

# Discussion of Sesia et al. and the Knockoff Framework

BY JONATHAN D. ROSENBLATT

*Dept. of Industrial Engineering and Management,  
 Ben Gurion University of the Negev, Israel.*

johnros@bgu.ac.il

5

AND YA'ACOV RITOV

*Department of Statistics, University of Michigan, USA*

yritov@umich.edu

AND JELLE J. GOEMAN

*Department of Biomedical Data Sciences, Leiden University Medical Center, The Netherlands.*

J.J.Goeman@lumc.nl

10

## 1. ON THE MOTIVATION

The authors of Sesia et al. (2018) set out to design a procedure for variable selection with provable finite sample statistical guarantees. The *knockoff+* algorithm proposed by Sesia et al. (2018), provably controls the FDR of conditionally independent variables. Denoting with  $X := (X_j)_{j=1}^p$  and  $Y$  the predictor and outcome variables, respectively. The *false discovery proportion*, a.k.a. the *false selection proportion*, is defined as  $FDP := V/R$  where  $R$  is the number of variables selected, and  $V$  is the number of falsely selected. A false selection defined by Sesia et al. (2018) to be a selected  $X_j$  where  $Y|X_{-j}$  is independent of  $X_j$ . The knockoff algorithm of Sesia et al. (2018) provably controls the  $FDR := \mathbb{E}[FDP]$ , at some user selected magnitude,  $q$ .

15

20

The fundamental idea of the method is to generate another set of variables, the *knockoffs*, that are exchangeable with the original  $X_j$ s, only that conditionally on  $X$  they are uncorrelated with  $Y$ . The method then proceeds to compute a test statistic that captures the difference between the apparent dependence between  $Y$  and  $X_j$  and to that of  $Y$  with the knockoff  $\tilde{X}_j$ .

Crucially for our discussion: (1) The FDR is an expectation with respect to variability in  $X$  and  $Y$ , i.e., it is a *random design* guarantee. (2) The procedure is *model-free*, i.e. *non-parametric*, with respect to prediction, in that nothing is assumed on the form of  $Y|X$ . (3) The framework assumes full knowledge of  $F_X$ , i.e., the joint distribution of predictors, marginalized over  $Y$ . In application it is based on the assumption that  $F_X$  can be estimated well enough to be assumed known. (4) The method aims at good variable selection and not prediction, estimation nor ranking.

25

30

We think of the method in Sesia et al. (2018) as an adaptation of Candès et al. (2018) to genome-wide association studies (GWAS). The differences between the two: (1) The non-null variables in Candès et al. (2018) are those that belong to the minimal set that renders all others independent, i.e., the non-null is the Markov-Blanket of  $X$  on  $Y$ . In Sesia et al. (2018), the non-null variables are those with non-null partial correlations. Under some minimal conditions, that exclude perfectly dependent variables, these two sets are the same. (2) In Candès et al. (2018)  $F_X$  is multivariate Gaussian, while in Sesia et al. (2018), a hidden Markov model (HMM). An important contribution in each paper is an algorithm for sampling knockoffs from the assumed model.

35

## 2. ON THE PROBLEM SETUP

The problem setup in Candès et al. (2018) and Sesia et al. (2018) deals with random design inference, in a non-parametric model. We find this to be a very useful setup for screening problems: It is consistent with the random designs typically found in observational studies, and it avoids the very-useful-yet-controversial linearity assumption.

A more surprising component of the problem setup is the knowledge of  $F_X$ , i.e., the joint distribution of predictors, marginalized over  $Y$ . Many authors would consider this an Oracle assumption, and given the difficulty of estimating joint distributions, an unrealistic one. GWAS represents a situation in which much is known about  $F_X$ , e.g. from the HapMap project (International HapMap Consortium, 2003). Other areas of potential application include Semi-Supervised Machine Learning. In general high-dimensional data settings, however,  $F_X$  has so many parameters that estimating it may turn out to be a more difficult problem than estimating the conditional distribution of  $Y|X$ . In all cases, robustness to errors in  $F_X$  is crucial for the practical usefulness of the method, and we are happy to see the promising preliminary results of Candès et al. (2018). On the other hand, the pruning step in the GWAS example of Sesia et al. (2018) suggests that strong dependencies may adversely affect the performance of the method; both with respect to power, and with respect to the tails of the FDP distribution. We expect this matter to be further investigated in the future.

## 3. KNOCKOFFS AS PSEUDO-VARIABLES

The idea of augmenting design matrices with random variables is not new. It has been suggested many times, for the purposes of prediction, variable ranking, large-sample support recovery, etc. Tusher et al. (2001) have already proposed the idea of permuting the original variables for FDR control on selected variables. While intuitive and elegant, their algorithm did not have any provable guarantees, and implies marginal nulls and not conditional. Some more algorithms adding “fake”, “phony”, “probes” or “pseudo variables”, are reviewed in Guyon & Elisseeff (2003).

Perhaps the most similar work is that of Wu et al. (2007), which not only propose adding “pseudo-variables” for the purpose of estimating the variable selection FDR, but also require two conditions very similar to the knockoff conditions. Wu et al. (2007) require that: “(A1) real unimportant variables and phony unimportant variables have the same probability of being selected on average”, and “(A2) real important variables have the same probability of being selected whether or not phony variables are present”. These two conditions cannot be satisfied according to Wu et al. (2007), but they are clearly related to the *pairwise exchangeability* and *nullity condition* in Sesia et al. (2018) and Candès et al. (2018). One may thus view the two knockoff conditions as a satisfiable version of Wu’s A1 and A2. To the credit of Wu et al. (2007) we quote their insights, which already hint at what will be later formalized in the knockoff conditions: “Permutation produces pseudovariables that when appended to the real data create what are essentially matched pairs. To each real variable there corresponds a pseudo variable with identical sample moments and also with preservation of correlations”.

## 4. CONDITIONAL RESAMPLING

The invariance property in Sesia et al. (2018) describes so-called *null-invariant transformations* (Goeman & Solari, 2010): the joint distribution of the augmented data  $(Y, X, \tilde{X})$  is invariant under the transformation  $\text{swap}(J)$  provided that  $J$  is a set of true nulls. Known null-invariant transformations, e.g. permutations and rotations (Langsrud, 2005), existed so far only for null hypotheses about marginal association. E.g. Tusher et al. (2001) use a knockoff-like framework to test marginal, not conditional, nulls. This method was recently proven to control tail probabilities of the FDP (Hemerik & Goeman, 2018). Which null hypothesis is to be preferred? Conditional or marginal? Think of two correlated SNPs. One causal and the other is not. In a screening experiment such as GWAS, we would like both to be selected. It thus seems to us that for the purpose of screening, marginal nulls, not conditional, are desirable.

Knockoffs represent, as far as we know, the first null-invariant transformations with respect to conditional nulls rather than marginal nulls. Seeing knockoffs as null-invariants opens the way for their more classical use, e.g. using multiple random knockoffs for controlling familywise error in the manner of Westfall & Young (1993). Since familywise error control is the norm in the field of GWAS, the latter would be a worthwhile extension. An alternative approach is resampling from the conditional distribution of  $X_j$  given  $X_{-j}$ . This approach was found to be powerful in (Candès et al., 2018, Fig.3), but dismissed there for computational reasons; it was not considered again in Sesia et al. (2018). This may be a missed opportunity. Conditional resampling can be done easily in the context of hidden Markov models, where, for example, if  $Y \in \{0, 1\}$ , one can easily test for conditional independence in the  $2 \times L$  table of  $Y$  and  $X_j$ . More generally, resampling-based methods such as Westfall & Young (1993) are computationally quite efficient, requiring only a small multiple of  $\alpha^{-1}$  resampled data sets for powerful FWER control at level  $\alpha$ . Alternatively, applying early stopping (e.g. Jiang & Salzman, 2012), and tail resampling (e.g. Yu et al., 2011) can be used for additional computational efficiency.

## 5. ON SAMPLING FROM HMMs

In Sesia et al. (2018) the explanatory variables are considered as HMM: there is a latent Markov process  $Z$ , such that the observations  $X_1, \dots, X_p$  are independent given  $Z$  and  $\mathcal{L}(X_j|Z) = \mathcal{L}(X_j|Z_j)$ . The knockoff construction algorithm is based on sampling from the Markov process  $X \rightarrow Z \rightarrow \tilde{Z} \rightarrow \tilde{X}$ , so that  $(\tilde{Z}, \tilde{X})$  is a knockoff of  $(Z, X)$ . Since we only need that  $\tilde{X}$  be a knockoff of  $X$ , then it is not clear to us why the extra step of a random sample of  $\tilde{Z}$  is needed? Can we not sample directly  $X \rightarrow Z \rightarrow \tilde{X}$ ?

Is knowledge of  $F_X$  realistic for HMMs? Any process can be approximated as a weak limit of a *non-stationary* HMM. The difficulty is, of course, with the number of needed parameters. For  $K$  hidden states, and  $L$  output values, then the number of parameters is approximately  $pK(K + L - 2)$ . With  $n$  and  $p$  as in the examples, the demand that the estimated model will be considered as “known”, enforces quite as small  $K$ , and hence the HMM assumption is restrictive. An HMM is not Markovian, but it is exponentially mixing. Practically speaking,  $K$  small does not enable medial range non-monotone dependency, as one may find in GWAS due to Linkage-Disequilibrium. If  $n$  forces  $K$  small, dependencies are local, so that marginal and conditional nulls are the same, except locally. One may thus test for independence between  $Y$  and  $X_j$  and some of its neighbors, while ignoring distant  $j$ ’s.

## 6. POWER CONSIDERATIONS AND ADVERSARIAL EXAMPLES

Sesia et al. (2018) and Candès et al. (2018) provide compelling power simulations. Below we try to design some adversarial examples, where knockoffs may not be favorable. From the practitioners’ view, a “knockoff handbook” will certainly be beneficial.

A non-favorable setup for knockoff is the following. We consider Gaussian  $X$ s which are highly correlated. High correlation between the  $X$  enforces a high correlation between  $X$  and  $\tilde{X}$ . In this setup, the knockoff framework may be worse than the simplest regression with FWER control. Why? Because the power gain due to FDR control is attenuated by the power cost of the tremendous generality of the knockoff assumptions.

More formally, consider  $|\mathcal{S}| \ll p$  and  $\mathcal{N}(0, 1)$  statistics under the null. Rejecting  $q|\mathcal{S}|$  true nulls on the average means rejecting variables with test statistics greater than  $\sqrt{2 \log(p/q|\mathcal{S}|)}$  on  $Z$  scale. A naive Bonferroni correction (which we do not advocate) implies a  $\Phi^{-1}(1 - q/p) \approx \sqrt{2 \log(p/q)}$  rejection cutoff. The ratio between the cutoff of some FDR controlling procedure, and FWER control with Bonferroni, is thus  $\{1 + \log(|\mathcal{S}|)/\log(p/q)\}^{1/2}$ . If  $|\mathcal{S}| \ll p$ , then FDR has more power than FWER, but not by much. We now need to show that the price of provable FDR control of the knockoffs, may be quite large compared to classical techniques.

We observe  $Y = \beta^T X + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . The predictors are correlated in pairs:  $X \in \mathbb{R}^{2p}$ ,  $X := (W_j^T)_{j=1}^p; W_j \sim \mathcal{N}_2(0, \Sigma)$ , where  $\Sigma = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ . Candès et al. (2018) show that the knockoffs can be

given by  $(W_j, \tilde{W}_j)^T$  multivariate normal with  $Cov[W_j, \tilde{W}_j] = \Sigma - \lambda I$  for  $0 < \lambda < 2(1 - \rho)$ . We fit the model  $Y = \sum W_j^T \gamma_j + \sum \tilde{W}_j^T \tilde{\gamma}_j + \epsilon$ . Now, by the structure  $Var[(X, \tilde{X})]$ , we see that an efficient estimation of  $\gamma_1 - \tilde{\gamma}_1, \dots, \gamma_p - \tilde{\gamma}_p$  and  $\gamma_1 + \tilde{\gamma}_1, \dots, \gamma_p + \tilde{\gamma}_p$  can be achieved by estimating each of them separately, in particular  $\hat{\gamma}_j - \hat{\tilde{\gamma}}_j \sim \mathcal{N}_2(\gamma_j, 2n^{-1}\sigma^2\lambda^{-1}I)$ . Compare this to the estimate of the given linear model  $Y = \sum W_j^T \gamma_j + \epsilon$ , where  $\hat{\gamma}_j \sim \mathcal{N}_2(\gamma_j, n^{-1}\Sigma^{-1})$ . If  $\rho \approx 1$ , and  $j$  is a null variable then the tails of  $|\hat{\beta}_j| - |\hat{\tilde{\beta}}_j|$  behave like the tails of  $\hat{\beta}_j - \hat{\tilde{\beta}}_j$ . We see thus, that using the knockoff methods increases the effect size needed to achieve a given power by a factor of  $(2(1 - \rho^2)/\lambda)^{1/2} > (1 + \rho)^{1/2}$ .

Combining with the fact that in our setup the FDR and FWER rejection cutoffs are similar, we conclude this is clearly an unfavorable setup for knockoffs. Note that the regular estimator was based only on the known distribution of  $X$  and not on the validity model of  $Y$  given  $X$ .

Although this white noise toy model was described in terms of the Gaussian linear model, the idea is probably valid in much more generality. The fundamental observation being that a high correlation within  $X$  implies high correlation between  $X$  and  $\tilde{X}$ , thus between  $T_j$  and  $\tilde{T}_j$ , even though they should be consistent for different values. This may explain why in Sec.7.1, Sesia et al. (2018, Sec.7.1) the authors apply a pruning pre-processing stage, which reduces correlations.

Another way to design an adversarial example, is simply letting the tails of the null test statistics to be heavier than the alternative. The selection cutoff will be driven by the few selected null variables with heavier tails, at the risk of masking non-nulls. This is perhaps a less realistic setup than our previous, but will clearly hurt the power of the knockoff procedure.

A standard FDR procedure should have significant level  $q$  if all variables are null. The suggested knock-off methods seem to have a problematic behavior if  $q|S| < 1$ . Two methods were suggested in Barber & Candès (2015) and Candès et al. (2018). The regular knockoff method has in this situation a significant level of approximately  $1/2$ , while the knockoff+ method has a negligible power of approximately  $2^{|S|-1/q}$  no matter how strong is the signal: we should have at least  $q^{-1} - |S|$  nulls in the right tail before the first null appears on the left.

## 7. FUTURE RESEARCH

We find the knockoff framework to be quite exciting. Not because it offers solutions to all possible difficulties, but on the contrary: because it sets the stage for many important research questions. Some of the following questions are already acknowledged in Candès et al. (2018, Sec.7.2): Are error guarantees **robust** to misspecification of  $F_X$ ? What test statistics have more **power**? Knockoffs are not uniquely defined so how to best generate them? Do methods defined for other null-invariants, such as permutations, generalize to knockoffs easily? How to sample knockoffs efficiently? Does screening with knockoffs have more power than the linear model (even if miss-specified)? What other algorithms are possible assuming  $F_X$  known? What can we borrow from the semi-supervised learning literature to the knockoff setup? Are the conditional nulls of Sesia et al. (2018) preferable over marginal nulls such as in Tusher et al. (2001)?

Dai & Barber (2016), Janson & Su (2016), Chen et al. (2017b), Chen et al. (2017a), and others, have already started to explore and extend the knockoff framework of Barber & Candès (2015). We expect many such explorations in the upcoming future.

## ACKNOWLEDGEMENT

The authors thank Dr. Aldo Solari, Dr. Livio Finos, Prof. Yuekai Sun, and the students in the University of Michigan Stats 710 class for fruitful discussions leading to this manuscript, while not necessarily agreeing with its conclusions.

## REFERENCES

- BARBER, R. F. & CANDÈS, E. J. (2015). Controlling the false discovery rate via knockoffs. *The Annals of Statistics* **43**, 2055–2085. 00123. 175
- CANDÈS, E., FAN, Y., JANSON, L. & LV, J. (2018). Panning for gold: ‘model-X’ knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **80**, 551–577. 00000. 180
- CHEN, J., HOU, A. & HOU, T. Y. (2017a). A Pseudo Knockoff Filter for Correlated Features. *arXiv:1708.09305 [math, stat]* 00000. 180
- CHEN, J., HOU, A. & HOU, T. Y. (2017b). Some Analysis of the Knockoff Filter and its Variants. *arXiv:1706.03400 [math, stat]* 00001. 180
- DAI, R. & BARBER, R. (2016). The knockoff filter for FDR control in group-sparse and multitask regression. In *International Conference on Machine Learning*. 00008.
- GOEMAN, J. J. & SOLARI, A. (2010). The sequential rejection principle of familywise error control. *The Annals of Statistics*, 3782–3810. 185
- GUYON, I. & ELISSEFF, A. (2003). An Introduction to Variable and Feature Selection. *Journal of Machine Learning Research* **3**, 1157–1182. 11418.
- HEMERIK, J. & GOEMAN, J. J. (2018). False discovery proportion estimation by permutations: confidence for significance analysis of microarrays. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **80**, 137–155. 190
- INTERNATIONAL HAPMAP CONSORTIUM (2003). The international hapmap project. *Nature* **426**, 789.
- JANSON, L. & SU, W. (2016). Familywise error rate control via knockoffs. *Electronic Journal of Statistics* **10**, 960–975. 00011.
- JIANG, H. & SALZMAN, J. (2012). Statistical properties of an early stopping rule for resampling-based multiple testing. *Biometrika* **99**, 973–980. 195
- LANGSRUD, Ø. (2005). Rotation tests. *Statistics and computing* **15**, 53–60.
- SEsia, M., SABATTI, C. & CANDÈS, E. J. (2018). Gene hunting with hidden markov model knockoffs. *Biometrika* 00000.
- TUSHER, V. G., TIBSHIRANI, R. & CHU, G. (2001). Significance analysis of microarrays applied to the ionizing radiation response. *Proceedings of the National Academy of Sciences* **98**, 5116–5121. 12117. 200
- WESTFALL, P. & YOUNG, S. (1993). *Resampling-based multiple testing: Examples and methods for p-value adjustment*. Wiley-Interscience.
- WU, Y., BOOS, D. D. & STEFANSKI, L. A. (2007). Controlling Variable Selection by the Addition of Pseudovariables. *Journal of the American Statistical Association* **102**, 235–243. 205
- YU, K., LIANG, F., CIAMPA, J. & CHATTERJEE, N. (2011). Efficient p-value evaluation for resampling-based tests. *Biostatistics* **12**, 582–593.

[Received y yd YandY. Editorial decision on y yd YandY]