11) Check for updates

Comments on "A Scale-Free Approach for False Discovery Rate Control in Generalized Linear Models"

Sai Li•, Yisha Yaob, and Cu n-H ui Zha n g<

• institute of Statistics and Big Data, Renmin University of China, Beijing, China; bDepartment of Biostatistics, Yale University, New Haven, CT; < Department of Statistic, sRutgers University, Piscataway, NJ

ARTICLE HI STORY

Received May 2023 Accepted June 2023

We congratulate Chenguang Dai, Buyu Lin, Xin Xing and Jun Liu on their development of interesting methods for asymptotic FDR control and related theory in generalized linear models. We appreciate this opportunity to comment on their thought-provoking paper.

1. Introduction

Consider testing null hypotheses Hi: 0i = 0 with statistics Zj, $j \in [p]$. Let So = $U \in [p]: 0i = 0$ } be the index set for the true nulls. Suppose it is sensible to reject Hj when Zj is large. Let $\Leftrightarrow t(x) = I(x)$:::: t} denote the threshold test For tests with a common threshold level t, $R1(X[pJ) = L \Rightarrow 1 < Pt(Xj)$ is the total number of rejected hypotheses and $R_{0,1}(X[p]) = LjESo (1)$ ((x, y)) the total number of false rejections. Benjamini and Hochberg (1995) advocated to measure the Type-I error by FDR in multiple testing and proved that for independent test statistics (x, y) with a known common continuous null survival function G_0 , the adaptive threshold level

$$T = \inf \{ t : R1(\mathbb{Z}[p) :::: Go(t)p/q \}$$
 (1.1)

of the global test for the intersection null $_{1}Hj$ in Simes (1986) also controls the FDR:

FDR= IE[
$$\frac{\text{Ro.T}(Z[p])}{\text{I V RT}(Z[p])}$$
] < $\frac{1}{D}$ (1.2)

where p_0 = IS ol-Benjamini and Yekutieli (2001) relaxed the mutual independence assumption on the test statistics to the positive regression dependency on each one from the subset So (PRDS) and the strict null assumption to IF{Z :::: t} :::; G_0 (t -) with G_0 (t -) = gk/p fork E[p], E[p].

Barber and Candes (2015) introduced knockoffs to expand the realm of non-asymptotic FDR control. Statistics Z, \ldots, Z are knockoffs of $Z1, \ldots, Z_n$, if $(Zj, Zj, j \in B, Zk > Z'k'k \in Be)$ has the same joint distribution as its B-swapped version $(Zj, zj, j \in B, Zk, Z1c, k \in eB)$ for every null subset $B = S_0$. Let Z, \ldots, Z be

knoc koffs of $Z_1, ..., Zp$,

$$S_j = \text{sgn}(Z_j - Z_j), \ \sqrt[3]{40} = s_j w(Z_j, Z_j), \ (1.3)$$

with a function satisfying/w(x, x) = fw(x, x) :::: O.T he y proposed to use the threshold test $\phi_1(W)$ at the adaptive threshold level

$$T(W[p]) = \inf_{\{t: \begin{cases} 8 + RtC_{-}W[p], \\ -1 \text{ v } R1(W[p]) \end{cases}} t \cdot \begin{cases} 8 + RtC_{-}W[p], \\ -1 \text{ v } R1(W[p]) \end{cases} q, \}$$
 (1.4)

and pro ved FDR::; q for the rule with $\delta = 1$ and a related error bound for $\delta = 0$. A crucial feature of the knockoffs is the sign symmetry $\{sj,j \in So\} \sim un \text{ iform}\{-1,1\}$ So in the null set. The knockoff approach has been further developed in Candes et al. (2018) and Huang and Janson (2020).

The beautyof the above results is the non-asymptotic nature of the FDR guarantee. Still, the PRDS assumption and the knowledge of the null survival function, or the knockoffs, may not be available in practice.

Following Xing, Zhao, and Liu (2021) and Dai et al. (2022), Dai et al. (2023) advocates the use of the thresholding level (1.4) on mirror statistics. Sufficient conditions are developed for approximate FDR control in moderate- and high-dimensional generalized linear models (GLM). Mirror statistics can be constructed by data splitting (DS) when asymptotically normal estimates of 0j are available. Let ℓ and ℓ be two independent estimators of 9IP ℓ mirror statistics are defined by

$$S_j = sgn(BJ > 1)sgn (0]2 >), M_j = sj M(el J1 > 1, 10]2 > 1), (1.5)$$

with a function s atisfying/M(x, x) = /M(x,x):::: 0asthe/w(·, •) in (1.3). Suppose

$$\ddot{O}j^{(k)}/an^{D} N(\mu,n,j,1), k=1,2,$$
 (1.6)

for some unknownparameters an and /1,nJ satisfying $sgn(\mu,nJ) = sgn(0)$. For example, in Proposition 3.1 of Dai et al. (2023), $0i = \bullet j/3r$, an = a*/.../n and /1,n,j = a*0i/an. Similar to the test

statistic Wj in (1.3), Mj is approximately marginally symmetric when 0i = 0 and likely to be positive and large when 0i = ifa = 0. Th is motivates the use of the threshold rule <pt on the mirror statistics at $t = T_8(M[b])$ in (1.4). In an asymptotic analysis, a condition on the level of dependence can be imposed instead of the independence of the signs $\{S_i, j \in S_0\}$ as required in the non-asymptotic FDR control with knockoffs.

Compared with the knockoff approach, the mirror statistics are more readily available with sufficiently large samples. Compared with asymptotic FDR cont rol with the Benjamini-Hochberg rule (1.1) based on (1.6), the mirror statistics approach does not require a specification of the scale o'n. Moreover, as mentioned in Dai et al. (2023), the asymptotic symmetry $IP'\{Mj > t\} :::; IP\{Mj :::: - t\}$ for $j \in S_0$, sufficient for the mirror statistic to make sense, may kick in with a smaller sample size requirement than the asymptotic normality (1.6).

2. Relationship of Mirror Statistics to Knockoffs

As functions of data, mirror statistics are equivalent to knock off test statistics through an algebraic transformation. Consider two-sided tests for simplicity. Given estimators e&/ and el;/, let

$$zj = |OI| > + eJ^2 > 1$$
; 2, $z; = |eJ'| > + eJ' > |I| 2$, $i \in [p1]$

We have $|Q1>1 \text{ y } 1Q2> | = \frac{1}{2} + \mathbf{Z}_1, |Q1>1 \text{ j } |Q_1^2> | = \text{ iz j } - \text{ z, i,}$ and sgn@1)sgn(ey>) = sgn(Zj - Z1) with the convention sgn(0) = 0. Thus, the mirror statistics can be written as

$$Mj = \operatorname{sgn}\{0J!' \mid \operatorname{sigl} \quad 2 \mid f M(/e]I \mid | e]2) \mid 1$$

$$= \operatorname{sgn}(Zj - Z;)fw(Zj,ZJ)$$

$$= wj \qquad (2.1)$$

with f w(x, x') = /M(X + x', Ix - x'). Con versely, M(x, x') =fw(|x + x'| | / 2, |x - x'| | / 2) given a choice of $fw(\cdot, \bullet)$. For example, Mj with/M(x, x') = (|x| + |x'|)/2 matches Wj with fw(x, x') =Jx J v Jx' J and Mj with / M(x, x') = Jxx' J matche s Wj with $fw(x,x') = |x|^2 \cdot x'^2 1;$

$$sgn(Zj - Z'.)(Zj v Z'.), \quad /M(x, x') = (Ix ! + lx' J) / 2,$$

$$M \cdot - \{I \text{ sgn}(Zj - z; \dot{j}) \text{ } IZJ - z; \dot{j}1, \quad /M(x, x') = Jxx'!.$$

It can be seen that the natural choice is M(x, x') = (|x| + |x'|) / 2in the first example.

The algebraci transformation of Wj to Mj is not guaranteed to offer new statistical insight, and vice versa. For example, when Wi are defined through the Lasso path in linear regression with non-orthogonaldesigns, the interpretation of the corresponding Mj is unclear.

When e & / and e l ; / are two independent copies of a Gaussian vector $N(O-n/1,n,\lceil p \rceil, :En)$,

$$Z[p]$$
" $N(O'n/Ln,[p],:En/2)/, Z[p]$ " $N(0,:En/2)/,$

and Z[b and z 1 are independent. Thus, the use of mirror statistics is equivalent to generating a copy of noise vector and treatingitasaknockoffof ZIP₁, with the correspondence between (ML) and f wC ·) in (2.1). For example, in linear regression with a design matrix X of rank p and Gaussian noise with a

known noise level o-, Gaussian mirror (Xing, Zhao, and Liu 2021) can be constructed in one shot by generating z 1 with $:En/2 = o^{-2}(XT \times)^{-1}$. Mirror statistics M[pJ are equivalent to test statistics W [p] based on knockoffs if and only if I:n is diagonal.

In an asymptotic analysis based on (1.6), the mirror statistic and knock off methods would have a symptotically the same FDR and power due to their algebraic equivalence (2.1) when the dependence between the test statistics has only an infinitesimal impact on the operating characteristics of the tests. This conditionon dependence is typically imposed in the literature through the sparsity of the correlation of the estimates in (1.6). In this asymptopia, an advantage of mirror statistics is their availability through DS.

Asymptotic FDR control in moderately high- and highdimensional settings has been considered in Liu(2013), Xia, Cai, and Li (2018), Javanmard and Javadi (2019), and Ma, Tony Cai, and Li (2021) with the Benjamini-Hochberg (BH) rule, and in Xing, Zhao, and Liu (2021) and Dai et al. (2022) with mirror statistics. The theoretical results in Dai et al. (2023) on datasplitting-based mirror statistics in GLM further develop this direction.

3. Scale-Free FDR Control

For asymptotic FDR control, the unknown scaling factor o-* in Sur and Candes (2019) and Proposition 3.1 of Dai et al. (2023) is no longer a sticking issue with DS. Under (1.6), scale-free asymptotic FDR control can be achieved by the BH rule with familiarStudent's t-statistics. Still, mirror statistics may relyless heavily on the asymptotic normality.

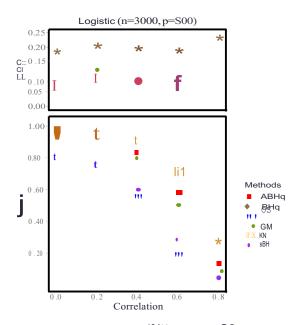
When the asymptotic normality holds with a common asymptotic variance as in (1.6) and pis large, the BH rule (1.1) can be used with the two-sided test statistics

$$Z_{j} = \frac{\left|\widehat{\theta}_{j}^{(1)} + \widehat{\theta}_{j}^{(2)}\right|}{\left(\sum_{i=1}^{p} \left(\widehat{\theta}_{i}^{(1)} - \widehat{\theta}_{i}^{(2)}\right)^{2} / (p-1)\right)^{1/2}}$$
(3.1)

and the standard absolute Gaussian survival function Go(t) = $2 < 1 \le 0$ We compare this procedure (3.1), denoted as "8"n- BH': with

DS and Gaussian mirror (GM) in Dai et al. (2023), the BH rule on the MLE (BHq), BH rule with adjusted p-values(ABHq) in Sur and Candes (2019), and model-X knockoff (KN) in Candes et al. (2018). Specifically, Figure 1 reports the simulation comparison among the above six procedures for feature selection in logistic regression in the same setting as in Figure 3 of Dai et al. (2023). Each point in Figure 1 represents the average of 50 independent replications with N(0, I:) iid designs, where :Eij = r-lij l, q = 0.1, n = 3000, p = 500, 11.B* llo = 50, and 1.Bl1 = signal for $j \in Sg$. It can be seen from these results that '7n-BH quite consistently exhibits FDR :::; $qp_0/p = 0.09$, slightly more power than DS, and less power than ABHq, GM and KN. We note that ABHq, GM and KN rely more heavilyon Gaussian

In the case of p > n, debiased Lasso exhibits heteroscedasticity (Candes et al. 2018; Dai et al. 2023). In the presence of heteroscedasticity,0 $\frac{1}{N}$ / o-n,j g $N(\mu, m, 1)$, scale if we asymptotic FDR control can be achieved with the absolute Cauchy



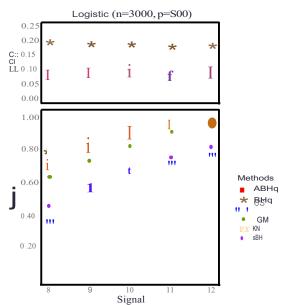


Figure 1. Left panel: the signal strength is fixed at /fJt I = 11 for /f E S8, and the correlations among the covariates vary. Right panel: the correlation is fixed at /f E S8, and the signal strength varies. Here "sBH" stands for an -BH.

ing $G_0(t) = 1 - (2/n)$ arctan(t). However, our simulation experiment demonstrates that a cost of this robustnessagainst heteroscedasticity is a significant loss of power.

4. Multiple Data Splitting

As knockoff and mirror statistics are randomized, it would be interesting to find methods to aggregate multiple randomized tests based on them for further improvement. See for example the discussion in section 7.2.3 of Candes et al. (2018). In Dai et al. (2022), the authors proposed a multiple data splitting (MDS) strategy to achieve this goal.

Similar to Dai et al. (2022), Proposition 3.3 seems to focus on special regimes aimed at the nearly complete FDR control in the form of FDR = o(1). This can be seen as follows.

Let $J = E^*[JU \in S]/(IS i \vee 1)]$ for a randomized selector $S \subset J$; p] where E^* denotes avera&ng over many copies θ ! random d S given data. If the selector S controls the FDP = $IS \cap S_0 I$ (1S1 V 1), I:tSo I: S q + op(1). Let $I(I):S \cdot \cdot \cdot :S I(p)$ be the ordered entries of I[p]. Define I(m) by

$$J(I) + \cdots +_{I(m)} : S q < I(I) + \cdots +_{I(m+I)}.$$
 (4.1)

Dai et al. (2023) proposes to test H_j by thresholding I_j at level $I_{(m)}$ and provides sufficient conditions for the asymptotic FDP control by such MDS schemes. The following lemma provides a sketch of a more explicit version of the proofs in Dai et al. (2022, 2023).

Lemma 1. Let $S1 = [p] \setminus So, P1 = P$ - Po, and I[p] be any vector of nonnegative statistics satisfying $I \neq I$ f = I and $I \neq I$ $f \in So_{I}$: $SI \neq I$ and L $I \neq I$ Let $I \neq I$ be as in (4.1). In the event where $I \neq I$ and L $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$ $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$ $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$ $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$ and L $I \neq I$ $I \neq I$

$$FDP = \frac{L \ iESo \ I\{I; > 1_{(m)}\}}{p - m}$$

where $E_i = E_i/(q + E_1)$.

In the nonsparse case, Proposition 3.3 of Dai et al. (2023) provides FDP = op(1) because the assumptions imply E₁ o(1), E₂ = o(1) and κ 1 is bounded away from zero in (4.2).

5. Fast Data Splitting Methods

In sparse GLM where Data = (X,y) E JRnx(p+I), a debiased estimator of $f3j^*$ can be written as """" (7J,uj,Data) where 7J is an initial estimate of /3* and u_j is an estimate of the direction u_j for the least favorable submodel for the estimation of f3l- f_1 and f_2 estimated by the Lasso, the computational cost of $f3/l_1$, given Data is

$$Cost(-"(.Bww, 'Data), j \in [p]) \times pCost(Lasso).$$

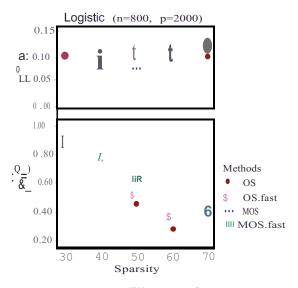
We ignore the computational cost of estimating scaling factors such as the $\not<$ and o_h in the discussion below (1.6) as they could be viewed as byproducts of the debiased Lasso program. Let Data(1,1, o_h), o_h o_h

Under proper sparsity and regularity conditions, the debiased Lasso provides

""",
$$Ut$$
", $Data$) - $id_{3}(f3^*, Uj, Data) = op(n 1^2)$.

so that the asymptotic normality is valid with a singleuj¹⁰(Data) foreach j based on the entire data in DS and MDS. We propose such procedures as fast DS and MDS algorithms in Table 1.

OS.fast needs to run Lasso p + 2 times while DS 2(p + 1) times. The computational cost of MOS.fast is of the order (T +



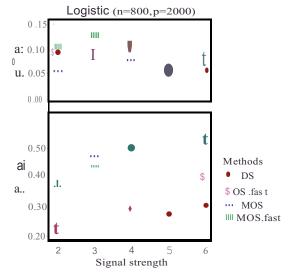


Figure 2. Left panel: the signal strength is 1/3/1 = 4 for j = 5, and the signal strength varies. Right panel: the sparsity is fixed at 60 and the signal strength varies.

Table 1. Fast DS and MOS algo rithms.

	DS.fast(Da ta)
1.	. Computeuj = u -""(Data),j $E[p]$.
2.	Generate Data <f>, Data <2> bysplitting data at random.</f>
3.	For $j \in [p]$, compute $M' \in \mathbb{R}$
	(0) , uj, Data (0) and respective scaling factors, $\mathbf{k} = 1, 2$.
4.	Compute 8 ¹), $\mathfrak{F}^{\{2\}}$ an $dM_j = M_j(\mathfrak{D}^{\{1\}})$, $\mathfrak{g}^{\{2\}}$, \mathfrak{e} in $d(F_{0r})$.

MDS.fast(Data,U[p]) with input U[p] as inStep 1 ofDS.fast

- 1. For $e \in en$, genera to Data <1,1, Data <2.1> by splitting data at random.
- 2. For $j \in [p]$, compute " ("fi'--"" (Data $\langle k,l \rangle$), $\langle Uj,Data \langle k,l \rangle$) and respective scaling factors, k=1,2
- 3. Compute $eJ^{1,1}$ $if^{2}l$) and $My \ge Mi[eJl,l]$ $8i^{2}$,l) if nd(Forl 1)
- 4. Compute $\Phi \& (M-1), s(I) \text{ and } 11) = IU E s(I > J; JW > I, JE [p], end(Fore)$
- 5. Compute/j = $I:t, U / T, j \in [p], a \cdot nd/(m)$ as in (4.1).
- 6. Compute a nd out put $/\{/>/(m)\}$, $j \in [p]$.

Compute and output $\diamondsuit \mathbb{Z}$ (M[p]),

p) x Cost(Lasso) while the cost of MDS is of the order $T \times p \times Cost(Lasso)$.

Figure 2 reports simulation performance of DS,OS.fast,MDS and MOS.fast in the settings of Figure 5 in Dai et al. (2023). Each point is averaged based on 50 independen t replications. The results demonstrate compara ble performance between DS and OS.fast and that between MDS and MOS.fast.

Supplementary Materials

In the supplement, we provide a proof of Lemma 1 and the R code for the simulation results reported in Figures 1 and 2.

Disclosure Statement

The authors report there are no competing interests to declare.

Funding

Li's research is supported in part by the National Natural Science Foundation of China (grant no. 12201630.) Yao's research is supported in part by

U.S. National Institutes of Health grants ROIHGO10171 and R01MH116527 and National Science Foundation grant DMS-2112711. Zhang's research is supported in part by National Science Foundation grants CCF-1934924, DMS2052949 and DMS2210850

References

Barber, R. F.,and Candes, E. J.(2015), "Controlling the False Discovery Rate via Knockoffs," *The Annals of statistics*, 43, 2055 - 2085. (1586]

Benjam in i, Y., and Hochberg, Y. (1995), "Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing," *Journal of the Royal Statistical Society*, Series B, 57, 289-300. [1586]

Benjamini, Y., and Yekutieli, D. (2001), "The Control of the False Discovery Rate in Multiple Testing Under Dependency;' *Annals of statistics*, 29, 1165-1188. (1586)

Candes, E., Fan, Y., Janson, L. and Lv, J. (2018), "Panning for Gold: Model-x Knockoffs for High-Dimensional Controlled Variable Selection," *Journal of the Royal Statistical Society*, Series B, 80, 551-577. [1586,1587,1588)

Dai, C., Lin, B., Xing, X., and Liu, J. S. (2022), "False Discovery RateContro 1 via Data Splitting," *Journal of the American Statistical Association*, 1-18, DOI: 10.1080/01621459.2022.2060113. [1586,1587,1588)

 - - (2023), "AScale-Free Approach for False Discovery RateControl in Generalized Linear Models; Journal of the American Statistical Association, 118,1551-1565. [15861587,1588,1589)

Huang, D., and Janson, L. (2020), "Relaxing the Assumptions of Knockoffs by Conditioning," *The Annals of Statistics*, 48, 3021-3042. (1586)

Javanmard, A., and Javadi, H. (2019), "False Discovery Rate Control via Debiased Lasso," *Electronic Journal of Statistics*, 13, 1212 - 1253. (1587)

Liu, W.(2013), "Gaussian Graphical Model Estimation with False Discovery Rate Control," *The Annals of Statistics*, **41**, 2948-2978. [1587]

Ma, R., Tony Cai, T., and Li, H. (2021), "Global and Simultaneous HypothesisTesting for High-Dimensional Logistic Regression Models," *Journal of the American Statistical Association*, 116, 984-998. [1587)

Simes, R. J. (1986),"An Improved Bonferroni Procedure for Multiple Tests of Significance; *Biometrika*, 73, 751-754. [1586)

Sur, P., and Candes, E.J. (2019), "A Modern Maximum-LikelihoodTheory for High-Dimensional Logistic Regression; Proceedings of the National Academy of Sciences, 116, 14516-14525. [1587)

Xia, Y., Cai, T. T., and Li, H. (2018), "Joint Testing and False Discovery Rate Control in High-Dimensional Multivariate Regression: *Biometrika*, 105, 249-269. [Is 87]

Xing, X., Zhao, Z., and Liu, J.S. (2021), "Controlling False Discovery Rate Using Gaussian Mirrors: Journal of the American Statistical Association, 118, 222-241. (1586, 1587)