

Generating knockoffs via conditional independence

Emanuela Dreassi

Università degli Studi di Firenze, Italy
e-mail: emanuela.dreassi@unifi.it

Fabrizio Leisen

King's College London, UK
e-mail: fabrizio.leisen@gmail.com

Luca Pratelli

Accademia Navale di Livorno, Italy
e-mail: Luca_pratelli@marina.difesa.it

Pietro Rigo*

Università di Bologna, Italy
e-mail: pietro.rigo@unibo.it

Abstract: Let X be a p -variate random vector and \tilde{X} a knockoff copy of X (in the sense of [9]). A new approach for constructing \tilde{X} (henceforth, NA) has been introduced in [8]. NA has essentially three advantages: (i) To build \tilde{X} is straightforward; (ii) The joint distribution of (X, \tilde{X}) can be written in closed form; (iii) \tilde{X} is often optimal under various criteria. However, for NA to apply, X_1, \dots, X_p should be conditionally independent given some random element Z . Our first result is that any probability measure μ on \mathbb{R}^p can be approximated by a probability measure μ_0 of the form

$$\mu_0(A_1 \times \dots \times A_p) = E \left\{ \prod_{i=1}^p P(X_i \in A_i | Z) \right\}.$$

The approximation is in total variation distance when μ is absolutely continuous, and an explicit formula for μ_0 is provided. If $X \sim \mu_0$, then X_1, \dots, X_p are conditionally independent. Hence, with a negligible error, one can assume $X \sim \mu_0$ and build \tilde{X} through NA. Our second result is a characterization of the knockoffs \tilde{X} obtained via NA. It is shown that \tilde{X} is of this type if and only if the pair (X, \tilde{X}) can be extended to an infinite sequence so as to satisfy certain invariance conditions. The basic tool for proving this fact is de Finetti's theorem for partially exchangeable sequences. In addition to the quoted results, an explicit formula for the conditional distribution of \tilde{X} given X is obtained in a few cases. In one of such cases, it is assumed $X_i \in \{0, 1\}$ for all i .

MSC2020 subject classifications: 62E10, 62H05, 60E05, 62J02.

*Corresponding author.

Keywords and phrases: Approximation, conditional independence, high-dimensional regression, knockoffs, multivariate dependence, partial exchangeability, variable selection.

Received June 2023.

1. Introduction

One of the main problems, both in statistics and machine learning, is to identify the explanatory variables which are to be discarded, for they don't have a meaningful effect on the response variable. To formalize, let X_1, \dots, X_p, Y be real random variables, where Y is regarded as the response variable and X_1, \dots, X_p as the explanatory variables. A *Markov blanket* is a minimal subset $\mathcal{S} \subset \{1, \dots, p\}$ such that

$$Y \perp\!\!\!\perp (X_i : i \notin \mathcal{S}) \mid (X_i : i \in \mathcal{S}).$$

Under mild conditions, a Markov blanket \mathcal{S} exists, is unique, and $\{1, \dots, p\} \setminus \mathcal{S}$ can be written as

$$\{1, \dots, p\} \setminus \mathcal{S} = \{i : Y \perp\!\!\!\perp X_i \mid (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_p)\};$$

see e.g. [9, p. 558] and [11, p. 8]. The problem mentioned above is to identify \mathcal{S} .

To any selection procedure concerned with this problem, we can associate the false discovery rate $E(\frac{|\hat{\mathcal{S}} \setminus \mathcal{S}|}{|\hat{\mathcal{S}}|})$, where $\hat{\mathcal{S}}$ denotes the estimate of \mathcal{S} provided by the procedure. As in the Neyman-Pearson theory, those selection procedures which take the false discovery rate under control worth special attention.

One such procedure has been introduced by Barber and Candès; see [2, 3, 5, 9, 14, 19]. Let

$$X = (X_1, \dots, X_p).$$

Roughly speaking, Barber and Candès' idea is to create an auxiliary vector

$$\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_p),$$

called a *knockoff copy* of X , which is able to capture the connections among X_1, \dots, X_p . Once \tilde{X} is given, each X_i is selected/discarded based on the comparison between it and \tilde{X}_i . Intuitively, \tilde{X}_i plays the role of a control for X_i , and X_i is selected if it appears to be considerably more associated with Y than its knockoff copy \tilde{X}_i . This procedure is a recent breakthrough as regards variable selection. In addition to take the false discovery rate under control, it has other merits. In particular, it works whatever the conditional distribution of Y given X . More precisely, for the knockoff procedure to apply, *one must assign $\mathcal{L}(X)$ but is not forced to specify $\mathcal{L}(Y \mid X)$* . (Here and in the sequel, for any random elements U and V , we denote by $\mathcal{L}(U)$ and $\mathcal{L}(U \mid V)$ the probability distribution of U and the conditional distribution of U given V , respectively).

Let us make precise the conditions required to \tilde{X} . For each $i \in \{1, \dots, p\}$ and each point $x \in \mathbb{R}^{2p}$, define $f_i(x) \in \mathbb{R}^{2p}$ by swapping x_i with x_{p+i} and

leaving all other coordinates of x fixed. Then, $f_i : \mathbb{R}^{2p} \rightarrow \mathbb{R}^{2p}$ is a permutation. For instance, for $p = 2$, one obtains $f_1(x) = (x_3, x_2, x_1, x_4)$ and $f_2(x) = (x_1, x_4, x_3, x_2)$. In this notation, \tilde{X} is a *knockoff copy of X* , or merely a *knockoff*, if

$$(i) \quad f_i(X, \tilde{X}) \sim (X, \tilde{X}) \text{ for each } i \in \{1, \dots, p\} \quad \text{and} \quad (ii) \quad \tilde{X} \perp\!\!\!\perp Y \mid X.$$

For the knockoff procedure to apply, one must select $\mathcal{L}(X)$ and construct \tilde{X} . However, obtaining \tilde{X} is not easy. Condition (ii) does not create any problems, for it is automatically true whenever \tilde{X} is built based only on X , neglecting any information about Y . On the contrary, condition (i) is quite difficult to be realized. Current tractable methods to achieve (i) require conditions on $\mathcal{L}(X)$. To our knowledge, such methods are available only when X is Gaussian [9], or the set of observed nodes in a hidden Markov model [19], or conditionally independent given some random element [8] and [14]. The third condition (conditional independence) is discussed in Section 1.1 and includes the other two as special cases. There are also some universal algorithms, such as the *Sequential Conditional Independent Pairs* [9] and the *Metropolized Knockoff Sampler* [5], which are virtually able to cover any choice of $\mathcal{L}(X)$. However, these algorithms do not provide a closed formula for \tilde{X} . More importantly, they are computationally intractable as soon as $\mathcal{L}(X)$ is complex; see [5] and [14]. As a matter of fact, they work effectively only for some choices of $\mathcal{L}(X)$ (such as graphical models) but not for all. A last remark is that, even if one succeeds to build \tilde{X} , the joint distribution of the pair (X, \tilde{X}) could be unknown. This is a further shortcoming. In fact, after observing $X = x$, it would be natural to sample a value \tilde{x} for \tilde{X} from the conditional distribution $\mathcal{L}(\tilde{X} \mid X = x)$. But this is impossible if $\mathcal{L}(\tilde{X} \mid X = x)$ is unknown.

In a nutshell, the above remarks may be summarized as follows. *If X is not conditionally independent* (in the sense of Section 1.1), then:

- How to build a reasonable knockoff \tilde{X} is unknown.
- The existing numerical algorithms are computationally heavy and may fail to work.
- Even if one succeeds to build \tilde{X} , the joint distribution of the pair (X, \tilde{X}) is unknown.

1.1. A new approach to knockoffs construction

As noted above, while powerful and effective, the knockoff procedure suffers from some shortcomings due to the difficulty of building a reasonable knockoff \tilde{X} . Such shortcomings are partially overcome by a new method for constructing \tilde{X} , based on conditional independence, introduced in [8]. Similar ideas were also previously developed in [14]. Another related reference is [4]. In this section, we recall the main features of this method.

Suppose that X_1, \dots, X_p are conditionally independent given some random element Z . Denote by Θ the set where Z takes values and by γ the probability

distribution of Z . Moreover, let \mathcal{B} be the Borel σ -field on \mathbb{R} and

$$P_i(A \mid \theta) = P(X_i \in A \mid Z = \theta) \quad \text{for all } i = 1, \dots, p, \theta \in \Theta \text{ and } A \in \mathcal{B}.$$

Note that γ is a probability measure on Θ and each $P_i(\cdot \mid \theta)$ is a probability measure on \mathbb{R} . Since X_1, \dots, X_p are conditionally independent given Z ,

$$\begin{aligned} P(X_1 \in A_1, \dots, X_p \in A_p) &= E \left\{ \prod_{i=1}^p P(X_i \in A_i \mid Z) \right\} \\ &= \int_{\Theta} \prod_{i=1}^p P_i(A_i \mid \theta) \gamma(d\theta) \quad \text{for all } A_1, \dots, A_p \in \mathcal{B}. \end{aligned} \quad (1)$$

Hence, one can define a probability measure λ on \mathbb{R}^{2p} as

$$\lambda(A_1 \times \dots \times A_p \times B_1 \times \dots \times B_p) = \int_{\Theta} \prod_{i=1}^p P_i(A_i \mid \theta) P_i(B_i \mid \theta) \gamma(d\theta)$$

where $A_i \in \mathcal{B}$ and $B_i \in \mathcal{B}$ for all i . In [8, Th. 12], it is shown that any p -variate random vector \tilde{X} such that

$$\mathcal{L}(X, \tilde{X}) = \lambda$$

is a knockoff copy of X .

Thus, arguing as above, not only one builds \tilde{X} in a straightforward way but also obtains the joint distribution of (X, \tilde{X}) , namely

$$\begin{aligned} P(X_1 \in A_1, \dots, X_p \in A_p, \tilde{X}_1 \in B_1, \dots, \tilde{X}_p \in B_p) \\ = \int_{\Theta} \prod_{i=1}^p P_i(A_i \mid \theta) P_i(B_i \mid \theta) \gamma(d\theta). \end{aligned} \quad (2)$$

The price to be paid is to assign $\mathcal{L}(X)$ so as to satisfy (1). (Recall that the choice of $\mathcal{L}(X)$ is a statistician's task). But this price is not expensive for two reasons. The first one is quite practical. The probability measures satisfying (1) are flexible enough to cover most real situations. Modeling X_1, \dots, X_p as conditionally independent (given some Z) is actually reasonable in a number of practical problems. The second reason is theoretical and is based on the results of this paper. Indeed, even if (1) fails, $\mathcal{L}(X)$ can be approximated arbitrarily well by probability measures satisfying (1); see Theorems 3 and 4 below.

The previous approach has two further advantages. First, \tilde{X} is often optimal under some criterions, such as mean absolute correlation and reconstructability. This is discussed in Example 1. However, we note by now that

$$\text{cov}(X_i, \tilde{X}_i) = 0 \quad \text{if } Z \text{ is such that } E(X_i \mid Z) = 0.$$

Second, even if it is not Bayesian from the conceptual point of view, the previous approach largely exploits Bayesian tools. Hence, to construct \tilde{X} and evaluate $\mathcal{L}(X, \tilde{X})$, all the Bayesian machinery can be recovered.

To illustrate, suppose that $P_i(\cdot | \theta)$ admits a density $f_i(\cdot | \theta)$ with respect to some dominating measure λ_i . For instance, λ_i could be Lebesgue measure or counting measure. Then, the joint densities of X and (X, \tilde{X}) are, respectively,

$$h(x) = h(x_1, \dots, x_p) = \int_{\Theta} \prod_{i=1}^p f_i(x_i | \theta) \gamma(d\theta) \quad \text{and}$$

$$f(x, \tilde{x}) = f(x_1, \dots, x_p, \tilde{x}_1, \dots, \tilde{x}_p) = \int_{\Theta} \prod_{i=1}^p f_i(x_i | \theta) f_i(\tilde{x}_i | \theta) \gamma(d\theta)$$

where x and \tilde{x} denote points of \mathbb{R}^p . In turn, assuming $h(x) > 0$ for the sake of simplicity, the conditional density of \tilde{X} given $X = x$ can be written as

$$\frac{f(x, \tilde{x})}{h(x)} = \frac{\int_{\Theta} \prod_{i=1}^p f_i(x_i | \theta) f_i(\tilde{x}_i | \theta) \gamma(d\theta)}{\int_{\Theta} \prod_{i=1}^p f_i(x_i | \theta) \gamma(d\theta)}.$$

Therefore, we have an explicit formula for $\mathcal{L}(\tilde{X} | X = x)$.

In the rest of this paper, to make the exposition easier, a knockoff \tilde{X} obtained as above (i.e., a knockoff \tilde{X} satisfying equation (2)) is said to be a *conditional independence knockoff* (CIK). To highlight the connection between \tilde{X} and X , we also say that \tilde{X} is the CIK of X .

1.2. Content of this paper

This paper is basically a follow up of [8]. It consists of two results, two examples, and a numerical experiment. The results are of the theoretical type. They aim to characterize the CIKs, to show that they can be applied to virtually any real situation, and to highlight some of their optimality properties. The examples provide an explicit formula for $\mathcal{L}(\tilde{X} | X = x)$ in two (meaningful) cases: mixtures of 2-valued (or 3-valued) distributions and mixtures of centered normal distributions. In particular, the first example deals with the case $X_i \in \{0, 1\}$ for all i . Such a case is important in applications, mainly in a genetic framework. Nevertheless, apart from our example, we are not aware of any theoretical investigation of this case. Finally, in the numerical experiment, the CIKs are tested against simulated and real data.

In the sequel, for any $d \geq 1$, a probability measure on \mathbb{R}^d is called *absolutely continuous* if it admits a density with respect to Lebesgue measure on \mathbb{R}^d . Moreover, \mathcal{P} is the class of all probability measures on \mathbb{R}^p and $\mathcal{P}_0 \subset \mathcal{P}$ is the subclass consisting of those $\mu_0 \in \mathcal{P}$ of the form

$$\mu_0(A_1 \times \dots \times A_p) = \int_{\Theta} \prod_{i=1}^p P_i(A_i | \theta) \gamma(d\theta),$$

for some choice of Θ , γ and $P_i(\cdot | \theta)$ such that $P_i(\cdot | \theta)$ is absolutely continuous for all i and θ .

We next briefly describe our two results. Moreover, by means of an example, we point out some optimality properties of the CIKs.

Our first result (henceforth, R1) is that, for all $\mu \in \mathcal{P}$ and $\epsilon > 0$, there is $\mu_0 \in \mathcal{P}_0$ such that

$$d_{BL}(\mu, \mu_0) < \epsilon \quad \text{and} \quad d_{TV}(\mu, \mu_0) < \epsilon \quad \text{if } \mu \text{ is absolutely continuous.}$$

In addition, an explicit formula for μ_0 is provided. Here, d_{BL} and d_{TV} are the bounded Lipschitz metric and the total variation metric, respectively. Their definitions are recalled in Section 2.

The motivation for R1 is that, to build a CIK, one needs $\mathcal{L}(X) \in \mathcal{P}_0$. This is not guaranteed, however, since the choice of $\mathcal{L}(X)$ is not subjected to any constraint. Hence, it is natural to investigate whether $\mathcal{L}(X)$ can be at least approximated by elements of \mathcal{P}_0 . Because of R1, this is actually true. Roughly speaking, R1 aims to support \mathcal{P}_0 by showing that its elements are (approximately) able to model any real situation.

In addition to the previous motivation, R1 has also some practical utility. Suppose $\mathcal{L}(X) = \mu$ for some $\mu \in \mathcal{P}$. To fix ideas, suppose μ is absolutely continuous. If μ is arbitrary, how to build a reasonable knockoff \tilde{X} is unknown. However, given $\epsilon > 0$, there is $\mu_0 \in \mathcal{P}_0$ such that $d_{TV}(\mu, \mu_0) < \epsilon$. Such a μ_0 can be built explicitly (recall that R1 provides an explicit formula for μ_0). Denote by T a p -variate random vector such that $\mathcal{L}(T) = \mu_0$. Since $\mu_0 \in \mathcal{P}_0$, the CIK \tilde{T} of T can be obtained straightforwardly. Then,

$$d_{TV}(\mathcal{L}(\tilde{X}), \mathcal{L}(\tilde{T})) = d_{TV}(\mu, \mu_0) < \epsilon$$

for *any* knockoff copy \tilde{X} of X . Hence, by the robustness properties of the knockoff procedure [3], \tilde{T} should be a reasonable approximation of \tilde{X} .

Our second result (henceforth, R2) is a characterization of the CIKs. Let \mathcal{K} denote the class of the CIKs, that is

$$\mathcal{K} = \{ \tilde{X} : \mathcal{L}(X, \tilde{X}) \text{ admits representation (2) for some } \Theta, \gamma \text{ and } P_i(\cdot | \theta) \}.$$

Moreover, for any knockoff \tilde{X} , say that (X, \tilde{X}) is *infinitely extendable* if there is an (infinite) sequence $V = (V_1, V_2, \dots)$ such that

- $(V_1, \dots, V_{2p}) \sim (X, \tilde{X})$;
- V satisfies the same invariance condition as (X, \tilde{X}) (this condition is formalized in Section 2.2).

Then, R2 states that

$$\tilde{X} \in \mathcal{K} \quad \Leftrightarrow \quad (X, \tilde{X}) \text{ is infinitely extendable.}$$

Hence, if (X, \tilde{X}) is required to be infinitely extendable, then X *must* be conditionally independent (given some Z) and \tilde{X} *must* be the CIK of X . The proof of R2 is based on de Finetti's theorem for partially exchangeable sequences.

Based on R2, a question is whether infinite extendability of (X, \tilde{X}) is a reasonable condition. To answer, two facts are to be stressed. Firstly, by de Finetti's theorem, infinite extendability of (X, \tilde{X}) essentially amounts to conditional independence of X and \tilde{X} . Secondly, for the knockoff procedure to have a low type II error rate, it is desirable that X and \tilde{X} are "as independent as possible"; see e.g. [9, p. 563] and [20]. Now, to have X and \tilde{X} as independent as possible, a reasonable strategy is to take X and \tilde{X} conditionally independent, or equivalently to require (X, \tilde{X}) to be infinitely extendable.

Example 1 (Optimality of the CIKs). Suppose $E(X_i^2) < \infty$ and $\text{var}(X_i) > 0$ for all i . Obviously, \tilde{X} should be selected so as to make the power of the knockoff procedure as high as possible. To this end, two criteria are to minimize the *mean absolute correlation*

$$\sum_{i=1}^p \left| \frac{\text{cov}(X_i, \tilde{X}_i)}{\text{var}(X_i)} \right|,$$

and to minimize the *reconstructability index*

$$\sum_{i=1}^p E\{\text{var}(X_i | L_i)\}^{-1} \quad \text{where } L_i = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_p, \tilde{X}_1, \dots, \tilde{X}_p).$$

The first criterion (mean absolute correlation) is quite popular in the machine learning community. At least in some cases, however, it is overcome by the second (reconstructability index); see [5] and [20]. Note also that

$$\begin{aligned} E\{\text{var}(X_i | L_i)\} &= E\{E(X_i^2 | L_i) - E(X_i | L_i)^2\} \\ &= E(X_i^2) - E\{E(X_i | L_i)^2\} \leq E(X_i^2). \end{aligned}$$

Suppose now that X_1, \dots, X_p are conditionally independent, given some random element Z , and \tilde{X} is the CIK of X . Suppose also that $E(X_i | Z) = 0$ a.s. for all i . Then,

$$\text{cov}(X_i, \tilde{X}_i) = E(X_i \tilde{X}_i) = E\{E(X_i \tilde{X}_i | Z)\} = E\{E(X_i | Z) E(\tilde{X}_i | Z)\} = 0.$$

Therefore, \tilde{X} is optimal under the first criterion. Moreover,

$$E(X_i | L_i) = E\{E(X_i | Z, L_i) | L_i\} = E\{E(X_i | Z) | L_i\} = 0 \quad \text{a.s.}$$

Hence, $E\{\text{var}(X_i | L_i)\} = E(X_i^2)$ and the reconstructability index attains its minimum value $\sum_{i=1}^p E(X_i^2)^{-1}$. Therefore, \tilde{X} is optimal under the second criterion as well.

2. Theoretical results

We first recall some (well known) definitions. A function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is said to be *Lipschitz* if there is a constant $b \geq 0$ such that

$$|f(x) - f(y)| \leq b \|x - y\| \quad \text{for all } x, y \in \mathbb{R}^p,$$

where $\|\cdot\|$ is the Euclidean norm. In this case, we also say that f is b -Lipschitz or that b is a Lipschitz constant for f .

We remind that \mathcal{P} denotes the class of all probability measures on \mathbb{R}^p . The *bounded Lipschitz metric* d_{BL} and the *total variation metric* d_{TV} are two distances on \mathcal{P} . If $\mu, \nu \in \mathcal{P}$, they are defined as

$$d_{BL}(\mu, \nu) = \sup_g \left| \int_{\mathbb{R}^p} g d\mu - \int_{\mathbb{R}^p} g d\nu \right| \quad \text{and} \quad d_{TV}(\mu, \nu) = \sup_A |\mu(A) - \nu(A)|$$

where \sup_g is over the 1-Lipschitz functions $g : \mathbb{R}^p \rightarrow [-1, 1]$ and \sup_A is over the Borel subsets $A \subset \mathbb{R}^p$. Among other things, d_{BL} has the property that

$$\mu_n \rightarrow \mu \text{ weakly} \quad \Leftrightarrow \quad d_{BL}(\mu_n, \mu) \rightarrow 0$$

where $\mu_n, \mu \in \mathcal{P}$. We also note that d_{BL} and d_{TV} are connected through the inequality $d_{BL} \leq 2 d_{TV}$.

We next turn to our main results.

2.1. \mathcal{P}_0 is dense in \mathcal{P}

Let \mathcal{P}_0 be the class of those probability measures $\mu_0 \in \mathcal{P}$ which can be written as

$$\mu_0(A_1 \times \cdots \times A_p) = \int_{\Theta} \prod_{i=1}^p P_i(A_i | \theta) \gamma(d\theta), \quad A_1, \dots, A_p \in \mathcal{B},$$

for some choice of Θ , γ and $P_i(\cdot | \theta)$. To avoid trivialities, $P_i(\cdot | \theta)$ is assumed to be absolutely continuous for all $i = 1, \dots, p$ and $\theta \in \Theta$. The latter assumption is motivated by the next example.

Example 2 (Why $P_i(\cdot | \theta)$ absolutely continuous). Suppose

$$\Theta = \mathbb{R}^p, \quad \gamma = \mathcal{L}(X) \quad \text{and} \quad P_i(A | \theta) = 1_A(\theta_i) \quad (3)$$

where $\theta = (\theta_1, \dots, \theta_p)$ denotes a point of \mathbb{R}^p . Then,

$$P(X_1 \in A_1, \dots, X_p \in A_p) = \int_{\mathbb{R}^p} \prod_{i=1}^p 1_{A_i}(\theta_i) \gamma(d\theta) = \int_{\Theta} \prod_{i=1}^p P_i(A_i | \theta) \gamma(d\theta).$$

Hence, without some further constraint (such as $P_i(\cdot | \theta)$ absolutely continuous), one would obtain $\mathcal{P}_0 = \mathcal{P}$ with Θ , γ and $P_i(\cdot | \theta)$ as in (3). However, this is not practically useful. In fact, under (3), the CIK \tilde{X} of X is the trivial knockoff $\tilde{X} = X$, which is unsuitable to perform the knockoff procedure.

If $\mathcal{L}(X) \in \mathcal{P}_0$, it is straightforward to obtain the CIK \tilde{X} of X and to write $\mathcal{L}(X, \tilde{X})$ in closed form. But clearly it may be that $\mathcal{L}(X) \notin \mathcal{P}_0$. In this case, it is quite natural to investigate whether $\mathcal{L}(X)$ can be approximated by elements of \mathcal{P}_0 . This is actually possible and the approximation is very strong if $\mathcal{L}(X)$ is absolutely continuous.

Theorem 3. For all $\mu \in \mathcal{P}$ and $\epsilon > 0$, there is $\mu_0 \in \mathcal{P}_0$ such that $d_{BL}(\mu, \mu_0) < \epsilon$. In particular, one such μ_0 is

$$\mu_0(A) = \int_{\mathbb{R}^p} \mathcal{N}_p(x, cI)(A) \mu(dx) \quad \text{for all Borel sets } A \subset \mathbb{R}^p \quad (4)$$

where $c = \epsilon^2/2p$ and $\mathcal{N}_p(x, cI)$ denotes the Gaussian law on \mathbb{R}^p with mean x and covariance matrix cI , i.e.

$$\mathcal{N}_p(x, cI)(A) = (2\pi c)^{-p/2} \int_A \exp\left\{-\frac{\|y-x\|^2}{2c}\right\} dy.$$

Theorem 4. Suppose $\mu \in \mathcal{P}$ is absolutely continuous. Then, for each $\epsilon > 0$, there is $\mu_0 \in \mathcal{P}_0$ such that $d_{TV}(\mu, \mu_0) < \epsilon$. Moreover, if μ has a Lipschitz density, one such μ_0 can be defined by (4) with

$$c = \frac{1}{4p} \left(\frac{\epsilon}{bm(B)} \right)^2,$$

where b is a Lipschitz constant for the density of μ , m is the Lebesgue measure on \mathbb{R}^p and $B \subset \mathbb{R}^p$ is any Borel set satisfying $\mu(B^c) < \epsilon/2$ and $0 < m(B) < \infty$.

Theorems 3 and 4 are proved in the Supplementary Material.

It is worth noting that, in addition to (4), there are other laws $\mu_0 \in \mathcal{P}_0$ satisfying the inequalities $d_{BL}(\mu, \mu_0) < \epsilon$ or $d_{TV}(\mu, \mu_0) < \epsilon$. Moreover, in the second part of Theorem 4, the Lipschitz condition on the density of μ can be weakened at the price of making μ_0 slightly more involved.

The motivation of Theorems 3–4 has been mentioned in Section 1.2. In short, if $\mathcal{L}(X) \notin \mathcal{P}_0$, the CIK \tilde{X} of X cannot be built. However, Theorems 3–4 imply that $\mathcal{L}(X)$ can be approximated by elements of \mathcal{P}_0 . Hence, with a negligible error, it can be assumed $X \sim \mu_0$ and the CIK \tilde{X} of X can be easily obtained. This is our main motivation. However, Theorems 3–4 have a practical implication as well. Suppose $\mathcal{L}(X) = \mu$ for some $\mu \in \mathcal{P}$. To fix ideas, suppose μ is absolutely continuous with a Lipschitz density. Fix $\epsilon > 0$, define $\mu_0 \in \mathcal{P}_0$ as in Theorem 4, and call T a p -variate vector such that $\mathcal{L}(T) = \mu_0$. Since $\mu_0 \in \mathcal{P}_0$, the CIK \tilde{T} of T can be easily built. Moreover, given any knockoff copy \tilde{X} of X , since $\tilde{X} \sim X \sim \mu$ and $\tilde{T} \sim T \sim \mu_0$, Theorem 4 yields

$$d_{TV}(\mathcal{L}(\tilde{X}), \mathcal{L}(\tilde{T})) = d_{TV}(\mu, \mu_0) < \epsilon.$$

Therefore, by the robustness properties of the knockoff procedure [3], \tilde{T} is expected to be a reasonable approximation of \tilde{X} . (Obviously, the latter claim should be supported by a numerical comparison of the power and the false discovery rate corresponding to \tilde{X} and \tilde{T} . Such a comparison is not trivial, however, since \tilde{X} is unknown for arbitrary μ).

Finally, two remarks are in order. The first is summarized by the following lemma.

Lemma 5. Let \tilde{T} be the CIK of T , where $T \sim \mu_0$ with μ_0 given by (4). Then,

$$(T, \tilde{T}) \sim (L + M, L + N)$$

where L, M, N are independent, $L \sim \mu$ and $M \sim N \sim \mathcal{N}_p(0, cI)$.

Proof. For any Borel sets $A, B \subset \mathbb{R}^p$, one obtains

$$\begin{aligned} & P(L + M \in A, L + N \in B) \\ &= \int_{\mathbb{R}^p} P(x + M \in A, x + N \in B) \mu(dx) \\ &= \int_{\mathbb{R}^p} P(x + M \in A) P(x + N \in B) \mu(dx) \\ &= \int_{\mathbb{R}^p} \mathcal{N}_p(x, cI)(A) \mathcal{N}_p(x, cI)(B) \mu(dx) = P(T \in A, \tilde{T} \in B). \quad \square \end{aligned}$$

Lemma 5 makes clear the structure of $\mathcal{L}(T, \tilde{T})$ and may be useful for sampling from such distribution.

The second remark is that, if $\mathcal{L}(X, \tilde{X})$ is absolutely continuous and has a Lipschitz density, the pair (T, \tilde{T}) can be taken such that

$$d_{TV}(\mathcal{L}(X, \tilde{X}), \mathcal{L}(T, \tilde{T})) < \epsilon.$$

In the notation $\mu^* = \mathcal{L}(X, \tilde{X})$ and $\mu_0^* = \mathcal{L}(T, \tilde{T})$, it suffices to let

$$\mu_0^*(A) = \int_{\mathbb{R}^{2p}} \mathcal{N}_{2p}(x, cI)(A) \mu^*(dx) \quad \text{for all Borel sets } A \subset \mathbb{R}^{2p}$$

where c is a suitable constant. Thus, $\mathcal{L}(X, \tilde{X})$ can be approximated in total variation by $\mathcal{L}(T, \tilde{T})$ for any knockoff \tilde{X} which makes $\mathcal{L}(X, \tilde{X})$ absolutely continuous with a Lipschitz density. While this fact is theoretically meaningful and supports the CIKs further, the above formula for μ_0^* has little practical use, since μ^* is generally unknown (it is even unknown how to obtain \tilde{X}).

2.2. A characterization of the CIKs

Recall that

$$\mathcal{K} = \{ \tilde{X} : \mathcal{L}(X, \tilde{X}) \text{ admits representation (2) for some } \Theta, \gamma \text{ and } P_i(\cdot | \theta) \}$$

is the class of the CIKs of X . Such a \mathcal{K} does not include all possible knockoffs. Here is a trivial example.

Example 6 (Not every knockoff is a CIK). Suppose that X_1, \dots, X_p are i.i.d. with $P(X_1 = 0) = P(X_1 = 1) = 1/2$. In this case, it would be natural to take \tilde{X} as an independent copy of X . But suppose we let

$$\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_p) = (1 - X_1, \dots, 1 - X_p).$$

Then, for all $a, b \in \{0, 1\}^p$,

$$P(X = a, \tilde{X} = b) = P(X = a) \quad \text{if } b_i = 1 - a_i \text{ for each } i = 1, \dots, p$$

while $P(X = a, \tilde{X} = b) = 0$ otherwise. Based on this fact, it is straightforward to verify that \tilde{X} is a knockoff copy of X . However, since $X_i^2 = X_i$, one obtains

$$\begin{aligned} \text{cov}(X_i, \tilde{X}_i) &= E\{X_i(1 - X_i)\} - E(X_i)^2 \\ &= E(X_i) - E(X_i^2) - E(X_i)^2 = -E(X_i)^2 < 0. \end{aligned}$$

Now, if $\tilde{X} \in \mathcal{K}$, Jensen's inequality implies $\text{cov}(X_i, \tilde{X}_i) \geq 0$. Hence, $\tilde{X} \notin \mathcal{K}$.

Based on Example 6, a question is how to identify the members of \mathcal{K} among all possible knockoffs \tilde{X} . To answer this question, we recall that (X, \tilde{X}) is said to be infinitely extendable if there exists an (infinite) sequence $V = (V_1, V_2, \dots)$ such that $(V_1, \dots, V_{2p}) \sim (X, \tilde{X})$ and V satisfies the same invariance condition as (X, \tilde{X}) . Formally, the latter request should be meant as follows. Given three integers i, j, k with $1 \leq i \leq p$ and $j, k \geq 0$, define a new sequence $V^* = (V_1^*, V_2^*, \dots)$ by swapping V_{kp+i} with V_{jp+i} and leaving all other elements of V fixed, that is,

$$V_{kp+i}^* = V_{jp+i}, \quad V_{jp+i}^* = V_{kp+i}, \quad V_r^* = V_r \quad \text{if } r \notin \{kp+i, jp+i\}.$$

Then, V is required to satisfy

$$V^* \sim V \quad \text{for all } 1 \leq i \leq p \text{ and } j, k \geq 0. \quad (5)$$

Condition (5) is nothing but a form of partial exchangeability; see [1] and [10]. In fact, the main tool for proving the next result is de Finetti's theorem for partially exchangeable sequences.

Theorem 7. *Let \tilde{X} be a knockoff copy of X . Then, $\tilde{X} \in \mathcal{K}$ if and only if (X, \tilde{X}) is infinitely extendable.*

The essence of Theorem 7 is that, if (X, \tilde{X}) is required to be infinitely extendable, then X *must* be conditionally independent (given some Z) and \tilde{X} *must* be the CIK of X . One reason for requiring infinite extendability has been given in Section 1.2. Essentially, infinite extendability of (X, \tilde{X}) amounts to conditional independence between X and \tilde{X} , which in turn implies optimality of \tilde{X} under various criteria for increasing the power of the knockoff procedure; see Example 1.

Proof of Theorem 7. Suppose $\tilde{X} \in \mathcal{K}$, that is, $\mathcal{L}(X, \tilde{X})$ admits representation (2) for some Θ, γ and $P_i(\cdot | \theta)$. For all $A_1, A_2, \dots \in \mathcal{B}$, define

$$P(V_1 \in A_1, V_2 \in A_2, \dots) = \int_{\Theta} \prod_{k=0}^{\infty} \prod_{j=1}^p P_j(A_{kp+j} | \theta) \gamma(d\theta).$$

Then, $V = (V_1, V_2, \dots)$ is an infinite sequence satisfying condition (5). Moreover, by (2), one obtains

$$\begin{aligned} & P(V_1 \in A_1, \dots, V_p \in A_p, V_{p+1} \in B_1, \dots, V_{2p} \in B_p) \\ &= \int_{\Theta} P_1(A_1 | \theta) \dots P_p(A_p | \theta) P_1(B_1 | \theta) \dots P_p(B_p | \theta) \gamma(d\theta) \\ &= P(X_1 \in A_1, \dots, X_p \in A_p, \tilde{X}_1 \in B_1, \dots, \tilde{X}_p \in B_p) \end{aligned}$$

whenever $A_i, B_i \in \mathcal{B}$ for each i . Hence, (X, \tilde{X}) is infinitely extendable. Conversely, suppose (X, \tilde{X}) is infinitely extendable and take an infinite sequence $V = (V_1, V_2, \dots)$ satisfying condition (5) and $(V_1, \dots, V_{2p}) \sim (X, \tilde{X})$. Let \mathcal{Q} denote the set of all probability measures on \mathbb{R} . By (5), V is partially exchangeable; see e.g. [1]. Hence, by de Finetti's theorem, there is a probability measure π on \mathcal{Q}^p such that

$$P(V_1 \in A_1, V_2 \in A_2, \dots) = \int_{\mathcal{Q}^p} \prod_{k=0}^{\infty} \prod_{j=1}^p q_j(A_{kp+j}) \pi(dq_1, \dots, dq_p)$$

for all $A_1, A_2, \dots \in \mathcal{B}$; see [1] again. (Such a π is usually called the de Finetti's measure of V). In particular,

$$\begin{aligned} & P(X_1 \in A_1, \dots, X_p \in A_p, \tilde{X}_1 \in B_1, \dots, \tilde{X}_p \in B_p) \\ &= P(V_1 \in A_1, \dots, V_p \in A_p, V_{p+1} \in B_1, \dots, V_{2p} \in B_p) \\ &= \int_{\mathcal{Q}^p} q_1(A_1) \dots q_p(A_p) q_1(B_1) \dots q_p(B_p) \pi(dq_1, \dots, dq_p). \end{aligned}$$

Thus, to conclude the proof, it suffices to let

$$\Theta = \mathcal{Q}^p, \quad \gamma = \pi, \quad P_i(\cdot | \theta) = q_i(\cdot) \text{ for all } \theta = (q_1, \dots, q_p) \in \mathcal{Q}^p. \quad \square$$

3. 2-valued and 3-valued covariates

In applications, an important special case is $X_i \in \{0, 1\}$. In a genetic framework, for instance, $X_i = 0$ or $X_i = 1$ according to whether the i -th gene is absent or present. Another meaningful case is $X_i \in \{0, 1, 2\}$, where $X_i = 2$ can be given various interpretations. For instance, $X_i = 2$ could mean that the absence/presence of the i -th gene cannot be established. Despite their practical significance, to our knowledge, these cases have not received much attention, from the theoretical point of view, to date. In this section, we try to fill this gap. We aim to build a CIK \tilde{X} when X is a vector of 2-valued or 3-valued random variables.

There are obviously various cases. For instance, some covariates are 2-valued, other 3-valued, and the remaining ones have a continuous distribution function. Here, we only focus on two extreme situations: either all covariates are 2-valued or all are 3-valued.

Suppose first $X_i \in \{0, 1\}$ for all i . To build a CIK, X_1, \dots, X_p must be conditionally independent given a random parameter θ . Here, it is natural to let $\theta = (\theta_1, \dots, \theta_p)$ with $\theta_i \in (0, 1)$ regarded as the (random) probability of the event $\{X_i = 1\}$. Accordingly, X_1, \dots, X_p are assumed conditionally independent given $\theta = (\theta_1, \dots, \theta_p)$ with $P(X_i = 1 \mid \theta) = \theta_i$. In this case,

$$P(X = x) = P(X_1 = x_1, \dots, X_p = x_p) = \int_{(0,1)^p} \prod_{i=1}^p \theta_i^{x_i} (1 - \theta_i)^{1-x_i} \gamma(d\theta)$$

where γ denotes the probability distribution of θ and

$$x = (x_1, \dots, x_p) \in \{0, 1\}^p.$$

To be concrete, we also assume

$$\theta = \lambda c$$

where $\lambda \in (0, 1)$ is a random scalar and $c = (c_1, \dots, c_p) \in (0, 1)^p$ a vector of known constants. Moreover, we take λ uniformly distributed on $(0, 1)$ and we let

$$s = \sum_{i=1}^p x_i, \quad S = \{i : x_i = 1\}, \quad d_0 = 1, \quad d_j = \sum_{\substack{i_1 < \dots < i_j \\ i_1, \dots, i_j \notin S}} c_{i_1} c_{i_2} \dots c_{i_j}.$$

Then, after some algebra, one obtains

$$P(X = x) = \prod_{i \in S} c_i \int_0^1 \lambda^s \prod_{i \notin S} (1 - \lambda c_i) d\lambda = \prod_{i \in S} c_i \sum_{j=0}^{p-s} (-1)^j \frac{d_j}{j + s + 1}.$$

Similarly, we can evaluate $P(X = x, \tilde{X} = \tilde{x})$ where

$$\tilde{x} = (\tilde{x}_1, \dots, \tilde{x}_p) \in \{0, 1\}^p.$$

To this end, define

$$t = \sum_{i=1}^p \tilde{x}_i, \quad T = \{i : \tilde{x}_i = 1\}, \quad e_0 = 1, \quad e_j = \sum_{\substack{i_1 < \dots < i_j \\ i_1, \dots, i_j \notin T}} c_{i_1} c_{i_2} \dots c_{i_j}.$$

Then,

$$\begin{aligned} P(X = x, \tilde{X} = \tilde{x}) &= \prod_{i \in S} c_i \prod_{i \in T} c_i \int_0^1 \lambda^{s+t} \prod_{i \notin S} (1 - \lambda c_i) \prod_{i \notin T} (1 - \lambda c_i) d\lambda \\ &= \prod_{i \in S} c_i \prod_{i \in T} c_i \sum_{j=0}^{p-s} \sum_{k=0}^{p-t} (-1)^{j+k} \frac{d_j e_k}{j + k + s + t + 1}. \end{aligned}$$

Finally,

$$\begin{aligned} P(\tilde{X} = \tilde{x} \mid X = x) &= \frac{P(X = x, \tilde{X} = \tilde{x})}{P(X = x)} \\ &= \prod_{i \in T} c_i \left(\sum_{j=0}^{p-s} (-1)^j \frac{d_j}{j+s+1} \right)^{-1} \sum_{j=0}^{p-s} \sum_{k=0}^{p-t} (-1)^{j+k} \frac{d_j e_k}{j+k+s+t+1}. \end{aligned}$$

We now have an explicit formula for $\mathcal{L}(\tilde{X} \mid X = x)$. In a sense, this is the best we can do. In fact, after observing $X = x$, a value \tilde{x} for \tilde{X} can be drawn directly from $\mathcal{L}(\tilde{X} \mid X = x)$.

Next, suppose that $X_i \in \{0, 1, 2\}$ for all i . To deal with this case, we assume X_1, \dots, X_p conditionally independent given λ with

$$P(X_i = 0 \mid \lambda) = \lambda(1 - c_i), \quad P(X_i = 1 \mid \lambda) = \lambda c_i, \quad P(X_i = 2 \mid \lambda) = 1 - \lambda,$$

where $\lambda \in (0, 1)$ is a random scalar and $c_i \in (0, 1)$ a fixed known constant. We give λ a beta distribution with parameters $a > 0$ and $b > 0$. Moreover, for all

$$x = (x_1, \dots, x_p) \in \{0, 1, 2\}^p \quad \text{and} \quad \tilde{x} = (\tilde{x}_1, \dots, \tilde{x}_p) \in \{0, 1, 2\}^p,$$

we let

$$\begin{aligned} S_0 &= \{i : x_i = 0\}, \quad S_1 = \{i : x_i = 1\}, \quad T_0 = \{i : \tilde{x}_i = 0\}, \quad T_1 = \{i : \tilde{x}_i = 1\}, \\ m_2 &= \text{card} \{i : x_i = 2\} \quad \text{and} \quad n_2 = \text{card} \{i : \tilde{x}_i = 2\}. \end{aligned}$$

Then,

$$\begin{aligned} P(X = x) &= \int_0^1 \prod_i P(X_i = x_i \mid \lambda) \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \lambda^{a-1} (1-\lambda)^{b-1} d\lambda \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \prod_{i \in S_1} c_i \prod_{i \in S_0} (1-c_i) \int_0^1 \lambda^{a+p-m_2-1} (1-\lambda)^{b+m_2-1} d\lambda \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+p-m_2)\Gamma(b+m_2)}{\Gamma(a+b+p)} \prod_{i \in S_1} c_i \prod_{i \in S_0} (1-c_i) \end{aligned}$$

and

$$\begin{aligned} P(X = x, \tilde{X} = \tilde{x}) &= \int_0^1 \prod_i P(X_i = x_i \mid \lambda) P(X_i = \tilde{x}_i \mid \lambda) \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \times \\ &\quad \times \lambda^{a-1} (1-\lambda)^{b-1} d\lambda = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+2p-m_2-n_2)\Gamma(b+m_2+n_2)}{\Gamma(a+b+2p)} \times \\ &\quad \times \prod_{i \in S_1} c_i \prod_{i \in T_1} c_i \prod_{i \in S_0} (1-c_i) \prod_{i \in T_0} (1-c_i). \end{aligned}$$

Hence, even in this case, we have an explicit formula for $\mathcal{L}(\tilde{X} | X = x)$, that is

$$\begin{aligned} P(\tilde{X} = \tilde{x} | X = x) &= \frac{P(X = x, \tilde{X} = \tilde{x})}{P(X = x)} \\ &= \frac{\Gamma(a + 2p - m_2 - n_2) \Gamma(b + m_2 + n_2) \Gamma(a + b + p)}{\Gamma(a + b + 2p) \Gamma(a + p - m_2) \Gamma(b + m_2)} \prod_{i \in T_1} c_i \prod_{i \in T_0} (1 - c_i). \end{aligned}$$

4. Mixtures of centered normal distributions

In this section, $\theta = (\theta_1, \dots, \theta_p)$ is a vector of strictly positive random variables and X_1, \dots, X_p are conditionally independent given θ with

$$X_i | \theta \sim \mathcal{N}_1(0, \theta_i) \quad \text{for each } i = 1, \dots, p.$$

Mixtures of centered normal distributions allow to model various real situations while preserving some properties of the Gaussian laws. For this reason, they are quite popular in applications; see e.g. [15] and references therein. Among other things, they are involved in Bayesian inference for logistic models [16] and they arise as the limit laws in the CLT for exchangeable random variables [7, Sect. 3]. A further motivation for this type of data is that $E(X_i | \theta) = 0$. Hence, the CIKs are optimal and in particular

$$\text{cov}(X_i, \tilde{X}_i) = 0 \quad \text{for all } i;$$

see Example 1.

To build a CIK, a “prior” γ on $\Theta = (0, \infty)^p$ is to be selected. Quite surprisingly, to our knowledge, the choice of γ seems to be almost neglected in the Bayesian literature (apart from the special case $\theta_1 = \dots = \theta_p$); see e.g. [13]. We next propose two choices of γ . As in Section 3, we let $\theta = \lambda c$ where $\lambda > 0$ is a scalar and $c = (c_1, \dots, c_p)$ a vector such that $c_i > 0$ for all i .

4.1. First choice of γ

We first assume that λ is random but c is not. Equivalently, we suppose that the ratios $\theta_i/\theta_j = c_i/c_j$ are non-random and known. While simple, this assumption makes sense in various applications, for instance in a financial framework.

The random variable λ is given an inverse Gamma distribution with parameters $a > 0$ and $b > 0$, that is, λ has density $\psi(x) = b^a \Gamma(a)^{-1} x^{-a-1} \exp(-b/x)$ for $x > 0$. In this case, the density of (X, \tilde{X}) is

$$f(x, \tilde{x}) = \int_0^\infty \prod_{i=1}^p f_i(x_i | \lambda, c) f_i(\tilde{x}_i | \lambda, c) \psi(\lambda) d\lambda$$

where x and \tilde{x} are points of \mathbb{R}^p and $f_i(\cdot | \lambda, c)$ is the density of $\mathcal{N}_1(0, c_i \lambda)$. Hence,

$$f(x, \tilde{x}) = \frac{b^a}{(2\pi)^p \Gamma(a) \prod_{i=1}^p c_i} \int_0^\infty \lambda^{-a-p-1} \exp\left\{-\frac{1}{\lambda} \left(b + \sum_{i=1}^p \frac{x_i^2 + \tilde{x}_i^2}{2c_i}\right)\right\} d\lambda$$

$$= \frac{b^a \Gamma(a+p)}{(2\pi)^p \Gamma(a) \prod_{i=1}^p c_i} \left(b + \sum_{i=1}^p \frac{x_i^2 + \tilde{x}_i^2}{2c_i} \right)^{-(a+p)}.$$

Similarly, the density of X is

$$\begin{aligned} h(x) &= \int_0^\infty \prod_{i=1}^p f_i(x_i | \lambda, c) \psi(\lambda) d\lambda \\ &= \frac{b^a \Gamma(a+p/2)}{(2\pi)^{p/2} \Gamma(a) \sqrt{\prod_{i=1}^p c_i}} \left(b + \sum_{i=1}^p \frac{x_i^2}{2c_i} \right)^{-(a+p/2)}. \end{aligned}$$

It is worth noting that f and h are densities of Student's- t distributions. We recall that the m -variate Student's- t distribution with k degrees of freedom is the absolutely continuous distribution on \mathbb{R}^m with density

$$\varphi(x) = \frac{\Gamma[(m+k)/2]}{\Gamma(k/2) (k\pi)^{m/2} \sqrt{\det \Sigma}} (1 + (1/k) x^T \Sigma^{-1} x)^{-(m+k)/2} \quad \text{for each } x \in \mathbb{R}^m,$$

where Σ is a symmetric positive definite $m \times m$ matrix. Hence, one obtains $\varphi = f$ if $m = 2p$, $k = 2a$ and $\Sigma = b a^{-1} \text{diag}(c_1, \dots, c_p, c_1, \dots, c_p)$ and $\varphi = h$ if $m = p$, $k = 2a$ and $\Sigma = b a^{-1} \text{diag}(c_1, \dots, c_p)$.

Finally, the conditional density of \tilde{X} given $X = x$ can be written as

$$g(\tilde{x} | x) = \frac{f(x, \tilde{x})}{h(x)} = \frac{\Gamma(a+p)}{(2\pi)^{p/2} \Gamma(a+p/2) \sqrt{\prod_{i=1}^p c_i}} \frac{(b + \sum_{i=1}^p \frac{x_i^2}{2c_i})^{a+p/2}}{(b + \sum_{i=1}^p \frac{x_i^2 + \tilde{x}_i^2}{2c_i})^{a+p}}.$$

Once again, $g(\cdot | x)$ is the density of a Student's- t distribution (with parameters depending on x). To see this, it suffices to let $m = p$, $k = 2a + p$, and

$$\Sigma = \frac{2}{2a+p} \left(b + \sum_{i=1}^p \frac{x_i^2}{2c_i} \right) \text{diag}(c_1, \dots, c_p).$$

Thus, we have an explicit formula for $g(\cdot | x)$ and this is quite useful in applications. A numerical example is in Section 5.

4.2. Second choice of γ

Suppose now that c is random and independent of λ . Let c be given an absolutely continuous distribution with density q . Then, f , h and g turn into

$$\begin{aligned} f(x, \tilde{x}) &= \int_{(0, \infty)^p} \int_0^{+\infty} \prod_{i=1}^p f_i(x_i | \lambda, c) f_i(\tilde{x}_i | \lambda, c) \psi(\lambda) q(c) d\lambda dc, \\ h(x) &= \int_{(0, \infty)^p} \int_0^{+\infty} \prod_{i=1}^p f_i(x_i | \lambda, c) \psi(\lambda) q(c) d\lambda dc \quad \text{and} \quad g(\tilde{x} | x) = \frac{f(x, \tilde{x})}{h(x)}. \end{aligned}$$

As an example, c_1, \dots, c_p could be taken i.i.d. according to a uniform distribution on some bounded interval $B \subset (0, \infty)$, i.e.,

$$q(c) = \frac{1}{\text{length}(B)^p} \prod_{i=1}^p 1_B(c_i).$$

In general, the above integrals cannot be explicitly evaluated. Hence, sampling from $g(\cdot | x)$ is not easy, but it is still possible by computational methods. For instance, we could proceed as follows. Since $g(\cdot | x)$ is proportional to $f(x, \cdot)$, we focus on $f(x, \cdot)$. Then, to sample from $f(x, \cdot)$, we adopt a data augmentation strategy where λ and c are treated as auxiliary variables. The idea is to consider the density function

$$g^*(\tilde{x}, \lambda, c | x) \propto \prod_{i=1}^p f_i(x_i | \lambda, c) f_i(\tilde{x}_i | \lambda, c) \psi(\lambda) q(c)$$

and perform a Gibbs sampling on the variables (\tilde{x}, λ, c) . We conclude this section by listing the full conditional distributions required to run the algorithm.

- Let $\tilde{x}_{-i} = (\tilde{x}_1, \dots, \tilde{x}_{i-1}, \tilde{x}_{i+1}, \dots, \tilde{x}_p)$. The full conditional distribution of $(\tilde{x}_i | \tilde{x}_{-i}, \lambda, c)$ is proportional to $f_i(\cdot | \lambda, c)$. This means that \tilde{x}_i can be sampled from a centered normal distribution with variance λc_i .
- The full conditional distribution of $(\lambda | \tilde{x}, c)$ is proportional to

$$\frac{\psi(\lambda)}{\lambda^p} \exp \left\{ -\frac{1}{\lambda} \sum_{i=1}^p \frac{x_i^2 + \tilde{x}_i^2}{2c_i} \right\}.$$

Hence, since ψ is the inverse gamma density with parameters a and b , the full conditional of λ is still an inverse gamma with parameters

$$a^* = a + p \quad \text{and} \quad b^* = b + \sum_{i=1}^p \frac{x_i^2 + \tilde{x}_i^2}{2c_i}.$$

Obviously, λ could be also given a different distribution. In this case, the corresponding full conditional is probably more involved, but one may use a metropolis within Gibbs step.

- Let $c_{-i} = (c_1, \dots, c_{i-1}, c_{i+1}, \dots, c_p)$. The full conditional distribution of $(c_i | \tilde{x}, \lambda, c_{-i})$ is proportional to

$$\frac{q(c)}{c_i} \exp \left\{ -\frac{1}{c_i} \frac{x_i^2 + \tilde{x}_i^2}{2\lambda} \right\}.$$

Sampling from the above is not straightforward and may require a metropolis within Gibbs step.

5. A numerical experiment

In this section, the CIKs are tested numerically against both simulated and real data. To this end, X is assumed to be as in Section 4.1. Hence, $\mathcal{L}(X)$ is a mixture of centered normal distributions and $\theta = \lambda c$, where the scalar λ has an inverse gamma distribution with parameters a and b while $c = (c_1, \dots, c_p)$ is a vector of strictly positive constants.

To learn something about the impact of the parameters, the experiment has been repeated for various choices of a , b and c . The obtained results are quite stable with respect to a and b but exhibit a notable variability with respect to c . In the sequel, a and b have been selected so as to control the mean and the variance of λ (which hold $b/(a-1)$ and $\frac{b^2}{(a-1)^2(a-2)}$, respectively, for $a > 2$). In case of real data (Section 5.2) a and b have been also tuned based on the observed value of X . The choice of c is certainly more delicate. As in Section 4.2, one option could be modeling c as a random vector (rather than a fixed vector). For instance, c_1, \dots, c_p could be i.i.d, according to a uniform distribution on some interval, and independent of λ . However, in this section, c is taken to be non-random. This choice has essentially three motivations. First, it may be convenient in real problems, in order to account for the different roles of the various covariates. Second, it is practically simpler since computational methods are not required. Third, if c is non-random, a direct comparison with Section 6.3 of [18] is easier.

One more remark is in order. To compare different knockoff procedures, three popular criteria are the power, the false discovery rate, and the observed correlations between the X_i and their knockoffs \tilde{X}_i . However, as regards the CIKs of Section 4.1, the third criterion is superfluous, since $\text{cov}(X_i, \tilde{X}_i) = 0$ for all i . Indeed, under the third criterion (as well as under the reconstructability criterion), the CIKs of Section 4.1 are superior to any other knockoff procedure; see Example 1.

5.1. Simulated data

According to the usual format (see e.g. [9] and [18]) the simulation experiment has been performed as follows.

- A subset $I \subset \{1, \dots, p\}$ such that $|I| = 60$ has been randomly selected and the coefficients β_1, \dots, β_p have been defined as

$$\beta_i = 0 \text{ if } i \notin I \quad \text{and} \quad \beta_i = \frac{u}{\sqrt{n}} \text{ if } i \in I.$$

Here, n is a positive integer and $u > 0$ a parameter called *signal amplitude*.

- n i.i.d. observations

$$X^{(j)} = (X_{1j}, \dots, X_{pj}), \quad j = 1, \dots, n,$$

have been generated from a p -variate Student's- t distribution with $2a$ degrees of freedom and matrix $\Sigma = b a^{-1} \text{diag}(c_1, \dots, c_p)$. Given $X^{(j)}$, the

corresponding response variable $Y^{(j)}$ has been defined as

$$Y^{(j)} = \sum_{i=1}^p \beta_i X_{ij} + e_j$$

where e_1, \dots, e_n are i.i.d. standard normal errors.

- For each $j = 1, \dots, n$, we sampled m CIKs, say $\tilde{X}^{(1,j)}, \dots, \tilde{X}^{(m,j)}$, from the conditional distribution of $\tilde{X}^{(j)}$ given $X^{(j)} = x^{(j)}$ where $x^{(j)}$ is the observed value of $X^{(j)}$. Precisely, for each $k = 1, \dots, m$, the value of $\tilde{X}^{(k,j)}$ was sampled from the p -variate Student's- t distribution with $2a + p$ degrees of freedom and matrix

$$\Sigma = \frac{2}{2a + p} \left(b + \sum_{i=1}^p \frac{x_{ij}^2}{2c_i} \right) \text{diag}(c_1, \dots, c_p).$$

- For each $k = 1, \dots, m$, the knockoff selection procedure has been applied to the data

$$\{Y^{(j)}, X^{(j)}, \tilde{X}^{(k,j)} : j = 1, \dots, n\}$$

so as to calculate the power and the false discovery rate, say $pow(k)$ and $fdr(k)$. To do this, we exploited the R-cran package `knockoff`:

<https://cran.r-project.org/web/packages/knockoff/index.html>

This package is based on the comparison between the lasso coefficient estimates of each covariate and its knockoff.

- The final outputs are the arithmetic means of the powers and the false discovery rates, i.e.,

$$pow = (1/m) \sum_{k=1}^m pow(k) \quad \text{and} \quad fdr = (1/m) \sum_{k=1}^m fdr(k).$$

To run the simulation experiment, we took $m = n = p = 1000$ and a theoretical value of the false discovery rate equal to 0.1. As already noted, the experiment has been repeated for various choices of the parameters a, b, c, u . Overall, the results have been quite stable with respect to all parameters but c . The specific results reported here correspond to $a = 6, b = 10, c_i = i$ and $u = 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 1, 1.5, 2, 2.5, 3$.

The observed results, in terms of pow and fdr , are summarized in Figure 1. The performance of the CIKs appears to be excellent, even if it slightly gets worse for small values of the amplitude u . It is worth noting that, as regards the power, the behavior of the CIKs is even optimal. This was quite expected, however, because of the optimality of the CIKs discussed in Example 1.

5.2. Real data

We next turn to real data. In this case, the CIKs can be compared with some other knockoff procedures, namely: The Benjamin and Hochberg method [6],

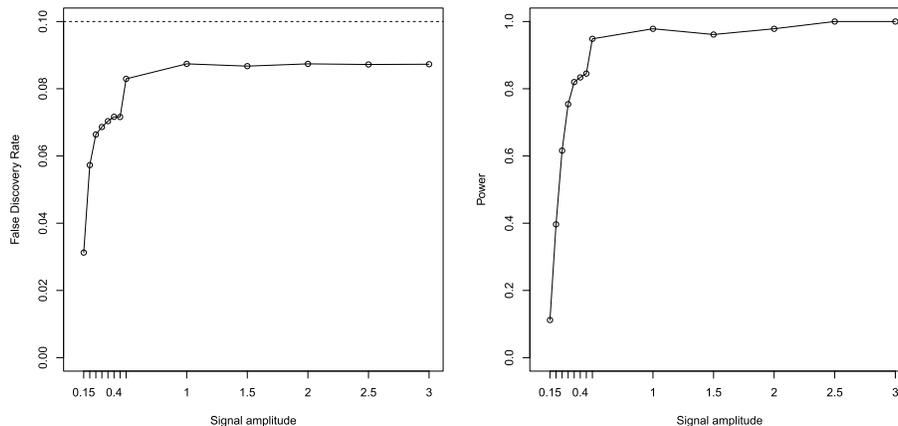


FIG 1. Results from the simulation experiment: False Discovery Rate (left) and Power (right) performances for the CIKs with Signal amplitude equal to 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 1, 1.5, 2, 2.5, 3

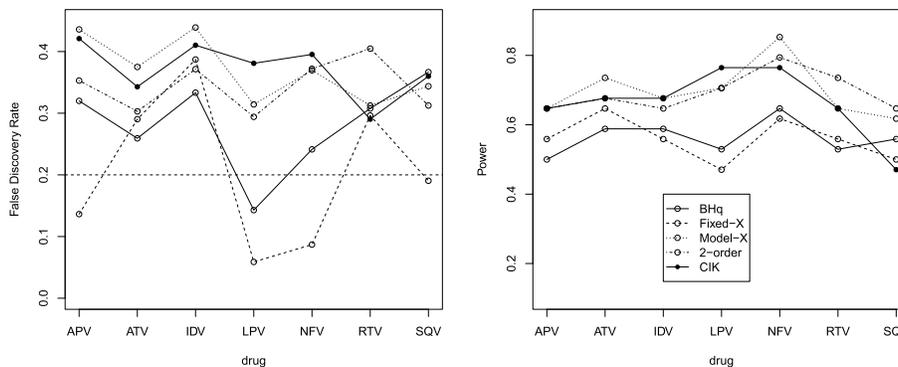


FIG 2. Results from the real example: False Discovery Rate (left) and Power (right) performances across methods and drugs.

denoted by BHq; The fixed X knockoff [2], denoted by Fixed- X ; The model- X Gaussian knockoff [9], denoted by Model- X ; The second-order knockoff [9, 18], denoted by Second-order. The comparison is based on the power, the false discovery rate, and the number of false and true discoveries. The results reported here correspond to $a = 4$, $b = 3$ and $c_i = i$.

We focus on the human immunodeficiency virus type 1 (HIV-1) dataset [17], which has been used in several papers on the knockoff procedure; see e.g. [2, 18]. The dimension of our dataset is $n = 846$ and $p = 341$, where n denotes the number of observations. The knockoff filter is applied to detect the mutations associated with drug resistance. In fact, the HIV-1 dataset provides drug resistance measurements. Furthermore, it includes genotype information from samples of

HIV-1, with separate data sets for resistance to protease inhibitors, nucleoside reverses transcriptase inhibitors, and non-nucleoside RT inhibitors. We deal with resistance to protease inhibitors, and we analyze separately the following drugs: amprenavir (APV), atazanavir (ATV), indinavir (IDV), lopinavir (LPV), nelfinavir (NFV), ritonavir (RTV) and saquinavir (SQV).

Figure 2 summarizes the performances of the five methods across different drugs in terms of power and false discovery rate. It turns out that, in most cases, the CIKs are performing well. Compared to the other procedures, the CIKs are performing better in terms of power for APV, IDV and LPV whilst are performing worse for SQV. In terms of false discovery rate, the CIKs perform better than others for RTV whilst are performing worse for LPV, NFV and SQV. Figure 3 shows the performances of the five methods for each drug related to their discoveries. We note that the number of true discoveries with the CIKs is higher compared to BHq and Fixed-X for all the drugs and similarly to Second-order and Model-X. We also highlight the performance of the CIKs in RTV with respect to the other methods.

To sum up, though the CIKs are not the best, they guarantee a good balance between power and false discovery rate and its performance is analogous to that of the other methods. For instance, as regards APV, ATV, IDV, LPV, NFV, the CIKs have a similar number of true discoveries with respect to Second-order and X-Model but also a fewer number of false discoveries.

Supplementary Material

We close the paper by proving Theorems 3 and 4. For $c > 0$ and $z \in \mathbb{R}^p$, we denote by $\phi_z(\cdot)$ the density function of $\mathcal{N}_p(z, cI)$, i.e.

$$\phi_z(x) = (2\pi c)^{-p/2} \exp\left\{-\frac{\|x-z\|^2}{2c}\right\} \quad \text{for all } x \in \mathbb{R}^p.$$

Proof of Theorem 3. Define $c = \epsilon^2/2p$ and

$$\mu_0(A) = \int_{\mathbb{R}^p} \mathcal{N}_p(x, cI)(A) \mu(dx) \quad \text{for all Borel sets } A \subset \mathbb{R}^p.$$

To see that $\mu_0 \in \mathcal{P}_0$, it suffices to let

$$\Theta = \mathbb{R}^p, \quad \gamma = \mu \quad \text{and} \quad P_i(\cdot | \theta) = \mathcal{N}_1(\theta_i, c),$$

where θ_i denotes the i -th coordinate of $\theta = (\theta_1, \dots, \theta_p) \in \mathbb{R}^p$. Obviously, $P_i(\cdot | \theta)$ is absolutely continuous. Moreover, since

$$\mathcal{N}_p(x, cI) = \mathcal{N}_1(x_1, c) \times \dots \times \mathcal{N}_1(x_p, c) \quad \text{for all } x = (x_1, \dots, x_p) \in \mathbb{R}^p,$$

one obtains

$$\mu_0(A_1 \times \dots \times A_p) = \int_{\Theta} \prod_{i=1}^p P_i(A_i | \theta) \gamma(d\theta) \quad \text{for all } A_1, \dots, A_p \in \mathcal{B}.$$

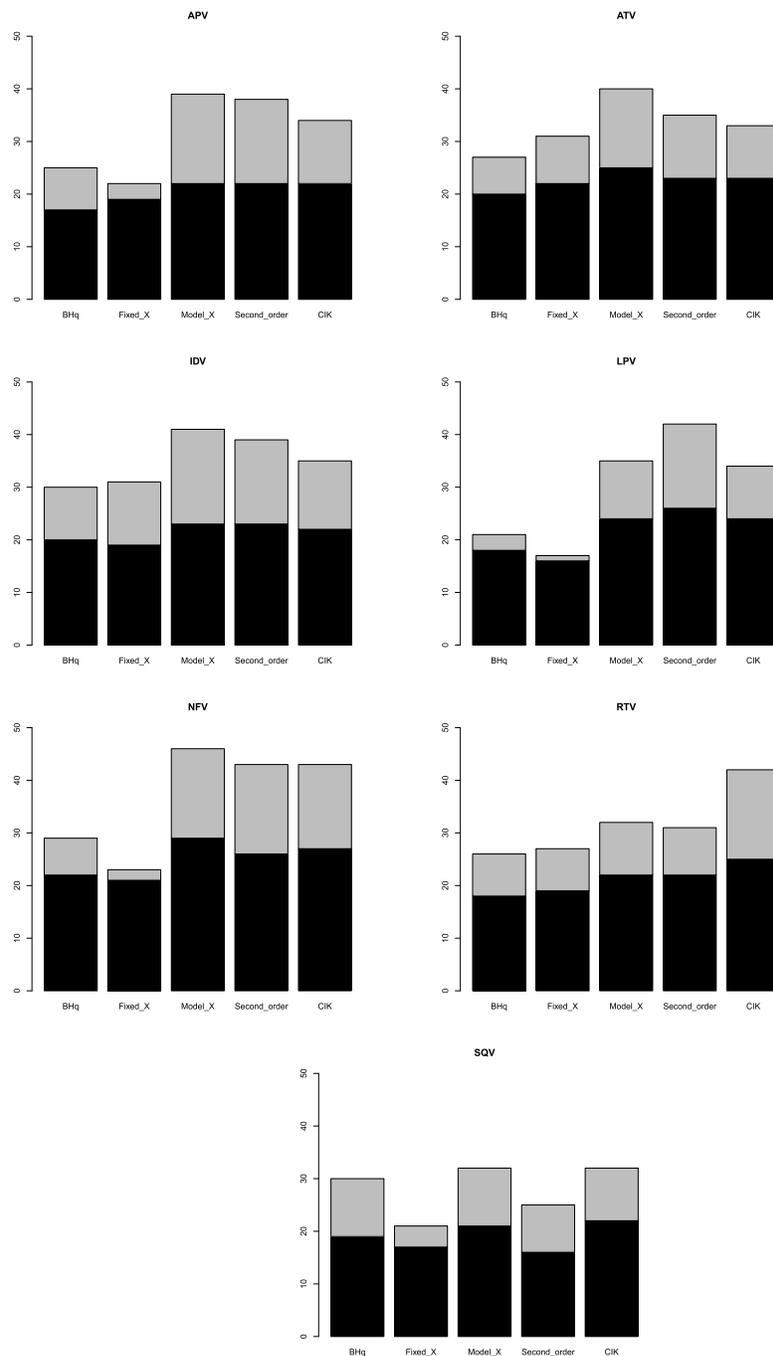


FIG 3. Comparison of different knockoff filters in terms of discoveries for each type of drug. The black part denotes the true discoveries whilst the gray part denotes the false discoveries.

We next prove $d_{BL}(\mu, \mu_0) < \epsilon$. Fix a 1-Lipschitz function $g : \mathbb{R}^p \rightarrow [-1, 1]$. Then,

$$\int_{\mathbb{R}^p} g d\mu_0 = \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} g(y) \phi_x(y) dy \mu(dx).$$

Since g is 1-Lipschitz and

$$\int_{\mathbb{R}^p} \|y - x\|^2 \phi_x(y) dy = pc,$$

it follows that

$$\begin{aligned} \left| \int_{\mathbb{R}^p} g d\mu_0 - \int_{\mathbb{R}^p} g d\mu \right| &= \left| \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} \{g(y) - g(x)\} \phi_x(y) dy \mu(dx) \right| \\ &\leq \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} |g(y) - g(x)| \phi_x(y) dy \mu(dx) \\ &\leq \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} \|y - x\| \phi_x(y) dy \mu(dx) \\ &\leq \int_{\mathbb{R}^p} \sqrt{\int_{\mathbb{R}^p} \|y - x\|^2 \phi_x(y) dy} \mu(dx) = \sqrt{pc} = \frac{\epsilon}{\sqrt{2}}. \end{aligned}$$

Therefore,

$$d_{BL}(\mu, \mu_0) = \sup_g \left| \int_{\mathbb{R}^p} g d\mu_0 - \int_{\mathbb{R}^p} g d\mu \right| \leq \frac{\epsilon}{\sqrt{2}} < \epsilon. \quad \square$$

Proof of Theorem 4. Let m denote the Lebesgue measure on \mathbb{R}^p . Suppose μ is absolutely continuous and denote by f a density of μ (with respect to m). Given $\epsilon > 0$, there is a function f_0 on \mathbb{R}^p such that:

- f_0 is a probability density (with respect to m);
- $\int_{\mathbb{R}^p} |f(x) - f_0(x)| dx < \epsilon$;
- f_0 is of the form $f_0 = \sum_{j=1}^k a_j 1_{R_j}$, where k is a positive integer, $a_j > 0$ a constant, and R_j a bounded rectangle, i.e.

$$R_j = I_{1j} \times \cdots \times I_{pj}$$

where I_{ij} is a bounded interval of the real line for each $i = 1, \dots, p$;

see e.g. Theorem (2.41) of [12, p. 69]. Define μ_0 as the probability measure on \mathbb{R}^p with density f_0 . Since μ and μ_0 are both absolutely continuous,

$$d_{TV}(\mu, \mu_0) = (1/2) \int_{\mathbb{R}^p} |f(x) - f_0(x)| dx < \epsilon/2 < \epsilon.$$

Moreover, μ_0 can be written as

$$\mu_0 = \sum_{j=1}^k a_j m(R_j) (\mathcal{U}_{1j} \times \cdots \times \mathcal{U}_{pj})$$

where \mathcal{U}_{ij} is the uniform distribution on the interval I_{ij} . Hence, letting

$$\Theta = \{1, \dots, k\}, \quad \gamma\{\theta\} = a_\theta m(R_\theta) \quad \text{and} \quad P_i(\cdot | \theta) = \mathcal{U}_{i\theta},$$

one obtains

$$\mu_0(A_1 \times \dots \times A_p) = \sum_{\theta=1}^k a_\theta m(R_\theta) \prod_{i=1}^p \mathcal{U}_{i\theta}(A_i) = \int_{\Theta} \prod_{i=1}^p P_i(A_i | \theta) \gamma(d\theta)$$

for all $A_1, \dots, A_p \in \mathcal{B}$. Hence, $\mu_0 \in \mathcal{P}_0$.

This proves the first part of the Theorem. To prove the second part, suppose f is Lipschitz and define μ_0 by (4) with $c = \frac{1}{4p} \left(\frac{\epsilon}{b m(B)}\right)^2$, where b is a Lipschitz constant for f and $B \subset \mathbb{R}^p$ a Borel set satisfying $\mu(B^c) < \epsilon/2$ and $0 < m(B) < \infty$. Since $\mu_0 \in \mathcal{P}_0$, as shown in the proof of Theorem 3, we have only to prove that $d_{TV}(\mu, \mu_0) < \epsilon$. The density f_0 of μ_0 can be written as

$$f_0(x) = \int_{\mathbb{R}^p} \phi_x(y) f(y) dy.$$

Therefore,

$$\begin{aligned} d_{TV}(\mu, \mu_0) &= \int_{\mathbb{R}^p} (f(x) - f_0(x))^+ dx \\ &\leq \int_{B^c} f(x) dx + \int_B |f(x) - f_0(x)| dx \\ &= \mu(B^c) + \int_B \left| \int_{\mathbb{R}^p} \{f(x) - f(y)\} \phi_x(y) dy \right| dx \\ &\leq \mu(B^c) + \int_B \int_{\mathbb{R}^p} |f(x) - f(y)| \phi_x(y) dy dx \\ &\leq \mu(B^c) + b \int_B \int_{\mathbb{R}^p} \|y - x\| \phi_x(y) dy dx \\ &\leq \mu(B^c) + b \int_B \sqrt{\int_{\mathbb{R}^p} \|y - x\|^2 \phi_x(y) dy} dx \\ &= \mu(B^c) + b \int_B \sqrt{p c} dx \\ &= \mu(B^c) + b m(B) \sqrt{p c} \\ &= \mu(B^c) + (\epsilon/2) < \epsilon \end{aligned}$$

where the last inequality is because $\mu(B^c) < \epsilon/2$. This concludes the proof. \square

Acknowledgments

We are grateful to Guido Consonni for a very useful conversation.

References

- [1] ALDOUS D.J. (1985). Exchangeability and related topics, in: *Ecole de Probabilites de Saint-Flour XIII*, Lect. Notes in Math., vol. 1117, Springer, Berlin. [MR0883646](#)
- [2] BARBER R.F., CANDES E.J. (2015). Controlling the false discovery rate via knockoffs, *Ann. Statist.*, **43**, 2055–2085. [MR3375876](#)
- [3] BARBER R.F., CANDES E.J., SAMWORTH R.J. (2020). Robust inference with knockoffs, *Ann. Statist.*, **48**, 1409–1431. [MR4124328](#)
- [4] BARBER R.F., JANSON L. (2022). Testing goodness-of-fit and conditional independence with approximate co-sufficient sampling, *Ann. Statist.*, **50**, 2514–2544. [MR4500617](#)
- [5] BATES S., CANDES E.J., JANSON L., WANG W. (2021). Metropolized knockoff sampling, *J.A.S.A.*, **116**, 1413–1427. [MR4309282](#)
- [6] BENJAMIN Y., HOCHBERG Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing, *J. R. Statist. Soc. B*, **57**, 289–300. [MR1325392](#)
- [7] BERTI P., PRATELLI L., RIGO P. (2004). Limit theorems for a class of identically distributed random variables, *Ann. Probab.*, **32**, 2029–2052. [MR2073184](#)
- [8] BERTI P., DREASSI E., LEISEN F., PRATELLI L., RIGO P. (2023). New perspectives on knockoffs construction, *J. Stat. Plan. Inference*, **223**, 1–14. [MR4467688](#)
- [9] CANDES E.J., FAN Y., JANSON L., LV J. (2018). Panning for gold: ‘model- X ’ knockoffs for high dimensional controlled variable selection, *J. R. Statist. Soc. B*, **80**, 551–577. [MR3798878](#)
- [10] DIACONIS P., FREEDMAN D. (1980). De Finetti’s theorem for Markov chains, *Ann. Probab.*, **8**, 115–130. [MR0556418](#)
- [11] EDWARDS D. (2000). *Introduction to Graphical Modelling*, Springer, New York. [MR1880319](#)
- [12] FOLLAND G.B. (1984). *Real Analysis*, Wiley, New York. [MR0767633](#)
- [13] GELMAN A., CARLIN J.B., STERN H.S., DUNSON D.B., VEHTARI A., RUBIN D.B. (2013). *Bayesian Data Analysis*, 3rd edition, Chapman and Hall/CRC Texts in Statistical Science, Boca Raton. [MR3235677](#)
- [14] GIMENEZ J.R., GHORBANI A., ZOU J. (2019). Knockoffs for the mass: new feature importance statistics with false discovery guarantees, in: *Proc. of the 22nd Interna. Conf. on Artificial Intelligence and Statistics 2019*, Naha, Okinawa, Japan, PMLR, vol. 89.
- [15] HINTZ E., HOFERT M., LEMIEUX C. (2021). Normal variance mixtures: distribution, density and parameter estimation, *Computat. Stat. Data Anal.*, **157**, 107175. [MR4204413](#)
- [16] POLSON N.G., SCOTT J.G., WINDLE J. (2013). Bayesian inference for logistic models using Pólya-Gamma latent variables, *J.A.S.A.*, **108**, 1339–1349. [MR3174712](#)
- [17] RHEE S.Y., TAYLOR J., WADHERA G., BEN-HUR A., BRUTLAG D.L., SHAFER R.W. (2006). Genotypic predictors of human immunodeficiency

- virus type 1 drug resistance, *Proc. Natl. Acad. Sci. USA*, **103**, 17355–17360.
- [18] ROMANO Y., SESIA M., CANDÉS E.J. (2020). Deep knockoffs, *J.A.S.A.*, **115**, 1861–1872. [MR4189763](#)
- [19] SESIA M., SABATTI C., CANDÉS E.J. (2019). Gene hunting with hidden Markov model knockoffs, *Biometrika*, **106**, 1–18. [MR3912377](#)
- [20] SPECTOR A., JANSON L. (2022). Powerful knockoffs via minimizing reconstructability, *Ann. Statist.*, **50**, 252–276. [MR4382016](#)