Benchmarking Debiasing Methods for LLM-based Parameter Estimates

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

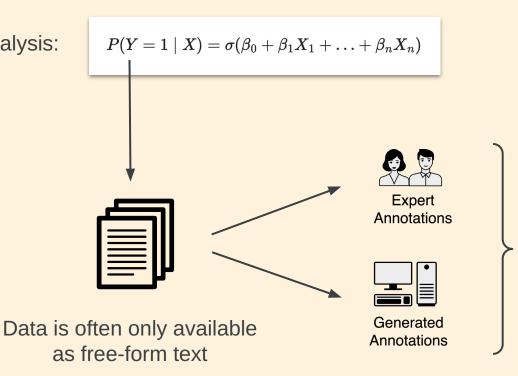






Motivation

Statistical analysis:



Debiasing methods:

Use small subset of expert annotations to bias-correct model from generated annotations.

Prediction-Powered Inference (PPI)*

Correct imputation estimate with rectifier:

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

Compute rectifier from gradients:

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

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:

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

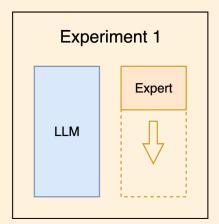
$$\mathbb{E}[ilde{y_i} - y_i \mid x_i] = 0$$

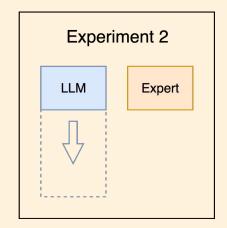
Research Questions

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

RQ2: What are the performance differences between debiasing methods?

Experiments

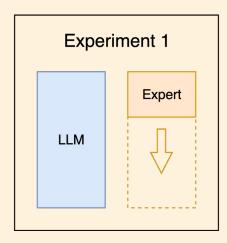


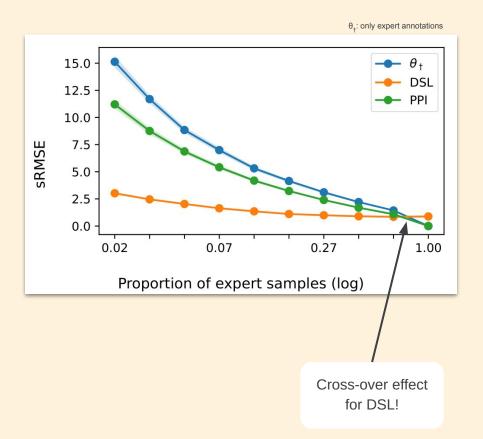


For both experiments:

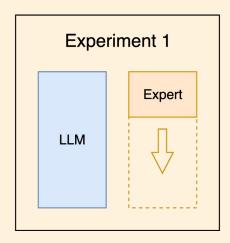
- Logistic regression
 - 4 discrete input features
 - Binary output
- 4 different LLMs
- 4 different datasets
- Evaluation: standardised RMSE with a reference model

Results

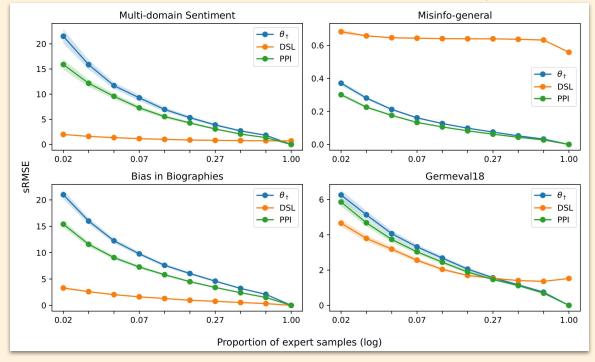




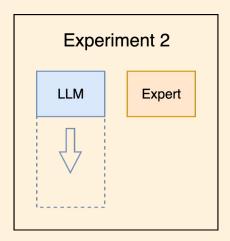
Results

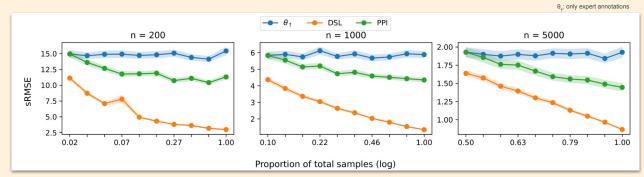


 θ_{t} : only expert annotations



Results





Conclusions

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.

Thank you for listening!





