

XAI in Epidemiology - Exploring Interacting Causes of a Health Outcome with CoOL (Causes of Outcome Learning)

Andreas Rieckmann¹, Piotr Dworzynski¹, Leila Arras², Sebastian Lapuschkin², Wojciech Samek², Onyebuchi A. Arah³, Naja H. Rod¹, Claus T. Ekstrøm⁴

¹Section of Epidemiology, Department of Public Health & Novo Nordisk Foundation, University of Copenhagen, Denmark

²Department of Artificial Intelligence, Fraunhofer Heinrich Hertz Institute & BIFOLD, Berlin, Germany

³Department of Epidemiology, University of California & Department of Statistics, UCLA College, USA

⁴Section of Biostatistics, Department of Public Health, University of Copenhagen, Denmark

TASK

Discover (interacting) exposures causing a (health) outcome

Tabular data:

Binary Exposures: X_i

- sex (male/female)
- age category
- taking drug A
- taking drug B
- physically active
- smoking
- high BMI
- high blood pressure
- exposed to pollutant C
- ...

Questions:

- Which exposures cause an increased risk?
- Do they act alone or in synergy?

Health Outcome: Y

Disease (1=Yes or 0=No)

Note: a similar task might also occur in other domains (social sciences, environment, public policy, finance, marketing, industry...)

Synergy/Interaction example:

		$X_1=0$	$X_1=1$
		No Asbestos \bar{A}	Asbestos A
$X_2=0$	No Smoking \bar{S}	0.11 %	0.67 %
$X_2=1$	Smoking S	0.95 %	4.50 %

Lung cancer risk

Interaction Coefficient: $IC = (R_{AS} - R_{\bar{A}\bar{S}}) - [(R_{AS} - R_{\bar{A}S}) + (R_{AS} - R_{A\bar{S}})]$

$$IC = R_{AS} - R_{\bar{A}\bar{S}} - R_{\bar{A}S} - R_{A\bar{S}}$$

$$= 4.5 - 0.67 - 0.95 + 0.11 = 2.99$$

- if $IC > 0$ positive interaction
- if $IC < 0$ negative interaction
- if $IC = 0$ no interaction

Can we discover such interactions with machine learning?

STANDARD SOLUTION

"Linear" regression

Model:

$$P(Y=1|X_1=x_1, X_2=x_2) = c_0 + c_1 \cdot x_1 + c_2 \cdot x_2 + c_3 \cdot x_1 \cdot x_2$$

regression terms = all possible combinations of input variables!

$$c_0 = P(Y=1|X_1=0, X_2=0) = \text{baseline risk}$$

$$c_1 = P(Y=1|X_1=1, X_2=0) - P(Y=1|X_1=0, X_2=0) = \text{risk due to } X_1 \text{ alone}$$

$$c_2 = P(Y=1|X_1=0, X_2=1) - P(Y=1|X_1=0, X_2=0) = \text{risk due to } X_2 \text{ alone}$$

$$c_3 = P(Y=1|X_1=1, X_2=1) - P(Y=1|X_1=1, X_2=0) - P(Y=1|X_1=0, X_2=1) + P(Y=1|X_1=0, X_2=0)$$

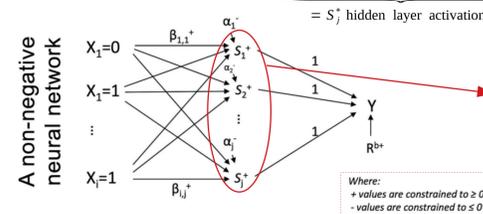
= risk due to interaction between X_1 and X_2 = IC

- Limitations:
- 2^N regression terms (11 binary exposures \rightarrow 2048 terms!)
 - overfitting
 - requires large dataset
 - computationally challenging
 - hard to interpret variables occurring in multiple interaction terms
 - misleading results even if p-value of regression coefficients is low

OUR APPROACH

Neural network + XAI

$$P(Y=1|X_1, X_2, \dots) = \sum_j (\text{ReLU}(\sum_i X_i \cdot \beta_{i,j}^+ + \alpha_j^-)) + R^{b+}$$



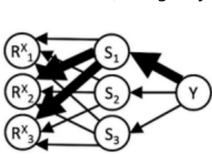
Non-linear hidden layer can modelize any higher-order interaction (no more need to test all terms explicitly!), as well as standalone exposure effects

- Stochastic Gradient Descent
- Squared prediction-error loss $(Y_{true} - \hat{P}_{Model}(Y|X))^2$
- Regularization through squared L2-norm penalty on weights $\|\beta\|_2^2$
- Initialize baseline risk parameter with mean risk $E[Y_{true}]$

METHOD

Discover high-risk subgroups

Step 1: Decompose the prediction for each individual into risk contributions per exposure (using Layer-wise Relevance Propagation) [Bach et al. 2015]

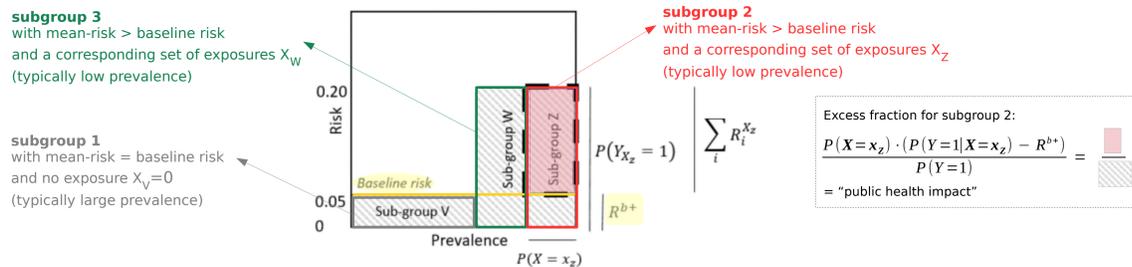


- Output layer: $R_{total} = P_{Model}(Y=1|X) - R^{b+}$
- Hidden layer: $R_j = \frac{S_j}{\sum_j S_j} R_{total}$
- Input layer: $R_i = \sum_j \frac{X_i \cdot \beta_{i,j}^+}{\sum_i X_i \cdot \beta_{i,j}^+} R_j$

Step 2: Cluster individuals based on their risk contributions (using Manhattan distances and the Ward's method for hierarchical clustering) [Strauss et al. 2017]

id	R^{b+}	R^{Man}	R^{Woman}	R^{DrugA}	$R^{No drugA}$	R^{DrugB}	$R^{No drugB}$
1	0.05	0	0	0	0	0	0
4	0.05	0	0	0	0	0	0
2	0.05	0.075	0	0.075	0	0	0
5	0.05	0.075	0	0.075	0	0	0
3	0.05	0	0.075	0	0	0.075	0
6	0.05	0	0.075	0	0	0.075	0

Result 1: Plot of prevalence and mean-risk per subgroup



Result 2: Table of mean-risk contributions per exposure for each subgroup

Mean risk contributions by sub-group (Standard deviation)	Baseline risk	sex_0	sex_1	drug_a_0	drug_a_1	drug_b_0	drug_b_1
subgroup 1	4.6% (0.0%) [4.6%]						
subgroup 2	4.6% (0.0%) [4.6%]		7.7% (0.0%) [0.0%]				7.7% (0.0%) [0.0%]
subgroup 3	4.6% (0.0%) [4.6%]	8% (0.1%) [0.0%]			7.7% (0.0%) [0.0%]		

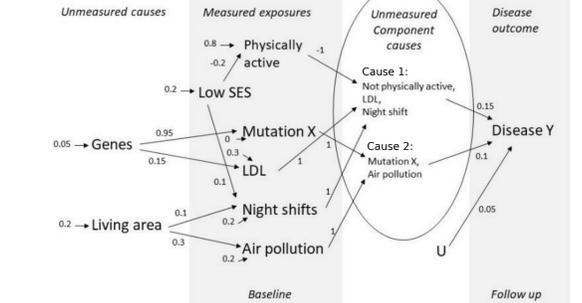
How to test exposures for synergy?

if "contribution of exposure X_i within subgroup" > "contribution of exposure X_i with other exposures set to 0" \Rightarrow then synergy of X_i with other exposures in the subgroup

EXPERIMENTS

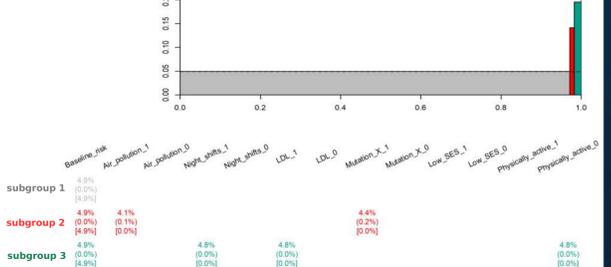
Simulation

Synthetic data:

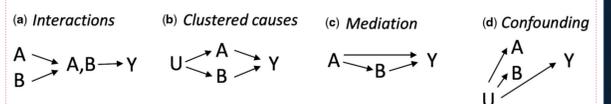


Ground truth risks:
 $P(Y) = 5.4\%$ = mean-risk in population
 $P(Y|U) = 5\%$ = baseline risk
 $P(\text{Cause 1}) = 1.8\%$ = prevalence of subgroup with Cause 1
 $P(Y|\text{Cause 1}) = 15\%$ = mean-risk for this subgroup
 $P(\text{Cause 2}) = 1.2\%$ = prevalence of subgroup with Cause 2
 $P(Y|\text{Cause 2}) = 10\%$ = mean-risk for this subgroup

CoOL Results:



More simulations and robustness checks are available in the Supplementary Material of the paper with below causal structures:



Publication:

Rieckmann et al. Causes of Outcome Learning: a causal inference-inspired machine learning approach to disentangling common combinations of potential causes of a health outcome. International Journal of Epidemiology, 2022

- Code: <https://cran.r-project.org/package=CoOL>
- Website: <https://www.causesofoutcomelearning.org>

