

A theoretical treatment of
conditional independence testing
under Model-X

Eugene Katsevich

August 3, 2020

Joint work with Aaditya Ramdas (CMU)



Preprint available at arxiv.org/abs/2005.05506

Conditional independence testing under Model-X

For random variables $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \in \mathbb{R}^{1+1+p}$, would like to test

$$H_0 : \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}$$

based on a sample $(X, Y, Z) = \{(X_i, Y_i, Z_i) : i = 1, \dots, n\}$.

Conditional independence testing under Model-X

For random variables $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \in \mathbb{R}^{1+1+p}$, would like to test

$$H_0 : \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}$$

based on a sample $(X, Y, Z) = \{(X_i, Y_i, Z_i) : i = 1, \dots, n\}$.

Model-X assumption (Candès et al., 2018)

$$f_{\mathbf{X}|\mathbf{Z}} = f_{\mathbf{X}|\mathbf{Z}}^* \text{ for known } f^*$$

Model-X methodologies

MX knockoffs and the conditional randomization test (CRT) proposed by Candès et al., 2018.

Algorithm: Conditional Randomization Test

Data: Samples (X_i, Y_i, Z_i) , test statistic $T(X, Y, Z)$, MX $f_{\mathbf{X}|Z}^*$

Compute $T(X, Y, Z)$;

for $b = 1, \dots, B$ **do**

 | Resample \tilde{X}_i^b from $f_{\mathbf{X}|Z}^*$ and recompute $T(\tilde{X}^b, Y, Z)$;

end

Return: $p_{\text{CRT}} = \frac{1}{1+B} \left(1 + \sum_{b=1}^B \mathbb{1}(T(\tilde{X}^b, Y, Z) \geq T(X, Y, Z)) \right)$

A variety of knockoffs and CRT extensions are now available.

Common themes in MX methodology

- T usually based on a statistical machine learning method
- Performance of the ML method impacts the power of the test
- Randomness in X used for inference, conditioning on Y and Z

Common themes in MX methodology

- T usually based on a statistical machine learning method
- Performance of the ML method impacts the power of the test
- Randomness in X used for inference, conditioning on Y and Z

Goal of this talk:

Develop a quantitative understanding of these themes.

Most powerful CRT against a point alternative

Given alternative distributions \bar{f}_Z and $\bar{f}_{Y|X,Z}$, consider testing

$$H_0 : (\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \sim f_Z f_{\mathbf{X}|Z}^* f_{\mathbf{Y}|Z} \quad \text{for some } f_Z, f_{\mathbf{Y}|Z}$$

$$H_1 : (\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \sim \bar{f}_Z f_{\mathbf{X}|Z}^* \bar{f}_{\mathbf{Y}|X,Z}.$$

Composite null prevents directly deducing the most powerful test.

Most powerful CRT against a point alternative

Given alternative distributions \bar{f}_Z and $\bar{f}_{Y|X,Z}$, consider testing

$$H_0 : (\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \sim f_Z f_{\mathbf{X}|Z}^* f_{\mathbf{Y}|Z} \quad \text{for some } f_Z, f_{\mathbf{Y}|Z}$$

$$H_1 : (\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \sim \bar{f}_Z f_{\mathbf{X}|Z}^* \bar{f}_{\mathbf{Y}|X,Z}.$$

Composite null prevents directly deducing the most powerful test.

But the CRT ϕ_T is also a *conditionally valid* test:

$$\sup_{H_0} \mathbb{E}[\phi_T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) | \mathbf{Y} = y, \mathbf{Z} = z] \leq \alpha \quad \text{for all } y, z.$$

What is the most powerful *conditionally valid* test?

Conditioning reduces a composite null to a point null

Fix realizations Y_i and Z_i for each i . Then,

$$\text{Under } H_0, \quad X_i | Y_i, Z_i \stackrel{\text{ind}}{\sim} f_{X_i | Z_i}^*;$$

$$\text{Under } H_1, \quad X_i | Y_i, Z_i \stackrel{\text{ind}}{\sim} f_{X_i | Z_i}^* \frac{\bar{f}_{Y_i | X_i, Z_i}}{f_{Y_i | Z_i}}.$$

So, conditioning on Y, Z gives a simple hypothesis testing problem.

Neyman-Pearson gives the most powerful CRT

By NP, most powerful conditionally valid test rejects for large

$$T^{\text{opt}}(X, Y, Z) = \prod_{i=1}^n \frac{\bar{f}_{Y_i|X_i, Z_i}}{\bar{f}_{Y_i|Z_i}} \propto \prod_{i=1}^n \bar{f}_{Y_i|X_i, Z_i}.$$

Neyman-Pearson gives the most powerful CRT

By NP, most powerful conditionally valid test rejects for large

$$T^{\text{opt}}(X, Y, Z) = \prod_{i=1}^n \frac{\bar{f}_{Y_i|X_i, Z_i}}{\bar{f}_{Y_i|Z_i}} \propto \prod_{i=1}^n \bar{f}_{Y_i|X_i, Z_i}.$$

Theorem¹

CRT based on T^{opt} is the most powerful conditionally valid test against $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \sim \bar{f}_Z f_{\mathbf{X}|Z}^* \bar{f}_{\mathbf{Y}|\mathbf{X}, \mathbf{Z}}$.

¹All results stated informally

Neyman-Pearson gives the most powerful CRT

By NP, most powerful conditionally valid test rejects for large

$$T^{\text{opt}}(X, Y, Z) = \prod_{i=1}^n \frac{\bar{f}_{Y_i|X_i, Z_i}}{\bar{f}_{Y_i|Z_i}} \propto \prod_{i=1}^n \bar{f}_{Y_i|X_i, Z_i}.$$

Theorem¹

CRT based on T^{opt} is the most powerful conditionally valid test against $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \sim \bar{f}_Z f_{\mathbf{X}^*|Z} \bar{f}_{\mathbf{Y}|\mathbf{X}, Z}$.

The knockoff statistic $T^{\text{opt}}([X, \tilde{X}], Y) = \prod_{i=1}^n \bar{f}_{Y_i|X_i}$ maximizes the probability $\mathbb{P}[T([X, \tilde{X}], Y) > T([X, \tilde{X}]_{\text{swap}(j)}, Y)]$.

¹All results stated informally

Neyman-Pearson gives the most powerful CRT

By NP, most powerful conditionally valid test rejects for large

$$T^{\text{opt}}(X, Y, Z) = \prod_{i=1}^n \frac{\bar{f}_{Y_i|X_i, Z_i}}{\bar{f}_{Y_i|Z_i}} \propto \prod_{i=1}^n \bar{f}_{Y_i|X_i, Z_i}.$$

Theorem¹

CRT based on T^{opt} is the most powerful conditionally valid test against $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \sim \bar{f}_{\mathbf{Z}} f_{\mathbf{X}^*|\mathbf{Z}} \bar{f}_{\mathbf{Y}|\mathbf{X}, \mathbf{Z}}$.

The knockoff statistic $T^{\text{opt}}([X, \tilde{X}], Y) = \prod_{i=1}^n \bar{f}_{Y_i|X_i}$ maximizes the probability $\mathbb{P}[T([X, \tilde{X}], Y) > T([X, \tilde{X}]_{\text{swap}(j)}, Y)]$.

ML methods T used in practice learn approximations to $\bar{f}_{\mathbf{Y}|\mathbf{X}, \mathbf{Z}}$.

¹All results stated informally

Connections

- *Least squares:* If $\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}\gamma + \epsilon$, the optimal CRT statistic is $\|\mathbf{Y} - \mathbf{X}\hat{\beta} - \mathbf{Z}\hat{\gamma}\|^2 - \|\mathbf{Y} - \mathbf{Z}\hat{\gamma}\|^2$, akin to the OLS F -statistic.

Connections

- *Least squares*: If $\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}\gamma + \epsilon$, the optimal CRT statistic is $\|\mathbf{Y} - \mathbf{X}\hat{\beta} - \mathbf{Z}\hat{\gamma}\|^2 - \|\mathbf{Y} - \mathbf{Z}\hat{\gamma}\|^2$, akin to the OLS F -statistic.
- *Unbiased testing*: In parametric families with nuisance params, MP unbiased test is MP test conditional on nuisance sufficient statistic.

Connections

- *Least squares*: If $\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}\gamma + \epsilon$, the optimal CRT statistic is $\|\mathbf{Y} - \mathbf{X}\beta - \mathbf{Z}\gamma\|^2 - \|\mathbf{Y} - \mathbf{Z}\gamma\|^2$, akin to the OLS F -statistic.
- *Unbiased testing*: In parametric families with nuisance params, MP unbiased test is MP test conditional on nuisance sufficient statistic.
- *Holdout randomization test*²: $\hat{f}_{\mathbf{Y}|\mathbf{X},\mathbf{Z}}$ learned on a training set and CRT based on loss $\sum_i \log \hat{f}_{Y_i|X_i,Z_i}$ run on a test set.

²Tansey et al., 2018, Bates et al., 2020

Impact of ML prediction error on CRT power

Impact of ML prediction error on CRT power

Suppose that

$$\mathbf{Y} = \mathbf{X}\beta + g(\mathbf{Z}) + \epsilon, \quad \epsilon \sim N(0, \sigma^2),$$

and we estimate \hat{g} based on a separate training set (like HRT).

How does test error in \hat{g} impact the power of the CRT based on \hat{g} ?

Impact of ML prediction error on CRT power

Suppose that

$$\mathbf{Y} = \mathbf{X}\beta + g(\mathbf{Z}) + \epsilon, \quad \epsilon \sim N(0, \sigma^2),$$

and we estimate \hat{g} based on a separate training set (like HRT).

How does test error in \hat{g} impact the power of the CRT based on \hat{g} ?

In particular, consider the CRT based on³

$$T(X, Y, Z) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu_i)(Y_i - \hat{g}(Z_i)),$$

where $\mu_i = \mathbb{E}[X_i | Z_i]$.

³Related to the *generalized covariance measure* of Shah and Peters (2020); studied in the double robustness literature, e.g. Chernozhukov et al. (2018).

Prediction error impacts asymptotic efficiency of the CRT

For simplicity, suppose $\text{Var}[\mathbf{X}|\mathbf{Z}] = s^2$ a.s. for some $s^2 > 0$.

Define the test error $\mathcal{E} = \mathbb{E}[(\hat{g}(\mathbf{Z}) - g(\mathbf{Z}))^2]$.

Fixing dimension and training set, let $n \rightarrow \infty$ and $\beta_n = \frac{h}{\sqrt{n}}$.

Prediction error impacts asymptotic efficiency of the CRT

For simplicity, suppose $\text{Var}[\mathbf{X}|\mathbf{Z}] = s^2$ a.s. for some $s^2 > 0$.

Define the test error $\mathcal{E} = \mathbb{E}[(\hat{g}(\mathbf{Z}) - g(\mathbf{Z}))^2]$.

Fixing dimension and training set, let $n \rightarrow \infty$ and $\beta_n = \frac{h}{\sqrt{n}}$.

Theorem

Under the MX assumption and bounded fourth moments,

$$\mathbb{E}[\phi_T(X, Y, Z)|Y, Z] \rightarrow \Phi\left(z_\alpha + \frac{hs}{\sqrt{\sigma^2 + \mathcal{E}}}\right),$$

almost surely in Y, Z .

More connections to linear regression

Note that

$$\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z}) = \mathbf{X}\beta + (g(\mathbf{Z}) - \widehat{\mathbf{g}}(\mathbf{Z}) + \epsilon) = \mathbf{X}\beta + \epsilon'; \quad \epsilon' \sim (0, \sigma^2 + \mathcal{E}).$$

Estimation error in $\widehat{\mathbf{g}}$ inflates the noise level by \mathcal{E} .

More connections to linear regression

Note that

$$\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z}) = \mathbf{X}\beta + (g(\mathbf{Z}) - \widehat{\mathbf{g}}(\mathbf{Z}) + \epsilon) = \mathbf{X}\beta + \epsilon'; \quad \epsilon' \sim (0, \sigma^2 + \mathcal{E}).$$

Estimation error in $\widehat{\mathbf{g}}$ inflates the noise level by \mathcal{E} .

$$T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \frac{1}{\sqrt{n}}(\mathbf{X} - \mu)^T(\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z})) \text{ is unnormalized OLS } t\text{-stat.}$$

More connections to linear regression

Note that

$$\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z}) = \mathbf{X}\beta + (g(\mathbf{Z}) - \widehat{\mathbf{g}}(\mathbf{Z}) + \epsilon) = \mathbf{X}\beta + \epsilon'; \quad \epsilon' \sim (0, \sigma^2 + \mathcal{E}).$$

Estimation error in $\widehat{\mathbf{g}}$ inflates the noise level by \mathcal{E} .

$$T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \frac{1}{\sqrt{n}}(\mathbf{X} - \mu)^T(\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z})) \text{ is unnormalized OLS } t\text{-stat.}$$

Under the null ($\beta = 0$),

$$T(\widetilde{\mathbf{X}}, \mathbf{Y}, \mathbf{Z}) | \mathbf{Y}, \mathbf{Z} \rightarrow N(0, s^2(\sigma^2 + \mathcal{E}))$$

More connections to linear regression

Note that

$$\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z}) = \mathbf{X}\beta + (\mathbf{g}(\mathbf{Z}) - \widehat{\mathbf{g}}(\mathbf{Z}) + \boldsymbol{\epsilon}) = \mathbf{X}\beta + \boldsymbol{\epsilon}'; \quad \boldsymbol{\epsilon}' \sim (0, \sigma^2 + \mathcal{E}).$$

Estimation error in $\widehat{\mathbf{g}}$ inflates the noise level by \mathcal{E} .

$$T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \frac{1}{\sqrt{n}}(\mathbf{X} - \boldsymbol{\mu})^T(\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z})) \text{ is unnormalized OLS } t\text{-stat.}$$

Under the null ($\beta = 0$),

$$T(\widetilde{\mathbf{X}}, \mathbf{Y}, \mathbf{Z}) | \mathbf{Y}, \mathbf{Z} \rightarrow N(0, s^2(\sigma^2 + \mathcal{E}))$$

$$T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) | \mathbf{X}, \mathbf{Z} \rightarrow N(0, s^2(\sigma^2 + \mathcal{E})).$$

More connections to linear regression

Note that

$$\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z}) = \mathbf{X}\beta + (\mathbf{g}(\mathbf{Z}) - \widehat{\mathbf{g}}(\mathbf{Z}) + \epsilon) = \mathbf{X}\beta + \epsilon'; \quad \epsilon' \sim (0, \sigma^2 + \mathcal{E}).$$

Estimation error in $\widehat{\mathbf{g}}$ inflates the noise level by \mathcal{E} .

$$T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \frac{1}{\sqrt{n}}(\mathbf{X} - \mu)^T(\mathbf{Y} - \widehat{\mathbf{g}}(\mathbf{Z})) \text{ is unnormalized OLS } t\text{-stat.}$$

Under the null ($\beta = 0$),

$$T(\widetilde{\mathbf{X}}, \mathbf{Y}, \mathbf{Z}) | \mathbf{Y}, \mathbf{Z} \rightarrow N(0, s^2(\sigma^2 + \mathcal{E}))$$

$$T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) | \mathbf{X}, \mathbf{Z} \rightarrow N(0, s^2(\sigma^2 + \mathcal{E})).$$

If the semiparametric model is true, CRT resampling distribution asymptotically equivalent to OLS null distribution.

Asymptotic validity under a weaker assumption than MX

Asymptotic validity under a weaker assumption than MX

Under the null,

$$\text{Var}[T_n | Y, Z] = \frac{1}{n} \sum_{i=1}^n \text{Var}[X_i | Z_i] (Y_i - \hat{g}(Z_i))^2 = S_n^2,$$

with S_n^2 known. We can show that, almost surely in (Y, Z) ,

$$\mathcal{L}(S_n^{-1} T_n | Y, Z) \rightarrow N(0, 1).$$

(S_n, T_n) only involves first and second moments $\mathbb{E}[\mathbf{X} | \mathbf{Z}]$, $\text{Var}[\mathbf{X} | \mathbf{Z}]$.

Asymptotic validity under a weaker assumption than MX

This observation motivates the following:

Definition (MX(2) assumption)

$(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ is such that $\mathbb{E}[\mathbf{X}|\mathbf{Z}] = \mu(\mathbf{Z})$ and $\text{Var}[\mathbf{X}|\mathbf{Z}] = s^2(\mathbf{Z})$,
for known functions $\mu(\cdot)$ and $s^2(\cdot)$.

Asymptotic validity under a weaker assumption than MX

This observation motivates the following:

Definition (MX(2) assumption)

$(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ is such that $\mathbb{E}[\mathbf{X}|\mathbf{Z}] = \mu(\mathbf{Z})$ and $\text{Var}[\mathbf{X}|\mathbf{Z}] = s^2(\mathbf{Z})$, for known functions $\mu(\cdot)$ and $s^2(\cdot)$.

There is an asymptotically valid conditional independence test that does not require the MX assumption or resampling:

Theorem

Under MX(2), $\phi = \mathbb{1}(S_n^{-1} T_n > z_{1-\alpha})$ has uniform asympt. level α .

Asymptotic validity under a weaker assumption than MX

This observation motivates the following:

Definition (MX(2) assumption)

$(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ is such that $\mathbb{E}[\mathbf{X}|\mathbf{Z}] = \mu(\mathbf{Z})$ and $\text{Var}[\mathbf{X}|\mathbf{Z}] = s^2(\mathbf{Z})$, for known functions $\mu(\cdot)$ and $s^2(\cdot)$.

There is an asymptotically valid conditional independence test that does not require the MX assumption or resampling:

Theorem

Under MX(2), $\phi = \mathbb{1}(S_n^{-1} T_n > z_{1-\alpha})$ has uniform asympt. level α .

Related to the double robustness literature, but the latter focuses on approximating the first moments $\mathbb{E}[\mathbf{X}|\mathbf{Z}]$ and $\mathbb{E}[\mathbf{Y}|\mathbf{Z}]$.

Connections to causal inference

The MX setting is like that of a randomized experiment:
 \mathbf{X} is the treatment; \mathbf{Y} is the response; \mathbf{Z} are the covariates.

Instead of complete randomization, the treatment \mathbf{X} is assigned to units based on the covariates \mathbf{Z} using a known mechanism $f_{\mathbf{X}|\mathbf{Z}}^*$.

Even in the absence of confounding, adjusting for covariates known to reduce variance in estimates of causal effect.

Connections to causal inference

Non-asymptotic tests based on resampling \mathbf{X} go back to Fisher (1935) and Rosenbaum (1984). Both treat \mathbf{Y}, \mathbf{Z} as fixed.

Asymptotic “superpopulation” approach (e.g. Robins et al., 1992) treats \mathbf{Y} as random, focused on semiparametric models.

Current work reinforces close links between the two approaches; see also discussion in Rosenbaum (2002).

Summary

In this talk, we

- Identified the CRT most powerful against point alternatives;

Summary

In this talk, we

- Identified the CRT most powerful against point alternatives;
- Expressed CRT's asymptotic power in terms of ML test error;

Summary

In this talk, we

- Identified the CRT most powerful against point alternatives;
- Expressed CRT's asymptotic power in terms of ML test error;
- Weakened the MX assumption, retaining asymptotic validity;

Summary

In this talk, we

- Identified the CRT most powerful against point alternatives;
- Expressed CRT's asymptotic power in terms of ML test error;
- Weakened the MX assumption, retaining asymptotic validity;
- Drew some connections between MX and causal inference.

Future work

Many questions still remain open:

- Are all valid tests under MX also conditionally valid? If not, are all optimal tests conditionally valid?

Future work

Many questions still remain open:

- Are all valid tests under MX also conditionally valid? If not, are all optimal tests conditionally valid?
- Optimality statements against composite alternatives?

Future work

Many questions still remain open:

- Are all valid tests under MX also conditionally valid? If not, are all optimal tests conditionally valid?
- Optimality statements against composite alternatives?
- Extensions of power results to high dimensions? Some results for lasso test statistics available for knockoffs⁴ and for CRT⁵.

⁴Weinstein et al. (2017, 2020), Fan et al. (2019), Liu and Rigollet (2019)

⁵Celentano et al. (2020)

Future work

Many questions still remain open:

- Are all valid tests under MX also conditionally valid? If not, are all optimal tests conditionally valid?
- Optimality statements against composite alternatives?
- Extensions of power results to high dimensions? Some results for lasso test statistics available for knockoffs⁴ and for CRT⁵.
- Further connections with causal inference and with existing asymptotic (doubly robust, semiparametric) inference?

⁴Weinstein et al. (2017, 2020), Fan et al. (2019), Liu and Rigollet (2019)

⁵Celentano et al. (2020)

Future work

Many questions still remain open:

- Are all valid tests under MX also conditionally valid? If not, are all optimal tests conditionally valid?
- Optimality statements against composite alternatives?
- Extensions of power results to high dimensions? Some results for lasso test statistics available for knockoffs⁴ and for CRT⁵.
- Further connections with causal inference and with existing asymptotic (doubly robust, semiparametric) inference?

These lines of inquiry can improve our understanding of MX methodologies and help guide their development in the future.

⁴Weinstein et al. (2017, 2020), Fan et al. (2019), Liu and Rigollet (2019)

⁵Celentano et al. (2020)