

DINFK

Detecting when the available data does not allow reliable inference

Fanny Yang, Statistical Machine Learning Group Department of Computer Science @ETH Zurich

joint work with students Alex Tifrea, Eric Stavarache, Piersilvio de Bartolomeis, Javier A. Martinez, Konstantin Donhauser













individual test sample with new disease/class







individual test sample with new disease/class







I. A lower bound for hidden confounding using randomized control trials

joint work with Piersilvio de Bartolomeis, Javier Abad Martinez, Konstantin Donhauser

work in progress















Definition of "strong" hidden confounding
Approach: How to detect it using RCT?

Potential outcome framework



Potential outcome framework



Potential outcome framework



Potential outcome framework



Potential outcome framework



Potential outcome framework



Additional assumptions:

- Transportability of CATE, i.e. $\mathbb{E}_{\mathbb{P}^{os}}[Y(1) Y(0) \mid X] = \mathbb{E}_{\mathbb{P}^{rct}}[Y(1) Y(0) \mid X]$
- Support inclusion $supp(\mathbb{P}^{rct}) \subseteq supp(\mathbb{P}^{os})$

Additional assumptions:

- Transportability of CATE, i.e. $\mathbb{E}_{\mathbb{P}^{os}}[Y(1) Y(0) \mid X] = \mathbb{E}_{\mathbb{P}^{rct}}[Y(1) Y(0) \mid X]$
- Support inclusion $supp(\mathbb{P}^{rct}) \subseteq supp(\mathbb{P}^{os})$

Definitions:

• \mathbb{P}^{os} satisfies MSM(Γ) if $\Gamma^{-1} \leq \frac{\mathbb{P}^{os}(T=1|X,U)}{\mathbb{P}^{os}(T=0|X,U)} / \frac{\mathbb{P}^{os}(T=1|X)}{\mathbb{P}^{os}(T=0|X)} \leq \Gamma$ almost surely (Tan-O6)

Additional assumptions:

- Transportability of CATE, i.e. $\mathbb{E}_{\mathbb{P}^{os}}[Y(1) Y(0) \mid X] = \mathbb{E}_{\mathbb{P}^{rct}}[Y(1) Y(0) \mid X]$
- Support inclusion $supp(\mathbb{P}^{rct}) \subseteq supp(\mathbb{P}^{os})$



Additional assumptions:

- Transportability of CATE, i.e. $\mathbb{E}_{\mathbb{P}^{os}}[Y(1) Y(0) \mid X] = \mathbb{E}_{\mathbb{P}^{rct}}[Y(1) Y(0) \mid X]$
- Support inclusion $supp(\mathbb{P}^{rct}) \subseteq supp(\mathbb{P}^{os})$



• true confounding strength $\Gamma^*(\mathbb{P}^{os})$: The smallest Γ for which \mathbb{P}^{os} satisfies $MSM(\Gamma)$

Additional assumptions:

- Transportability of CATE, i.e. $\mathbb{E}_{\mathbb{P}^{os}}[Y(1) Y(0) \mid X] = \mathbb{E}_{\mathbb{P}^{rct}}[Y(1) Y(0) \mid X]$
- Support inclusion $supp(\mathbb{P}^{rct}) \subseteq supp(\mathbb{P}^{os})$



• true confounding strength $\Gamma^*(\mathbb{P}^{os})$: The smallest Γ for which \mathbb{P}^{os} satisfies $MSM(\Gamma)$

Scenarios we want to detect: when true confounding Γ^* of \mathbb{P}^{os} is too large

Definition of "strong" hidden confounding
Approach: How to detect it using RCT?

Our plug-and-play approach for desired significance α :

1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$

- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$

- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$
- Test $H_0(\Gamma) \Longrightarrow$ test whether $\mu \in [\mu_{\Gamma}^-, \mu_{\Gamma}^+]$ with ATE $\mu = \mathbb{E}_{\mathbb{P}}[Y(1) Y(0)]$ and

ATE sensitivity bounds
$$\mu_{\Gamma}^{-} = \inf_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)], \quad \mu_{\Gamma}^{+} = \sup_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)]$$

- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$
- Test $H_0(\Gamma) \Longrightarrow$ test whether $\mu \in [\mu_{\Gamma}^-, \mu_{\Gamma}^+]$ with ATE $\mu = \mathbb{E}_{\mathbb{P}}[Y(1) Y(0)]$ and

ATE sensitivity bounds
$$\mu_{\Gamma}^{-} = \inf_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)], \quad \mu_{\Gamma}^{+} = \sup_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)]$$

all full distributions that yield observed $\mathbb{P}_{X,Y,T}^{os}$ and satisfy MSM(Γ)

- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$
- Test $H_0(\Gamma) \Longrightarrow$ test whether $\mu \in [\mu_{\Gamma}^-, \mu_{\Gamma}^+]$ with ATE $\mu = \mathbb{E}_{\mathbb{P}}[Y(1) Y(0)]$ and

ATE sensitivity bounds
$$\mu_{\Gamma}^{-} = \inf_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)], \ \mu_{\Gamma}^{+} = \sup_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)]$$



- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$
- Test $H_0(\Gamma) \Longrightarrow$ test whether $\mu \in [\mu_{\Gamma}^-, \mu_{\Gamma}^+]$ with ATE $\mu = \mathbb{E}_{\mathbb{P}}[Y(1) Y(0)]$ and

ATE sensitivity bounds
$$\mu_{\Gamma}^{-} = \inf_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)], \ \mu_{\Gamma}^{+} = \sup_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)]$$



- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$
- Test $H_0(\Gamma) \Longrightarrow$ test whether $\mu \in [\mu_{\Gamma}^-, \mu_{\Gamma}^+]$ with ATE $\mu = \mathbb{E}_{\mathbb{P}}[Y(1) Y(0)]$ and

ATE sensitivity bounds
$$\mu_{\Gamma}^{-} = \inf_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)], \ \mu_{\Gamma}^{+} = \sup_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)]$$



- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$
- Test $H_0(\Gamma) \Longrightarrow$ test whether $\mu \in [\mu_{\Gamma}^-, \mu_{\Gamma}^+]$ with ATE $\mu = \mathbb{E}_{\mathbb{P}}[Y(1) Y(0)]$ and

ATE sensitivity bounds
$$\mu_{\Gamma}^{-} = \inf_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)], \ \mu_{\Gamma}^{+} = \sup_{\widetilde{\mathbb{P}} \in P_{\Gamma}(\mathbb{P}_{X,Y,T}^{os})} \mathbb{E}_{\widetilde{\mathbb{P}}}[Y(1) - Y(0)]$$



Our plug-and-play approach for desired significance α :

- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$



probability of rejection over 20 runs on semi-synthetic data

Our plug-and-play approach for desired significance α :

- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$



probability of rejection over 20 runs on semi-synthetic data
Our paradigm: finding a lower bound

Our plug-and-play approach for desired significance α :

- 1. Test $\phi_{\alpha}(\Gamma)$ of the null $H_0(\Gamma)$: \mathbb{P}^{os} satisfies $MSM(\Gamma) \iff \Gamma^* \leq \Gamma$
- 2. Report $\hat{\Gamma}_{LB} = \inf \{ \Gamma: \phi_{\alpha}(\Gamma) = 0 \}$ and flag if $\hat{\Gamma}_{LB} > \Gamma_{\text{thresh}}$



- probability of rejection over 20 runs on semi-synthetic data
- asymptotically $\mathbb{P}(\widehat{\Gamma}_{LB} > \Gamma^*) \leq \alpha$ implied by $\mathbb{P}(\phi_a(\Gamma^*) = 1) \leq \alpha$

- without RCT and using sensitivity bounds
 - $_{\circ}$ quantification via critical value $\widehat{\Gamma}_{ct}$ that changes causal conclusions

e.g. vanderWeele-Ding-17, Jin-Ren-Candes-23 etc.

- without RCT and using sensitivity bounds
 - quantification via critical value Γ̂_{ct} that changes causal conclusions
 e.g. vanderWeele-Ding-17, Jin-Ren-Candes-23 etc.
 but: unclear relation to true Γ*

- without RCT and using sensitivity bounds
 - \circ quantification via critical value $\hat{\Gamma}_{ct}$ that changes causal conclusions e.g. vanderWeele-Ding-17, Jin-Ren-Candes-23 etc.

but: unclear relation to true Γ^*

can test joint null hypothesis ATE(obs. study) > 0 and MSM(Γ) holds
 e.g. Yadlowsky-Namkoong-Basu-Duchi-Tian-22, Jin-Ren-Candes-23

- without RCT and using sensitivity bounds
 - \circ quantification via critical value $\hat{\Gamma}_{ct}$ that changes causal conclusions e.g. vanderWeele-Ding-17, Jin-Ren-Candes-23 etc.

but: unclear relation to true Γ^*

can test joint null hypothesis ATE(obs. study) > 0 and MSM(Γ) holds
 e.g. Yadlowsky-Namkoong-Basu-Duchi-Tian-22, Jin-Ren-Candes-23

but: rejection only means either MSM(Γ) assumption wrong or ATE ≤ 0

- without RCT and using sensitivity bounds
 - \circ quantification via critical value $\hat{\Gamma}_{ct}$ that changes causal conclusions e.g. vanderWeele-Ding-17, Jin-Ren-Candes-23 etc.

but: unclear relation to true Γ^*

can test joint null hypothesis ATE(obs. study) > 0 and MSM(Γ) holds
 e.g. Yadlowsky-Namkoong-Basu-Duchi-Tian-22, Jin-Ren-Candes-23

but: rejection only means either MSM(Γ) assumption wrong or ATE ≤ 0

• with RCT:

 $_{\circ}\,$ binary test for existence of confounding with H_0: $\Gamma^{\star}>1$

e.g. Viele et al '14, Hussein-Oberst-Shih-Sontag '22

Previous paradigms that can be used for detection

- without RCT and using sensitivity bounds
 - $\circ\,$ quantification via critical gamma value $\widehat{\Gamma}_{ct}\,$ that changes causal conclusions



• with RCT:

• binary te e.g. Viele e \rightarrow flag even if Γ^* small Sontag '22

Evaluation on real-world data (WHI)

- Randomized trial and observational study run by the NHLBI (1993-2005)
- Treatment: hormone replacement therapy
- Outcomes: coronary heart disease

Evaluation on real-world data (WHI)

- Randomized trial and observational study run by the NHLBI (1993-2005)
- Treatment: hormone replacement therapy
- Outcomes: coronary heart disease
- hidden confounder (revealed later): start of treatment



Evaluation on real-world data (WHI)

- Randomized trial and observational study run by the NHLBI (1993-2005)
- Treatment: hormone replacement therapy
- Outcomes: coronary heart disease
- hidden confounder (revealed later): start of treatment

	Coronary heart disease	
treated	as trial start	ted before trial
$\hat{\Gamma}_{CT}$	1.017	1.164
$\hat{\Gamma}_{LB}$	1.009	1.224
ψ_{bin}	1	1
ψ_{sens}	0	1



OS

rct

start of trial

duration of treatment

- Compute $\hat{\Gamma}_{CT}$ that changes ATE sign and compare let "expert" assess "likeliness"
- ψ_{bin} : tests for existence, e.g. check $\hat{\Gamma}_{LB} > 1$
- ψ_{sens} (ours): check whether too large $\hat{\Gamma}_{LB} > \hat{\Gamma}_{CT}$



Current and future work

Higher power using

- kernelized test as opposed to averaging
- non-"adversarial" sensitivity model

Current and future work

Higher power using

- kernelized test as opposed to averaging
- non-"adversarial" sensitivity model

Extended applicability:

- multiple observational studies (no RCT)
- Automatic detection of hidden confounders from set of features





II. Semi-supervised novelty detection using ensembles with regularized disagreement

joint work with Alexandru Tifrea, Eric Stavarache

published at UAI '22









Novelty detection method tells user that software doesn't "know enough" to predict new point



Novelty detection method tells user that software doesn't "know enough" to predict new point



Novelty detection method tells user that software doesn't "know enough" to predict new point



Definition: Points we can't make inference on
 Approach: How to detect those samples?

What's "novel" to a trained model?

"novel" / o.o.d. points: test points $x \in X$ the model cannot reliably predict.

What's "novel" to a trained model?

"novel" / o.o.d. points: test points $x \in X$ the model cannot reliably predict.

First: which points $x \in X$ can a model predict "reliably" in an unseen test set?

• i.d. generalization from finite samples (traditional learning theory) and

What's "novel" to a trained model?

"novel" / o.o.d. points: test points $x \in X$ the model cannot reliably predict.

First: which points $x \in X$ can a model predict "reliably" in an unseen test set?

- i.d. generalization from finite samples (traditional learning theory) and
- o.o.d. generalization (extrapolatable from training distribution) -

depends on test shift & model complexity



- True classifierTraining support P
 - 🗱 Unlabeled test data

Extrapolatable given training distribution + linear ground truth:

Points $x \in X$ where the set of all linear Bayes optimal classifiers agree on



True classifierTraining support P

🙁 Unlabeled test data

Extrapolatable given training distribution + linear ground truth:

Points $x \in X$ where the set of all linear Bayes optimal classifiers agree on



Extrapolatable given training distribution + linear ground truth:

Points $x \in X$ where the set of all linear Bayes optimal classifiers agree on



Goal now: how to output green area

Definition: Points we can't make inference on
 Approach: How to detect those samples?

- with good validation accuracy on old classes
- but different predictions outside of training distribution

- with good validation accuracy on old classes
- but different predictions outside of training distribution
- \rightarrow flag all points where the models disagree as "novel"

- with good validation accuracy on old classes
- but different predictions outside of training distribution
- \rightarrow flag all points where the models disagree as "novel"





- with good validation accuracy on old classes
- but different predictions outside of training distribution
- \rightarrow flag all points where the models disagree as "novel"







Classifier II

- with good validation accuracy on old classes
- but different predictions outside of training distribution
- \rightarrow flag all points where the models disagree as "novel"





Key for "good performance": Complexity of ensemble models being only as large as needed



Key for "good performance": Complexity of ensemble models being only as large as needed



Key for "good performance": Complexity of ensemble models being only as large as needed



Key for "good performance": Complexity of ensemble models being only as large as needed



Idea for right amount of disagreement: maximize disagreement s.t. validation error of all models small

"regularization"
Key for our improvement: Regularized disagreement

Key for "good performance": Complexity of ensemble models being only as large as needed



Idea for right amount of disagreement: maximize disagreement s.t. validation error of all models small

"regularization"

using unlabeled test data

Key for our improvement: Regularized disagreement

Key for "good performance": Complexity of ensemble models being only as large as needed



Idea for right amount of disagreement: maximize disagreement s.t. validation error of all models small

"regularization"

using unlabeled test data

using labeled training data

The near OOD problem on images with DNN

CIFAR-10

Chest X-Ray & retinal datasets



The near OOD problem on images with DNN

Chest X-Ray & retinal datasets







- "Hidden yet quantifiable: A lower bound for confounding strength using randomized trials" by Piersilvio De Bartolomeis*, Javier Abad*, Konstantin Donhauser, FY, arxiv preprint
- "Semi-supervised novelty detection using ensembles with regularized disagreement" by Alexandru Țifrea, Eric Stavarache, and FY, (UAI), 2022

• • labeled training points

- training distribution
- novel classes
- unlabeled test points
 from both old & new classes



• • labeled training points

- training distribution
- novel classes
- unlabeled test points
 from both old & new classes



- • labeled training points
 - training distribution
 - novel classes
- unlabeled test points
 from both old & new classes



• Artificially label all unlabeled test data with one label

- • labeled training points
 - training distribution
 - novel classes
- unlabeled test points
 from both old & new classes



• Artificially label all unlabeled test data with one label





Regularizing disagreement using labeled data



Current and future work

Non-adversarial confounding



Discussion of the paradigm

- Propose two tests $\phi(\Gamma)$ based on (C)ATE sensitivity analysis intervals
 - obs: estimate mu with importance weighting rct, then ATE sensitivity
 valid when ATE bounds are asymptotically normal
 - rct: estimate mu on rct, then CATE sensitivity on obs -> average on rct valid when CATE sensitivity bounds converge at a $1/\sqrt{n}$ rate and $n_{rct} \ll n_{os}$