
A Perspective on False Discovery Rate Control via Knockoffs

Anonymous Author(s)

Affiliation

Address

email

Abstract

The knockoff filter introduced by Barber and Candès 2016 is an elegant framework for controlling the false discovery rate in variable selection. Yet, there is no conclusive result on the power (type II error rate) analysis, or how to choose the knockoff generation method, even in the Gaussian setting. When the predictors are i.i.d. Gaussian, it is known that as the signal to noise ratio tend to infinity, the power of the knockoff filter tends to 1 under any fixed FDR budget. However, when the predictors have a general covariance structure Σ , it is not obvious that one can define an analogous notion of the signal to noise ratio. We introduce the notion of *effective signal deficiency* (ESD) as any functional of Σ , such that the power tend to 1 *if and only if* this functional tends to 0 (under given noise level, sparsity ratio, and sampling rate). We then study the ESD for Lasso and the knockoff filter with different knockoff constructions, assuming the correctness of the replica method prediction for Lasso. As a baseline for comparison, we show that using Lasso with an oracle for choosing the threshold that gives the correct FDR, the ESD tends to 0 if and only if the empirical distribution of the diagonals of the precision matrix $\mathbf{P} := \Sigma^{-1}$ converges to 0 in distribution. In other words, the ESD can be taken as $\|(P_{jj})_{j=1}^p\|_{LP} := \inf \left\{ \epsilon > 0: \frac{1}{p} |\{P_{jj} \geq \epsilon\}| \leq \epsilon \right\}$. For the knockoff filter, if $\underline{\mathbf{P}}$ is the $2p \times 2p$ precision matrix for the predictors and knockoff variables, we show that the ESD is $\|(\underline{P}_{jj})_{j=1}^{2p}\|_{LP}$. We then find more explicit formulae for various specific knockoff constructions. We introduce the *conditional independence knockoff*, which always exists for Gaussian tree graphical models (or when the graph is sufficiently sparse), and show that its ESD is $\|(\Sigma_{jj} P_{jj}^2)_{j=1}^p\|_{LP}$. In contrast, for the equi-knockoffs in the literature, the ESD can achieve $\lambda_{max}(\mathbf{P})$, which is prohibitive when a small set of predictors are highly correlated.

1 Introduction

Modern large-scale data analysis often concerns the problem of finding a small set of highly informative predictors, among a larger set often of size comparable or larger than the number of observations. Examples include selecting a few genes related to a certain disease, or discovering a number of demographic attributes linked to the crime rates in a community. False-discovery rate (FDR) control, popularized by Benjamini and Hochberg's [BH95], has become a now-standard criterion for the type I errors in such large-scale hypothesis testing problems. Under orthogonal designs and assuming that the p -values under the null hypothesis are known, the Benjamini-Hochberg method is guaranteed to bound the FDR below any desired threshold ([BH95][STS04]). Recently, the knockoff filter [BC15][CFJL18] has emerged as a competitive approach for FDR control, which extends to setting beyond orthogonal designs and known p -values under null, and has demonstrated great empirical success [KS][SKB⁺].

37 The knockoff filter is build on the Lasso estimator. While it is possible to perform variable selection
 38 by a simple threshold test for the Lasso coefficients, it may not be feasible to determine such a
 39 threshold that gives rise to the desired FDR. The idea of the knockoff filter is to generate knockoff
 40 (fake) variables that has the same distribution as the true ones, but conditionally independent of the
 41 observations, and then regress the observations on both the true variables and the knockoff variables.
 42 Roughly speaking, one can then determine a threshold with the desired FDR guarantee, by leveraging
 43 on the Lasso coefficients for the knockoff variables, knowing that they are nulls.

44 It is known that any construction of the knockoffs that satisfies a certain exchangeability condition is
 45 guaranteed to control the FDR [BC15, Theorem 1]. Such a construction is not unique, and we would
 46 like to choose one with a high power (type II error rate). Heuristically, we would like the knockoffs
 47 to be as uncorrelated with the true predictors as possible, while the exchangeability condition is
 48 satisfied. Various algorithms for generating the knockoffs have been proposed in different settings,
 49 which typically involves solving an optimization that minimizes a heuristically chosen correlation
 50 measure [BC15][CFJL18][RSC]. Knockoff constructions with analytic expressions are rare (with
 51 the exception of the equi-knockoff [BC15] and metropolized knockoff sampling [BCJW19]). Partly
 52 due to this, analytical studies of the power of the knockoff filter has been very limited, even in the
 53 Gaussian setting. In the special case where the predictors are independent, one can generate the
 54 knockoffs simply independent of the true predictors, in which case [WBC17] has shown that the
 55 power tends to 1 as the signal to noise ratio tends to infinity (under a fixed sampling rate), by
 56 leveraging results on the Lasso statistics [BM12][SBC17]. For the case of correlated predictors,
 57 [FDLL19] proved a lower bound on the power, where the limiting (sample size $n \rightarrow \infty$) power is
 58 bounded below in terms of the number of predictors p and extremal eigenvalues of the covariance
 59 matrix of the true and knockoff variables. The assumption of bounded eigenvalues may not appear
 60 to capture the crux of the matter in certain scenarios. For example, if all predictors are independent
 61 except that two of them are always equal, then the minimum eigenvalue of the covariance matrix is
 62 zero, yet the FDR and the power is almost unchanged as we experimentally observed for the knockoff
 63 filter.

64 In this paper, we again consider the knockoff filter in the Gaussian case, but ask the following question:
 65 for a fixed sampling rate n/p and given a sequence of predictor covariance matrices $(\Sigma^{(p)})_{p=1}^\infty$, is
 66 power of the knockoff filter tending to 1 under any fixed FDR budget? Using the results on the Lasso
 67 statistics in [JM14], we find that the answer essentially depends on whether the empirical distribution
 68 of the diagonals of $\underline{\mathbf{P}} := \underline{\Sigma}^{-1}$ converges to 0, where $\underline{\Sigma}$ denotes the covariance matrix of the true and
 69 knockoff variables. Note that $\underline{\Sigma}^{-1}$ depends on the method of generating knockoffs, and hence this
 70 observation can be useful in the comparison of various knockoff constructions; an explicit evaluation
 71 function will be provided in (13), which we call *effective signal deficiency*.

72 A second contribution is to propose a new rule of generating the knockoffs, called *conditionally*
 73 *independent knockoffs* (CIK), which possesses both simple analytic expressions and excellent ex-
 74 perimental performance. CIK does not exist for all Σ , but we show its existence for tree graphical
 75 models or other sufficiently sparse graphs. Note that in practice, the so-called model-X knockoff filter
 76 requires the knowledge of Σ , an estimation of which is often prohibitive except when the graph has
 77 sparse or tree structures. CIK has simple explicit expressions of the effective signal deficiency for tree
 78 models, since the empirical distribution of the diagonals of Σ^{-1} is the same as that of $(P_{jj}^2 \Sigma_{jj})_{j=1}^p$.
 79 We remark that CIK is different than *metropolized knockoff sampling* studied in [BCJW19] (originally
 80 appeared in [CFJL18, Section 3.4.1]), even in the case of Gaussian Markov chains. The latter exists
 81 for generic distributions and is computationally efficient for Markov chains.

82 2 Preliminaries on the knockoff filter

83 Notation: $[n] := \{1, \dots, n\}$. We use boldface such as $\mathbf{X} := (X_{ij})_{i \in [n], j \in [p]}$ and $\mathbf{Y} := (Y_i)_{i \in [n]} =$
 84 Y^n to denote matrix and vectors. $\|\theta\|_0$ and $\|\theta\|_1$ denote the standard ℓ_0 and ℓ_1 norms of vectors.
 85 $\text{diag}(\mathbf{s})$ is a diagonal matrix when \mathbf{s} is a vector, and $\text{diag}(\mathbf{P})$ denotes the vector of diagonal entries
 86 when \mathbf{P} is a matrix. In discussions of knockoffs, we use the underline to indicate instances in the
 87 extended cases with knockoff variables, e.g., $\underline{\theta}$ denotes a $2p$ -vector when θ is a p -vector. $Q(r)$, $r \in \mathbb{R}$
 88 denotes the Gaussian tail probability. $\mathbf{A} \preceq \mathbf{B}$ means that the matrix $\mathbf{B} - \mathbf{A}$ is positive semidefinite.

89 Suppose that the true observation model is $Y = \sum_{j=1}^p \theta_{0,j} X_j + N$, where $\theta_0 \in \mathbb{R}^p$ are the unknown
 90 parameters, $X^p := (X_j)_{j=1}^p$ are the (observable) predictors, and N is the noise. We adopt the model-

91 X framework [CFJL18], where X^p are assumed to be random variables with known distribution. This
 92 is a “semi-supervised learning” setting where a large number of unlabelled samples are available for
 93 estimating the distribution of X^p . Knowing a sample size n number of observations and predictor
 94 values, the knockoff filter aims to determine the set of active predictor, $\{j: \theta_{0,j} \neq 0\}$, while
 95 controlling the false discovery rate (FDR)

$$FDR := \mathbb{E} \left[\frac{|\mathcal{H}_0 \cap \hat{\mathcal{H}}_1|}{|\hat{\mathcal{H}}_1|} \right] \quad (1)$$

96 below a given threshold. Here, $H_0 := \{j: \theta_{0,j} = 0\}$ and $\hat{\mathcal{H}}_1$ denotes the set of selected predictors.
 97 The method is to generate knockoff variables $\tilde{X}_1, \dots, \tilde{X}_p$, with the property that

$$(X^p, \tilde{X}^p)_{\text{swap}(S)} = (X^p, \tilde{X}^p) \quad (2)$$

98 in distribution, for any set $S \in \{1, \dots, p\}$. The swap operation means switching the true and
 99 knockoff coordinates with indices in S ; for example, if $p = 2$ and $S = \{1\}$, then $(X^2, \tilde{X}^2)_{\text{swap}(S)} =$
 100 $(\tilde{X}_1, X_2, X_1, \tilde{X}_2)$.

101 Recall that we use the underline to indicate instances in the extended cases with knockoff variables.
 102 For example, $\underline{\theta}$ is a $2p$ -vector, $\underline{\Sigma}$ is a $2p \times 2p$ matrix. The the knockoff filter performs the following:
 103 regress Y on $[X^p, \tilde{X}^p]$, let $\hat{\underline{\theta}}_1, \dots, \hat{\underline{\theta}}_{2p}$ be the Lasso coefficients, and put

$$W_j := |\hat{\underline{\theta}}_j| - |\hat{\underline{\theta}}_{j+p}|, \quad (3)$$

104 $j = 1, \dots, p$. Choose the data dependent threshold $T > 0$ by the following rule

$$T := \min \left\{ t \in \mathcal{W}: \frac{|\{j: W_j \leq -t\}|}{|\{j: W_j \geq t\}| \vee 1} \leq q \right\} \quad (4)$$

105 where $\mathcal{W} := \{|W_j|: j = 1, \dots, p\} \setminus \{0\}$, and q equals the given FDR budget. Then select j for
 106 which $W_j > T$ as the active predictors. It is shown in [BC15, Theorem 1] that this procedure bounds¹
 107 FDR below q . More generally, the FDR is controlled below q as long as $(W_j)_{j=1}^p$ depends on $(\underline{\mathbf{X}}, \mathbf{Y})$
 108 only through $(\underline{\mathbf{X}}^\top \underline{\mathbf{X}}, \underline{\mathbf{X}}^\top \mathbf{Y})$, and satisfies the antisymmetry property [BC15, Section 2.2].

109 For Gaussian X^p , note that the exchangeability condition implies that the covariance of (X^p, \tilde{X}^p)
 110 has the form

$$\underline{\Sigma} = \begin{bmatrix} \Sigma & \Sigma - \text{diag}(\mathbf{s}) \\ \Sigma - \text{diag}(\mathbf{s}) & \Sigma \end{bmatrix}. \quad (5)$$

111 As observed in [BC15], positive semi-definiteness of this matrix is equivalent to

$$\text{diag}(\mathbf{s}) \succeq \mathbf{0}, \quad (6)$$

$$2\Sigma - \text{diag}(\mathbf{s}) \succeq \mathbf{0}. \quad (7)$$

112 Previous methods for generating the knockoffs (computing \mathbf{s}) include the following [CFJL18]:

- 113 • The equi-knockoffs construction chooses s_1, \dots, s_p all equal. Note that the maximum of
 114 such value compatible with (7) is $2\lambda_{\min}(\Sigma)$. [CFJL18] assumed the normalization $\Sigma_{jj} = 1$,
 115 $j = 1, \dots, p$, and recommended choosing

$$s_j = 2\lambda_{\min}(\Sigma) \wedge 1, \quad (8)$$

116 with the goal of minimizing the correlation between X_j and \tilde{X}_j .

- 117 • The semidefinite program (SDP) construction solves the following

$$\text{minimize} \quad \sum_{j=1}^p |\Sigma_{jj} - s_j| \quad (9)$$

$$\text{s.t.} \quad 0 \leq s_j \leq \Sigma_{jj}, \quad \text{diag}(\mathbf{s}) \preceq 2\Sigma. \quad (10)$$

¹While [BC15] discusses fixed knockoff whereas the present paper mainly concerns the model-X knockoff, the proof in [BC15] still works in the model-X case (see the explanation in [CFJL18]).

- The approximate semidefinite program (ASDP) construction first solves (9) and (10) with Σ replaced by a certain block diagonal matrix, returning a vector \hat{s} . Then the final result is chosen as $s = \gamma \hat{s}$, where γ is the maximum scalar that fulfils (7).

We do not discuss other knockoff constructions, such as the exact construction [CFJL18, Section 3.4.1] and deep knockoff [RSC], are mostly targeting at general non-Gaussian distributions.

Define

$$POWER := \mathbb{E} \left[\frac{|\hat{\mathcal{H}}_1 \cap \mathcal{H}_1|}{|\mathcal{H}_1|} \right] \quad (11)$$

where $\mathcal{H}_1 := \{j: \theta_{0,j} \neq 0\}$. Previously, [JM14] performed power analysis for i.i.d. design in the linear (fixed n/p) regime. [FDLL19] performed power analysis for a general Σ , and showed a consistency result as $n \rightarrow \infty$.

3 Summary of the main results

The present paper is interested in the linear (fixed n/p) regime and general Σ . At first glance, it is not even obvious that any meaningful result can be said for a general sequence $(\Sigma^{(p)})_{p=1}^\infty$. A starting point is the observation that, under mild conditions, the empirical distribution of the errors in the regression coefficients divided by $(\Sigma^{-1})_{jj}^{1/2}$ is asymptotically normal. This has been hinted or has explicitly appeared in various literature on regression problems, e.g. [JM14][EKBB⁺13, Lemma 1].

More formally, we consider the asymptotic setting where the sequence of instances $\{(\theta_0^{(p)}, \Sigma^{(p)})\}_{p \geq 1}$ has a *standard distributional limit* in the sense of [JM14, Definition 4.1] (reproduced in Definition 3). This assumption implies that the empirical distribution of $\{(\theta_{0,j}, \underline{P}_{jj})\}_{j=1}^p$ converges weakly to some probability measure on \mathbb{R}^2 , and it is believed to be not much stronger than that, although a more precise characterization of such sequences remains an outstanding question [JM14]. Assuming the correctness of the replica method calculations, [JM14, Replica Method Claim 4.6] provided mild conditions under which the standard distributional limit exists. In particular, as observed by [JM14], those conditions are satisfied for block diagonal covariance matrices in which the empirical distribution of the block converges. In other words, in a “direct sum” version of the problem where we have a fixed Σ but have k independent copies of those predictors, and let n also grow linearly in k , the assumption for the replica method claim in [JM14, Replica Method Claim 4.6] is always satisfied. We remark that in the regime of vanishing $\|\theta_0\|_0/p$, there are also rigorous (without appealing to the replica method) results showing that the weak convergence of the distribution of $\{(\theta_{0,j}, \underline{P}_{jj})\}_{j=1}^p$ is essentially sufficient for the existence of a standard distributional limit ([JM14, Theorem 4.5]), although the present paper does not concern that regime. We introduce:

Definition 1 (Effective signal deficiency). Given any sampling rate² $\delta := n/p > 1$, noise level $N_i \sim \mathcal{N}(0, \sigma^2)$ (under the model $Y_i = \sum_{j=1}^p \theta_{0,j} X_{ij} + N_i, i = 1, \dots, n$), and a variable selection algorithm, define the *effective signal deficiency* (ESD) as any function of a sequence $(\theta_0^{(p)}, \Sigma^{(p)})_{p \geq 1}$ with a standard distributional limit, such that the following property holds: for any $\epsilon > 0$, there exists $\epsilon' > 0$ such that $ESD < \epsilon'$ ensures that $\limsup_{p \rightarrow \infty} \max\{FDR^{(p)}, 1 - POWER^{(p)}\} < \epsilon$.

We are often interested in settings where θ_0 has given sparsity level and bounds on the amplitude of the nonzero coefficients, so effectively ESD is a function of the sequence $(\Sigma^{(p)})_{p \geq 1}$. Also note that by definition, ESD is not unique, and our goal is to find simple representations of the equivalent class. ESD is a potentially useful concept in comparing or evaluating different ways of generating knockoff matrices. As an analogy, think of the various notions of convergences of probability measures. A sequence of probability measures may converge in one topology but not in another. Similarly, one may cook up different functionals of the covariance matrix, such as $\lim_{p \rightarrow \infty} p \operatorname{Tr}^{-1}(\Sigma)$ and $\lim_{p \rightarrow \infty} p \operatorname{Tr}(\Sigma^{-1})$, which both intuitively characterize some sort of signal deficiency since they tend to be small when the signal gets stronger. However, they are not equivalent, and the second convergence to 0 is *stronger* in the sense that the first must vanish when the second vanishes. ESD is

²We assume $\delta > 1$ for convenience so that the parameter τ in the replica analysis can be bounded independently of Σ . However the result in [JM14, Definition 4.1] applies to any $\delta > 0$.

intended to be the correct notion of “convergence” that characterizes FDR tending to 0 and power tending to 1.

Of course, by definition it is not obvious that a succinct expression of such an effective signal deficiency exists. Remarkably, we find that the effective signal deficiency can be characterized by the convergence of certain empirical distribution derived from Σ . The effective signal deficiency for various (old and new) algorithms as follows:

- Not using knockoffs: one may use Lasso to regress \mathbf{Y} on \mathbf{X} , and obtain $\hat{\theta}$:

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\theta\|^2 + \lambda \|\theta\|_1 \right\} \quad (12)$$

The parameter λ can be chosen as any fixed positive number independent of p . Instead of a direct threshold test on $\hat{\theta}$, we compute an “unbiased version” $\hat{\theta}^u$ (defined in (15)) as in [JM14] for simplicity of the analysis, and pass a threshold to select non-nulls. Suppose that there is an oracle telling how to pick the threshold to make FDR at the desired level. We show that ESD for this oracle algorithm is the limit $p \rightarrow \infty$ of

$$\|(P_{jj})_{j=1}^p\|_{LP} := \inf \left\{ \epsilon > 0 : \frac{1}{p} |\{P_{jj} \geq \epsilon\}| \leq \epsilon \right\}. \quad (13)$$

The assumption of the standard distributional limit ensures the weak convergence of the empirical distribution of $(P_{jj})_{j=1}^p$, and hence the convergence of (13). For simplicity, we may simply say $\|(P_{jj})_{j=1}^p\|_{LP}$ is ESD without mentioning the limit in p , when there is no confusion. In other words, we defined $\|(P_{jj})_{j=1}^p\|_{LP}$ as the distance between the empirical distance of $\operatorname{diag}(\mathbf{P})$ and the delta measure at 0, under the Lévy-Prokhorov metric³ (we are abusing the notation of norms even though this is not a norm). Note that $\|\cdot\|_{LP}$ can be replaced by any metric compatible with the weak convergence topology.

- General knockoff: for a general (potentially non-analytic) knockoff construction, it seems hopeless to find simple expressions of ESD in terms of Σ . Nevertheless, if $(\theta_0^{(p)}, \Sigma^{(p)})$ has a standard distributional limit, we can express ESD as $\|(P_{jj})_{j=1}^{2p}\|_{LP}$, where we recall that \mathbf{P} is the extended precision matrix including the knockoff variables. We next find more explicit expressions in terms of \mathbf{P} for specific constructions:
- Equi-knockoff: we show that ESD is at least $\lambda_{\max}(\mathbf{P})$. This is also achievable by choosing $s_j = \lambda_{\min}(\Sigma)$ (note that this is slightly different than (8) in [BC15][CFJL18]).
- We introduce a new method for generating the knockoff matrix, called *conditional independence knockoff*. If the Gaussian graphical model is from a tree, the conditional independence knockoff always exists, and the ESD is $\|P_{jj}^2 \Sigma_{jj}\|_{LP}$.

As noted in [FDLL19], although knockoff filter has the advantage of controlling FDR, it usually has a lower power than Lasso with oracle threshold. We use oracle threshold Lasso as a baseline for comparison, and indeed its ESD is smaller than that of other algorithms.

The last knockoff construction, conditional independence knockoff, appears to be new. It is both analytically simple and empirically competitive. Comparing equi- and conditional independence knockoff: the latter is more robust, since having a small fraction of j with large $P_{jj}^2 \Sigma_{jj}$ does not increase $\|\cdot\|_{LP}$ much. For example, if the p and $p-1$ th predictors are equal, then the ESD for conditional independence knockoff almost does not change, but equi-knockoff completely fails. While the solution in the (approximate) semidefinite knockoff is not analytically simple, empirically we find that the conditional independence knockoff usually has similar or improved performance.

4 Baseline: Lasso with oracle threshold

Before analyzing any algorithm, let us observe the following converse bound, which is information-theoretic (i.e. not limited to Lasso or any particular algorithm). This result lower bounds the effective

³Generally, the Lévy-Prokhorov distance between two probability measures μ and ν is defined as $\inf\{\epsilon > 0 | \mu(A) \leq \nu(A^\epsilon) + \epsilon, \nu(A) \leq \mu(A^\epsilon) + \epsilon, \forall A\}$, where A^ϵ denotes the ϵ -neighborhood of A .

signal deficiency (ESD) by $\|(P_{jj})_{j=1}^p\|_{LP}$, where $\|\cdot\|_{LP}$ was defined in (13). Intuitively, the result comes from the fact that the conditional variance of X_j given $X_{\setminus j}$ is P_{jj}^{-1} .

Proposition 2 (Converse). *Fix $\alpha \in (0, 1)$, $n, p \in \mathbb{N}$, and let Σ be the covariance matrix of the Gaussian predictors. Let $\theta_1, \dots, \theta_p$ be i.i.d. $\text{Ber}(\alpha)$. Assume that the noise variance (for each sample) is n . Assume that there exists an algorithm satisfying $FDR \leq q$, $POWER \geq 1 - \epsilon$. Then for $n, p \geq N(\epsilon, \alpha)$ large enough, we have*

$$\|(P_{jj})_{j=1}^p\|_{LP} \leq \max \left\{ \frac{1.1}{\left(2Q^{-1}(\sqrt{\frac{q}{1-\alpha}} + \sqrt{2\epsilon})\right)^2}, \frac{q}{1-\alpha} + \sqrt{\epsilon} \right\}, \quad (14)$$

where $Q(a) :=$ denotes the standard Gaussian tail probability. In particular, $\max\{\epsilon, q\} \rightarrow 0$ implies that $\|(P_{jj})_{j=1}^p\|_{LP} \rightarrow 0$.

We next show that $\|(P_{jj})_{j=1}^p\|_{LP}$ is also an achievable ESD, and in fact achievable by appropriately using Lasso. More precisely, for a sequence of instances having a standard distributional limit defined as follows (introduced by [JM14, Definition 4.1]), $\lim_{p \rightarrow \infty} \|(P_{jj})_{j=1}^p\|_{LP}$ is an achievable ESD, where the existence of the limit is ensured by the standard distributional limit assumption.

Definition 3 (Standard distributional limit). A sequence $\{(\Sigma^{(p)}, \theta_0^{(p)}, m^{(p)}, n^{(p)}, \sigma^{(p)})\}_{p \geq 1}$ is said to have a *standard distributional limit* if there exists $\tau \neq 0$ and potentially random $d \in \mathbb{R}$ such that $\{\theta_{0,j}, (\hat{\theta}_j^u - \theta_{0,j})/\tau, (\Sigma^{-1})_{jj}\}_{j=1}^m$ converges almost surely to a probability measure ν on \mathbb{R}^3 . Here, $\Sigma^{(p)}$ is $m^{(p)} \times m^{(p)}$ matrix, $\hat{\theta}^u$ is defined in terms of the Lasso estimator:

$$\hat{\theta}^u := \hat{\theta} + \frac{d}{n} \Sigma^{-1} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X} \hat{\theta}), \quad (15)$$

ν is the probability distribution of $(\Theta_0, \Upsilon^{1/2} Z, \Upsilon)$, where $Z \sim \mathcal{N}(0, 1)$, and Θ_0 and Υ are some random variables independent of Z .

As mentioned in [JM14], characterizing instances having a standard distributional limit is highly nontrivial. Yet, at least, the definition is non-empty since it contains the case of standard Gaussian design. Moreover, a non-rigorous replica argument indicates that the standard distributional limit exists as long as a certain functional defined on \mathbb{R}^2 has a differentiable limit [JM14, Replica Method Claim 4.6], which is always satisfied for block diagonal Σ where the empirical distribution of the blocks converges.

We now consider a simple variable selection algorithm based on the Lasso estimator (without using knockoffs), where indices j for which $|\hat{\theta}_j^u|$ exceeds a certain threshold are selected (see definition of $\hat{\theta}^u$ in (15)). In practice, the knockoff filter has the advantage of controlling FDR. However, if Lasso is used with the right threshold giving the correct FDR, then Lasso has higher power than the knockoff filter (see the discussion in [FDLL19]).

Proposition 4 (Lasso achievability). *Let $\{(\Sigma^{(p)}, \theta_0^{(p)}, n^{(p)}, \sigma^{(p)})\}_{p \geq 1}$ be any sequence having a standard distributional limit, where*

$$|\theta_{0,j}| \geq 1, \quad \forall j: \theta_{0,j} \neq 0, \quad (16)$$

$$\limsup_{p \rightarrow \infty} \|\theta_0^{(p)}\|_1/p = \beta < \infty, \quad (17)$$

$$\lim_{p \rightarrow \infty} \|\theta_0^{(p)}\|_0/p = \alpha, \quad (18)$$

$$\lim_{p \rightarrow \infty} n^{(p)}/p = \delta, \quad (19)$$

$$\sigma^{(p)} = \sqrt{n} \sigma_0, \quad (20)$$

$$\lim_{p \rightarrow \infty} \|(P_{jj})_{j=1}^p\|_{LP} = L. \quad (21)$$

Then using a threshold test for $\hat{\theta}^u$ defined in (15), where the Lasso parameter $\lambda > 0$ is any number independent of p , one can bound both $FDR^{(p)}$ and $1 - POWER^{(p)}$ by $f_{\alpha, \beta, \delta, \sigma_0, \lambda}(L)$ almost surely for large enough p , where $f_{\alpha, \beta, \delta, \sigma_0, \lambda}(\cdot)$ is a function that vanishes at the origin for any fixed $\alpha, \beta, \delta, \sigma_0, \lambda$. For explicit bounds, see (22) and (23).

240 Explicitly, we have the following bounds almost surely,

$$\limsup_{p \rightarrow \infty} FDR^{(p)} \leq 2Q\left(\frac{1}{2\tau\sqrt{L}}\right) + L; \quad (22)$$

$$\liminf_{p \rightarrow \infty} POWER^{(p)} \geq 1 - \frac{1}{\alpha} \left[2Q\left(\frac{1}{2\tau\sqrt{L}}\right) + L \right], \quad (23)$$

241 where τ is from the definition of the standard distributional limit, and is bounded by

$$\tau^2 \leq \frac{2\delta\sigma_0^2}{\delta-1} + \frac{2\lambda\delta(\delta+1)\beta}{(\delta-1)^3}. \quad (24)$$

242 5 Results for general knockoff matrices

243 Given Σ , let $\underline{\Sigma}$ be the extended $2p \times 2p$ covariance matrix for the true predictors and their knockoffs.
 244 Let $\underline{\theta}_0$ be the $2p$ -vector where the indices corresponding to the knockoff variables ($j = p+1, \dots, 2p$)
 245 are 0. Consider the procedure of the knockoff filter described in Section 2, with a slight tweak: define
 246 $W_j := |\hat{\theta}_j^u| - |\hat{\theta}_{j+p}^u|$, instead of (3), where unbiased regression coefficients $\hat{\theta}^u$ is defined analogous
 247 to (15). This definition of W_j still fulfills the sufficiency and antisymmetry condition in [BC15,
 248 Section 2.2], so FDR can still be controlled. This change allows us to perform analysis using results
 249 in [JM14]. We also assume that the Lasso parameter λ is an arbitrary number independent of p . Then,
 250 assuming the existence of standard distributional limit, we show that $\|(\underline{P}_{jj})_{j=1}^{2p}\|_{LP}$ is ESD.

251 **Lemma 5.** Let $\{(\underline{\Sigma}^{(p)}, \underline{\theta}_0^{(p)}, 2p, n^{(p)}, \sigma^{(p)})\}_{p \geq 1}$ be a sequence having a standard distribution limit.
 252 Suppose that (16)-(20) still applies, with $\delta > 2$. Suppose that

$$q > \frac{2Q(1/3\tau\sqrt{L}) + L}{\alpha/2 - 4Q(1/3\tau\sqrt{L}) - 2L} \quad (25)$$

253 where τ is from the definition of the standard distributional limit, which is bounded as in (24). Let
 254 $L := \lim_{p \rightarrow \infty} \|(\underline{P}_{jj})_{j=1}^{2p}\|_{LP}$. Then almost surely, running the knockoff filter with FDR budget at q
 255 achieves power

$$\liminf_{p \rightarrow \infty} POWER^{(p)} \geq 1 - \frac{8}{\alpha} Q(1/3\tau\sqrt{L}) - \frac{4L}{\alpha}. \quad (26)$$

256 In particular, the asymptotic power tends to 1 as $L \rightarrow 0$.

257 6 Conditional independence knockoff and ESD

258 We introduce the *conditional independence knockoff*, where X_j and \tilde{X}_j are conditionally independent
 259 of $X_{\setminus j}$, for each $j = 1, \dots, p$. This condition implies that

$$s_j = \text{Var}(X_j | X_{\setminus j}) = P_{jj}^{-1}, \quad j = 1, \dots, p, \quad (27)$$

260 where s_1, \dots, s_p are as defined in (5). However such an \mathbf{s} may violate the positive semidefinite
 261 assumption for the joint covariance matrix (example with $p = 3$ exists). Yet, interestingly, we find
 262 that in the case of tree graphical model, this construction always exists. In many practical scenarios,
 263 the predictors X^p comes from a tree graphical model, and we can estimate the underlying graph using
 264 the Chow-Liu algorithm [CL68].

265 **Theorem 6.** $\underline{\Sigma}$ defined in (5) is positive semidefinite with \mathbf{s} defined in (27), if either of the following
 266 is satisfied: 1) Σ is the covariance matrix of a tree graphical model; 2) \mathbf{P} is diagonally dominant.

267 Either condition in the theorem intuitively imposes that the graph is sparse. In practice, Σ needs to be
 268 estimated, which is generally only feasible with some sparse structure (e.g. via graphical lasso).

269 Assuming the existence of a standard distributional limit, we have the following results:

270 **Theorem 7.** For tree graphical models, assuming that the algorithm is knockoff filter using the
 271 conditional independence knockoff, $\|(P_{jj}^2 \Sigma^{jj})_{j=1}^p\|_{LP}$ is an effective signal deficiency.

272 **Theorem 8.** The effective signal deficiency for equi-knockoff with $s_j = \lambda_{\min}(\Sigma)$, $j = 1, \dots, p$ is
 273 $\lambda_{\min}(\Sigma)$.

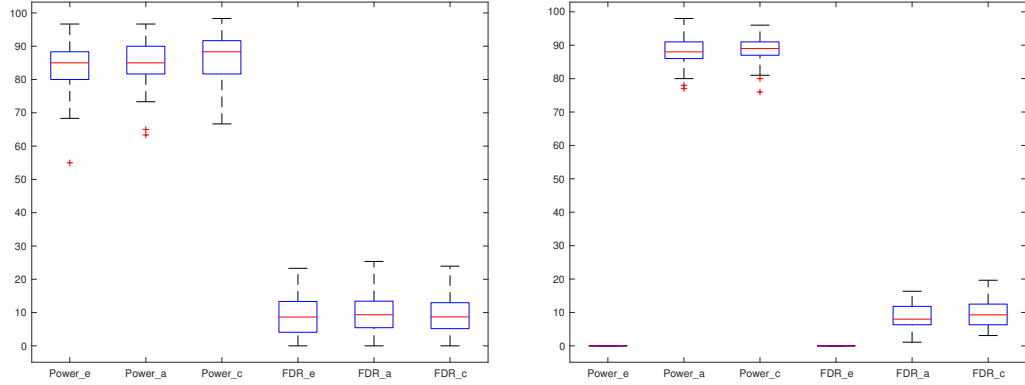


Figure 1: Left: Binary tree, equal correlations. Extension _e, _a, _c refers to equi-knockoffs, asdp knockoffs, and conditional independence knockoff, respectively. Right: Markov chain, randomly chosen correlation strengths.

7 Experimental results

We first consider the setting where the conditional independence graph X_1, \dots, X_p forms a binary tree, in which $X_1, \dots, X_p \sim \mathcal{N}(0, 1)$. The correlations between adjacent nodes are all equal to 0.5. Choose $k = 100$ out of $p = 1000$ indices uniformly at random as the support of θ , and set $\theta_j = 4.5$ for j in the support. Generate $n = 1000$ samples $Y_i = \mathbf{X}_i\theta + N_i$ where $N_i \sim \mathcal{N}(0, n)$.

Figure 1 Left shows the box plots of the power and FDR for the knockoff filter with three different knockoff constructions, equi-knockoff, sd-knockoff, and conditional independence knockoff. The FDR is controlled at the target $q = 0.1$ in all three cases. The powers are similar, but the rough trend is $POWER_e < POWER_s < POWER_c$. We then compare the effective signal deficiency. Note that in the current setting, $\text{Var}(\underline{X}_j | \underline{X}_{\setminus j}) \leq 1$, and hence $\underline{P}_{jj} \geq 1$, for each $j = 1, \dots, 2p$, and we always have $\|(\underline{P}_{jj})_{j=1}^{2p}\|_{LP} = 1$ by the definition (13), which cannot reveal any useful information for comparison. To resolve this, we can scale down \underline{P}_{jj} by a common factor before computing the LP norms, noting that such a scaled version of the LP norm is still a valid effective signal deficiency (in the same equivalence class). Lacking a systematic way of choosing such a scaling factor, heuristically we choose it so that the LP norms for the three algorithms are all “bounded away from 0 and 1”. We find that $\|(\underline{P}_{e,jj})_{j=1}^{2p}/2000\|_{LP} = 0.5010$, $\|(\underline{P}_{s,jj})_{j=1}^{2p}/2000\|_{LP} = 0.0480$, and $\|(\underline{P}_{c,jj})_{j=1}^{2p}/2000\|_{LP} = 0.0025$, and their ordering matches the ordering of the powers.

In the previous example, the simplest equi-knockoff has a highly competitive performance. However, this is an artifact of the fact that the data covariance is highly structured (i.e., correlations are all the same). If the correlations have high fluctuations, and in particular, a small number of node pairs are highly correlated, then the equi-knockoff has a much worse performance. This is demonstrated in the next example. Consider the setting where X_1, \dots, X_p forms a Markov chain, in which $X_1, \dots, X_p \sim \mathcal{N}(0, 1)$. The correlation between X_j and X_{j+1} is $\rho_j := G_j 1\{|G_j| \leq 1\}$, where $G_j \sim \mathcal{N}(0, 0.25)$, $j = 1, \dots, p - 1$ are chosen independently. Choose $k = 100$ out of $p = 1000$ indices uniformly at random as the support of θ , and set $\theta_j = 4.5$ for j in the support. Generate $n = 1200$ samples $Y_i = \mathbf{X}_i\theta + N_i$ where $N_i \sim \mathcal{N}(0, 0.49n)$.

Figure 1 Right shows the box plots of the power and FDR for the knockoff filter with three different knockoff constructions. The target FDR $q = 0.1$. Since the correlations are now chosen randomly, with high probability there exist highly nodes, and hence $\lambda_{\min}(\Sigma)$ can be very small, in which case the equi-knockoff performs poorly. (Figure 1 shows the case where the correlations are truncated between $[-1, 1]$. If we truncate the correlation to a smaller interval around 0, we can observe that $POWER_e$ goes up). However $POWER_c$ is similar to $POWER_s$, with the median of the former slightly higher. To compare the ESD, first scale down \underline{P}_{jj} by a heuristically chosen factor. We find $\|(\underline{P}_{e,jj})_{j=1}^{2p}/100\|_{LP} = 0.9995$, $\|(\underline{P}_{s,jj})_{j=1}^{2p}/100\|_{LP} = 0.8660$, and $\|(\underline{P}_{c,jj})_{j=1}^{2p}/100\|_{LP} = 0.1075$, and their ordering matches the ordering of the powers of the three knockoff constructions.

References

- [BC15] Rina Foygel Barber and Emmanuel J. Candès. Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5):2055–2085, 2015.
- [BCJW19] Stephen Bates, Emmanuel Candès, Lucas Janson, and Wenshuo Wang. Metropolized knockoff sampling. *arXiv preprint arXiv:1903.00434*, 2019.
- [BH95] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
- [BM12] Mohsen Bayati and Andrea Montanari. The lasso risk for Gaussian matrices. *IEEE Transactions on Information Theory*, 58(4):1997–2017, 2012.
- [CFJL18] Emmanuel Candès, Yingying Fan, Lucas Janson, and Jinchi Lv. Panning for gold: ‘model-x’ knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(3):551–577, 2018.
- [CL68] C Chow and Cong Liu. Approximating discrete probability distributions with dependence trees. *IEEE transactions on Information Theory*, 14(3):462–467, 1968.
- [EKBB⁺13] Nouredine El Karoui, Derek Bean, Peter J Bickel, Chinghay Lim, and Bin Yu. On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences*, 110(36):14557–14562, 2013.
- [FDLL19] Yingying Fan, Emre Demirkaya, Gaorong Li, and Jinchi Lv. Rank: large-scale inference with graphical nonlinear knockoffs. *Journal of the American Statistical Association*, pages 1–43, 2019.
- [JM14] Adel Javanmard and Andrea Montanari. Hypothesis testing in high-dimensional regression under the gaussian random design model: Asymptotic theory. *IEEE Transactions on Information Theory*, 60(10):6522–6554, 2014.
- [KS] Eugene Katsevich and Chiara Sabatti. Multilayer knockoff filter: Controlled variable selection at multiple resolutions. *arXiv:1706.09375 (2017)*.
- [RSC] Yaniv Romano, Matteo Sesia, and Emmanuel Candès. Deep knockoffs. *arXiv:1811.06687 (2018)*.
- [SBC17] Weijie Su, Malgorzata Bogdan, and Emmanuel Candès. False discoveries occur early on the lasso path. *The Annals of Statistics*, 45(5):2133–2150, 2017.
- [SKB⁺] Matteo Sesia, Eugene Katsevich, Stephen Bates, Emmanuel Candès, and Chiara Sabatti. Multi-resolution localization of causal variants across the genome. *bioRxiv (2019)*.
- [STS04] John D Storey, Jonathan E Taylor, and David Siegmund. Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(1):187–205, 2004.
- [WBC17] Asaf Weinstein, Rina Barber, and Emmanuel Candès. A power and prediction analysis for knockoffs with lasso statistics. *arXiv preprint arXiv:1712.06465*, 2017.

A Proof of Proposition 2

Proof. We can prove such an impossibility bound even assuming that when deciding whether j is null, one has the full information of $\theta_{0,\setminus j}$. Then the problem becomes testing a single parameter θ_j from the effective observation $\theta_j(X_j - \mathbb{E}[X_j|X_{\setminus j}]) + W$ where the signal $X_j - \mathbb{E}[X_j|X_{\setminus j}]$ has variance $[(\Sigma^{-1})_{jj}]^{-1}$.

353 Using the fact that $|\{j: \hat{\theta}_j \neq 0\}| \leq p$, we obtain from the definition of FDR that $\frac{1}{p} \sum_{j=1}^p \mathbb{P}[\theta_j =$
 354 $0, \hat{\theta}_j \neq 0] \leq q$, and hence

$$\frac{1}{p} \sum_{j=1}^p \mathbb{P}[\hat{\theta}_j \neq 0 | \theta_j = 0] \leq \frac{q}{1-\alpha}. \quad (28)$$

355 By the Markov inequality we have $\frac{1}{q} |\mathcal{S}_1| \leq \sqrt{\frac{q}{1-\alpha}}$ where we defined

$$\mathcal{S}_1 := \left\{ j: \mathbb{P}[\hat{\theta}_j \neq 0 | \theta_j = 0] > \sqrt{\frac{q}{1-\alpha}} \right\}. \quad (29)$$

356 For each $j \in [p] \setminus \mathcal{S}_1$, we have

$$\sqrt{\frac{q}{1-\alpha}} > \mathbb{P}[\hat{\theta}_j \neq 0 | \theta_j = 0] \quad (30)$$

$$= \mathbb{P} \left[\hat{\theta}_j \neq 0 \left| \frac{1}{n} \|\bar{\mathbf{X}}_j\|^2 P_{jj} \leq 1.1, \theta_j = 0 \right. \right] \cdot \mathbb{P} \left[\frac{1}{n} \|\bar{\mathbf{X}}_j\|^2 P_{jj} \leq 1.1 \right] \quad (31)$$

$$\geq \mathbb{P} \left[\hat{\theta}_j \neq 0 \left| \frac{1}{n} \|\bar{\mathbf{X}}_j\|^2 P_{jj} \leq 1.1, \theta_j = 0 \right. \right] \cdot (1 - o_n(1)) \quad (32)$$

357 where we defined $\bar{\mathbf{X}}_j := \mathbf{X}_j - \mathbb{E}[\mathbf{X}_j | \mathbf{X}_{\setminus j}]$; (31) used the independence of $\bar{\mathbf{X}}_j$ and θ_j ; (32) used the
 358 concentration of the χ^2 distribution.

359 Let us turn to Type-II error:

$$1 - \epsilon \leq \mathbb{E} \left[\frac{|\{j: \theta_j \neq 0, \hat{\theta}_j \neq 0\}|}{1 \vee |\{j: \theta_j \neq 0\}|} \right] \quad (33)$$

$$\leq \mathbb{E} \left[\frac{|\{j: \theta_j \neq 0, \hat{\theta}_j \neq 0\}|}{1 \vee |\{j: \theta_j \neq 0\}|} \cdot 1_{\{|\{j: \theta_j \neq 0\}| \geq (1-\epsilon)\alpha p\}} \right] \\ + \mathbb{P}[|\{j: \theta_j \neq 0\}| < (1-\epsilon)\alpha p] \quad (34)$$

$$\leq \mathbb{E} \left[\frac{|\{j: \theta_j \neq 0, \hat{\theta}_j \neq 0\}|}{(1-\epsilon)\alpha p} \right] + \mathbb{P}[|\{j: \theta_j \neq 0\}| < (1-\epsilon)\alpha p]. \quad (35)$$

360 Therefore,⁴

$$\frac{1}{p} \sum_{j=1}^p \mathbb{P}[\hat{\theta}_j \neq 0 | \theta_j \neq 0] = \frac{1}{\alpha p} \sum_{j=1}^p \mathbb{P}[\hat{\theta}_j \neq 0, \theta_j \neq 0] \quad (36)$$

$$\geq (1-\epsilon)(1-\epsilon - o_p(1; \epsilon, \alpha)) \quad (37)$$

$$\geq 1 - 2\epsilon - o_p(1; \epsilon, \alpha). \quad (38)$$

361 By the Markov inequality, $\frac{1}{p} |\mathcal{S}_2| \leq \sqrt{2\epsilon + o_p(1; \epsilon, \alpha)}$, where we defined

$$\mathcal{S}_2 := \left\{ j: \mathbb{P}[\hat{\theta}_j = 0, \theta_j \neq 0] > \sqrt{2\epsilon + o_p(1; \epsilon, \alpha)} \right\}. \quad (39)$$

362 For $j \in [p] \setminus \mathcal{S}_2$, we have

$$\sqrt{2\epsilon + o_p(1; \epsilon, \alpha)} \geq \mathbb{P}[\hat{\theta}_j = 0, \theta_j \neq 0] \quad (40)$$

$$= \mathbb{P} \left[\hat{\theta}_j = 0 \left| \frac{1}{n} \|\bar{\mathbf{X}}_j\|^2 P_{jj} \leq 1.1, \theta_j \neq 0 \right. \right] \cdot \mathbb{P} \left[\frac{1}{n} \|\bar{\mathbf{X}}_j\|^2 P_{jj} \leq 1.1 \right] \quad (41)$$

$$\geq \mathbb{P} \left[\hat{\theta}_j = 0 \left| \frac{1}{n} \|\bar{\mathbf{X}}_j\|^2 P_{jj} \leq 1.1, \theta_j \neq 0 \right. \right] \cdot (1 - o_n(1)). \quad (42)$$

⁴ $o_p(1; \epsilon, p)$ means a sequence indexed by p , which vanishes for any fixed ϵ, p .

363 Using Neyman-Pearson's lemma, we can lower bound $\mathbb{P} \left[\hat{\theta}_j \neq 0 \mid \frac{1}{n} \|\bar{\mathbf{X}}_j\|^2 P_{jj} \leq 1.1, \theta_j = 0 \right] +$
 364 $\mathbb{P} \left[\hat{\theta}_j = 0 \mid \frac{1}{n} \|\bar{\mathbf{X}}_j\|^2 P_{jj} \leq 1.1, \theta_j \neq 0 \right]$ by

$$1 - \frac{1}{2} |\mu_1 - \mu_2| \geq 1 - \sqrt{\frac{1}{2} D(\mu_1 \| \mu_2)} \quad (43)$$

$$\geq 2Q \left(\frac{1}{2} \sqrt{\frac{1.1}{P_{jj}}} \right) \quad (44)$$

365 where μ_1 and μ_2 are one dimensional Gaussian distributions with the same variance n but differ in
 366 mean by $\sqrt{1.1n/P_{jj}}$. But, for $j \in [p] \setminus (\mathcal{S}_1 \cup \mathcal{S}_2)$, we can upper bound it by

$$\phi := \frac{1}{1 - o_n(1)} \left(\sqrt{\frac{q}{1 - \alpha}} + \sqrt{2\epsilon + o_p(1; \epsilon, \alpha)} \right), \quad (45)$$

367 and hence from (44) and (45),

$$P_{jj} \leq \frac{1.1}{(2Q^{-1}(\phi/2))^2}. \quad (46)$$

368 Thus

$$\|(P_{jj})_{j=1}^p\|_{LP} \leq \max \left\{ \frac{1.1}{(2Q^{-1}(\phi/2))^2}, \frac{1}{p} |\mathcal{S}_1 \cup \mathcal{S}_2| \right\} \quad (47)$$

$$\leq \max \left\{ \frac{1.1}{(2Q^{-1}(\phi/2))^2}, \frac{q}{1 - \alpha} + \sqrt{2\epsilon + o_p(1; \epsilon, \alpha)} \right\}. \quad (48)$$

369 \square

370 B Proof of Proposition 4

371 *Proof.* The main work is to show that τ defined in the standard distributional limit is bounded
 372 independently of Σ . Recall from [JM14, (37)] that τ satisfies the equation

$$\tau^2 = \sigma_0^2 + \frac{1}{\delta} \lim_{p \rightarrow \infty} \frac{1}{p} \mathbb{E}[\|\eta_{1/d}(\theta_0 + \tau \Sigma^{-1/2} \mathbf{Z}) - \theta_0\|_{\Sigma}^2] \quad (49)$$

373 where $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\|\mathbf{y}\|_{\Sigma} := \sqrt{\mathbf{y}^T \Sigma \mathbf{y}}$, $1/d = 1 - \|\hat{\theta}\|_0/n \geq 1 - 1/\delta$. and the proxy operator is
 374 defined by

$$\eta_{1/d}(\mathbf{y}) := \operatorname{argmin}_{\theta \in \mathbb{R}^p} \left\{ \frac{1}{2d} \|\theta - \mathbf{y}\|_{\Sigma}^2 + \lambda \|\theta\|_1 \right\}. \quad (50)$$

375 We note that the proxy operator $\eta_{1/d}$ is non-expansive in $\|\cdot\|_{\Sigma}$. Indeed, consider arbitrary $\mathbf{y}_1, \mathbf{y}_2$,
 376 and let θ_1, θ_2 be such that $0 \in \Sigma(\theta_k - \mathbf{y}_k) + \partial \mathcal{L}(\theta_k)$, $k = 1, 2$, where $\partial \mathcal{L}$ denotes the subgradient
 377 of the convex functional $\lambda d \|\cdot\|_1$. We then have

$$\Sigma(\mathbf{y}_1 - \mathbf{y}_2) \in \Sigma(\theta_1 - \theta_2) + \partial \mathcal{L}(\theta_1) - \partial \mathcal{L}(\theta_2), \quad (51)$$

378 and hence there exist $G_1 \in \partial \mathcal{L}(\theta_1)$ and $G_2 \in \partial \mathcal{L}(\theta_2)$ such that $\|\theta_1 - \theta_2\|_{\Sigma}^2 \leq \|\theta_1 - \theta_2\|_{\Sigma}^2 +$
 379 $\langle G_1 - G_2, \theta_1 - \theta_2 \rangle = \langle \theta_1 - \theta_2, \mathbf{y}_1 - \mathbf{y}_2 \rangle_{\Sigma} \leq \|\theta_1 - \theta_2\|_{\Sigma} \cdot \|\mathbf{y}_1 - \mathbf{y}_2\|_{\Sigma}$, where we used the
 380 convexity of $\mathcal{L}(\cdot)$. This shows the non-expansiveness of $\eta_{1/d}$. We now upper bound the right side of
 381 (49) by noting that

$$\begin{aligned} & \|\eta_{1/d}(\theta_0 + \tau \Sigma^{-1/2} \mathbf{Z}) - \theta_0\|_{\Sigma}^2 \\ & \leq \frac{1 + \delta}{2} \|\eta_{1/d}(\theta_0 + \tau \Sigma^{-1/2} \mathbf{Z}) - \eta_{1/d}(\theta_0)\|_{\Sigma}^2 + \frac{\delta + 1}{\delta - 1} \|\eta_{1/d}(\theta_0) - \theta_0\|_{\Sigma}^2 \end{aligned} \quad (52)$$

$$\leq \frac{1 + \delta}{2} \|\tau \Sigma^{-1/2} \mathbf{Z}\|_{\Sigma}^2 + \frac{\delta + 1}{\delta - 1} \|\eta_{1/d}(\theta_0) - \theta_0\|_{\Sigma}^2 \quad (53)$$

$$\leq \frac{1 + \delta}{2} \tau^2 p + \frac{\delta + 1}{\delta - 1} (\|\eta_{1/d}(\theta_0) - \theta_0\|_{\Sigma}^2 + 2d\lambda \|\eta_{1/d}(\theta_0)\|_1) \quad (54)$$

$$\leq \frac{1 + \delta}{2} \tau^2 p + \frac{\delta + 1}{\delta - 1} \cdot 2\lambda d \|\theta_0\|_1. \quad (55)$$

Therefore by (49),

$$\tau^2 \leq \frac{2\delta\sigma_0^2}{\delta-1} + \frac{2\lambda\delta(\delta+1)\beta}{(\delta-1)^3}. \quad (56)$$

Suppose that the algorithm selects j such that $|\hat{\theta}_j^u| < t$ as nulls, for some threshold $t \in (0, 1)$. Let $\mathcal{H}_0 := \{j: \theta_{0,j} = 0\}$ and $\mathcal{H}_1 := [p] \setminus \mathcal{H}_0$. We have

$$|\{j \in \mathcal{H}_0, |\hat{\theta}_j^u| \geq t\}| \leq \inf_{s>0} \left[|\{j \in \mathcal{H}_0: |\hat{\theta}_j^u| \geq \tau s P_{jj}^{1/2}\}| + |\{j \in \mathcal{H}_0: \tau s P_{jj}^{1/2} \geq t\}| \right]. \quad (57)$$

For fixed s independent of p , by the definition of standard distributional limit, we have, almost surely,

$$\limsup_{p \rightarrow \infty} \frac{1}{p} |\{j \in \mathcal{H}_0: |\hat{\theta}_j^u| \geq \tau s P_{jj}^{1/2}\}| \leq \mathbb{P}[|Z| \Upsilon^{1/2} \geq s \Upsilon^{1/2}] \quad (58)$$

$$= \mathbb{P}[|Z| \geq s] \quad (59)$$

where $Z \sim \mathcal{N}(0, 1)$ and Υ is a random variable whose distribution is the weak limit of the empirical distribution of $(P_{jj})_{j=1}^p$ (from the definition of the standard empirical distribution). Moreover, if $s \leq \frac{t}{\tau\sqrt{L}}$, then

$$\limsup_{p \rightarrow \infty} \frac{1}{p} |\{\tau s P_{jj}^{1/2} \geq t\}| \leq \mathbb{P}\left[\Upsilon \geq \left(\frac{t}{s\tau}\right)^2\right] \quad (60)$$

$$\leq L \quad (61)$$

almost surely, where $L := \lim_{p \rightarrow \infty} \|(P_{jj}^{(p)})_{j=1}^p\|_{LP}$. Substituting into (57), we obtain

$$\limsup_{p \rightarrow \infty} \frac{1}{p} |\{j \in \mathcal{H}_0, |\hat{\theta}_j^u| \geq t\}| \leq \mathbb{P}\left[|Z| \geq \frac{t}{\tau\sqrt{L}}\right] + L. \quad (62)$$

By the same arguments, we also have

$$\limsup_{p \rightarrow \infty} \frac{1}{p} |\{j \in \mathcal{H}_1, |\hat{\theta}_j^u| \leq t\}| \leq \limsup_{p \rightarrow \infty} \frac{1}{p} |\{j \in \mathcal{H}_1, |\hat{\theta}_j^u - \theta_{0,j}| \geq 1 - t\}| \quad (63)$$

$$\leq \mathbb{P}\left[|Z| \geq \frac{1-t}{\tau\sqrt{L}}\right] + L \quad (64)$$

almost surely. Choosing $t = 1/2$ and noting that $|\mathcal{H}_1|/p \rightarrow \alpha$ shows that we can bound

$$\limsup_{p \rightarrow \infty} FDR^{(p)} \leq 2Q\left(\frac{1}{2\tau\sqrt{L}}\right) + L \quad (65)$$

and

$$\liminf_{p \rightarrow \infty} POWER^{(p)} \geq 1 - \frac{1}{\alpha} \left[2Q\left(\frac{1}{2\tau\sqrt{L}}\right) + L \right]. \quad (66)$$

□

C Proof of Lemma 5

Proof. According to the definition of the standard distributional limit, there exists $\tau \neq 0$ such that with probability 1, the empirical distribution of $\left\{ \left(\frac{\hat{\theta}_j^u - \theta_{0,j}}{\tau}, (\underline{\Sigma}^{-1})_{jj} \right) \right\}_{j=1}^{2p}$ (which is random since \mathbf{Y} and \mathbf{X} are random) convergences weakly to the distribution of $(\Upsilon^{1/2}Z, \Upsilon)$ where $Z \sim \mathcal{N}(0, 1)$ is independent of Υ .

Since for any number t' , $W_j := |\hat{\theta}_j^u| - |\hat{\theta}_{j+d}^u| \leq -t'$ implies $|\hat{\theta}_{j+d}^u| \geq t'$, by the same steps up to (62), we have

$$\limsup_{p \rightarrow \infty} \frac{1}{2p} |\{j \in [p]: W_j \leq -t'\}| \leq 2Q(t'/\tau\sqrt{L}) + L. \quad (67)$$

401 But

$$|\{j \in \mathcal{H}_1 : W_j \geq t'\}| \geq |\mathcal{H}_1| - |\{j \in \mathcal{H}_1 : W_j \leq t'\}| \quad (68)$$

$$\geq |\mathcal{H}_1| - |\{j \in \mathcal{H}_1 : |\hat{\theta}_j^u| \leq 2t'\}| - |\{j \in \mathcal{H}_1 : |\hat{\theta}_{j+d}^u| \geq t'\}| \quad (69)$$

$$\geq |\mathcal{H}_1| - |\{j \in [p] : |\hat{\theta}_j^u - \theta_{0,j}| \geq 1 - 2t'\}| - |\{j \in [p] : |\hat{\theta}_{j+d}^u| \geq t'\}|, \quad (70)$$

402 where $\mathcal{H}_1 := \{j : \theta_{0,j} \neq 0\}$. Now again using the same steps up to (62), we conclude that almost
403 surely,

$$\liminf_{p \rightarrow \infty} \frac{1}{2p} |\{j \in \mathcal{H}_1 : W_j \geq t'\}| \geq \alpha/2 - 2Q((1-t')/\tau\sqrt{L}) - 2Q(t'/\tau\sqrt{L}) - 2L. \quad (71)$$

404 Since

$$T := \min \left\{ t : \frac{|\{j \in [p] : W_j \leq -t\}|}{|\{j \in [p] : W_j \geq t\}| \vee 1} \leq q \right\} \quad (72)$$

405 and we chose

$$\frac{2Q(t'/\tau\sqrt{L}) + L}{\alpha/2 - 2Q((1-t')/\tau\sqrt{L}) - 2Q(t'/\tau\sqrt{L}) - 2L} < q, \quad (73)$$

406 we see that almost surely, $T \leq t'$ for p large enough. Thus the number of true positives using the data
407 dependent threshold T is larger than the number of true positives using the threshold t' . The claim
408 follows by choosing $t' = 1/3$ and using (71). \square

409 D Proofs in Section 6

410 *Proof of Theorem 6.* From linear algebra, we see that a necessary and sufficient condition such that
411 (27) fulfills the positive semidefinite condition for the joint covariance matrix is that

$$2 \operatorname{diag}(\operatorname{diag}(\Sigma^{-1})) - \Sigma^{-1} \succeq 0 \quad (74)$$

412 In other words, we want the precision matrix to maintain p.s.d. after flipping the signs of the
413 off-diagonals. This is true in the diagonally dominant case.

414 Using Hammersley theorem we know that the nonzero patten of the precision matrix (inverse of the
415 covariance matrix) corresponds to the connectivity graph of the graphical model, which is a tree in
416 the current case. The claim then follows from Lemma 9 below. \square

417 **Lemma 9.** *\mathbf{P} is a square matrix and the nonzero pattern of \mathbf{P} corresponds to a forest (a union of*
418 *trees), then \mathbf{P} and $2 \operatorname{diag}(\mathbf{P}) - \mathbf{P}$ have the same set of eigenvalues.*

419 *Proof.* Assume without loss of generality that the first entry corresponds to a leaf and the second entry
420 corresponds to its unique neighbor. We can expand the determinant to check that the characteristic
421 polynomial satisfies

$$\det(\lambda \mathbf{I}_p - \mathbf{P}) = (\lambda - P_{11}) \det(\lambda \mathbf{I}_{p-1} - P_{[2:n] \times [2:p]}) - P_{12}^2 \det(\lambda \mathbf{I}_{p-2} - P_{[3:p] \times [3:p]}) \quad (75)$$

422 where p is the size of \mathbf{P} , and $P_{[2:p] \times [2:p]}$ denotes the principle submatrix of \mathbf{P} consisting of entries of
423 \mathbf{P} with indices in $\{2, \dots, p\} \times \{2, \dots, p\}$. Note that $P_{[2:p] \times [2:p]}$ and $P_{[3:p] \times [3:p]}$ also correspond to
424 forrests. By induction, we see that the off-diagonal coefficients enter the characteristic polynomial
425 only through their squares. In other words, the characteristic polynomial is unchanged after flipping
426 the signs of the off-diagonals. \square

427 *Proof of Theorem 7.* The $\{1, \dots, p\} \times \{1, \dots, p\}$ -submatrix of the precision matrix satisfies

$$(\underline{\mathbf{P}}_{[p] \times [p]})^{-1} = 2 \operatorname{diag}(\mathbf{s}) - \operatorname{diag}(\mathbf{s}) \Sigma^{-1} \operatorname{diag}(\mathbf{s}) \quad (76)$$

$$= 2 \operatorname{diag}^{-1}(\mathbf{P}) - \operatorname{diag}^{-1}(\mathbf{P}) \mathbf{P} \operatorname{diag}^{-1}(\mathbf{P}) \quad (77)$$

428 where $s_j = P_{jj}^{-1}$, $j = 1, \dots, p$ in the case of conditional expectation knockoff. Note that
429 $2 \operatorname{diag}^{-1}(\mathbf{P}) - \operatorname{diag}^{-1}(\mathbf{P}) \mathbf{P} \operatorname{diag}^{-1}(\mathbf{P})$ and $\operatorname{diag}^{-1}(\mathbf{P}) \mathbf{P} \operatorname{diag}^{-1}(\mathbf{P})$ have the same diagonals, but

430 the off-diagonals are of the opposite signs and equal absolute values. When \mathbf{P} is assumed to be
 431 associated with a tree, these two matrices have the same spectral, and in particular, have the same
 432 determinant. By the same reasoning, all their principal minors have the same determinant. Therefore
 433 the

$$\text{diag}(\mathbf{P}_{[p] \times [p]}) = \text{diag}^{-1}(\mathbf{P}) \text{diag}(\mathbf{P}^{-1}) \text{diag}^{-1}(\mathbf{P}) \quad (78)$$

$$= \text{diag}^{-1}(\mathbf{P}) \text{diag}(\mathbf{\Sigma}) \text{diag}^{-1}(\mathbf{P}). \quad (79)$$

434 □

435 *Proof of Theorem 8.* Recall that $(\mathbf{P}_{[p] \times [p]})^{-1} = 2 \text{diag}(\mathbf{s}) - \text{diag}(\mathbf{s}) \mathbf{\Sigma}^{-1} \text{diag}(\mathbf{s})$. In the case of
 436 equi-knockoff, one selects $s_j = \lambda_{\min}(\mathbf{\Sigma})$, and we have

$$\lambda_{\min}(\mathbf{\Sigma}) \mathbf{I} \preceq 2 \text{diag}(\mathbf{s}) - \text{diag}(\mathbf{s}) \mathbf{\Sigma}^{-1} \text{diag}(\mathbf{s}) \preceq \lambda_{\min}(\mathbf{\Sigma}) \mathbf{I}. \quad (80)$$

437 □

438 **E Notes on the experiments**

439 Our code is built upon the knockoff software on Emmanuel Candès's website,

440 <https://web.stanford.edu/group/candes/knockoffs/software/knockoffs/>

441 with the slight modification that W_j is computed using the unbiased coefficients (see Section 5). We
 442 hope to post the details of the changes of the codes and our simulation codes at the time of final
 443 submission.

444 It is worth mentioning that the code chooses the Lasso parameter λ via cross validation, whereas our
 445 theoretical analysis chooses any λ independent of p .