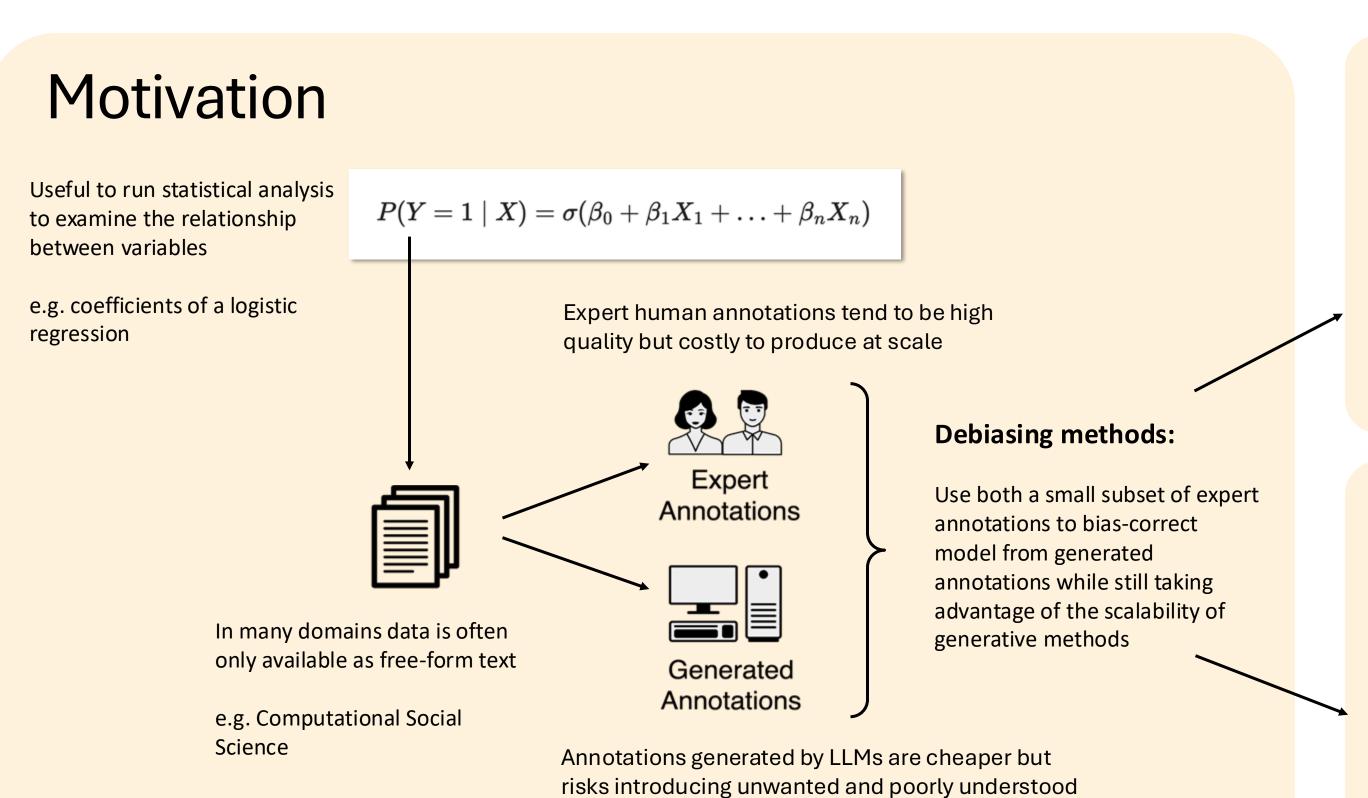
Benchmarking Debiasing Methods for LLM-based Parameter Estimates



Nicolas Audinet de Pieuchon, Adel Daoud, Connor T. Jerzak, Moa Johansson, Richard Johansson





Prediction-Powered Inference (PPI)*

Correct imputation estimate with rectifier:

$$\hat{ heta} = ilde{ heta} - r_{ heta}$$

Compute rectifier from gradients:

$$r_{ heta} = \mathbb{E}[
abla_{ heta}\ell_{ heta}(x_i,y_i) -
abla_{ heta}\ell_{ heta}(x_i,\hat{y}_i)]$$

*Angelopoulos, Anastasios N., et al. "Prediction-powered inference." Science 382.6671 (2023): 669-674

Design-based Supervised Learning (DSL)[†]

Assume known expert labelling distribution:

$$\pi(x_i,\hat{y}_i) = P(R_i = 1 \mid x_i,\hat{y}_i)$$

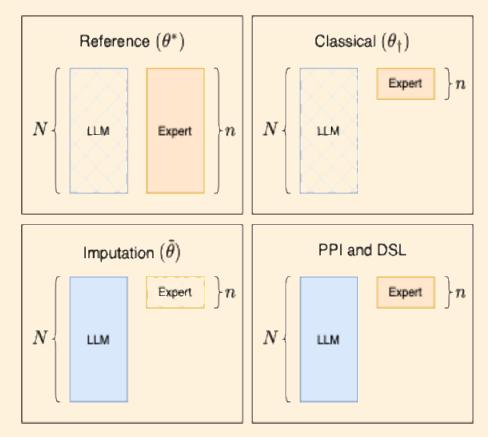
Double-robust estimate:

$$ilde{y_i} = \hat{y_i} - rac{R_i}{\pi_i}(\hat{y_i} - y_i)$$

 $\mathbb{E}[ilde{y_i} - y_i \mid x_i] = 0$

†Egami, Naoki, et al. "Using large language model annotations for the social sciences: A general framework of using predicted variables in downstream analyses." Preprint from November 17 (2024): 2024.

Setup



For all experiments:

- Logistic regression
 - 4 discrete input features
 - Binary output
- 4 different LLMs
- 4 different datasets
- Evaluation: standardised RMSE with a reference model

Results

biases into the downstream analysis.\

RQ1: When is it preferable to use debiasing methods over just the expert annotations?

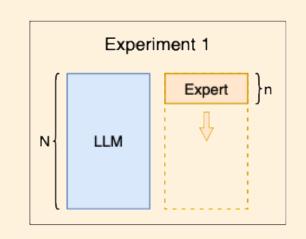
Both DSL and PPI strictly outperform using only expert annotations.

RQ2: What are the performance differences between debiasing methods?

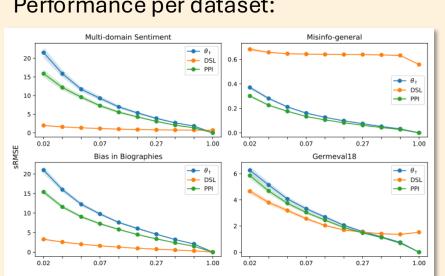
DSL tends to outperform PPI, but performance is more variable.

Experiment 1

Vary the proportion of samples also annotated by human experts:

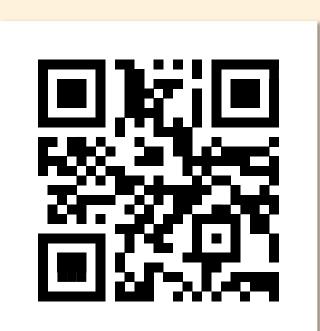


15.0 -0.07 0.27 Proportion of expert samples (log) Performance per dataset:



Given a fixed number of samples, how many of them should I annotate by hand to observe significant benefits from the debiasing methods?

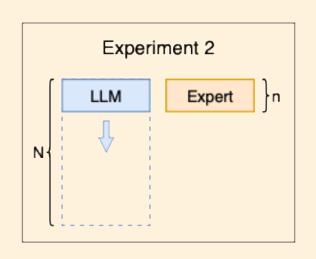
Scan to read the full paper!

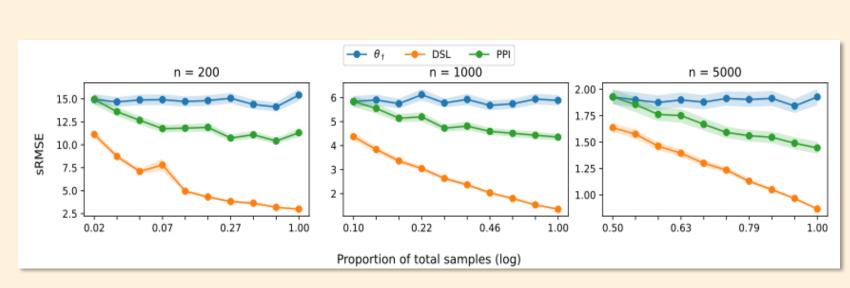




Experiment 2

Vary the number of samples with generated annotations:





Given a fixed number of expert-annotated samples, how many additional samples should I annotate with LLMs? Are there diminishing returns?