# Searching for Consistent Associations with a Multi-Environment Knockoff Filter

Shuangning Li

Stanford University

Deep Learning for Genetics, Genomics and Metagenomics

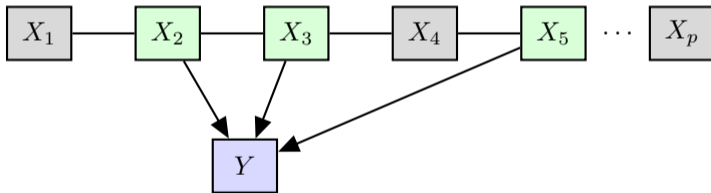# Joint work with



M. Sesia



Y. Romano



E. Candès



C. Sabatti

Li, Sesia, Romano, Candès & Sabatti. **Searching for Consistent Associations with a Multi-Environment Knockoff Filter.** *Biometrika*, 2021.
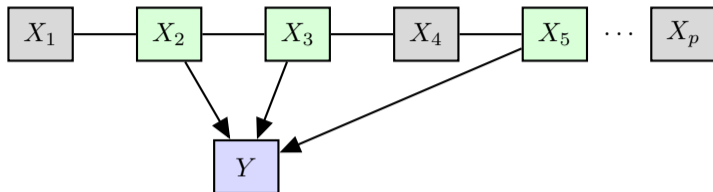
# Introduction

- In Genome-Wide Association Studies (GWAS), geneticists have measured hundreds of thousands of genetic variants and wish to know which of these influence a trait.
  E.g. What are the genes that influence **height**?



- In standard analysis, focus variables $X_j$'s that are associated with $Y$.
  The results are very far from identifying "causal" genetic variants.

# Conditional Independence
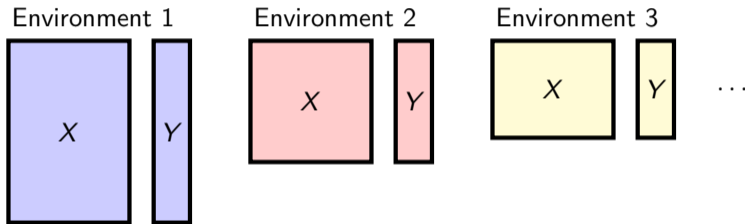


▶ Better goal: test for conditional independence

$$H_j : X_j \perp\!\!\!\perp Y \mid X_{-j}.$$

If $H_j$ is true, the $j$-th variable does not provide information about the response $Y$ beyond what is already provided by all the other variables.

▶ Control the false discovery rate

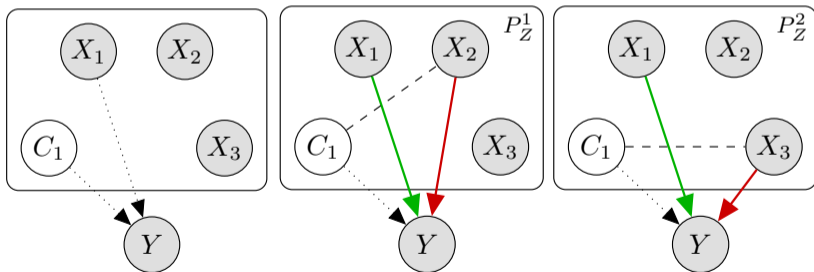$$\text{FDR} = \mathbb{E}\left[\frac{\#\text{ false positives}}{\#\text{ selections}}\right].$$

# Consistence across environments



Environment 1    Environment 2    Environment 3

$X$  $Y$     $X$  $Y$     $X$  $Y$   $\cdots$

▶ We say a variable $j$ is null in environment $e$ if $X_j^e \perp\!\!\!\perp Y^e | X_{-j}^e$. We would like to find variables that are non-null in **all** environments.

▶ In other words, now a variable is null for "consistent independence hypothesis", if it is null in at least one environments.

# Unobserved Confounder

- Shaded nodes are observed; "white" nodes are not.
- Dotted arrows represent true causal model connecting $Z$ to $Y$
- Broken lines identify correlations across variables
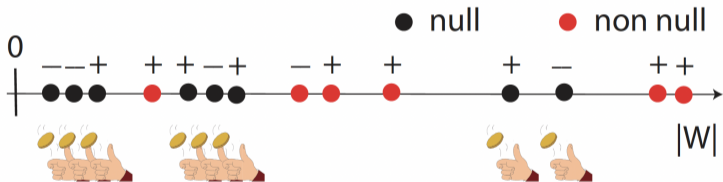- Filled arrows indicated detected conditional association

# Knockoffs

▶ The method of **knockoffs** (Barber and Candès, 2015; Candès et al., 2018) allows one to test the conditional independence hypothesis and provably controls the FDR.

▶ Construct knockoffs $\Rightarrow$ Get important statistics $\Rightarrow$ Report selected set
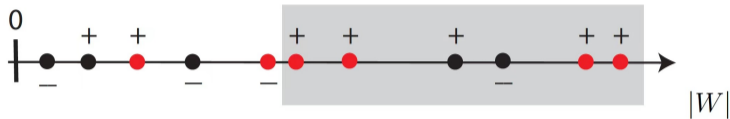
# Knockoffs

- We compute importance statistics $W$ from machine learning algorithms.
- Large $W_j$ says that variable $j$ appears important.
- Conditional on $|W|$, signs of null $W_j$'s are i.i.d. coin flips — crucial for FDR control!



- FDR is still controlled if the signs are $\leq$ coin flips (More conservative!)

# Knockoffs



▶ Let
$$\tau = \min\left\{ t : \widehat{\mathrm{FDP}}(t) = \frac{1+|\mathcal{S}^-(t)|}{1 \vee |\mathcal{S}^+(t)|} \le q \right\}$$

▶ Report $\hat{\mathcal{S}} = \{W_j \ge \tau\}$.

▶ Then false discovery rate is controlled.

$$\mathbb{E}\left[\frac{\# \text{ false positives}}{\# \text{ selections}}\right] \le q$$

# Multi-Environment Knockoff Filter (simple version!)

1. Compute importance statistics for each environment $\{W^e\}_{e=1}^{E}$ as usual.
2. Merge the statistics as follows

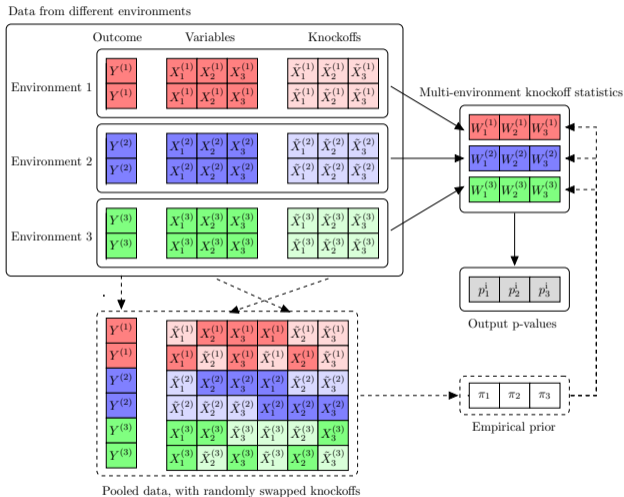$$\text{sign}\,(W_j) = \min_{e}\,\text{sign}\,(W_j^e)$$

$$|W_j| = f\,(|W_j^1|, |W_j^2|, \ldots, |W_j^E|)$$

   For example,

$$|W_j| = \prod_{e=1}^{E} |W_j^e|$$

3. Report a selected set based on $W$.

▶ If a variable is "null", then it's null in at least one environment. Say it's null in envir 1. Then $\text{sign}\,(W_j) = \min_e\,\text{sign}\,(W_j^e) \leq \text{sign}\,(W_j^1) \leq$ coin flips.
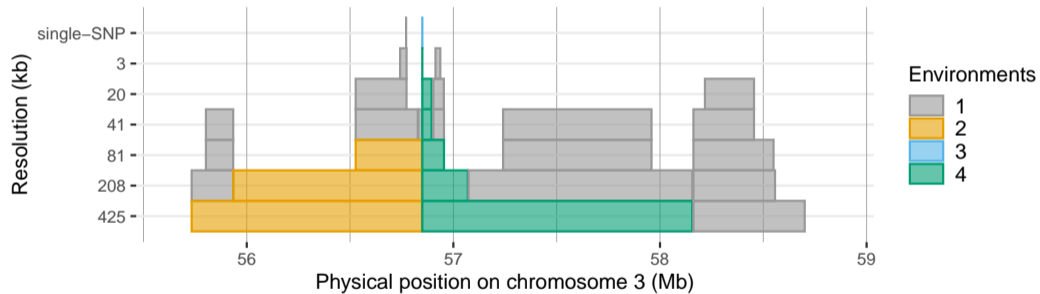
# Multi-Environment Knockoff Filter (complicated version!)

# Partial Conjunction

▶ If some of the environments have small sample size, then power can be low.
▶ Weaker goal: find variables that are non-null in $\geq r$ environments

# UK biobank data analysis

| Environment | Sample size | Self-reported ancestries |
|---|---|---|
| African | 7,623 | "African", "Caribbean", "Any other black background", "Black or Black British" |
| Asian | 3,284 | "Asian or Asian British", "Chinese", "Any other Asian background" |
| British | 429,934 | "British" |
| European | 28,994 | "Any other white background", "Irish", "White" |
| Indian | 7,628 | "Indian", "Pakistani", "Bangladeshi" |

# UK biobank data analysis

# References

Barber, R. F. and Candès, E. (2015). Controlling the false discovery rate via knockoffs. *Ann. Stat.*, 43(5):2055–2085.

Candès, E., Fan, Y., Janson, L., and Lv, J. (2018). Panning for gold: "model-X" knockoffs for high dimensional controlled variable selection. *J. R. Stat. Soc. B*, 80(3):551–577.

Thank you!