

A Counterfactual Perspective of **Heritability, Explainability, and ANOVA**

Qingyuan Zhao

Statistical Laboratory, University of Cambridge

January 7, 2026 @ UCL Gatsby Unit

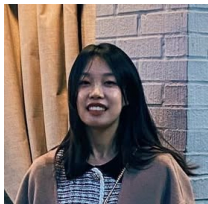
This talk

Ongoing project, two working papers:

1. **Heritability: A Counterfactual Perspective.** (Joint work with Haochen Lei, Jieru Shi, Hongyuan Cao. Available on my website.)
2. **Counterfactual Explainability and Analysis of Variance.** (Joint work with Zijun Gao. Available on arXiv:2411.01625.)



Haochen Lei
Florida State



Jieru Shi
UCL



Hongyuan Cao
Florida State



Zijun Gao
USC Marshall

Motivation

We wanted to answer a simple question:

Is there a good notion of **causal variable importance**?

This question is shaped by frustration from three ends:

1. The **causal inference** literature is obsessed with estimating the effect of a specific intervention (particularly the $ATE = E[Y(1) - Y(0)]$).
2. The **global sensitivity analysis** literature (e.g. functional ANOVA) assumes independent inputs and deterministic outputs.
3. The **genetics** literature is dominated by using linear models to estimate heritability.

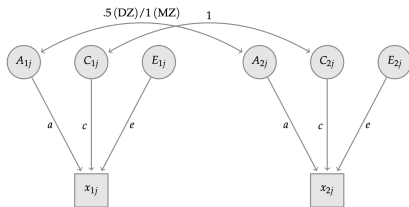
An exception

- ▶ Twin study is a classical method for estimating heritability. **Falconer's formula:**

$$h^2_{\text{twin}} = 2(\rho_{\text{MZ}} - \rho_{\text{DZ}}).$$

$\rho_{\text{MZ}}/\rho_{\text{DZ}}$ is the correlation of a trait (e.g. height) between monozygotic/dizygotic twins.
(The study of correlation between relatives goes back to Fisher's classic 1918 paper.)

- ▶ This is often justified using the **ACE model**: $Y = A + C + E$.
- ▶ Kohler et al. (2011)¹ formalized this as a linear **structural equation model**:



¹Social science methods for twins data: Integrating causality, endowments, and heritability.
Biodemography and Social Biology 2011.

Our first idea

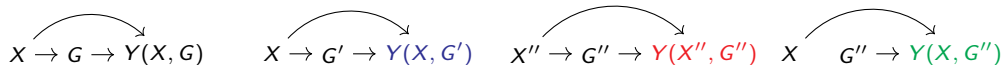
- ▶ We made real progress after thinking about the “**fraternal twin counterfactual**”.
- ▶ Let $Y(G)$ be the potential outcome of someone with genotype G , and $Y(G')$ be the counterfactual of a hypothetical twin with genotype G' that is an i.i.d. copy of G .
- ▶ We can define **counterfactual heritability** as

$$\begin{aligned}\xi &= \frac{\text{Var}(Y(G) - Y(G'))}{2 \text{Var}(Y(G))} = \frac{\text{Var}(Y(G, E) - Y(G', E))}{2 \text{Var}(Y(G, E))}. \quad (\text{Environment } E \text{ is shared.}) \\ &= 1 - \text{Cor}(Y(G), Y(G')).\end{aligned}$$

- ▶ Not hard to see that under the ACE model, $\xi = h_{\text{twin}}^2/2$. The factor $1/2$ is because ξ only measure heritability in **one generation**.
- ▶ In many way, this is a nice definition. A practical challenge is that, without further strong assumptions (e.g. additivity in the ACE model), the value of ξ **cannot be uniquely determined** even with infinite amount of data.

Distinguishing notions of heritability by the counterfactual comparison

- ▶ Geneticists talk about **broad-sense heritability**, **narrow-sense heritability**, and **SNP heritability**. Often only defined under specific models and are tricky to interpret.
- ▶ Let X be parent genotypes and G be children genotypes.



(a) Potential outcomes. (b) Fraternal counterfactual. (c) Unrelated counterfactual. (d) Adopted counterfactual.

- (b) **Fraternal**: “one-generation” heritability $\propto \text{Var}\{Y(X, G) - Y(X, G')\}$.
 - (c) **Unrelated**: “infinite-generation” heritability $\propto \text{Var}\{Y(X, G) - Y(X', G'')\}$.
 - (d) **Adopted**: heritability $\propto \text{Var}\{Y(X, G) - Y(X, G'')\}$.
- ▶ Implicit in relatedness disequilibrium regression (Young et al., 2018)² but not always ≤ 1 .

²Relatedness disequilibrium regression estimates heritability without environmental bias. *Nature Genetics*.

Partial identification (Theorem 2 in Lei, Shi, Cao, Zhao)

Suppose $G \perp\!\!\!\perp \{Y(g) : g \in \mathcal{G}\} \mid X$. Then

$$\xi'_I \leq \xi_I \leq \xi \leq \min\{\xi_u, \xi'_u\},$$

where

$$\begin{aligned}\xi'_I &= \frac{E[\text{Var}\{E(Y \mid G, X) \mid X\}]}{\text{Var}(Y)} = 1 - \frac{E(\text{Var}(Y \mid G, X)) + \text{Var}(E(Y \mid X))}{\text{Var}(Y)}, \\ \xi_I &= \frac{E\left[\left\{F_{G,X}^{-1}(U) - F_{G',X}^{-1}(U)\right\}^2\right]}{2 \text{Var}(Y)}, \quad \xi_u = \frac{E\left[\left\{F_{G,X}^{-1}(U) - F_{G',X}^{-1}(1-U)\right\}^2\right]}{2 \text{Var}(Y)}, \\ \xi'_u &= \frac{E(\text{Var}(Y \mid X))}{\text{Var}(Y)} = 1 - \frac{\text{Var}(E(Y \mid X))}{\text{Var}(Y)}.\end{aligned}$$

- ▶ $F_{G,X}^{-1}(\cdot)$ is the conditional quantile function of Y given G, X . $U \sim \text{Unif}[0, 1]$.
- ▶ The lower bound ξ_I is tight. The upper bound ξ_u is tight when G is binary but generally not very useful.

Comparison with other notions of heritability

- ▶ Let $H^2 = \frac{\text{Var}(E(Y | G))}{\text{Var}(Y)}$ and $h^2 = \frac{\text{Var}(G^T \theta)}{\text{Var}(Y)}$ (θ is the least squares projection).
- ▶ $G_1, G_2 \sim \text{Bern}(0.5)$, $X, E_1, E_2 \sim N(0, 0.25)$, **all independent**.

$Y(g) =$	Narrow h^2	Broad H^2	Counterfactual ξ	Lower ξ'_l	Tight lower ξ_l	Upper ξ'_u
$\beta_1 g_1 + \beta_2 g_2 + E_1$	$\frac{\beta_1^2 + \beta_2^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{\beta_1^2 + \beta_2^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{\beta_1^2 + \beta_2^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{\beta_1^2 + \beta_2^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{\beta_1^2 + \beta_2^2}{1 + \beta_1^2 + \beta_2^2}$	1
$\beta g_1 g_2 + E_1$	$\frac{2\beta^2}{4 + 3\beta^2}$	$\frac{3\beta^2}{4 + 3\beta^2}$	$\frac{3\beta^2}{4 + 3\beta^2}$	$\frac{3\beta^2}{4 + 3\beta^2}$	$\frac{3\beta^2}{4 + 3\beta^2}$	1
$\beta g_1 E_2 + E_1$	0	0	$\frac{\beta^2}{4 + 2\beta^2}$	0	$\frac{(\sqrt{1 + \beta^2} - 1)^2}{4 + 2\beta^2}$	1
$\beta_1 g_1 + \beta_2 X + E_1$	$\frac{\beta_1^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{\beta_1^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{\beta_1^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{\beta_1^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{\beta_1^2}{1 + \beta_1^2 + \beta_2^2}$	$\frac{1 + \beta_1^2}{1 + \beta_1^2 + \beta_2^2}$
$\beta g_1 X + E_1$	0	0	$\frac{\beta^2}{4 + 2\beta^2}$	$\frac{\beta^2}{4 + 2\beta^2}$	$\frac{\beta^2}{4 + 2\beta^2}$	$\frac{4 + \beta^2}{4 + 2\beta^2}$

Our second idea

- ▶ Let $\xi(G)$ be the counterfactual heritability above. Can similarly define $\xi(E)$ and $\xi(G \vee E)$.
- ▶ Row 3 above shows that G receives “credit” from the interaction term $G \times E$.
- ▶ Can we define **interaction explainability** using the inclusion-exclusion principle as

$$\xi(G \wedge E) = \xi(G) + \xi(E) - \xi(G \vee E)?$$

- ▶ It turns out that this can be rewritten as

$$\xi(G \wedge E) = \frac{\text{Var}(Y(\textcolor{blue}{G}, \textcolor{blue}{E}) - Y(\textcolor{red}{G}', \textcolor{blue}{E}) - Y(\textcolor{blue}{G}, \textcolor{red}{E}') + Y(\textcolor{red}{G}', \textcolor{red}{E}'))}{4 \text{Var}(Y)} \geq 0.$$

The **second-order finite difference** has indeed been used to describe interaction.

- ▶ We wrote a conference paper based on this and shared it with Professor Art Owen. He pointed out that this is the **superset importance** introduced by Giles Hooker³ with a variance-based formula in Liu and Owen (2006)⁴.

³Discovering additive structure in black box functions. KDD 2004.

⁴Estimating mean dimensionality of analysis of variance decompositions. JASA 2006.

Related literature: Functional ANOVA

- ▶ Let f be a function of independent random variables W_1, \dots, W_K .
- ▶ We can always write $f(W)$ as

$$f(W) = \sum_{S \subseteq [K]} f_S(W_S),$$

where the terms can be obtained inductively by (let $f_\emptyset(W) = E[f(W)]$)

$$f_S(w_S) = E \left[f(W) - \sum_{S' \subset S} f_{S'}(W) \mid W_S = w_S \right], \quad \text{for } S \subseteq [K].$$

- ▶ The terms are orthogonal: $E[f_S(W_S)f_{S'}(W_{S'})] = 0$ for all different $S, S' \subseteq [K]$.
- ▶ So we have the **functional ANOVA** decomposition:

$$\text{var}(f(W)) = \sum_{S \subseteq [K]} \sigma_S^2, \quad \text{where } \sigma_S^2 := \text{var}(f_S(W_S)).$$

- ▶ This goes back to Hoeffding (1948, *AOMS*).

Related literature: global sensitivity analysis

Some notions of variable importance for a subset $\mathcal{S} \subseteq [K]$:

- ▶ **Sobol's lower and upper sensitivity indices:** $\underline{\tau}_{\mathcal{S}}^2 = \sum_{\mathcal{S}' \subseteq \mathcal{S}} \sigma_{\mathcal{S}'}^2$, and $\bar{\tau}_{\mathcal{S}}^2 = \sum_{\mathcal{S}' \cap \mathcal{S} \neq \emptyset} \sigma_{\mathcal{S}'}^2$.
- ▶ **Super-set importance:** $\bar{\sigma}_{\mathcal{S}}^2 = \sum_{\mathcal{S}' \supseteq \mathcal{S}} \sigma_{\mathcal{S}'}^2$.
- ▶ **Shapley value** (not going into details...).

Pick-freeze method

The following formulas are used to accelerate GSA:

$$\underline{\tau}_{\mathcal{S}}^2 = \text{Cov}(f(W), f(W_{\mathcal{S}}, W'_{-\mathcal{S}})) \quad \text{and} \quad \bar{\tau}_{\mathcal{S}}^2 = \frac{1}{2} \mathbb{E}[\{f(W) - f(W'_{\mathcal{S}}, W_{-\mathcal{S}})\}^2],$$

$$\bar{\sigma}_{\mathcal{S}}^2 = 2^{-|\mathcal{S}|} \text{var}\{I_{\mathcal{S}}(W, W')\}, \quad \text{for all } \mathcal{S} \subseteq [K],$$

where $I_{\mathcal{S}}(w, w') = \sum_{\mathcal{S}' \subseteq \mathcal{S}} (-1)^{|\mathcal{S}| - |\mathcal{S}'|} f(w'_{\mathcal{S}'}, w_{-\mathcal{S}'})$ is an interaction contrast (which forms the

anchored decomposition, not going into details...)

Unified perspective: explainability as a probability measure

- The **explanation algebra** $\mathcal{E}(W)$ is the Boolean algebra generated by W_1, \dots, W_K using conjunction \vee , disjunction \wedge , and negation \neg .

Theorem 1 in Gao and Zhao

Let $\xi_1, \xi_2, \xi_3, \xi_4$ be any probability measures on $\mathcal{E}(W)$ such that, for all $\mathcal{S} \subseteq [K]$,

$$\xi_1((\wedge_{k \in \mathcal{S}} W_k) \wedge (\wedge_{k \notin \mathcal{S}} \neg W_k)) = \frac{\sigma_{\mathcal{S}}^2}{\text{var}(f(W))},$$

$$\xi_2(\neg(\vee_{k \notin \mathcal{S}} W_k)) = \frac{\underline{\tau}_{\mathcal{S}}^2}{\text{var}(f(W))},$$

$$\xi_3(\vee_{k \in \mathcal{S}} W_k) = \frac{\bar{\tau}_{\mathcal{S}}^2}{\text{var}(f(W))},$$

$$\xi_4(\wedge_{k \in \mathcal{S}} W_k) = \frac{\bar{\sigma}_{\mathcal{S}}^2}{\text{var}(f(W))}.$$

Then $\xi_1 = \xi_2 = \xi_3 = \xi_4$.

Implications

This is essentially a re-formulation of functiona ANOVA, but we can use the familiar probability theory to derive many implications:

1. $\xi(\vee_{k \in \mathcal{S}} W_k) = 1$.
2. $\xi(\vee_{k \in \mathcal{S}'} W_k) \geq \xi(\vee_{k \in \mathcal{S}} W_k)$ and $\xi(\wedge_{k \in \mathcal{S}'} W_k) \leq \xi(\wedge_{k \in \mathcal{S}} W_k)$ for any $\mathcal{S} \subseteq \mathcal{S}' \subseteq [K]$.
3. $\xi(W_1) + \dots \xi(W_K) \geq \xi(\vee_{k \in [K]} W_k)$.

Remarks

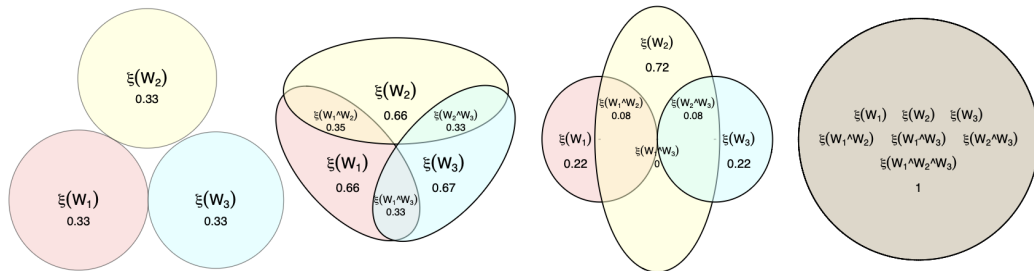
- ▶ 2.2 is what motivated Hooker's **superset-importance** (name is unhelpful...).
- ▶ 3 is the **Efron-Stein inequality**

$$\text{Var}(f(W)) \leq \sum_{k=1}^K \text{Var}(f(W) - f(W'_k, W_{-k})).$$

This can be improved by the Boole-Bonferroni inequalities.

- ▶ For practice, the important thing is that we can use **Venn's diagram** to visualize ξ .

Examples



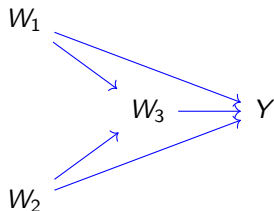
Left to right:

1. Linear: $f(w) = w_1 + w_2 + w_3$;
2. Quadratic: $f(w) = w_1 w_2 + w_1 w_3 + w_2 w_3$;
3. Single-layer NN (sigmoid activation): $f(w) = (1 + e^{10w_1 + 10w_2})^{-1} + (1 + e^{10w_2 + 10w_3})^{-1}$;
4. Multilinear monomial: $f(w) = w_1 w_2 w_3$.

Third idea

- ▶ How to extend this to dependent explanations?
- ▶ Our idea is to use a directed acyclic graph (DAG) to model their causal dependence and **define “fraternal counterfactuals” by resampling the intrinsic noise.**

Example



Pearl's NPSEM-IE assumes (in potential outcomes terms)

$$W_1 = f_1(E_1),$$

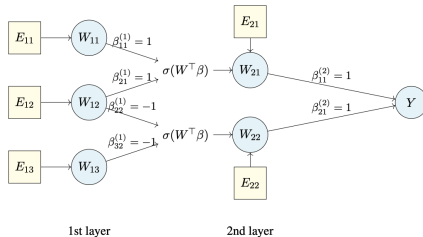
$$W_2 = f_2(E_2),$$

$$W_3(w_1, w_2) = f_3(w_1, w_2, E_3),$$

$$Y(w_1, w_2, w_3) = f_4(w_1, w_2, w_3, E_4),$$

- ▶ By recursive substitution, we can write Y as a function of the intrinsic noises E_1, E_2, E_3, E_4 , which are assumed to be independent in Pearl's model.
- ▶ Janzing et al. (2024, AISTATS) suggested the same idea for causal extensions of general dependence measures (e.g. mutual information) in a more conceptual paper.

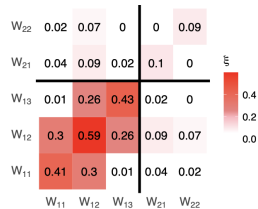
Toy example



(a) Network structure.

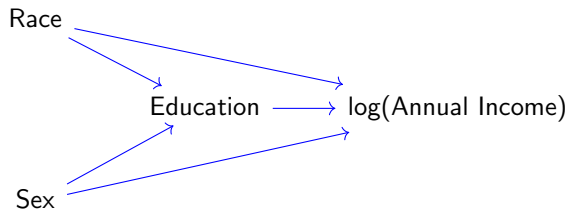


(b) Equal noise variance.



(c) Second layer noise variance decreased.

Real data example: Explaining income inequality



- ▶ Datasets: UCI Adult (1994) and ACSIncome (2018).
- ▶ We assumed the basic potential outcomes are **comonotone** (\Rightarrow **point identification**) and sampled the counterfactuals after estimating the conditional distributions using XGBoost.

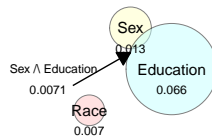
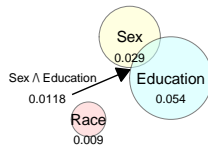
Results

Age

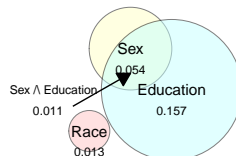
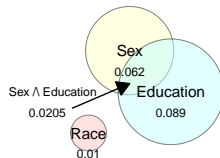
Income in 1994

Income in 2018

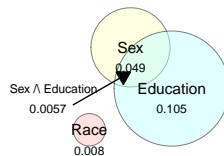
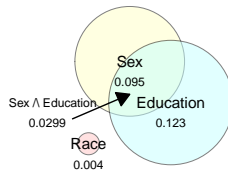
[25, 30)



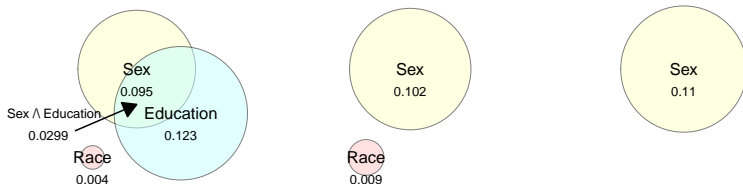
[40, 45)



Over 60



Consistency across graphs



- ▶ Counterfactual explainability generally depends on the causal DAG but has certain consistency properties with respect to finer mechanistic explanations (see the paper for some formal results).
- ▶ In this example, Sex should have the same total explainability regardless of whether Race and Education are included.

Conclusion

The papers have more theorems (inc. an axiomatization of ξ), examples, and discussion.

Why is counterfactual heritability/explainability a good idea?

In theory

It provides a counterfactual extension to functional ANOVA that can be applied with dependent explanations and non-deterministic outputs. Has certain consistency properties.

In practice

Uses contextual mechanistic information. Easy to visualize and interpret.

Thank you!