

Sensitivity Analysis with the R^2 -Calculus

Tobias Freidling Qingyuan Zhao
Statistical Laboratory, University of Cambridge

Sensitivity Analysis as Optimisation Problem

Model

- Model for observed and unobserved variables, O and U : $(O, U) \sim \mathbb{P}_{\theta, \psi}$
- θ parametrises the observable and ψ the unobservable aspects of $\mathbb{P}_{\theta, \psi}$.
- Quantity of interest: $g_n(\hat{\theta}, \psi)$, e.g. a point estimate or confidence interval for a causal effect.

Sensitivity Analysis

- Causal assumptions like 'no unmeasured confounding' correspond to fixing the value of ψ .
- Sensitivity analysis considers a set of plausible values Ψ instead.
- The range of values of g_n is the solution of

$$\max / \min g_n(\hat{\theta}, \psi) \quad \text{subject to } \psi \in \Psi.$$

Challenges

- Identification of $g_n(\hat{\theta}, \psi)$ in terms of the sensitivity parameters ψ
- Translating domain knowledge into the constraints $\psi \in \Psi$

Using interpretable R^2 -values helps practitioners to express their beliefs about the unmeasured confounder. The rules of the R^2 -calculus translate these into Ψ .

R^2 -Calculus

The R^2 -calculus assembles algebraic rules for R^2 -values and correlations as a coherent system of their own.

Definitions Let $Y \in \mathbb{R}^n$, $X \in \mathbb{R}^{n \times p}$ and $W \in \mathbb{R}^{n \times q}$ be n i.i.d. samples.

- R^2 -value and partial R^2 -value:

$$R_{Y \sim X}^2 := 1 - \frac{\text{var}(Y - X\hat{\beta}_X)}{\text{var}(Y)}, \quad R_{Y \sim X|W}^2 := \frac{R_{Y \sim X+W}^2 - R_{Y \sim W}^2}{1 - R_{Y \sim W}^2}.$$

- R -value: $R_{Y \sim X} := \text{corr}(Y, X)$ for $X \in \mathbb{R}^n$.

- f -value: $f_{Y \sim X} := R_{Y \sim X} / \sqrt{1 - R_{Y \sim X}^2}$.

Some Calculation Rules

- Decomposition of unexplained variance:

$$1 - R_{Y \sim X+W}^2 = (1 - R_{Y \sim X|W}^2)(1 - R_{Y \sim W}^2)$$

- Recursive partial correlation formula:

$$R_{Y \sim X|W} = \frac{R_{Y \sim X} - R_{Y \sim W} R_{X \sim W}}{\sqrt{1 - R_{Y \sim W}^2} \sqrt{1 - R_{X \sim W}^2}}, \quad \text{for } X, W \in \mathbb{R}^n$$

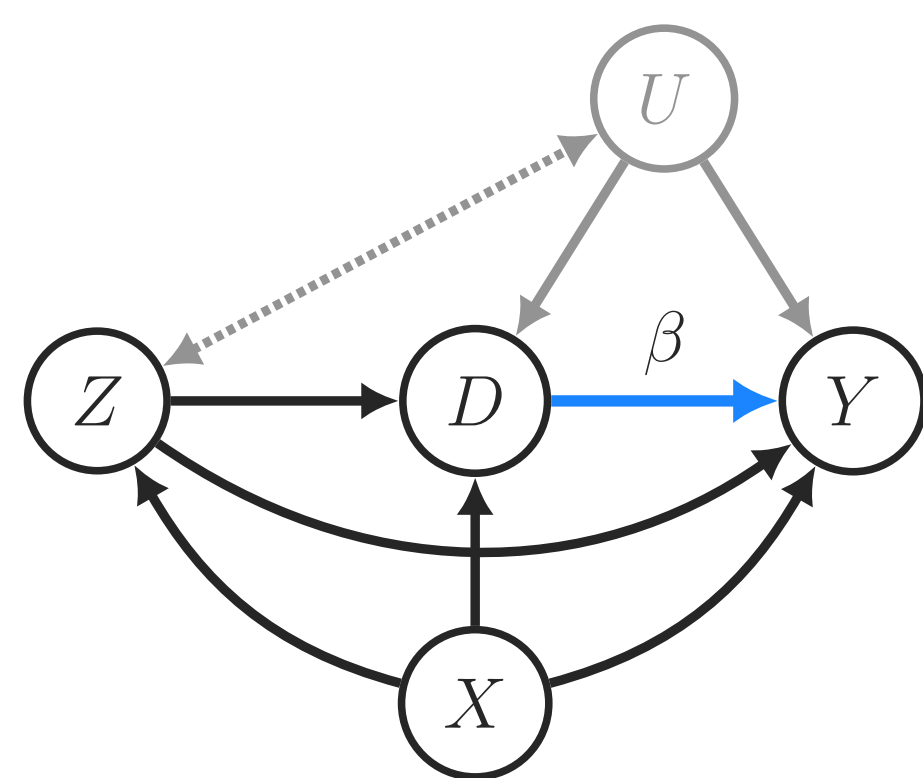
Bias of the k -class estimator

Linear model with continuous outcome $Y \in \mathbb{R}^n$, observed variables $D \in \mathbb{R}^n$, $Z \in \mathbb{R}^n$, $X \in \mathbb{R}^{n \times p}$ and unmeasured confounder $U \in \mathbb{R}^n$:

$$Y := \nu + D\beta + Z\delta + X\xi + U\lambda + \varepsilon.$$

We estimate β , the causal effect of D on Y , with the k -class estimator

$$\hat{\beta}_k = \frac{\text{cov}(D^{\perp X}, Y^{\perp X}) - k \text{cov}(D^{\perp Z, X}, Y^{\perp Z, X})}{\text{var}(D^{\perp X}) - k \text{var}(D^{\perp Z, X})}.$$



It interpolates between the OLS-estimator ($k \rightarrow -\infty$) and the IV-estimator ($k = 1$) with instrumental variable Z .

Via the R^2 -calculus and a technical result in [1], we get:

$$\hat{\beta}_k - \hat{\beta}_{\text{OR}} = \left[\frac{f_{Y \sim Z|X, D} R_{D \sim Z|X} + R_{Y \sim U|X, Z, D} f_{D \sim U|X, Z}}{1 - k + k R_{D \sim Z|X}^2} \right] \frac{\text{sd}(Y^{\perp X, Z, D})}{\text{sd}(D^{\perp X, Z})}.$$

Similarly, we can identify the ends of a $1 - \alpha$ confidence interval. This extends previous work on the OLS-estimator [2].

We choose the sensitivity parameters $\psi = (R_{D \sim U|X, Z}, R_{Y \sim U|X, Z, D})$.

References

- Carrie A. Hosman, Ben B. Hansen, and Paul W. Holland. The sensitivity of linear regression coefficients' confidence limits to the omission of a confounder. *The Annals of Applied Statistics*, 4(2):849 – 870, 2010.
- Carlos Cinelli and Chad Hazlett. Making sense of sensitivity: extending omitted variable bias. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(1):39–67, 2020.
- David Card. Using Geographic Variation in College Proximity to Estimate the Return to Schooling. Technical Report w4483, National Bureau of Economic Research, 1993.
- Tobias Freidling and Qingyuan Zhao. Sensitivity Analysis with the R^2 -calculus. *Working Paper (available upon request)*, 2022.

Specifying Bounds Ψ

Assume that $X = (\tilde{X}, \hat{X})$ can be separated such that $\hat{X} \perp U | \tilde{X}$ holds.

Linear Regression

- Choose constraints from the $U \rightarrow D$ and $U \rightarrow Y$ block of the table below.
- Range bounds: direct specification of ψ , e.g. $R_{D \sim U|X, Z} \in [-0.3, 0.2]$.
- Comparative bounds: benchmarking against observed covariates, e.g. $R_{D \sim U|X, Z}^2 \leq 0.5 R_{D \sim \hat{X}_j | \tilde{X}, \hat{X}_j, Z}^2$ means that U can explain at most half as much variability in D than \hat{X}_j can after partialling out $(\tilde{X}, \hat{X}_j, Z)$.

Comparative bounds are translated into constraints on ψ with the R^2 -calculus.

Instrumental Variable

- Choose any constraints from the table.
- R -values that parametrise the direct effect of Z on Y and the correlation between Z and U are connected to ψ via

$$f_{Y \sim Z|X, U, D} \sqrt{1 - R_{Y \sim U|X, D, Z}^2} = f_{Y \sim Z|X, D} \sqrt{1 - R_{Z \sim U|X, D}^2} - R_{Y \sim U|X, D, Z} R_{Z \sim U|X, D},$$

$$f_{Z \sim U|X, D} \sqrt{1 - R_{D \sim U|X, Z}^2} = f_{Z \sim U|X} \sqrt{1 - R_{D \sim Z|X}^2} - R_{D \sim Z|X} R_{D \sim U|X, Z}.$$

- Add these equations as constraints to the optimisation problem.

$U \rightarrow D$	1. $R_{D \sim U X, Z} \in [B_l^D, B_u^D]$ 2. $R_{D \sim U X, Z}^2 \leq \eta_D R_{D \sim \hat{X}_j \tilde{X}, \hat{X}_j, Z}^2$
$U \rightarrow Y$	1. $R_{Y \sim U X, Z, D} \in [B_l^Y, B_u^Y]$ 2. $R_{Y \sim U X, Z, D}^2 \leq \eta_Y R_{Y \sim \hat{X}_j \tilde{X}, \hat{X}_j, Z, D}^2$ 3. $R_{Y \sim U X, Z, D}^2 \leq \eta_Y R_{Y \sim \hat{X}_j \tilde{X}, \hat{X}_j, Z, D}^2$
$U \leftrightarrow Z$	1. $R_{Z \sim U X} \in [B_l^Z, B_u^Z]$ 2. $R_{Z \sim U X, Z}^2 \leq \eta_Z R_{Z \sim \hat{X}_j \tilde{X}, \hat{X}_j, Z}^2$
$Z \rightarrow Y$	1. $R_{Y \sim Z X, U, D} \in [B_l^{YZ}, B_u^{YZ}]$ 2. $R_{Y \sim Z X, U, D}^2 \leq \eta_{YZ} R_{Y \sim \hat{X}_j \tilde{X}, \hat{X}_j, U, D}^2$

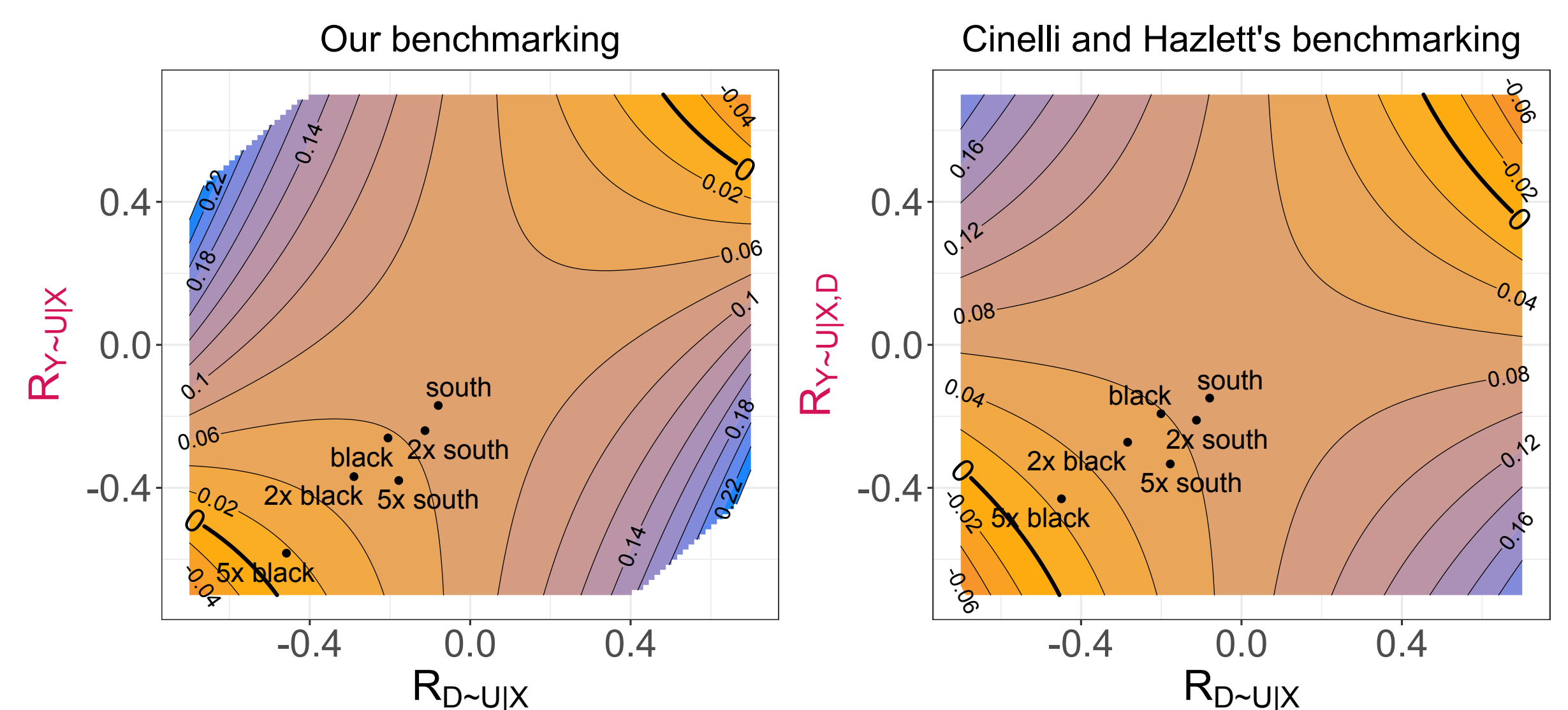
Table Specify interpretable bounds and use the R^2 -calculus to translate them into the constraints Ψ of the optimisation problem.

Insights

Data example: Inference on the causal effect of education on earnings, cf. [3].

Linear Regression

Lower end of the 95% confidence interval for different values of the sensitivity parameters with comparison points:



Instrumental Variable

We use college proximity as instrument for education and set the bounds

$$R_{D \sim U|X, Z} \in [-0.9, 0.9], \quad R_{Y \sim U|X, Z, D}^2 \leq 5 R_{Y \sim X_j | X_{-j}, Z, D}^2,$$

where X_j is an indicator for being black. For different IV-related bounds, we get

