j} \big| \, - \, U_j \, \big) \, \, \big| \, C_t \quad , \end{array}$$

where  $\Phi(x) = 1 - \Phi(x)$ . The second to last inequality is due to

$$\frac{x}{x^2+1}\varphi(x)<\, \Phi(x), \ \ \text{for all} \ x>\, 0.$$

On the other hand,

$$X = Pr(W_j > t) = E \Phi(P/|T_{1j}| - U_j) |C_t| \{1 + o(1)\}.$$

Therefore, by Markov inequality and similar arguments in the proof of Lemma S.2, the assertion holds if  $c_{np}a_{np}s_n\bar{q_n}$   $\sqrt[4]{n\log p}\log \bar{q_n}$   $h_n\to 0$ . Note that  $h_n$  can be made arbitrarily small as long as  $h_n\to\infty$  as  $n\to\infty$ , from which we completes the proof.

In the next lemma, we obtain the approximation result for those j with relatively large  $\mu_i$ .

Lemma S.5 Suppose Conditions 1 2 3 4 and 7 hold and

 $c_{np}a_{np}s_n\bar{q}_n \sqrt[V]{n\log p}(\log \bar{q}_n)^{1+\gamma} \rightarrow 0$ . Then, for any M > 0,

$$\sup_{\substack{M \leq t \leq G^{-1}(\alpha\eta_n/q_{0n})\\ M \leq t \leq G^{-1}(\alpha\eta_n/q_{0n})}} \frac{P^{\frac{p}{j \boxtimes G^c} I(W_j \geq t)}}{P^{\frac{p}{j \boxtimes G^c} I(W_j \geq t)}} \quad 1 = o_p(1),$$

$$\sup_{\substack{M \leq t \leq G^{-1}(\alpha\eta_n/q_{0n})\\ M \leq t \leq G^{-1}(\alpha\eta_n/q_{0n})}} P^{\frac{p}{j \boxtimes G^c} I(W_j \leq -t)}_{j \boxtimes G^c I(W_j \leq -t)} \quad 1 = o_p(1).$$

Proof. Under the designed conditions, we have  $W_j = \oint W_j \{1 + o_p(1)\}$  for any  $j \supseteq G^c$  uniformly. Then the results follow.

# D Proofs of Additional Theoretical Results

# D.1 Proof of Lemma 1 (the coin-flip property under dependence)

Observe that  $W_j = \sqrt[4]{n_1 n_2} |\mathbf{b}_{1j}| \mathbf{b}_{2j} / \sigma_j^2 =: c_j \times |\mathbf{b}_{2j}|$ . Conditional on  $D_1$ , we have  $W_j \mid W_{-j} \mid \mathbb{R}$   $N(\mu_j \mid_{-j}, \sigma_{j}^2 \mid_{-j})$  with

$$\begin{split} \mu_j|_{-j} &= \text{Cov}(W_j, W_{-j}) \text{Var}(W_{-j})^{-1}(W_{-j} - EW_{-j}) \text{ and} \\ \\ \sigma_{j|-j}^2 &= \text{Var}(W_j) - \text{Cov}(W_j, W_{-j}) \{ \text{Var}(W_{-j}) \}^{-1} \text{Cov}(W_j, W_{-j})^{>}. \end{split}$$

For any k,l  $\mathbb{C}$  S, we have  $Cov(W_k, W_l) = c_k c_l Q_{kl}$ . Let  $C = diag\{c_1, \ldots, c_{q_n}\}$  and D = CQC.

Then  $Cov(W_j, W_{-j}) = D_{j,-j}$  and  $Var(W_{-j}) = D_{-j,-j}$ . So, we obtain that

$$\begin{split} & \text{Pr}(W_{j} > 0 \mid \mid W_{j} \mid, W_{-j}, D_{1}) \\ & = \frac{\varphi \cdot \frac{\mid W_{j} \mid -\mu_{j\mid -j}}{\sigma_{j\mid -j}}}{\varphi \cdot \frac{\mid W_{j} \mid -\mu_{j\mid -j}}{\sigma_{j\mid -j}}} + \varphi \cdot \frac{\mid W_{j} \mid +\mu_{j\mid -j}}{\sigma_{j\mid -j}} \\ & = \frac{\varphi \cdot \frac{\mid W_{j} \mid -D_{j,-j} D_{1,-j} D_{1,-j} (W_{-j} - EW_{-j})}{\varphi \cdot \frac{\mid W_{j} \mid -D_{j,-j} D_{-j,-j} D_{-j,-j} D_{-j,-j}}{\varphi \cdot \frac{\mid W_{j} \mid -D_{j,-j} D_{1,-j} D_{-j,-j} D_{$$

Denote  $Q_{-i,j} = 0$  the jth column of Q excluding  $Q_{ij}$ . Finally we have

$$Pr(W_{j} > 0 | |W_{j}|, W_{-j}) = E \{Pr(W_{j} > 0 | |W_{j}|, W_{-j}, D_{1}) | |W_{j}|, W_{-j}\}$$

$$= E \{\Delta_{i}(|W_{i}|, W_{-i}, D_{1}) | |W_{i}|, W_{-i}\} - 1/2.$$

It can be easily verified that if  $Q_{-j,j} = 0$ ,  $\Delta_j(|W_j|, W_{-j}, D_1) = 1/2$  and consequently  $\Delta_j = 0$ .

## D.2 Asymptotic results for R-SDA and two-sample SDA

The next result is a direct corollary of Theorem 2 which establishes the FDR control of the multi-splitting procedure R-SDA.

Corollary 1 Suppose Conditions  $\fbox{1}6$  hold. For any  $\alpha$   $\fbox{2}$  (0,1) and a given B, the FDR of the R-SDA method satisfies  $\limsup_{(n,p)\to\infty} \mathsf{FDR} \le \alpha$ .

As in (14), the FDP is controlled for each replication so is the FDP of R-SDA, resulting in the FDR control.

To establish the FDR control result of SDA procedure for the two-sample problem, we introduce a new sequence of independent random variables  $\{\xi_i\}$  defined as follows:

$$\xi_{i} - \omega = \begin{cases} n_{2}/n_{2}^{(1)}(\xi_{2i}^{(1)} - \mu^{(1)}); & 1 \leq i \leq n_{2}^{(1)}; \\ \frac{1}{2} - n_{2}/n_{2}^{(2)}(\xi_{2i+1}^{(2)} - \mu^{(2)}); & n_{2}^{(1)} + 1 \leq i \leq n_{2}. \end{cases}$$

Note that

$$\xi_2^{(1)} - \xi_2^{(2)} - \omega = \frac{1}{n_2^{(1)}} \hat{X}_{i=1}^{(1)} (\xi_{2i}^{(1)} - \mu^{(1)}) - \frac{1}{n_2^{(2)}} \hat{X}_{i=1}^{(2)} (\xi_{2i}^{(2)} - \mu^{(2)}) = \frac{1}{n_2} \hat{X}_{i=1}^{n_2} (\xi_i - \omega).$$

By the proofs for Theorem 2 if we replace  $\mu$  as  $\omega$  and set  $\Omega^{-1} = \Sigma^{(1)}/\% + \Sigma^{(2)}/(1-\%)$  with %=  $\lim n_i^{(1)}/n_i$ , Theorem 2 holds also for the two-sample problem.

Corollary 2 Suppose Conditions  $\boxed{16}$  hold. For any  $\alpha$   $\boxed{2}$  (0, 1) and 0 < % < 1, the FDR of the SDA for the two-sample problem satisfies  $\limsup_{(n,p)\to\infty} \mathsf{FDR} \le \alpha$ .

We want to emphasize that as long as Condition  $\square$  is satisfied, the above results hold for other choices of  $T_{1j}$  as discussed in Appendix A.2. For example, consider a hard-thresholding estimator  $\mathbf{p}_{1j} = \xi_{1j} \mathbb{I}(|\xi_{1j}| > c^p \overline{\log p/n})$  for some c > 0. We know that  $c_{np} = p \overline{\log p/n}$  if  $\xi_{ij}$ 's have uniformly bounded fourth moments.

#### D.3 Additional lemmas

The first one is the standard Bernstein's inequality.

Lemma S.6 (Bernstein's inequality) Let  $X_1, \ldots, X_n$  be independent centered random variables a.s. bounded by A <  $\infty$  in absolute value. Let  $\sigma^2 = n^{-1} \sum_{i=1}^{p} E(X_i^2)$ . Then for all x > 0,

$$Pr \int_{i=1}^{\eta} X_i \ge x \le exp - \frac{x^2}{2n\sigma^2 + 2Ax/3}$$
.

The second one is a moderate deviation result for the mean; See Petrov (2002).

Lemma S.7 (Moderate deviation for the independent sum) Suppose that  $X_1,\ldots,X_n$  are independent random variables with mean zero, satisfying  $E(|X_j|^{2+\delta})<\infty$  ( $j=1,2,\ldots$ ). Let  $B_n=P_{i=1}^n$   $E(X_i^2)$ . Then,

$$\frac{\Pr(P_{i=1}^{n}X_{i}>x^{\sqrt{B_{n}}})}{1-\Phi(x)}\to 1,$$

as  $n \to \infty$  uniformly in x in the domain  $0 \le x \le C\{2 \log(1/L_n)\}^{1/2}$ , where  $L_n = B_n^{-1-\delta/2} P_{i=1}^n E[X_i]^{2+\delta}$  and C is a positive constant satisfying the condition C < 1.

The next lemma establishes uniform bounds for  $\mathbf{b}_{2i}$ .

Lemma S.8 Suppose Conditions  $\square$  and  $\square$  hold. Then, as  $n \rightarrow \infty$ ,  $\Pr \sigma^{-1}$ 

$$|\mathbf{p}_{2j} - \mathbf{p}_{j}| > \sigma^{p} C \log \bar{q} / n_{2} |D_{j_{h}} = o(1/\bar{q}),$$

holds uniformly in S, where C > 4.

Proof. Write

$$\mathbf{p}_{2j} - \mu_j = n_{\bar{2}}^{1} \sum_{i=1}^{X^{n_2}} e_{\bar{j}} \quad X_{\bar{j}S} X_{2S} \quad {}^{-1} X_{\bar{j}S} \epsilon_i := n_{\bar{2}}^{1} \sum_{i=1}^{X^{n_2}} e_{\bar{j}}.$$

Let  $m_n$  =  $(n_2\bar{q}_n)^{1/\theta+\gamma}\,K_{n\,2}$  and note that

$$\begin{aligned} &_{ij} = _{ij}I(|_{ij}| \leq m_n) - E\{_{j}I(|_{j}| \leq m_n)\} + _{ij}I(|_{ij}| > m_n) - E\{_{j}I(|_{j}| > m_n)\} \\ &=: _{ij,1} + _{ij,2}. \end{aligned}$$

Conditioned on the first split D<sub>1</sub>,

$$\begin{array}{l} \text{Pr} ( \mid \sqrt[V]{\overline{n_2}} ( | \textbf{h}_{2j} - \mu_j ) \mid > \sigma_j x \ \text{ for some } j \mid D_1 ) \\ = \text{Pr} \quad & X^{D_2} \\ & \text{ij,1} + & X^{D_2} \\ & \text{ij,2} > & \sqrt[V]{\overline{n_2}} \sigma_j x \ \text{ for some } j \mid D_1 \\ & \text{if } & \text{ij,2} > & \sqrt[V]{\overline{n_2}} \sigma_j x \ \text{ for some } j \mid D_1 \\ & \text{if } & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if } \\ & \text{if } & \text{if } & \text{if } \\ & \text{if } & \text{if }$$

Here a is a small positive value.

Firstly consider the term  $P_1$ . Note that  $_{1j,1},\ldots,_{n_2j,1}$  are independent centered random variables a.s. bounded by  $2m_n$  in absolute value. Then the Bernstein inequality in Lemma 5.6 yields that

$$P_{1} \leq 2q_{n} \max_{j} \exp \left( -\frac{n_{2}\sigma_{j}^{2}(1-\frac{b}{2})^{2}}{2n_{2}E(\frac{1}{2}) + 22m_{n} \cdot n\overline{q}_{j}x(1-a)/3} \right).$$

Recall that  $_{ij,1} = _{ij}I(|_{ij}| \le m_n) - E[_{j}I(|_{j}| \le m_n)]$ . Thus

$$E(_{j},_{1}^{2}) = Var\{_{j}I(|_{j}| \le m_{n})\} \le E\{_{j}I(|_{j}| \le m_{n})\} \le E(_{j}) = Q_{jj}.$$

We then have:

$$P_{1} \leq 2q_{n} \max_{j} \exp \left(-\frac{n_{2}\sigma_{j}^{2}x^{2}(1-a)^{2}}{2n_{2}Q_{jj}+2\cdot 2m_{n}\cdot \sqrt[4]{n_{2}}\sigma_{j}x(1-a)/3}\right)$$

$$\leq 2\bar{q}_{n} \max_{j} \exp \left(-\frac{x^{2}(1-a)^{2}}{2\sigma^{2}+4(1-a)\sigma_{j}^{-1}xm_{n}/(3\sqrt[4]{n_{2}})}\right). \tag{S.7}$$

Next we turn to consider P2. First note that

$$P_{2} \leq Pr \max_{i=1}^{|X|^{2}} ||I(|_{ij}| > m_{n}) + \max_{j} n_{2}E\{|_{j}|I(|_{j}| > m_{n})\} > \sqrt[V]{n_{2}\sigma_{j}\overline{x}a} |D_{1}|$$

Further note that

$$E^{2}\{|_{j}|I(|_{j}| > m_{n})\} \le E(^{2})P_{j}r(|_{j}| > m_{n}) \le E(^{2})\frac{E(|_{j}|^{\theta})}{\int_{0}^{\infty} m^{\theta}}$$

We then conclude that

$$\max_{j} n_{2} E\{|_{j} | I(|_{j}| > m_{n})\} \leq \max_{j} n_{2} \frac{q}{E(^{2}) E(|_{j}^{\theta})} = o(^{\sqrt[4]{n_{2}}}).$$

From this, we then have

$$P_{2} \leq \Pr \max_{i=1}^{|X|^{2}} |_{ij} |I(|_{ij}| > m_{n}) > {}^{V}n_{2}\overline{\sigma_{j}}xa/2 |D_{1}|$$

$$\leq \Pr \max_{j} |_{ij}| > m_{n} \text{ for some } i |D_{1}|$$

$$\leq n_{2} \frac{E(kA(S)\varepsilon_{i}k_{\infty}^{\theta})}{m_{n}^{\theta}} = o(\bar{q}_{n}^{-1}). \tag{S.8}$$

Let  $x = \sigma^{\sqrt[4]{C} \log \overline{q_n}}$ . From the inequalities (5.6), (5.7), and (5.8), we conclude that

$$\begin{split} & \text{Pr}\,(\,|\,^{\sqrt[N]{n_2}}\,(\not\!\! h_{2j}-\mu_j)\,|\,>\,\sigma_j\,x \ \text{ for some } j \ |\, D_1) \\ & \leq \, 2\bar{q}_n \, \max_j \exp \ -\frac{x^2(1-a)^2}{2\,\sigma^2 +\, 4(1-a)\sigma_j^{-1}xm_n/(3\,\sqrt[N]{n_2})} \ +\, o(q_n^{-1}) = \, o(\bar{q}_n^{-1}). \end{split}$$

holds uniformly in S, where we use the condition  $m_n/p \sqrt{\log \bar{q}_n} = o(1)$  which is implied by Condition 3.

### E Additional Numerical Results

#### E.1 Estimated co<sup>v</sup>ariance structures

This section compares the methods mentioned in Section of the unknown covariance case. In practice, one should adopt the most appropriate estimator tailored to specific correlation structures. Specifically, we have used the method based on Cholesky decomposition in Bickel and Levina (2008), the POET method proposed by Fan et al. (2013), and the graphical lasso (Friedman et al., 2008) to estimate the unknown Structures (I)—(III), respectively.

Figure \$\sigma\$ follows the settings in Figure 4 (except that the covariance matrix or its inverse is estimated). Figures \$\sigma\$ uses the same settings as those in Figures \$\sigma\$ with estimated covariance matrix. We omit a detailed discussion as the observed patterns seem to be very similar to those in the known covariance case (except that the FDR control sometimes becomes less accurate due to the additional estimation errors). Our conclusions based on Figure \$\sigma\$ and Figures \$\sigma\$ remain essentially the same as before. Knockoff and R-SDA seem to be the only methods that can control the FDR reasonably well in all scenarios, with the R-SDA method having much higher power in most scenarios.

## E.2 Additional comparisons

Figure  $\boxed{\textbf{S5}}$  demonstrates the FDR and AP for various signal magnitude  $\mu$  under the compound symmetry error structure (II) and three error distributions, for known and unknown covariance structures, respectively.

## E.3 Boxplots of FDPs

When the noises are sampled from the multivariate normal distribution, Figure  $\frac{56}{56}$  shows the boxplot of the FDP and AP of the testing procedures for  $\pi_1 = 0.05$  and 0.2, while fixing  $(n, p, \alpha) = (90, 500, 0.2)$ . The signal magnitude  $\mu$  is adjusted according to the covariance structures so that the APs are in a similar range. While the BH is conservative with little power, the R-SDA outperforms

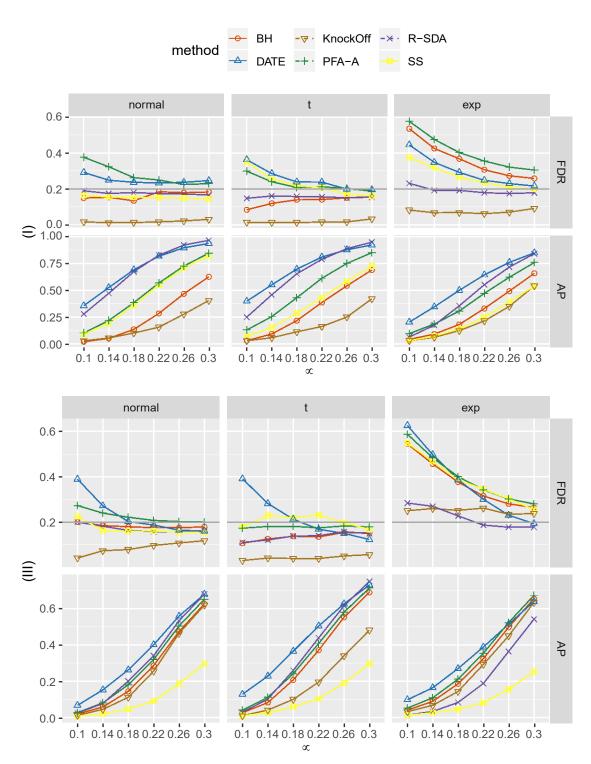


Figure S3: FDR and AP comparison for varying  $\mu$  in Settings (I) and (III) with estimated covariance matrix.

the DATE and PFA in the sense that it provides more accurate estimate of FDP and generally higher power. The conclusions are consistent for different choices of  $\pi_1$ , with narrower interquartile range of FDP and AP for larger  $\pi_1$ . As we can expect, the R-SDA has smaller variation than the single-splitting SDA.

#### E.4 The impact of the number of tests and sample sizes

We also conduct experiments by altering the number of tests p, while keeping  $(n, \pi_1, \alpha) = (90, 0.1, 0.2)$ . To make the AP comparable across p, the signal  $\mu$  is adjusted via  $\mu = C^p \log p/n$  with C depending on the covariance structures. The results are summarized in the top half of Figure 57. We can see that all methods have more accurate control of FDR as p increases, but the PFA<sub>A</sub> and DATE fail to control the FDR when p is small. To investigate the impact on sample sizes with unknown covariance, we set  $\mu = C^p \log(p)/n$ , fix  $(p, \pi_1, \alpha) = (500, 0.1, 0.2)$ , and consider the normal error. The results are summarized in the bottom half of Figure 57. We can see that that our R-SDA method is able to control the FDR and close to the nominal level regardless of the choice of n. Its superior performance relative to the other three methods is significant in some cases. Though all the methods exhibit steady AP pattern, the BH, PFA and DATE appear to need larger sample to achieve satisfactory FDR control than the R-SDA does. This again concurs with our theoretical result in Theorem 2 and demonstrates the advantage of using the nonparametric estimation of FDP in the SDA procedure.

### E.5 List of selected genes by different methods

Table S1 reports the list of 19 most differentially expressed probe sets obtained by the methods R-SDA, BH, SS, PFA-A and DATE in the real-data example.

Table S1: Differentially expressed probe sets in the B lineage ALL with BCR/ABL versus NEG molecular rearrangement, for five different multiple testing adjustment methods

R-SDA	ВН	SS	PFA	DATE
1635_at	1636_g_at	39730_at	1636_g_at	36502 <b>_</b> at
39730_at	39730 <b>_</b> at	39317_at	39730 <b>_</b> at	38385 <b>_</b> at
1636_g_at	1635_at	37027 <b>_</b> at	1635_at	40202 <b>_</b> at
36502_at	1674_at	38052_at	1674 <b>_</b> at	37403_at
37403_at	40504 <b>_</b> at	1635 <b>_</b> at	40202 <b>_</b> at	38052 <b>_</b> at
32134_at	40202 <b>_</b> at	1636 <b>_</b> g_at	37403 <b>_</b> at	33690_at
38052 <b>_</b> at	37015 <b>_</b> at	40202 <b>_</b> at	32434 <b>_</b> at	39317 <b>_</b> at
36821_at	37027 <b>_</b> at	34850_at	37014 <b>_</b> at	40876_at
38385 <b>_</b> at	32434_at	37403 <b>_</b> at	32979 <b>_</b> at	33440 <b>_</b> at
37027 <b>_</b> at	40167_s_at	37024 <b>_</b> at	1249 <b>_</b> at	1674 <b>_</b> at
1674 <b>_</b> at	40480_s_at	1249 <b>_</b> at	38111 <b>_</b> at	36908 <b>_</b> at
41872 <b>_</b> at	36591 <b>_</b> at	36802 <b>_</b> at	37015 <b>_</b> at	33774 <b>_</b> at
33440 <b>_</b> at	33774 <b>_</b> at	37025 <b>_</b> at	37147 <b>_</b> at	39730 <b>_</b> at
32434_at	37403 <b>_</b> at	32979_at	40504 <b>_</b> at	41592 <b>_</b> at
40876 <b>_</b> at	37014 <b>_</b> at	34870 <b>_</b> at	33440 <b>_</b> at	32134 <b>_</b> at
40202 <b>_</b> at	37363 <b>_</b> at	36502 <b>_</b> at	38112_g_at	39070 <b>_</b> at
39317 <b>_</b> at	34472 <b>_</b> at	33891 <b>_</b> at	36502 <b>_</b> at	37558 <b>_</b> at
32562 <b>_</b> at	32542 <b>_</b> at	34800 <b>_</b> at	31786 <b>_</b> at	33304 <b>_</b> at
34990 <b>_</b> at	39329 <b>_</b> at	36543_at	34850 <b>_</b> at	34180_at

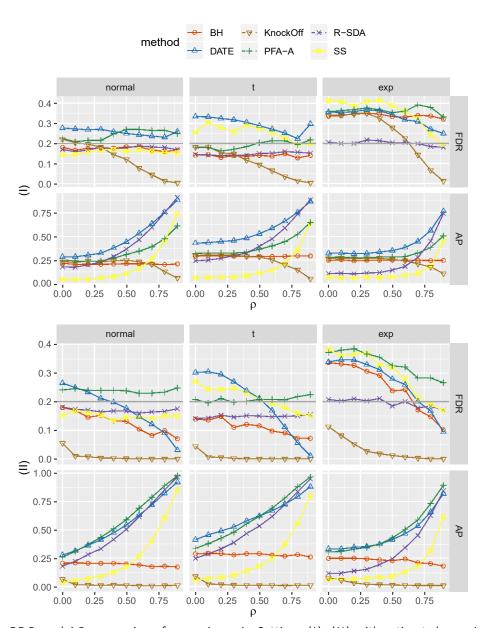


Figure S4: FDR and AP comparison for varying  $\rho$  in Settings (I)–(II) with estimated covariance matrix.

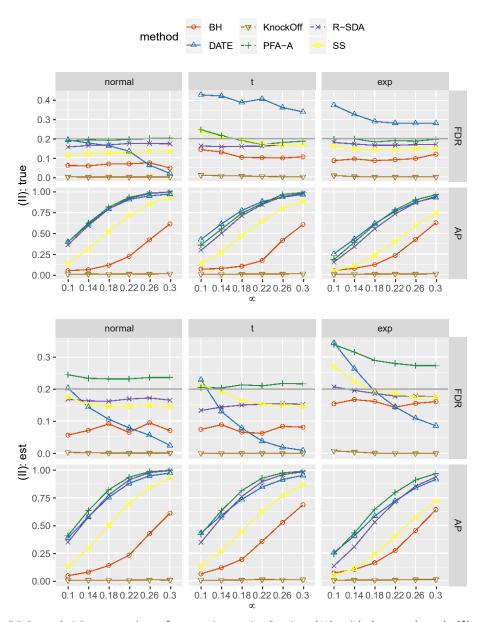


Figure S5: FDR and AP comparison for varying  $\mu$  in Setting (II) with known (top half) and unknown variances (bottom half).

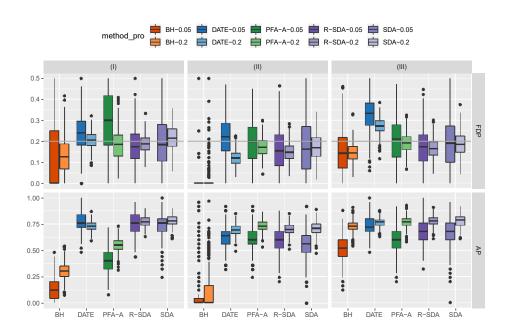


Figure S6: The boxplot of FDP and AP when the proportions of alternative are 0.05 and 0.2. The normal error is considered and  $(n, p, \alpha) = (90, 500, 0.2)$ . The signal strength  $\mu$  is set as 0.2, 0.15, 0.3 for the covariance structures (I)-(III), respectively.

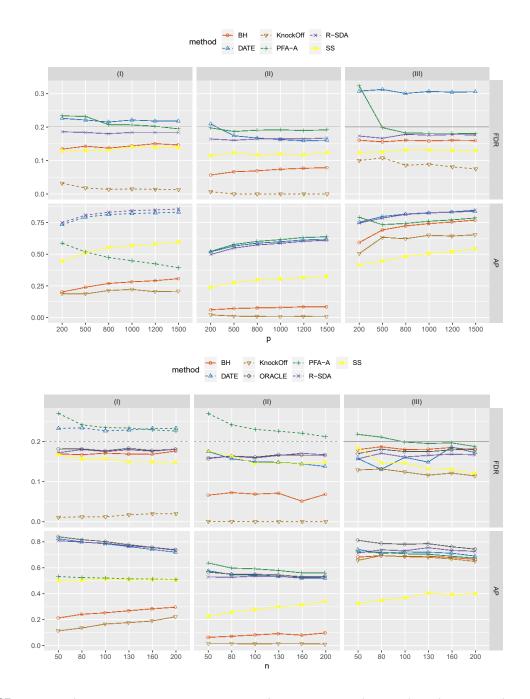


Figure S7: Top half: The empirical FDR and AP for varying p.  $(n,\pi_1,\alpha)$  = (90,0.1,0.2) and  $\mu_n$  =  $C^p \overline{\log(p)/n}$  with C = 0.8,0.5,1.2. Bottom half: The FDR and AP for varying n when the covariances are estimated.  $(p,\pi_1,\alpha)$  = (500,0.1,0.2) and  $\mu_n$  =  $C^p \overline{\log p/n}$  with C = 0.8,0.5,1.2.

