# Some Contributions to Modern Multiple Hypothesis Testing in High-dimension

Binh Nguyen

Inria Parietal, CEA, Laboratoire de Mathématiques d'Orsay, EDMH, Université Paris-Saclay

December 10, 2021

Ph.D. Advisors:   Sylvain Arlot *(Université Paris-Saclay)*
                  Bertrand Thirion *(INRIA Parietal, CEA)*

Ph.D. Committee:  Jelle Goeman *(Leiden University – Reviewer)*
                  Etienne Roquain *(LPSM, Sorbonne Université – Reviewer)*
                  Jeanette Mumford *(Stanford University)*
                  Claire Boyer *(LPSM, Sorbonne Université)*
                  Christophe Giraud *(Université Paris-Saclay)*

# Outline

Motivation

Aggregation of Multiple Knockoffs

A Conditional Randomization Test for High-dimensional Logistic Regression

Conclusions & Perspectives

# Outline

# Reproducibility Crisis: on Popular Media...

# Most Discoveries Might Be False (Ioannidis, 2005)

## Naive Hypothesis Testing

▶ $p = 100,000$ hypotheses (brain voxels), only $2,000$ are important.

▶ Testing at $5\%$ significant level, assume all important variables are selected:

$$\text{False Discovery Proportion} = \frac{5\% \times 98,000}{2000 + 5\% \times 98,000} \approx 70\%$$

# Most Discoveries Might Be False (Ioannidis, 2005)

## Naive Hypothesis Testing

▶ $p = 100,000$ hypotheses (brain voxels), only $2,000$ are important.

▶ Testing at $5\%$ significant level, assume all important variables are selected:

$$\text{False Discovery Proportion} = \frac{5\% \times 98,000}{2000 + 5\% \times 98,000} \approx 70\%$$

## False Discovery Rate (Benjamini and Hochberg, 1995)

▶ False Discovery Rate: the average number of *false discoveries* made among all discoveries.

▶ FDR control is less conservative than Family-Wise Error Rate control

# Marginal Inference

- $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$. Example: X is MRI data, y outcome
- Linear Model
$$y = X\boldsymbol{\beta}^0 + \sigma\boldsymbol{\xi},$$
  with $\sigma > 0$, $\boldsymbol{\xi} \sim \mathcal{N}(0, I_n)$
- Support set $\mathcal{S} \triangleq \left\{ j \in [p] \,\middle|\, \beta_i^0 \neq 0 \right\}$;
- Objective: find $\hat{\mathcal{S}} \subset \mathcal{S}$ as large as possible

# Marginal Inference

- $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$. Example: X is MRI data, y outcome
- Linear Model
$$y = X\boldsymbol{\beta}^0 + \sigma\boldsymbol{\xi},$$
  with $\sigma > 0$, $\boldsymbol{\xi} \sim \mathcal{N}(0, I_n)$
- Support set $\mathcal{S} \triangleq \left\{ j \in [p] \,\middle|\, \beta_i^0 \neq 0 \right\}$;
- Objective: find $\hat{\mathcal{S}} \subset \mathcal{S}$ as large as possible

## Marginal Testing

For each $j = 1, \ldots p$:

$$\text{(null) } \mathcal{H}_0^j : X_{*,j} \perp y \quad \text{vs.} \quad \text{(alternative) } \mathcal{H}_\alpha^j : X_{*,j} \not\perp y$$

$\longrightarrow$ FDR control: easy, solvable problem (Poldrack et al., 2012)

# Conditional Inference

## Conditional Independence Testing

Generalized Linear Model (GLM): $\mathrm{y} = g(\mathrm{X}\boldsymbol{\beta^0}) + \sigma\boldsymbol{\xi}$

Testing variable $j$ but also taking interaction with other variables $\mathrm{X}_{-j}$

(null) $\mathcal{H}_0^j : X_{*,j} \perp y \mid \mathrm{X}_{-j}$    vs.    (alternative) $\mathcal{H}_\alpha^j : X_{*,j} \not\perp y \mid \mathrm{X}_{-j}$,

or, equivalently

(null) $\mathcal{H}_0^j : \beta_j^0 = 0$    vs.    (alternative) $\mathcal{H}_\alpha^j : \beta_j^0 \neq 0$.

# FDR control with Conditional Inference

Conditional inference is challenging in <u>high-dimensional settings</u>: how to obtain statistical guarantee: p-value, confidence interval?

$$\longrightarrow \text{FDR controlling?}$$

[1] Barber and Candès (2015); Candès et al. (2018)

# FDR control with Conditional Inference

Conditional inference is challenging in <u>high-dimensional settings</u>: how to obtain statistical guarantee: p-value, confidence interval?

$$\longrightarrow \text{FDR controlling?}$$

## Knockoff Inference [1]

*State-of-the-art* in high-dimension conditional inference with guaranteed FDR control

---

[1] Barber and Candès (2015); Candès et al. (2018)

# Knockoff Inference

## Knockoff variables (Candès et al., 2018)

$\tilde{X} = (\tilde{x}_1, \ldots, \tilde{x}_p)$ is model-X knockoff variables of $X = (x_1, \ldots, x_p)$ iff:

1. For all subset $\mathcal{K} \subset \{1, \ldots, p\}$: $(X, \tilde{X})_{\mathrm{swap}(\mathcal{K})} \overset{d}{=} (X, \tilde{X})$
2. $\tilde{X} \perp y \mid X$



Knockoff variables: *noisy copies* of original variables

# Knockoff Inference

## Step 1 – Model-X Knockoff

Assuming distribution of $X$ is known, construct knockoff variables, concatenate $[\mathrm{X}, \tilde{\mathrm{X}}] \in \mathbb{R}^{n \times 2p}$

## Step 2

Calculate knockoff test-statistics W: *Lasso coefficient-difference*, obtain
$$\hat{\boldsymbol{\beta}} = \min_{\mathbf{w} \in \mathbb{R}^{2p}} \frac{1}{2} \|\mathbf{y} - [\mathrm{X}, \tilde{\mathrm{X}}]\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$$

then take the difference: $W_j = |\hat{\beta}_j(\lambda)| - |\hat{\beta}_{j+p}(\lambda)|$ for each $j$

# Knockoff Inference

## Step 3 – FDR control threshold

For given $t > 0$, False Discoveries Proportion can be estimated as:

$$\widehat{\text{FDP}}(t) = \frac{1 + \#\{j \in [p] \mid W_j \leq -t\}}{\#\{j \in [p] \mid W_j \geq t\} \vee 1}$$

then, for FDR level $\alpha \in (0, 1)$, calculate the threshold

$$\tau = \min\left\{ t > 0 \mid \widehat{\text{FDP}}(t) \leq \alpha \right\}$$

## Step 4

Select the variables: $\hat{S}(\tau) = \{j \in [p] \mid W_j \geq \tau\}$

# FDP estimation with Knockoff Statistic



0

Figure: Knockoff Statistic Distribution

# FDP estimation with Knockoff Statistic



$\hat{\mathcal{S}} = \{j : W_j \geq t\}$

Figure: Knockoff Statistic Distribution

# FDP estimation with Knockoff Statistic



Figure: Knockoff Statistic Distribution

# FDP estimation with Knockoff Statistic



Figure: Knockoff Statistic Distribution

Candès et al. (2018, Lemma 3.3): Under $\mathcal{H}_0^j : \beta_j^0 = 0$, the distribution of $W_j$ is symmetric around 0, *i.e.* $(W_j, -W_k)$ are exchangeable.

# Knockoff Inference: Theoretical Guarantee

**Theorem (Barber and Candès, 2015; Candès et al., 2018)**

$$\text{FDR}(\tau) = \mathbb{E}\left[\frac{|\hat{\mathcal{S}}(\tau) \cap \mathcal{S}^c|}{|\hat{\mathcal{S}}(\tau)| \vee 1}\right] \leq \alpha,$$

where $\mathcal{S}^c = [p] \backslash \mathcal{S}$: set of null index.

► Result is non-asymptotic.

► Model-X assumption: distribution of $X$ is known.

► Proof: using martingale theory (optional stopping time theorem).

# Knockoff Inference: Theoretical Guarantee

**Theorem (Barber and Candès, 2015; Candès et al., 2018)**

$$\mathrm{FDR}(\tau) = \mathbb{E}\left[\frac{|\hat{\mathcal{S}}(\tau) \cap \mathcal{S}^c|}{|\hat{\mathcal{S}}(\tau)| \vee 1}\right] \leq \alpha,$$

where $\mathcal{S}^c = [p] \backslash \mathcal{S}$: set of null index.

▶ Result is non-asymptotic.

▶ Model-X assumption: distribution of $X$ is known.

▶ Proof: using martingale theory (optional stopping time theorem).

⚠ Major issue: inference results are random.

# Demonstration: Instability of Knockoff Procedure

$$\mathbf{y} = \mathbf{X}\beta^{\mathbf{0}} + \sigma\xi$$

$\rho$    sparsity    snr

- ▶ $n = 500$ , $p = 1000$
- ▶ $\mathrm{X} \sim \mathcal{N}(0, \boldsymbol{\Sigma})$
- ▶ $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & \rho & \rho^2 & \dots & \rho^{p-1} \\ \rho^1 & 1 & \rho & \dots & \rho^{p-2} \\ \vdots & \dots & \ddots & \dots & \vdots \\ \rho^{p-2} & \rho^{p-3} & \dots & 1 & \rho \\ \rho^{p-1} & \rho^{p-2} & \rho^{p-3} & \dots & 1 \end{bmatrix}$ , with $\rho \in [0, 1)$
- ▶ $\xi \sim \mathcal{N}(0, \mathrm{I}_n)$
- ▶ $\mathtt{sparsity} = \dfrac{|\mathcal{S}|}{p}$

# Demonstration: Instability of Knockoff Procedure



Figure: 100 runs of knockoff inference on the <u>same simulated dataset</u>
n=500, p=1000, snr=3.0, $\rho = 0.7$, sparsity = 0.06

⚠ Large variance on both FDP and Power

# Outline

# Proposed Solution: Knockoff Statistics conversion



$$\widehat{\mathrm{FDP}}(t) = \frac{1 + \#\{j : W_j \le -t\}}{\#\{j : W_j \ge t\}}$$

$N^o$ False Positives $\approx$
$\#\{j : W_j \le -t\}$

$\hat{\mathcal{S}} = \{j : W_j \ge t\}$

-t      0      t

# Proposed Solution: Knockoff Statistics conversion



Introduce the intermediate p-values: convert Knockoff statistic $W_j$ to $\hat{p}_j$:

$$\hat{p}_j = \begin{cases} \dfrac{1 + \#\{k : W_k \leq -W_j\}}{p} & \text{if} \quad W_j > 0 \\ 1 \quad \text{if} \quad W_j \leq 0 \end{cases}$$

# AKO – Aggregation of Multiple Knockoffs

- ▶ Running multiple sampling of knockoffs, find knockoff statistics
- ▶ Convert knockoff statistics to intermediate p-values
- ▶ Quantile-aggregation of p-values (Meinshausen et al., 2009)

N., Chevalier, Thirion & Arlot (2020)

# AKO – Aggregation of Multiple Knockoffs

▶ Running multiple sampling of knockoffs, find knockoff statistics

▶ Convert knockoff statistics to intermediate p-values

▶ Quantile-aggregation of p-values (Meinshausen et al., 2009)

## Step 1: For $b = 1, 2, \ldots, B$:

▶ Run knockoff sampling, calculate test statistic $\{W_j^{(b)}\}_{j=1}^{p}$

▶ Convert the test statistic $W_j^{(b)}$ to $\hat{p}_j^{(b)}$:

$$\hat{p}_j^{(b)} = \begin{cases} \dfrac{1 + \#\{k : W_k^{(b)} \leq -W_j^{(b)}\}}{p} & \text{if} \quad W_j^{(b)} > 0 \\ 1 \quad \text{if} \quad W_j \leq 0 \end{cases}$$

N., Chevalier, Thirion & Arlot (2020)

# AKO – Aggregation of Multiple Knockoffs



## Step 2 – P-values Aggregation (Meinshausen et al., 2009)

$$\bar{p}_j = \min\left\{1, \gamma^{-1} q_\gamma(\hat{p}_j^{(b)})\right\} \quad \forall j \in [p]$$

For $\gamma \in (0, 1)$ with $q_\gamma(\cdot)$ the empirical $\gamma$-quantile function.

N., Chevalier, Thirion & Arlot (2020)

# AKO – Aggregation of Multiple Knockoffs

## Step 3 – FDR control with $\{\bar{p}_j\}_{j=1}^p$

▶ Order $\bar{p}_j$ ascendingly: $\bar{p}_{(1)} < \bar{p}_{(2)} \cdots < \bar{p}_{(p)}$

▶ Given FDR control level $\alpha \in (0, 1)$, find largest $k$ such that:
  ▶ $\bar{p}_{(k)} \leq k\alpha/p$ (Benjamini and Hochberg, 1995), or
  ▶ $\bar{p}_{(k)} \leq \dfrac{k\alpha}{p \sum_{i=1}^p 1/i}$ (Benjamini and Yekutieli, 2001)

  $\longrightarrow$ FDR threshold: $\tau = \bar{p}_{(k)}$

## Step 4 – Estimate $\hat{\mathcal{S}}$

▶ $\hat{\mathcal{S}}_{\text{AKO}} = \{j \in [p] \mid \bar{p}_j \leq \tau\}$

---

N., Chevalier, Thirion & Arlot (2020)

# Theoretical Results for AKO

**Assumption (Null Distribution of Knockoff Statistic)**

*The null knockoff statistics $(W_j)_{j \in S^c}$ are* i.i.d.

**Lemma**

*Under the above assumption, and furthermore assume $|S^c| \geq 2$, for all $j \in S^c$ the intermediate p-value $\hat{p}_j$ satisfies*

$$\forall t \in (0, 1) : \quad \mathbb{P}(\hat{p}_j \leq t) \leq \frac{p}{|S^c|} t$$

**Remark**

An improved version of Lemma 2, N., Chevalier, Thirion & Arlot (2020).

# Theoretical Results for AKO

## Theorem (Finite-sample guarantee of FDR control)

*Assuming the null knockoff statistics $(W_j)_{j \in \mathcal{S}^c}$ are i.i.d. , and $|\mathcal{S}^c| \geq 2$, then for an arbitrary number of samplings $B$, the output $\hat{\mathcal{S}}_{AKO}$ of Aggregation of Multiple Knockoff (AKO) controls FDR under predefined level $\alpha \in (0, 1)$, i.e.*

$$\mathbb{E}\left[\frac{|\hat{\mathcal{S}}_{AKO} \cap \mathcal{S}^c|}{|\hat{\mathcal{S}}_{AKO}| \vee 1}\right] \leq \alpha$$

## Remark

▶ An improved version of Theorem 1, N., Chevalier, Thirion & Arlot (2020).

▶ AKO with $B = 1$ is equivalent to KO.

# Experimental Results - Synthetic Data



Histogram of FDP & Power under the <u>same simulated dataset</u>:

▶ 2500 runs of Original Knockoff (KO – top)

▶ 100 runs of Aggregated Knockoff (AKO, $B = 25$ – bottom)

# Experimental Results - Synthetic Data

▶ Vary each of the three simulation parameters while keeping the others fixed

▶ Benchmarking methods:
  - ▶ *Ours: Aggregation of Multiple Knockoffs (AKO)*
  - ▶ Vanilla Knockoff (KO) (Barber and Candès, 2015; Candès et al., 2018)
  - ▶ Related knockoff aggregation methods: Holden and Helton (2018) (KO-HL), Emery and Keich (2019) (KO-EK), Gimenez and Zou (2019) (KO-GZ)
  - ▶ Debiased Lasso (DL-BH) (Javanmard and Javadi, 2019)

# Experimental Results - Synthetic Data

▶ Vary each of the three simulation parameters while keeping the others fixed



Figure: 100 runs with varying simulation parameters. Default: SNR = 3.0, $\rho = 0.5$, sparsity = 0.06. FDR is controlled at level $\alpha = 0.1$.

# Experimental Results - Brain Imaging

▶ Data: Human Connectome Project

▶ Objective: predict the experimental condition per task given brain activity

▶ $n = 900$ subjects, $p \approx 212000$

▶ Preprocessing: dimension reduction by clustering

$$p = 212000 \longrightarrow p = 1000$$



Figure: Detection of significant brain regions for HCP data – Emotion task (face vs. shape) (900 subjects)

▶ FDR control at $\alpha = 0.1$.

▶ Orange: brain areas with positive weight.

▶ Blue: brain areas with negative weight.

# Experimental Results - Brain Imaging



Figure: Jaccard index measuring the Jaccard similarity between the KO/AKO solutions and the Debiased Lasso (DL) solution over 7 tasks of HCP900.

# Outline

# Binary classification with logistic relationship

▶ *Binary* response vector $y \in \{0, 1\}^n$.

▶ Logistic relationship

$$\mathbb{P}(y_i = 1 \mid X_{i,*}) = \frac{1}{1 + \exp(-X_{i,*}^T \boldsymbol{\beta}^0)}.$$

▶ Estimate $\beta^0$ with Penalized Logistic Regression:

$$\hat{\boldsymbol{\beta}}^{\text{PEN}} = \text{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^{n} \log \left[ 1 + \exp(-y_i (X_{i,*}^T \boldsymbol{\beta})) \right] + \lambda \left\| \boldsymbol{\beta} \right\|_1.$$

# Penalized Logistic Regression

$$\hat{\boldsymbol{\beta}}^{\text{PEN}} = \text{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^{n} \log \left[ 1 + \exp(-y_i(\mathrm{X}_{i,*}^T \boldsymbol{\beta})) \right] + \lambda \left\| \boldsymbol{\beta} \right\|_1 .$$

▶ When $n < p$: hard problem (Sur and Candès, 2019; Zhao et al., 2020)
   $\longrightarrow$ P-value? Confidence interval? Conditional Independence Testing?
▶ Original Knockoff Inference: possible with $\ell_1$-logistic loss.

# Penalized Logistic Regression

$$\hat{\boldsymbol{\beta}}^{\text{PEN}} = \text{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^{n} \log \left[ 1 + \exp(-y_i (\text{X}_{i,*}^T \boldsymbol{\beta})) \right] + \lambda \left\| \boldsymbol{\beta} \right\|_1 .$$

▶ When $n < p$: hard problem (Sur and Candès, 2019; Zhao et al., 2020)
$\longrightarrow$ P-value? Confidence interval? Conditional Independence Testing?

▶ Original Knockoff Inference: possible with $\ell_1$-logistic loss.

## Conditional Randomization Test (CRT)

Candès et al. (2018): An alternative, more straight-forward method to knockoff inference.

# Conditional Randomization Test (CRT)

**Algorithm 1: Conditional Randomization Test**

1   INPUT dataset $(X, y)$, with $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$, number of sampling runs $B$, test statistic $T_j$, conditional distribution $P_{j|-j}$ for each $j = 1, \ldots, p$ ;

2   OUTPUT vector of p-values $\{\hat{p}_j\}_{j=1}^p$;

3   **for** $j = 1, 2, \ldots, p$ **do**

4     **for** $b = 1, 2, \ldots, B$ **do**

5       1. Generate $\tilde{X}_{*,j}^{(b)}$, a noisy variable from $P_{j|-j}$;

6       2. Compute test statistics $T_j$ for original variable and $\tilde{T}_j^{(b)}$ for noisy variables;

7     **end**

8     Compute the empirical p-value

$$\hat{p}_j = \frac{1 + \sum_{b=1}^{B} \mathbb{1}_{\{\tilde{T}_j^{(b)} \geq T_j\}}}{1 + B}$$

9   **end**

# Conditional Randomization Test (CRT)

⚠ Huge computational cost: $B$ inferences for *each* variable $j$
$\longrightarrow \mathcal{O}(Bp^4)$ with Lasso program to compute $T_j$

Distillation Conditional Randomization Test (Liu et al., 2020):
analytical formula for p-values

▶ Remove the multiple sampling of noisy variables.

▶ Pre-screening step: estimate $\hat{\mathcal{S}}^{\mathrm{SCREENING}} \subset [p]$, only calculate test-statistics inside this set.

# Distillation Conditional Randomization Test (dCRT)

## Algorithm 2: Lasso-dCRT (Liu et al., 2020)

1 INPUT dataset $(X, y)$, $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$;

2 OUTPUT vector of p-values $\{p_j\}_{j=1}^p$;

3 $\hat{\mathcal{S}}^{\text{SCREENING}} = \{j \in [p] \mid \hat{\beta}_j^{\text{PEN}} \neq 0\}$;

4 **for** $j \notin \hat{\mathcal{S}}^{SCREENING}$ **do**

5 $\quad \big| \quad p_j = 1$

6 **end**

7 **for** $j \in \hat{\mathcal{S}}^{SCREENING}$ **do**

8 $\quad \big| \quad$ 1. Distill info. of $X_{-j}$ to $X_{*,j}$ and $y$, obtain $\hat{\beta}^{d_{X_{*,j}}}$ and $\hat{\beta}^{d_{y,j}}$

9 $\quad \big| \quad$ 2. Obtain test statistic:

$$T_j = \sqrt{n} \; \frac{(y - X_{-j}\hat{\beta}^{d_{y,j}})^T (x_j - X_{-j}\hat{\beta}^{d_{X_{*,j}}})}{\|y - X_{-j}\hat{\beta}^{d_{y,j}}\|_2 \|X_{*,j} - X_{-j}\hat{\beta}^{d_{X_{*,j}}}\|_2}$$

$\quad \big| \quad$ 3. Compute (two-sided) p-value $p_j = 2[1 - \Phi(|T_j|)]$

10 **end**

# Distillation Operation

For each variable $j$, *remove* all the conditional information of the remaining variables $X_{-j}$ to $X_{*,j}$ and to $y$

## Lasso-Distillation

▶ $\hat{\boldsymbol{\beta}}^{d_y,j} = \text{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^{p-1}} \sum_{i=1}^n \log\left[1 + \exp(-y_i(X_{i,-j}^T \boldsymbol{\beta}))\right] + \lambda \|\boldsymbol{\beta}\|_1$

▶ $\hat{\boldsymbol{\beta}}^{d_{X_{*,j}}}(\lambda) = \text{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^{p-1}} \frac{1}{2} \|X_{*,j} - X_{-j}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$

## dCRT test statistics

$$T_j = \sqrt{n} \; \frac{(y - X_{-j}\hat{\boldsymbol{\beta}}^{d_y,j})^T (x_j - X_{-j}\hat{\boldsymbol{\beta}}^{d_{X_{*,j}}})}{\|y - X_{-j}\hat{\boldsymbol{\beta}}^{d_y,j}\|_2 \|X_{*,j} - X_{-j}\hat{\boldsymbol{\beta}}^{d_{X_{*,j}}}\|_2} \xrightarrow[n \to +\infty]{\mathcal{H}_0^j} \mathcal{N}(0,1) \,.$$

conditional to $y$ and $X_{-j}$

# Distillation Operator for Logistic Regression?

▶ Lasso-distillation in Liu et al. (2020): model misspecification with logistic relationship

▶ Demo:

$$\mathbf{y} = \texttt{logit}(\mathbf{X}\boldsymbol{\beta}^0 + \sigma\boldsymbol{\xi})$$

$\rho$

sparsity    snr

▶ 100 simulations, $p = 400$, $\mathrm{X} \sim \mathcal{N}(0, \boldsymbol{\Sigma})$ with $\boldsymbol{\Sigma}$ a Toeplitz matrix.

# Null distribution of dCRT test statistic



(a) n = 200          (b) n = 400          (c) n = 800

- ▶ QQ-Plot for one null dCRT statistic, 1000 samplings
- ▶ Fixed $p = 400$ varying, $n \in \{200, 400, 800\}$
- ▶ Theoretical quantile is of a standard Gaussian distribution

# Null distribution of dCRT test statistic



(a) n = 200      (b) n = 400      (c) n = 800

▶ QQ-Plot for one null dCRT statistic, 1000 samplings

▶ Fixed $p = 400$ varying, $n \in \{200, 400, 800\}$

▶ Theoretical quantile is of a standard Gaussian distribution

⚠ Null distribution is far from standard normal

# Adaptation of CRT to high-dim logistic regresssion

► Ning and Liu (2017): $T_j^{\text{decorr}}$ – decorrelating test-statistic $T_j$

► Finding $\hat{\boldsymbol{\beta}}^{d_y,j}$: find $\hat{\boldsymbol{\beta}}^{\text{PEN}}$, then omitting the $j$th coefficient

► Finding $\hat{\boldsymbol{\beta}}^{d_{\text{x}_{*,j}}}$: using weighted Lasso instead of standard Lasso.

# Adaptation of CRT to high-dim logistic regresssion

- Ning and Liu (2017): $T_j^{\texttt{decorr}}$ – decorrelating test-statistic $T_j$
- Finding $\hat{\boldsymbol{\beta}}^{d_y,j}$: find $\hat{\boldsymbol{\beta}}^{\texttt{PEN}}$, then omitting the $j$th coefficient
- Finding $\hat{\boldsymbol{\beta}}^{d_{\mathrm{X}_*,j}}$: using weighted Lasso instead of standard Lasso.

## Intuition: based on classical Rao's test score

$$\hat{\boldsymbol{\beta}}^{\texttt{PEN}} = \operatorname{argmin}_{\boldsymbol{\beta}\in\mathbb{R}^p} \underbrace{\sum_{i=1}^{n} \log\left[1 + \exp(-y_i(\mathrm{X}_{i,*}^T\boldsymbol{\beta}))\right]}_{\ell(\boldsymbol{\beta})} + \lambda\left\|\boldsymbol{\beta}\right\|_1$$

$$T_j^{\texttt{Rao}} = n^{1/2}\nabla_{\beta_j}\ell(\boldsymbol{\beta})\hat{\mathrm{I}}_{j|-j}^{-1/2}$$

- In high-dimension, $T_j^{\texttt{Rao}}$ is biased.
- The general formula of decorrelated test score $T_j^{\texttt{decorr}}$ is a debiased version of $T_j^{\texttt{Rao}}$.

# Proposed Solution: CRT-Logit

---

**Algorithm 3: CRT-logit**

---

1  INPUT dataset $(X, y)$, $X \in \mathbb{R}^{n \times p}, y \in \mathbb{R}^n$;

2  OUTPUT vector of p-values $\{p_j\}_{j=1}^p$;

3  $\hat{\boldsymbol{\beta}} \leftarrow$ `penalized_MLE`$(X, y)$; $\hat{\mathcal{S}}^{\texttt{screening}} \leftarrow \{j \in [p] \mid \hat{\beta}_j^{\texttt{MLE}} \neq 0\}$;

4  **for** $j \notin \hat{\mathcal{S}}^{screening}$ **do**

5     |  $p_j = 1$

6  **end**

7  **for** $j \in \hat{\mathcal{S}}^{screening}$ **do**

8     |  1. $\hat{\boldsymbol{\beta}}^{d_{X_{*,j}}} \leftarrow$ `scaled_lasso`$(X_{*,j}, X_{*,-j})$

9     |  2. $\hat{\boldsymbol{\beta}}^{d_{y,j}} \leftarrow (\hat{\beta}_1, \hat{\beta}_2, \ldots, \hat{\beta}_{j-1}, \hat{\beta}_{j+1}, \ldots, \hat{\beta}_p)$

10     |  3. $T_j^{\texttt{decorr}} \leftarrow$ `decorrelated_test_score`$(X, y)$

11     |  4. $p_j \leftarrow 2[1 - \Phi(|T_j^{\texttt{decorr}}|)]$

12  **end**

---

# Effectiveness of decorrelation on test statistics



(a) n = 200    (b) n = 400    (c) n = 800

# Simulation: Mildly High-dimensional Scenario



▶ 100 runs of simulations across varying parameters; FDR controlled $\alpha = 0.1$.

▶ Methods: Debiased Lasso (`dlasso`), model-X Knockoff (`KO-logit`), original dCRT (`dCRT`), our version of CRT (`CRT-logit`).

# Problem: Curse of Dimensionality



⚠ Failure of detecting variables when dimension grows large.

# Inference with Clusters of Variables

▶ Solution: Dimension reduction via spatially constrained clustering $p \longrightarrow C$ such that $C \ll p$: `cCRT-logit`



▶ Stabilize inference results with multiple clusterings + p-values aggregation (`cCRT-logit-agg`):

# Statistical inference with spatial tolerance

▶ Brain spatial organization: "close" voxels ↔ "close" weights



☐ Null weight voxels

■ Positive weight voxels

■ Negative weight voxels

▶ Spatial tolerance $\delta$ for false discoveries: $\mathrm{FDR}^\delta$



Declared Significant

True Positive

$\delta$

True Positive

# False Discovery Rate with spatial tolerance



- ▶ Distance between voxels: $d(j, k)$ for $(j, k) \in [p]^2$
- ▶ $\delta$-null region: $N^\delta = \left\{ j \in [p] \mid \forall k \in [p], d(j, k) \leq \delta \implies \beta_k^0 = 0 \right\}$

---

## FDP$^\delta$ and FDR$^\delta$

Given an estimation of the support $\hat{\mathcal{S}}$:

$$\text{FDP}^\delta = \frac{|\{N^\delta \cap \hat{\mathcal{S}}\}|}{|\hat{\mathcal{S}}| \vee 1}$$

$$\text{FDR}^\delta = \mathbb{E}[\text{FDP}^\delta]$$

# Theoretical Results for CRT-logit

Estimate support, for $\alpha \in (0,1)$:

▶ $\hat{\mathcal{S}}_{\texttt{cCRT-logit}} = \texttt{FDR\_control}(\{\hat{p}_j^{\,\texttt{cCRT-logit}}\}_{j=1}^p, \alpha)$

▶ $\hat{\mathcal{S}}_{\texttt{cCRT-logit-agg}} = \texttt{FDR\_control}(\{\hat{p}_j^{\,\texttt{cCRT-logit}}\}_{j=1}^p, \alpha)$

## Conjecture

If the clusters are independent, and all the clusters from all partitions considered have a diameter smaller than $\delta$, and the variables located between clusters are positively correlated, then, the output $\hat{\mathcal{S}}_{\texttt{cCRT-logit}}$ and $\hat{\mathcal{S}}_{\texttt{cCRT-logit-agg}}$ control $FDR^\delta$ under predefined level $\alpha \in (0,1)$, i.e.

$$\limsup_{n \to \infty} \mathbb{E}\left[\frac{|\hat{\mathcal{S}}_{\texttt{cCRT-logit}} \cap N^\delta|}{|\hat{\mathcal{S}}_{AKO}| \vee 1}\right] \leq \alpha$$

and

$$\limsup_{n \to \infty} \mathbb{E}\left[\frac{|\hat{\mathcal{S}}_{\texttt{cCRT-logit-agg}} \cap N^\delta|}{|\hat{\mathcal{S}}_{AKO}| \vee 1}\right] \leq \alpha$$

where $N^\delta$ is the $\delta$-null region defined above.

# Semi-simulated dataset (HCP 900)

Semi-simulated dataset:

▶ Use real data X (*e.g.* emotion task).

▶ build $\boldsymbol{\beta}^0$ independently from data of different task, *e.g.* X_motor_foot.

▶ Generate synthetic responses y from X_emotion and $\boldsymbol{\beta}^0$.



$$\mathbb{P}(\mathbf{y}_i = 1 \mid \mathbf{X}_{i,*}) = \frac{1}{1 + \exp(-\mathbf{X}_{i,*}\boldsymbol{\beta}^0 + \sigma\xi_i)}$$

# Semi-simulated dataset (HCP 900)



▶ FDR/Average Power of 50 runs of simulations on Human Brain Connectome dataset.

▶ Parameters: $n = 800$ (taken from 400 subjects), SNR $= 1.5$. FDR$^\delta$ is controlled at level $\alpha = 0.1$ and $\delta = 8$.

▶ Methods (clustering versions): Desparsified Lasso (`cdlasso`), model-X Knockoff (`cKO-logit`), original dCRT (`cdCRT`), our version of CRT (`cCRT-logit`) and the aggregation of CRT-logit across clusterings (`cCRT-logit-agg.`)

# Semi-simulated dataset (HCP 900)

# Related: Ensemble of Clustered Knockoffs

# Outline

# Summary

New procedures for statistical inference with high-dimensional data

## Aggregation of Multiple Knockoffs

► FDR control guarantee.
► Demonstrated empirically: more stable in inference results and higher statistical power.

## Conditional Randomization Test for high-dimensional logistic regression (CRT-logit)

► Reduce computational cost of original CRT.
► Ensemble of clusterings version works well in very high-dimension.

## Remark

Clustered version involves additional assumptions for statistical guarantee.

# Perspectives

▶ Formal statement and proof of the Conjecture on FDR control with CRT-logit

▶ Theoretical analysis of clustering inference with Knockoffs and CRT-logit: relaxing the assumption on independence of clusters.

▶ Applications for genomics data.

▶ Generative networks for knockoff variables generation.

# Perspectives

▶ Formal statement and proof of the Conjecture on FDR control with CRT-logit

▶ Theoretical analysis of clustering inference with Knockoffs and CRT-logit: relaxing the assumption on independence of clusters.

▶ Applications for genomics data.

▶ Generative networks for knockoff variables generation.

Thank you for listening!

# Second-order Model-X Knockoffs

Shares the first two moments - mean and covariance, *i.e.* :

$$\mathbb{E}[\tilde{X}] = \mathbb{E}[X], \quad \mathbb{E}[\tilde{X}^T \tilde{X}] = \Sigma \quad \text{and} \quad \mathbb{E}[\tilde{X}^T X] = \Sigma - \text{diag}\{s\}$$

<u>Additional assumption</u>: X has Gaussian design

$$\longrightarrow \tilde{x}_j \mid x_j \stackrel{d}{=} \mathcal{N}(\boldsymbol{\mu}, V)$$

$\longrightarrow$ Finding $\text{diag}\{s\}$ by:

▶ Semi-definite Programming (SDP)

▶ Approximate Semi-definite program (ASDP)

▶ Equi-correlated

# Knockoff Statistic

## Definition (Candès et al. (2018))

A knockoff statistic $W = \{W_j\}_{j \in [p]}$ is a measure of feature importance that satisfies the two following properties:

1. Depends only on $X, \tilde{X}$ and $y$

$$W = f(X, \tilde{X}, y), \text{ and}$$

2. Swapping the original variable column $x_j$ and its knockoff column $\tilde{x}_j$ will switch the sign of $W_j$ iff $j$ is in the support set $\mathcal{S}$:

$$W_j([X, \tilde{X}]_{swap(S)}, y) = \begin{cases} W_j([X, \tilde{X}], y) \text{ if } j \in \mathcal{S}^c \\ -W_j([X, \tilde{X}], y) \text{ if } j \in \mathcal{S} \end{cases}$$

# Theoretical Results for AKO

## Assumption (Null Distribution of Knockoff Statistic)

*Under the null hypothesis $H_{0,j} : \beta_j^0 = 0$, the Knockoff Statistics follow the same null distribution.*

## Lemma (Lemma 2 – N., Chevalier, Thirion, Arlot, 2020)

*Under the above assumption, and furthermore assume $|\mathcal{S}^c| \geq 2$, for all $j \in \mathcal{S}^c$ the intermediate p-value $\hat{p}_j$ satisfies*

$$\forall t \in (0, 1) : \quad \mathbb{P}(\hat{p}_j \leq t) \leq \frac{\kappa p}{|\mathcal{S}^c|} t$$

*where*

$$\kappa = \frac{\sqrt{22} - 2}{7\sqrt{22} - 32} \leq 3.24$$

# Theoretical Results for AKO

> **Theorem (Theorem 1 – N., Chevalier, Thirion, Arlot, 2020)**
>
> *(Finite-sample guarantee of FDR control)*
> *If, under the null hypothesis $H_{0,j} : \beta_j^0 = 0$, the Knockoff Statistics follow the same distribution, and if $|\mathcal{S}^c| \geq 2$, then for an arbitrary number of samplings $B$, the output $\hat{\mathcal{S}}_{AKO}$ of Aggregation of Multiple Knockoff (AKO) controls FDR under predefined level $\alpha \in (0, 1)$, i.e.*
>
> $$\mathbb{E}\left[\frac{|\hat{\mathcal{S}}_{AKO} \cap \mathcal{S}^c|}{|\hat{\mathcal{S}}_{AKO}| \vee 1}\right] \leq \kappa\alpha$$
>
> *where $\kappa = \dfrac{\sqrt{22} - 2}{7\sqrt{22} - 32} \leq 3.24$.*

# AKO extra results - Genome Wide Association Study

▶ Data: Flowering Phenotype of Arabidopsis Thaliana (FT_GH) – $n = 166, p = 9938$

▶ Objective: detect association of 174 candidate genes with phenotype FT_GH that dictates flowering time (Atwell et al., 2010).

▶ Preprocessing: dimension reduction following Slim et al. (2019)
$$p = 9938 \longrightarrow p = 1500.$$

| Method | Detected Genes |
|--------|----------------|
| AKO | AT2G21070, AT4G02780, AT5G47640 |
| KO | AT2G21070 |
| KO-GZ | AT2G21070 |
| DL-BH | — |

Table: List of detected genes associated with phenotype FT_GH.

From previous studies: AT2G21070 (Kim et al., 2008), AT4G02780 (Silverstone et al., 1998), AT5G47640 (Cai et al., 2007)

# Adaptation of CRT to high-dim logistic regresssion

- $T_j^{\text{decorr}}$: Decorrelating test-statistic $T_j$ (Ning and Liu, 2017)
- Finding $\hat{\boldsymbol{\beta}}^{d_y,j}$: find $\hat{\boldsymbol{\beta}}^{\text{PEN}}$, then omitting the $j$th coefficient, *i.e.*

$$\hat{\boldsymbol{\beta}}_j^{d_y,j} = (\hat{\beta}_1, \hat{\beta}_2, \ldots, \hat{\beta}_{j-1}, \hat{\beta}_{j+1}, \ldots, \hat{\beta}_p)$$

- Finding $\hat{\boldsymbol{\beta}}^{d_{\mathrm{X}*,j}}$: using weighted Lasso instead of standard Lasso.

$$\hat{\boldsymbol{\beta}}^{d_{\mathrm{X}*,j}} = \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^{p-1}} \frac{1}{n} \sum_{i=1}^{n} \frac{\exp{(\hat{\boldsymbol{\beta}}^T \mathbf{x}_i)}}{[1 + \exp{(\hat{\boldsymbol{\beta}}^T \mathbf{x}_i)}]^2} (x_{i,j} - \boldsymbol{\beta}^T \mathrm{X}_{-j})^2 + \lambda \left\| \boldsymbol{\beta} \right\|_1$$

# Adaptation of CRT to high-dim logistic regresssion

## Decorrelated test statistic

$T_j^{\mathbf{decorr}} =$

$$-(n\,\hat{\mathrm{I}}_{j|-j})^{-1/2} \sum_{i=1}^{n} \left[ y_i - \frac{1}{1 + \exp\left(-\mathrm{X}_{i,-j}\hat{\boldsymbol{\beta}}^{d_y,j}\right)} \right] \left[ \mathrm{X}_{i,j} - \mathrm{X}_{i,-j}^T \hat{\boldsymbol{\beta}}^{d_{\mathrm{X}_{*,j}}} \right],$$

where $\hat{\mathrm{I}}_{j|-j}$ is the estimated partial Fisher information:

$$\hat{\mathrm{I}}_{j|-j} = \frac{1}{n} \sum_{i=1}^{n} \frac{\exp\left(\hat{\boldsymbol{\beta}}\mathrm{X}_{i,*}\right)}{[1 + \exp\left(\hat{\boldsymbol{\beta}}\mathrm{X}_{i,*}\right)]^2} (\mathrm{X}_{i,j} - \hat{\boldsymbol{\beta}}^{d_x\,T}\mathrm{X}_{i,-j})^2\ \mathrm{X}_{i,j}.$$

## Asymptotic distribution (Ning and Liu, 2017)

$$T_j^{\mathbf{decorr}} \xrightarrow[n\to+\infty]{\mathcal{H}_0^j} \mathcal{N}(0,1)$$

# Assumptions in Clustered Inference

> **Spatial homogeneity with distance $\delta$**
>
> For all $(j, k) \in [p] \times [p]$, $d(j, k) \leq \delta$ implies that $\boldsymbol{\Sigma}_{j,k} \geq 0$, where $\boldsymbol{\Sigma}_{j,k} \triangleq \mathrm{Cov}(\mathbf{x}_j, \mathbf{x}_k)$.

> **Sparse-smooth with distance $\delta$**
>
> For all $(j, k) \in [p] \times [p]$, $d(j, k) \leq \delta$ implies that $\mathrm{sign}(\boldsymbol{\beta}_j^0) = \mathrm{sign}(\boldsymbol{\beta}_k^0)$.