Causality Basics

(Based on Peters et al., 2017, "Elements of Causal Inference: Foundations and Learning Algorithms") Chang Liu 2020.12.16

Content

- Causality and Structural Causal Models.
- Causal Inference.
- Causal Discovery.

Causality

• Formal definition of causality:

"two variables have a causal relation, if *intervening* the cause may change the effect, but not *vice versa*" [Pearl, 2009; Peters et al., 2017].

- Intervention: change the value of a variable by leveraging mechanisms and changing variables out of the considered system.
- Example: for the Altitude and average T emperature of a city, $A \rightarrow T$.
 - Running a huge heater (intv. T) does not lower A.
 - Raising the city by a huge elevator (intv. A) lowers T.
- Causality contains more information than observation (static data).
 - E.g., both p(A)p(T|A) ($A \rightarrow T$) and p(T)p(A|T) ($T \rightarrow A$) can describe the observed relation p(A,T).

Causality

The Independent Mechanism Principle.

• Principle 2.1 (Independent mechanisms)

The causal generative process of a system's variables is composed of *autonomous* modules that do not inform or influence each other.

In the probabilistic case, this means that the conditional distribution of each variable given its causes (i.e., its mechanism) does not inform or influence the other conditional distributions.



For describing causal relations.

- Definition 6.2 (Structural causal models) An SCM $\mathfrak{C} \coloneqq (\mathbf{S}, P_{\mathbf{N}})$ consists of a collection \mathbf{S} of d(structural) assignments, $X_j \coloneqq f_j(\operatorname{pa}_j, N_j), j = 1, ..., d$ where pa_j are called parents of X_j , and a joint distr. $P_{\mathbf{N}} = P_{N_1,...,N_d}$ over the noise variables, which we require to be jointly independent.
 - The joint indep. requirement on N comes from the Independent Mechanism Principle (the right block).
- Causal graph G: the Directed Acyclic Graph constructed based on the assignments.
- **Proposition 6.3 (Entailed distributions)** An SCM \mathfrak{C} defines a unique distribution over the variables $\mathbf{X} = (X_1, \dots, X_d)$ such that $X_j = f_j(\operatorname{pa}_j, N_j)$, $j = 1, \dots, d$, in distribution, called the entailed distribution $P_{\mathbf{X}}^{\mathfrak{C}}$ and sometimes write $P_{\mathbf{X}}$.
- Implications of an SCM:



Describing interventions using SCMs.

- "What if X was set to x?""What if patients took the treatment?"
- **Definition 6.8 (Intervention distribution)** Consider an SCM $\mathfrak{C} := (\mathbf{S}, P_{\mathbf{N}})$ and its entailed distribution $P_{\mathbf{X}}^{\mathfrak{C}}$. We *replace* one (or several) of the structural assignments to obtain a new SCM $\widetilde{\mathfrak{C}}$. Assume that we replace the assignment for X_k by: $X_k := \tilde{f}(\widetilde{pa}_k, \widetilde{N}_k)$, where all the new and old noises are required to be jointly independent.
 - Atomic/ideal, structural [Eberhardt and Scheines, 2007]/surgical [Pearl, 2009]/independent /deterministic [Korb et al., 2004] intervention: when $\tilde{f}(\tilde{pa}_k, \tilde{N}_k)$ puts a point mass on a value a. Simply write $P_{\mathbf{x}}^{\mathfrak{C}; do(X_k \coloneqq a)}$.

• **Soft intervention**: intervene with a distribution.



Describing interventions using SCMs.

- **Definition 6.12 (Total causal effect)** Given an SCM \mathfrak{C} , there is a total causal effect from X to Y, if X and Y are *dependent* in $P_{\mathbf{x}}^{\mathfrak{C};do(X:=\widetilde{N}_X)}$ for some r.v. \widetilde{N}_X .
- Proposition 6.13 (Total causal effects) Equivalent statements given an SCM &:
 - 1. There is a total causal effect from X to Y.
 - 2. There are $x^{(1)}$ and $x^{(2)}$ such that $P_v^{\mathfrak{C};do(X:=x^{(1)})} \neq P_v^{\mathfrak{C};do(X:=x^{(2)})}$.
 - 3. There is $x^{(1)}$ such that $P_Y^{\mathfrak{C};do(X:=x^{(1)})} \neq P_Y^{\mathfrak{C}}$.
 - 4. X and Y are dependent in $P_{X,Y}^{\mathfrak{C};do(X:=\widetilde{N}_X)}$ for any \widetilde{N}_X with full support.
- Proposition 6.14 (Graphical criteria for total causal effects) Consider SCM \mathfrak{C} with corresponding graph \mathcal{G} .
 - If there is no directed path from X to Y, then there is no total causal effect.
 - Sometimes there is a directed path but no total causal effect.

Describing counterfactuals using SCMs.

- "What if I took the treatment?"
 "What would the outcome be had I acted differently at that moment/situation?"
- Intervention {in a given situation/context} / {for a given individual/unit}.
 Intervention with a conditional.
- **Definition 6.17 (Counterfactuals)** Consider SCM $\mathfrak{C} \coloneqq (\mathbf{S}, P_{\mathbf{N}})$. Given some observations \mathbf{x} , define a *counterfactual* SCM by replacing the distribution of noise variables: $\mathfrak{C}_{\mathbf{X}=\mathbf{x}} \coloneqq (\mathbf{S}, P_{\mathbf{N}}^{\mathfrak{C}|\mathbf{X}=\mathbf{x}})$, where $P_{\mathbf{N}}^{\mathfrak{C}|\mathbf{X}=\mathbf{x}} \coloneqq P_{\mathbf{N}|\mathbf{X}=\mathbf{x}}$ (need not be jointly independent anymore).





Observational implications (SCM as a BayesNet); causality and correlation.

- **Definition 6.1 (d-separation)**. A path *P* is *d-separated* by a set of nodes **S**, iff. *P* either has an emitter ($\rightarrow X \rightarrow$ or $\leftarrow X \rightarrow$) in **S**, or a collider ($\rightarrow X \leftarrow$) that is not in **S** nor are its descendants.
- Definition 6.21 (Markov property) Given a DAG \mathcal{G} and a joint distribution P_X , Theorem 6.22 (Equivalence of this distribution is said to satisfy

(i) the **global Markov property** with respect to the DAG \mathcal{G} if

$$\mathbf{A} \perp\!\!\!\perp_{\mathcal{G}} \mathbf{B} \,|\, \mathbf{C} \,\Rightarrow\, \mathbf{A} \perp\!\!\!\perp \mathbf{B} \,|\, \mathbf{C}$$

for all disjoint vertex sets $\mathbf{A}, \mathbf{B}, \mathbf{C}$ (the symbol $\perp_{\mathcal{G}}$ denotes d-separation see Definition 6.1),

- (ii) the **local Markov property** with respect to the DAG \mathcal{G} if each variable is independent of its non-descendants given its parents, and
- (iii) the **Markov factorization property** with respect to the DAG \mathcal{G} if

$$p(\mathbf{x}) = p(x_1, \dots, x_d) = \prod_{j=1}^d p(x_j | \mathbf{pa}_j^{\mathcal{G}}).$$

- **Markov properties)**. When $P_{\mathbf{X}}$ has a density.
- (iii) is what we mean by using a DAG/BeyesNet to define a joint distribution (i.e., the entailed distribution by the SCM). Its equivalence to (i) means that we can read off conditional independence assertions of this joint distribution from the graph.

Observational implications (SCM as a BayesNet); causality and correlation.

- Confounding bias.
 - T and Y are correlated if X is not given.
 - E.g., *T* = chocolate consumption, *Y* = #Nobel prices, *X* = development status.
- Simpson's paradox.
 - *T*: treatment. *Y*: recovery. *X*: socio-economic status.
 - Is the treatment effective for recovery? p(Y = 1|T = 1) > p(Y = 1|T = 0), but $p^{do(T=1)}(Y = 1) < p^{do(T=0)}(Y = 1)$.
 - $p(Y|t) = \sum_{x} p(Y|x,t)p(x|t)$:

conditioning on T = 1 infers a good X thus a higher probability of recovery.





Observational implications (SCM as a BayesNet); causality and correlation.

- Selection bias.
 - *H* and *F* are correlated if secretly/unconsciously conditioned on *R*.
- Berkson's paradox.
 - "Why are handsome men such jerks?" [Ellenberg, 2014]
 - *H*= handsome, *F*= friendly, *R*= in a relation.
 - Single women often date men with R = 0. Such men tend to be either not H or not R.
- (<u>spurious correlations</u>).





• Example 6.19: Two SCMs that induce the same graph, observational distributions, and intervention distributions, could entail different counterfactual statements ("probabilistically and interventionally equivalent" but not "counterfactually equivalent").

Levels of model according to the prediction ability:

Level	Typical	Typical Questions	Examples
(Symbol)	Activity		
1. Association	Seeing	What is?	What does a symptom tell
P(y x)		How would seeing X	me about a disease?
		change my belief in Y ?	What does a survey tell us
			about the election results?
2. Intervention	Doing	What if?	What if I take aspirin, will
P(y do(x),z)		What if I do X?	my headache be cured?
123-18408. SHADAO 1810-00 - 681			What if we ban cigarettes?
3. Counterfactuals	Imagining,	Why?	Was it the aspirin that
$P(y_x x',y')$	Retrospection	Was it X that caused Y ?	stopped my headache?
		What if I had acted	Would Kennedy be alive
		differently?	had Oswald not shot him?
			What if I had not been
			smoking the past 2 years?

CAUSAL HIERARCHY THEOREM

 L_{1}, L_{2}, L_{3}

[Bareinboim, Correa, Ibeling, Icard, 2020]

Given that an SCM $M \rightarrow PCH$, we can show the following:

Theorem (CHT). With respect to Lebesgue measure over (a suitable encoding of L_3 -equivalence classes of) SCMs, the subset in which any PCH 'collapse' is measure zero.

Informally, for almost any SCM (i.e., almost any possible environment), the PCH does not collapse, i.e., the layers of the hierarchy remains distinct.

L₁

Corollary. To answer question at Layer i (about a certain interaction), one needs knowledge at layer i or higher.

- Given an SCM, infer the distribution/outcome of a variable, under an intervention or unit-level intervention (counterfactual).
 - Estimate causal effect:
 - Average treatment effect: $\mathbb{E}^{do(T=1)}[Y] \mathbb{E}^{do(T=0)}[Y]$.
 - Individual treatment effect: $\mathbb{E}^{do(T=1)}[Y|x] \mathbb{E}^{do(T=0)}[Y|x]$.
- Estimate the interventional distribution is the key task (counterfactual is the conditional in an intervened SCM).

- Truncated factorization [Pearl, 1993] / G-computation formula [Robins, 1986] / manipulation theorem [Spirtes et al., 2000]: $p^{\mathfrak{C};do(X_k:=\widetilde{N}_k)}(x_1, ..., x_d) = \widetilde{p}(x_k) \prod_{j \neq k} p^{\mathfrak{C}}(x_j | x_{pa_j})$, where \widetilde{p} is the density of \widetilde{N}_k .
 - Almost evident from the interventional SCM.
 - Relies on the autonomy/modularity of causality (principle of independence of mechanisms):

one causal mechanism does not influence or inform other causal mechanisms.

- Definition 6.38 (Valid adjustment set, v.a.s) Consider an SCM \mathfrak{C} over nodes V. Let $X, Y \in \mathbf{V}$ with $Y \notin pa_X$. We call a set $Z \subseteq \mathbf{V} \setminus \{X, Y\}$ a v.a.s for (X, Y) if $p^{\mathfrak{C}; do(X \coloneqq x)}(y) = \sum_{\mathbf{z}} p^{\mathfrak{C}}(y|x, \mathbf{z}) p^{\mathfrak{C}}(\mathbf{z})$.
 - Importance: we can then use the *observations* of variables in a v.a.s to estimate an *intervention* distribution.
 - By the previous result, pa_X is a v.a.s.



• Are there any other v.a.s's?

- **Proposition 6.41** (Valid adjustment sets) The following **Z**'s are v.a.s's for (X, Y) where $X, Y \in \mathbf{X}$ and $Y \notin pa_X$:
 - 1. "Parent adjustment": $\mathbf{Z} \coloneqq pa_X$;
 - 2. "Backdoor criterion" (Pearl's admissible/sufficient set): Any $\mathbf{Z} \subseteq \mathbf{X} \setminus \{X, Y\}$ with:
 - Z contains no descendant of X, and
 - **Z** blocks all paths from X to Y entering X through the backdoor $(X \leftarrow \cdots Y)$;
 - 3. "Toward necessity" (characterize all v.a.s) [Shpitser et al., 2010; Perkovic et al., 2015]: Any $\mathbf{Z} \subseteq \mathbf{X} \setminus \{X, Y\}$ with:
 - Z contains no nodes (X excluded) on a directed path from X to Y nor their descendants, and
 - **Z** blocks all non-directed paths from *X* to *Y*.



Figure 6.5: Only the path $X \leftarrow A \rightarrow K \rightarrow Y$ is a "backdoor path" from X to Y. The set $\mathbb{Z} = \{K\}$ satisfies the backdoor criterion (see Proposition 6.41 (ii)); but $\mathbb{Z} = \{F, C, K\}$ is also a valid adjustment set for (X, Y); see Proposition 6.41 (iii).

- 2. "Backdoor criterion" (Pearl's admissible/sufficient set): Any $\mathbf{Z} \subseteq \mathbf{X} \setminus \{X, Y\}$ with:
- Z contains no descendant of X, and
- **Z** blocks all paths from X to Y entering X through the backdoor $(X \leftarrow \cdots Y)$;
- 3. "Toward necessity" (characterize all v.a.s): Any $\mathbf{Z} \subseteq \mathbf{X} \setminus \{X, Y\}$ with:
- Z contains no nodes (X excluded) on a directed path from X to Y nor their descendants, and
- **Z** blocks all non-directed paths from *X* to *Y*.

Do-calculus: alternative way to estimate an *interventional* distribution from *observations*:

- The three rules of **do-calculus**: consider disjoint subsets **X**, **Y**, **Z**, **W** from DAG *G*.
 - 1. "Insertion/deletion of observations": $p^{(c;do(X := x))}(y|z, w) = p^{(c;do(X := x))}(y|w)$, if **Y** and **Z** are d-separated by **X**, **W** in a graph where incoming edges in **X** have been removed.
 - "Action-observation exchange": p^{C;do(X:=x,Z:=z)}(y|w) = p^{C;do(X:=x)}(y|z,w),
 if Y and Z are d-separated by X, W in a graph where incoming edges in X and outgoing edges from Z have been removed.
 - 3. "Insertion/deletion of actions": p^{C;do(X:=x,Z:=z)}(y|w) = p^{C;do(X:=x)}(y|w), if Y and Z are d-separated by X, W in a graph where incoming edges in X and Z(W) have been removed. Here, Z(W) is the subset of nodes in Z that are not ancestors of any node in W in a graph that is obtained from G after removing all edges into X.

Do-calculus: alternative way to estimate an *interventional* distribution from *observations*:

• Theorem 6.45 (Do-calculus).

- 1. The three rules are complete: all identifiable intv. distrs. can be computed by an iterative application of these three rules [Huang and Valtorta, 2006, Shpitser and Pearl, 2006] (more general than v.a.s).
- 2. There is an algorithm [Tian, 2002] that is guaranteed [Huang and Valtorta, 2006, Shpitser and Pearl, 2006] to find *all* identifiable intv. distrs.
- 3. There is a nec. & suff. graphical criterion for identifiability of intv. distrs. [Shpitser and Pearl, 2006, Corollary 3], based on so-called hedges [see also Huang and Valtorta, 2006].

Do-calculus: alternative way to estimate an *interventional* distribution from *observations*:

• Example 6.46 (Front-door adjustment).



- There is no v.a.s if we do not observe U.
- But $p^{\mathfrak{C};do(X:=x)}(y)$ is identifiable using the do-calculus if $p^{\mathfrak{C}}(z|x) > 0$: $p^{\mathfrak{C};do(X:=x)}(y) = \sum_{z} p^{\mathfrak{C}}(z|x) \sum_{x'} p^{\mathfrak{C}}(y|x',z) p^{\mathfrak{C}}(x').$
 - Observing Z in addition to X and \tilde{Y} reveals causal information nicely.
 - Causal relations can be explored by observing the "channel" Z that carries the "signal" from X to Y.

Learn a causal model (causal graph of an SCM) from data.

• In most cases, only observational data is available: Observational causal discovery.

Two main-stream methods:

- Conditional Independence (CI) / constraint based methods (PC algorithm).
- Directional methods (additive noise models).

- Assumption:
 - $P_{\mathbf{X}}$ is Markovian (makes CI assertions as the graph does, in one of the three equiv. forms) and faithful (makes no more CI assertions than the graph does).
 - The Markov equivalence class is then identifiable (Lemma 7.2; Spirtes et al., 2000).
 - **Definition 6.24 (Markov equivalence of graphs)**: Two DAGs are said to be Markov. equiv., if they have the same set of distributions Markovian w.r.t them. Equivalently, they have the same set of d-separation.
 - Lemma 6.25 (Graphical criteria for Markov equivalence; Verma and Pearl [1991]; Frydenberg [1990]) Two DAGs are Markov equiv., iff. they have the same skeleton and the same immoralities (v-structures).



- Instances:
 - Inductive causation (IC) algorithm [Pearl, 2009].
 - SGS algorithm (Spirtes, Glymour, Scheines) [Spirtes et al., 2000].
 - PC algorithm (Peter Spirtes & Clark Glymour) [Spirtes et al., 2000].

- Algorithm:
 - Estimation of Skeleton:
 - Lemma 7.8 [Verma and Pearl, 1991, Lemma 1] (i) Two nodes X, Y in a DAG are adjacent, iff. they cannot be d-separated by any subset $S \subseteq X \setminus \{X, Y\}$. (ii) If two nodes X, Y in a DAG are not adjacent, then they are d-separated by either pa_X or pa_Y .
 - IC/SGS: For each pair of nodes (*X*, *Y*), search through all such *S*'s. *X*, *Y* are adjacent iff. they are cond. dep. given any one of such *S*'s (i).
 - PC: efficient search for S. Start with a fully connected undirected graph, and traverse over such S's in an order of increasing size. For size k, one only has to go through S's that are subsets either of the neighbors of X or of the neighbors of Y (inverse negative of (ii)).

- Algorithm:
 - Orientation of Edges:
 - For a structure X Z Y with no edge betw. X, Y in the skeleton, let A be a set that d-separates X and Y. Then the structure can be oriented as $X \to Z \leftarrow Y$, iff. $Z \notin A$.
 - More edges can be oriented in order to, e.g., avoid cycles.
 - Meek's orientation rules [Meek, 1995]: a complete set of orientation rules.

Additive Noise Models (ANM).

• Philosophy:

Define a family of "particularly natural" conditionals, such that if a conditional can be described using this family, then the conditional in the opposite direction cannot be covered by this family.

- Determines a causal direction even for two variables with observational data.
- **Definition 7.3 (ANMs)** We call an SCM \mathfrak{C} an ANM if the structural assignments are of the form $X_j \coloneqq f_j(\operatorname{pa}_j) + N_j$, j = 1, ..., d. Further assume that f_j 's are differentiable and N_j 's have strictly positive densities.
 - Recall that in defining SCM, the N_j 's are required to be jointly independent, in respect to the independent mechanism principle.

Additive Noise Models (ANM).

• Philosophy:



Additive Noise Models (ANM).

• All the **non-identifiable cases of ANMs** [Zhang & Hyvarinen, 2009, Thm. 8; Peters et al., 2014, Prop. 23]:

Consider the bivariate ANM $X \to Y$, $Y \coloneqq f(X) + N$, where $N \perp X$ and f is three times differentiable. Assume that $f'(x)(\log p_N)''(y) = 0$ only at finitely many points (x, y). If there is a backward ANM $Y \to X$, then one of the followings must hold:

	X	Ν	f
(i)	Gaussian	Gaussian	linear
(ii)	log-mix-lin-exp	log-mix-lin-exp	linear
(iii) <i>,</i> (iv)	log-mix-lin-exp	one-sided asymptotically exponential, or generalized mixture of two exponentials.	strictly monotonic with $\lim_{x\to\infty \text{ or } -\infty} f'(x) = 0.$
(v)	generalized mixture of two exponentials	two-sided asymptotically exponential	strictly monotonic with $\lim_{x\to\infty \text{ or } -\infty} f'(x) = 0.$

Common choices of ANMs:

- In multivariate case, even linear Gaussian with equal noise variances is identifiable [Peters and Buhlmann, 2014].
- Linear Non-Gaussian Acyclic Models (LiNGAMs) [Thm. 7.6; Shimizu et al., 2006].
- Nonlinear Gaussian Additive Noise Models [Thm. 7.7; Peters et al., 2014, Cor. 31].

Learning causal graph with ANMs.

- Score-based methods.
 - Maximize likelihood using a greedy search technique (for optimization over the graph space).
 - For nonlinear Gaussian ANMs [Buhlmann et al., 2014]:
 - Given a current graph \mathcal{G} , regress each var. on its parents and obtain the score $\log p(\mathcal{D}|\mathcal{G})$ $\coloneqq \log p(\mathcal{D}|\mathcal{G}, \hat{\theta}) = -\sum_{j=1}^{d} \widehat{var}[R_j]$, where $R_j \coloneqq X_j - \widehat{f}_j(\operatorname{pa}_j)$ is the residual.
 - When computing the score of a neighboring graph, utilize the decomposition and update only the altered summand.
 - It can be shown to be consistent [Buhlmann et al., 2014].
 - If the noise cannot be assumed to be Gaussian, one can estimate the noise distr. [Nowzohour and Buhlmann, 2016] and obtain an entropy-like score.

Learning causal graph with ANMs.

• RESIT [Mooij et al., 2009; Peters et al., 2014]:

Algorithm 1 Regression with subsequent independence test (RESIT)	
1: Input: I.i.d. samples of a <i>p</i> -dimensional distribution on (X_1, \ldots, X_p) 2: $S := \{1, \ldots, p\}, \pi := []$ 3: PHASE 1: Determine topological order. 4: repeat 5: for $k \in S$ do 6: Regress X_k on $\{X_i\}_{i \in S \setminus \{k\}}$. 7: Measure dependence between residuals and $\{X_i\}_{i \in S \setminus \{k\}}$. 8: end for 9: Let k^* be the <i>k</i> with the weakest dependence. 10: $S := S \setminus \{k^*\}$	14: PHASE 2: Remove superfluous edges. 15: for $k \in \{2,, p\}$ do 16: for $\ell \in pa(\pi(k))$ do 17: Regress $X_{\pi(k)}$ on $\{X_i\}_{i \in pa(\pi(k)) \setminus \{\ell\}}$. 18: if residuals are independent of $\{X_i\}_{i \in \{\pi(1),,\pi(k-1)\}}$ then 19: $pa(\pi(k)) := pa(\pi(k)) \setminus \{\ell\}$ 20: end if 21: end for 22: end for 23: Output: $(pa(1),,pa(p))$
11: $pa(k^*) := S$ 12: $\pi := [k^*, \pi]$ (π will be the topological order, its last component b 13: until $\#S = 0$	being a sink)

Phase 1: Based on the fact that for each node X_i the corresp. N_i is indep. of all non-desc. of X_i . Particularly, for each sink node X_i , $N_i \perp (\mathbf{X} \setminus \{X_i\})$.

Phase 2: visit every node and eliminate incoming edges until the residuals are not indep. anymore.

Learning causal graph with ANMs.

- Independence-Based Score [Peters et al., 2014]:
 - RESIT is asymmetric, thus mistakes will propagate and accumulate over the iterations.
 - Intuition: (i) under G_0 , all the residuals are indep. of each other; (ii) causal minimality helps identifying a unique DAG:

$$\hat{\mathcal{G}} = \underset{\mathcal{G}}{\operatorname{argmin}} \sum_{i=1}^{p} \mathrm{DM}(\operatorname{res}_{i}^{\mathcal{G},\mathrm{RM}}, \operatorname{res}_{-i}^{\mathcal{G},\mathrm{RM}}) + \lambda \,\#(\mathrm{edges})\,.$$
(9)

- $\operatorname{res}_{i}^{\mathcal{G},\mathrm{RM}}$ are the residuals of node X_{i} when it is regressed on its parents specified in \mathcal{G} using regression method RM.
- DM denotes a dependence measure. The second term encourages causal minimality.
- ICA-based methods for LiNGAMs [Shimizu et al., 2006; 2011].

- CI-based methods.
 - More faithful to data.
 - Identif. only up to a Markov equiv. class.
 - CI test methods may be tricky.
- ANMs.
 - A relative strong assumption/belief (a family is particularly natural for modeling conditionals).
 - Identification within a Markov equiv. class.
- Leveraging intervened datasets.
 - Based on the invariance principle of causality.
 - Identification within a Markov equiv. class.
 - Stronger requirement on data.
 - [Peters et al., 2016; Romeijn & Williamson, 2018; Bengio et al., 2019].





Thanks!