

Ambiguity-aware Artificial Intelligence via Statistical Inference

Edgar Dobriban

Department of Statistics and Data Science, the Wharton School, University of Pennsylvania

August 12, 2024



Table of Contents

Motivation

Uncertainty quantification in AI/ML

AI Evaluation

Algorithmic Fairness

Context

- ▶ Artificial Intelligence (AI) is emerging as an important technology

Context

- ▶ Artificial Intelligence (AI) is emerging as an important technology
- ▶ Success stories: Large language models (ChatGPT, Claude), Image generation (Dall-E, Stable Diffusion), Protein structure prediction ...

Claude Q Visualizing Missing at Random Data ▾

ED Can you create a visualization to illustrate the notion of missing at random? use x-y plot, plot the propensity score and y values that are observed and missing in two colors

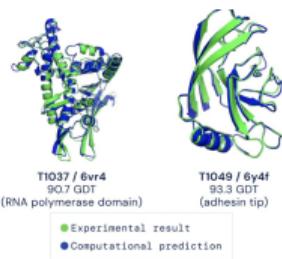
Certainly, I can create a visualization to illustrate the concept of Missing at Random (MAR) using an x-y plot with propensity scores and observed/missing y values in two colors. I'll use a React component to create this visualization.

</> Missing at Random Visualization
Click to open component

```
← Python Code for Propensity Score Plot with Horizontal Lines
plt.scatter(x, y_obs, c='blue', s=100, alpha=0.5, edgecolor='black', linewidth=1)
plt.scatter(x, y_miss, c='red', s=100, alpha=0.5, edgecolor='black', linewidth=1)

# Set Labels and title
plt.xlabel('x', fontsize=18)
plt.ylabel('Propensity Score', fontsize=18)
plt.title('Propensity Score with Horizontal Lines', fontsize=22, fontweight='bold')

# Set axis limits
plt.xlim(-5, 5)
plt.ylim(0, 1)
```



Goals of research

- ▶ What should statisticians do to maximize impact during this time?

Goals of research

- ▶ What should statisticians do to maximize impact during this time?
- ▶ What new methods should be developed?

Goals of research

- ▶ What should statisticians do to maximize impact during this time?
- ▶ What new methods should be developed?
- ▶ Address emerging areas of research where statistics could have an impact:
 - ▶ Uncertainty quantification (Ambiguity-awareness) in AI/ML

Goals of research

- ▶ What should statisticians do to maximize impact during this time?
- ▶ What new methods should be developed?
- ▶ Address emerging areas of research where statistics could have an impact:
 - ▶ Uncertainty quantification (Ambiguity-awareness) in AI/ML
 - ▶ AI evaluation

Goals of research

- ▶ What should statisticians do to maximize impact during this time?
- ▶ What new methods should be developed?
- ▶ Address emerging areas of research where statistics could have an impact:
 - ▶ Uncertainty quantification (Ambiguity-awareness) in AI/ML
 - ▶ AI evaluation
 - ▶ Algorithmic fairness
 - ▶ ...

Table of Contents

Motivation

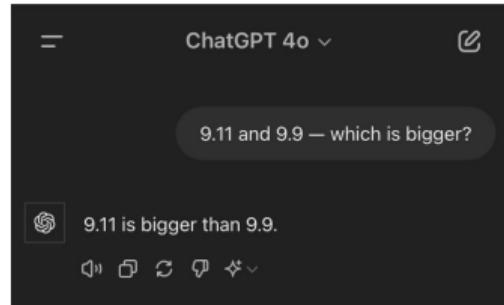
Uncertainty quantification in AI/ML

AI Evaluation

Algorithmic Fairness

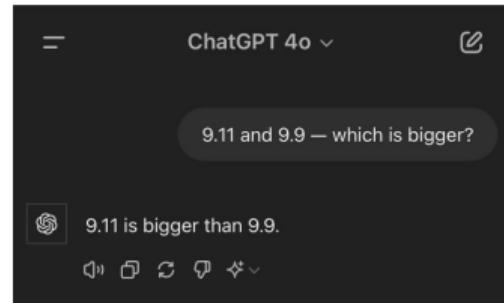
Uncertainty quantification in AI/ML

- ▶ AI systems can be wrong



Uncertainty quantification in AI/ML

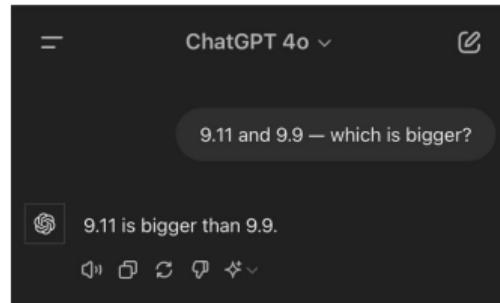
- ▶ AI systems can be wrong



- ▶ One emerging idea: Quantify uncertainty (and then perhaps refrain from generating when it is high)

Uncertainty quantification in AI/ML

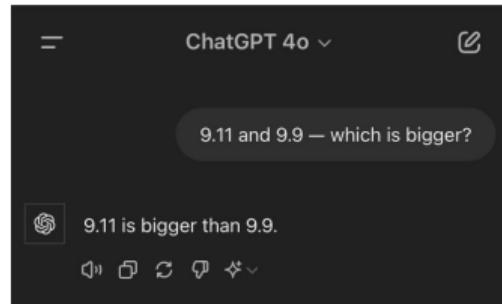
- ▶ AI systems can be wrong



- ▶ One emerging idea: Quantify uncertainty (and then perhaps refrain from generating when it is high)
- ▶ But what does it mean to quantify uncertainty? and what properties does good uncertainty quantification have?

Uncertainty quantification in AI/ML

- ▶ AI systems can be wrong



- ▶ One emerging idea: Quantify uncertainty (and then perhaps refrain from generating when it is high)
- ▶ But what does it mean to quantify uncertainty? and what properties does good uncertainty quantification have?

Uncertainty in Language Models: Assessment through Rank-Calibration

Xinmeng Huang^{*†}

Shuo Li^{*}

Mengxin Yu[†]

Matteo Sesia[‡]

Hamed Hassani[†]

Insup Lee[†]

Osbert Bastani^{§†}

Edgar Dobriban^{§†}

Example uncertainty measures

- ▶ Many uncertainty measures for LLMs have been discussed/proposed
 - ▶ Perplexity: $U(x, y) = \hat{p}(y|x)^{1/\text{len}(y)}$; [related to NLL – $\log \hat{p}(y|x)$]

Example uncertainty measures

- ▶ Many uncertainty measures for LLMs have been discussed/proposed
 - ▶ Perplexity: $U(x, y) = \hat{p}(y|x)^{1/\text{len}(y)}$; [related to NLL – $\log \hat{p}(y|x)$]
 - ▶ Semantic entropy: generate multiple y s, cluster them based on meaning, calculate entropy (Kuhn et al., 2023)

Answers to the question “What is the capital of France?”

Answer s	Likelihood $p(s x)$	Semantic likelihood $\sum_{s \in c} p(s x)$
Paris	0.5	0.9
It's Paris	0.4	
London	0.1	0.1
Entropy	0.94	0.33

Example uncertainty measures

- ▶ Many uncertainty measures for LLMs have been discussed/proposed
 - ▶ Perplexity: $U(x, y) = \hat{p}(y|x)^{1/\text{len}(y)}$; [related to NLL – $\log \hat{p}(y|x)$]
 - ▶ Semantic entropy: generate multiple ys, cluster them based on meaning, calculate entropy (Kuhn et al., 2023)

Answers to the question “What is the capital of France?”

Answer s	Likelihood $p(s x)$	Semantic likelihood $\sum_{s \in c} p(s x)$
Paris	0.5	0.9
It's Paris	0.4	
London	0.1	0.1
Entropy	0.94	0.33

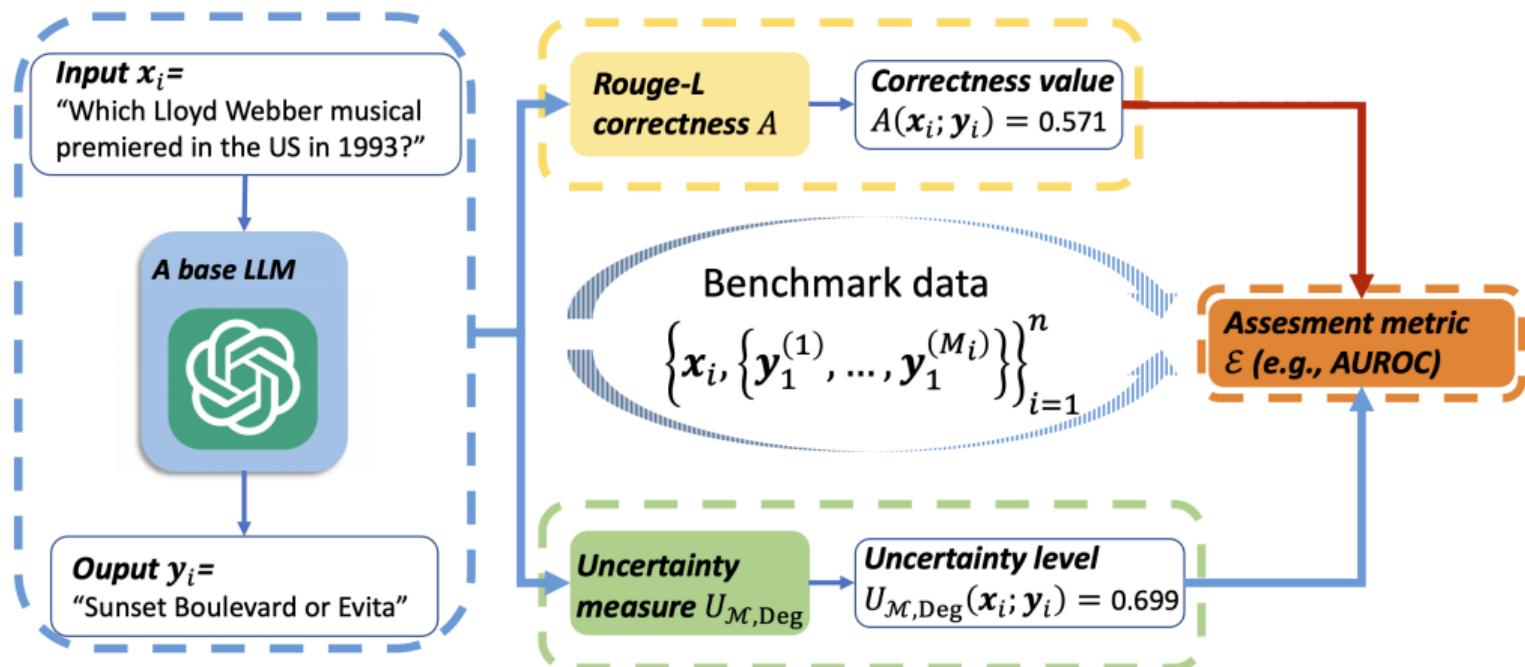
- ▶ Affinity graph: generate multiple ys, put pairwise similarities in a matrix, find eigenvalues/vectors (denoted EigV) (Lin et al., 2023)

Good uncertainty measures

- Our idea: *performance should decrease as a function of the uncertainty*

Good uncertainty measures

- Our idea: *performance should decrease as a function of the uncertainty*



Regression function and indication diagram

- ▶ Define the **regression function** $\text{reg}(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$, $u \mapsto \mathbb{E}_{\mathbf{x}, \mathbf{y}}[A(\mathbf{x}; \mathbf{y}) \mid U(\mathbf{x}; \mathbf{y}) = u]$, representing the *expected correctness level A given an uncertainty level U = u*.

Regression function and indication diagram

- ▶ Define the **regression function** $\text{reg}(\cdot) : \mathbb{R} \rightarrow \mathbb{R}, u \mapsto \mathbb{E}_{\mathbf{x}, \mathbf{y}}[A(\mathbf{x}; \mathbf{y}) \mid U(\mathbf{x}; \mathbf{y}) = u]$, representing the *expected correctness level A given an uncertainty level U = u*.
- ▶ **Indication diagram:** plot of estimated regression function as a fn. of U-percentiles.

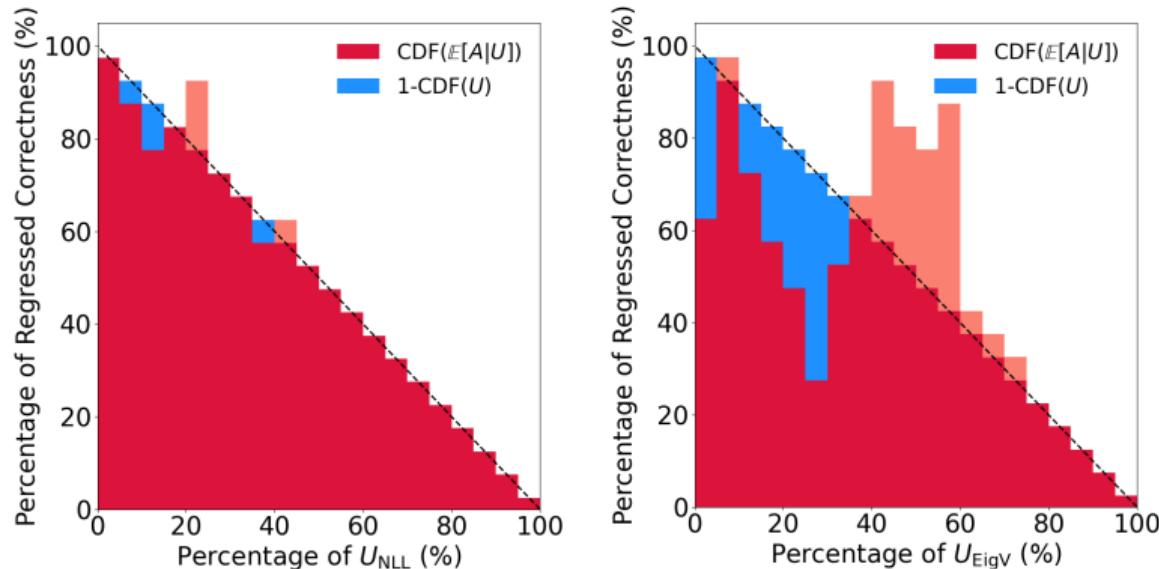


Figure: Indication diagrams for U_{NLL} (negative log-likelihood) and U_{EigV} , for the GPT-3.5-turbo model on the TriviaQA benchmark.

Rank-calibration

Definition (RANK-CALIBRATION)

An uncertainty measure U is *rank-calibrated* if reg is strictly *monotone decreasing*: on average, lower uncertainty implies higher generative quality.

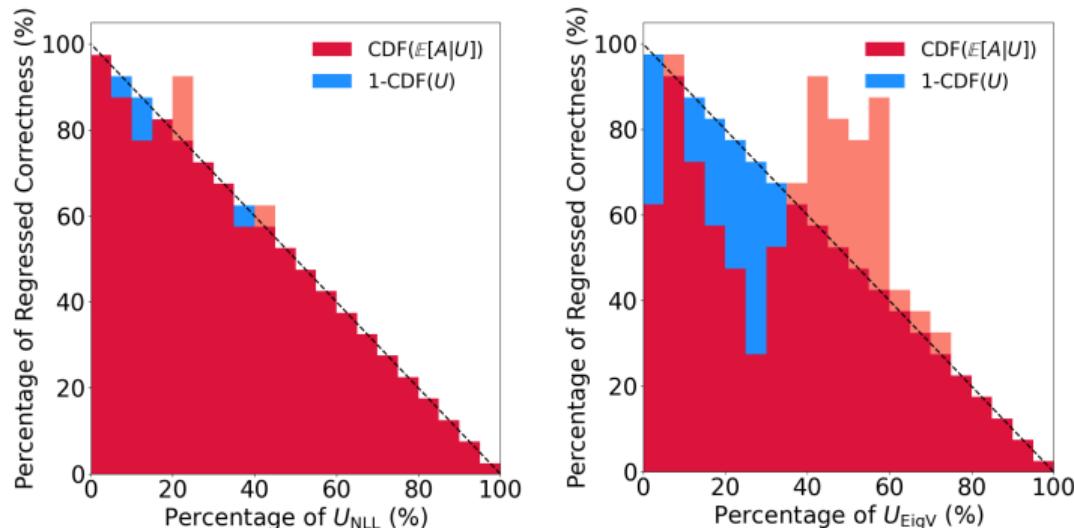


Figure: Indication diagrams for U_{NLL} (negative log-likelihood) and U_{EigV} , for the GPT-3.5-turbo model on the TriviaQA benchmark.

Rank-calibration error (RCE)

Rank-calibration implies $\mathbb{P}(U \leq u') = \mathbb{P}(\text{reg}(U) \geq \text{reg}(u'))$ for all $u' \in \text{dom}(U)$.

Rank-calibration error (RCE)

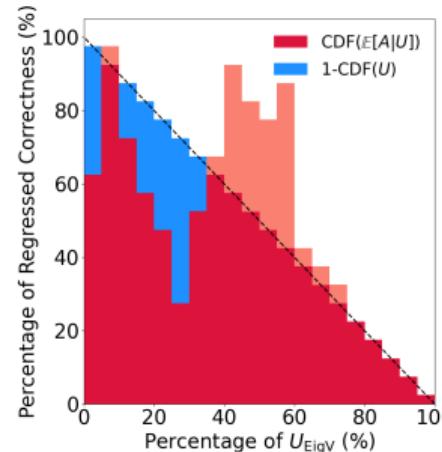
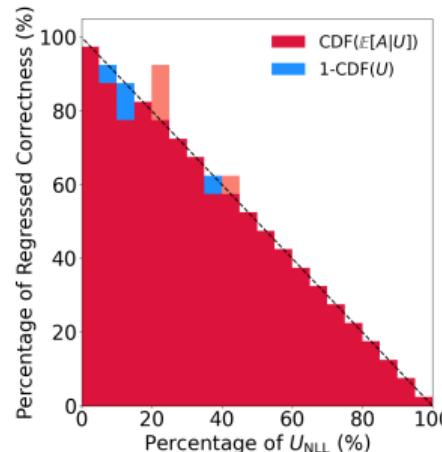
Rank-calibration implies $\mathbb{P}(U \leq u') = \mathbb{P}(\text{reg}(U) \geq \text{reg}(u'))$ for all $u' \in \text{dom}(U)$.

Definition (RANK-CALIBRATION ERROR)

The RCE of an uncertainty measure U is defined as

$$\mathbb{E}_{U'} [|\mathbb{P}_U(U \leq U') - \mathbb{P}_U(\text{reg}(U) \geq \text{reg}(U'))|],$$

where U' is an independent copy of U .



Empirical results

correctness	temperature	U_{Deg}	U_{EigV}	U_{NLL}	U_{SE}
bert	0.5	0.212 ± 0.040	0.212 ± 0.041	0.043 ± 0.006	0.052 ± 0.009
	1.0	0.129 ± 0.020	0.133 ± 0.020	0.039 ± 0.007	0.052 ± 0.012
	1.5	0.053 ± 0.011	0.074 ± 0.012	0.031 ± 0.007	0.081 ± 0.009
meteor	0.5	0.211 ± 0.045	0.208 ± 0.047	0.179 ± 0.021	0.234 ± 0.019
	1.0	0.131 ± 0.024	0.131 ± 0.022	0.146 ± 0.011	0.209 ± 0.012
	1.5	0.059 ± 0.011	0.077 ± 0.012	0.119 ± 0.010	0.176 ± 0.015
rougeL	0.5	0.210 ± 0.042	0.207 ± 0.041	0.041 ± 0.007	0.050 ± 0.008
	1.0	0.126 ± 0.019	0.129 ± 0.019	0.038 ± 0.007	0.059 ± 0.009
	1.5	0.059 ± 0.012	0.079 ± 0.011	0.034 ± 0.008	0.104 ± 0.007
rouge1	0.5	0.212 ± 0.043	0.209 ± 0.042	0.040 ± 0.007	0.050 ± 0.008
	1.0	0.126 ± 0.018	0.130 ± 0.021	0.039 ± 0.007	0.060 ± 0.009
	1.5	0.060 ± 0.011	0.078 ± 0.012	0.034 ± 0.008	0.105 ± 0.008

Table: RCE results for various experimental configurations.

Finding: models often have large rank calibration error (so uncertainty measures do not reflect performance).

Reducing RCE via re-calibration

- ▶ Re-calibrate uncertainty measure: change it to piece-wise constant estimate of regression function; Inspired by classical re-calibration (Mincer-Zamowitz, 1969)

Reducing RCE via re-calibration

- ▶ Re-calibrate uncertainty measure: change it to piece-wise constant estimate of regression function; Inspired by classical re-calibration (Mincer-Zamowitz, 1969)
- ▶ Improves RCE:

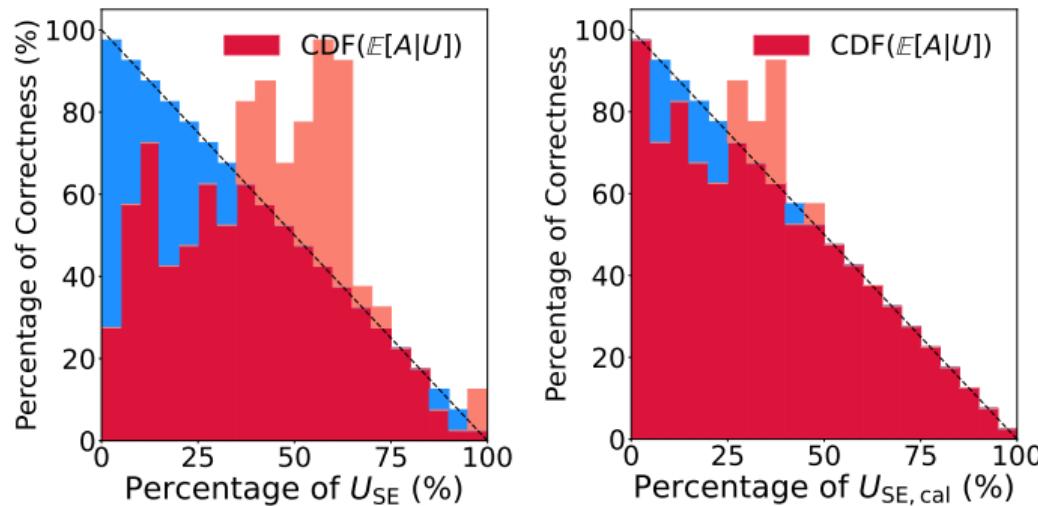


Figure: Indication diagrams of U_{SE} and $U_{SE,cal}$ (re-calibrated) for GPT-3.5-turbo (temperature 1.0) on TriviaQA with the Meteor correctness.

Summary for uncertainty quantification in AI/ML

- ▶ Proposed rank-calibration for LLMs: higher uncertainty should imply lower quality (not satisfied by default!)

Summary for uncertainty quantification in AI/ML

- ▶ Proposed rank-calibration for LLMs: higher uncertainty should imply lower quality (not satisfied by default!)
- ▶ Proposed non-parametric re-calibration to reduce RCE

Summary for uncertainty quantification in AI/ML

- ▶ Proposed rank-calibration for LLMs: higher uncertainty should imply lower quality (not satisfied by default!)
- ▶ Proposed non-parametric re-calibration to reduce RCE
- ▶ Showcased the utility of ideas from statistics

Uncertainty in Language Models: Assessment through Rank-Calibration

Xinmeng Huang^{*†}

Shuo Li^{*†}

Mengxin Yu[†]

Matteo Sesia[‡]

Hamed Hassani[†]

Insup Lee[†]

Osbert Bastani^{§†}

Edgar Dobriban^{§†}

Table of Contents

Motivation

Uncertainty quantification in AI/ML

AI Evaluation

Algorithmic Fairness

AI evaluation

- ▶ Evaluating AI can be hard

AI evaluation

- ▶ Evaluating AI can be hard
- ▶ Say want to measure *reasoning ability*:
 - ▶ Collect test data on reasoning problems (e.g., high school math)
 - ▶ Evaluate error rate

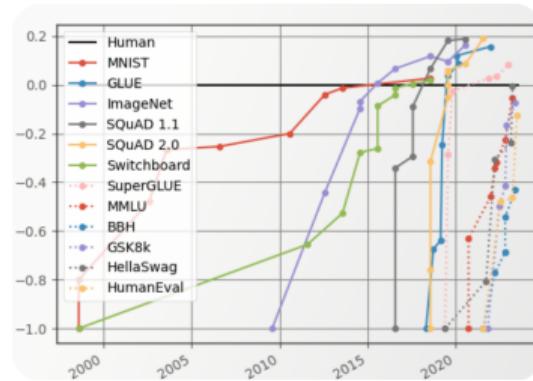


Figure: Plotting progress in AI (Kiel et al. 2023)

AI evaluation

- ▶ Evaluating AI can be hard
- ▶ Say want to measure *reasoning ability*:
 - ▶ Collect test data on reasoning problems (e.g., high school math)
 - ▶ Evaluate error rate

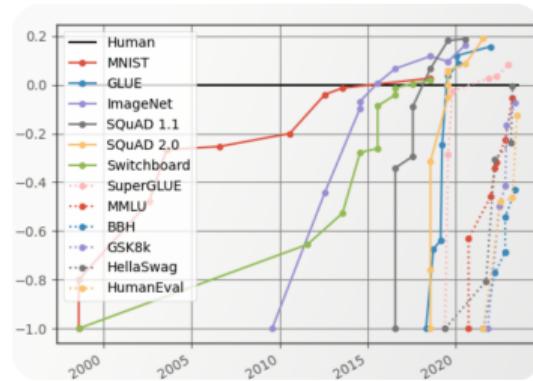


Figure: Plotting progress in AI (Kiel et al. 2023)

- ▶ Many challenges
 - ▶ Test data is necessarily arbitrary and limited
 - ▶ Correctness evaluation can be ambiguous (often use other LLMs!)
 - ▶ Evaluating the biggest models can be expensive (lots of API calls)

A statistical approach to efficient AI evaluation

tinyBenchmarks: evaluating LLMs with fewer examples

Felipe Maia Polo¹ Lucas Weber² Leshem Choshen^{3,4} Yuekai Sun¹ Gongjun Xu¹ Mikhail Yurochkin^{3,5}

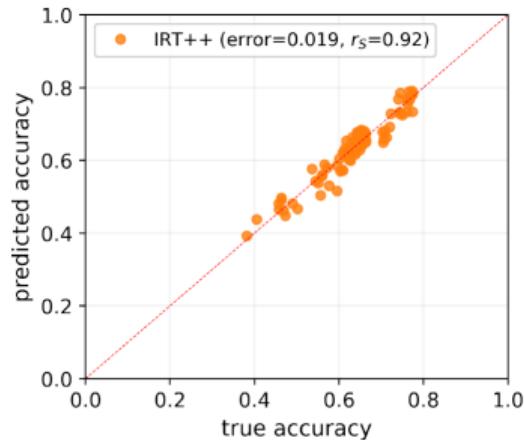


Figure 1. Estimating accuracy on MMLU (true accuracy) using 100 curated examples (predicted accuracy). IRT++, our best-performing evaluation strategy, predicts the accuracy of recent LLMs released between December 30th and January 18th within 1.9% of their true accuracy on all of MMLU (14K examples).

Item response theory (IRT) for AI evaluation

4. Better performance estimation with IRT

4.1. The IRT model

The two-parameter multidimensional IRT model assumes that the probability of the LLM j getting example i correctly is given by

$$p_{il} \triangleq \mathbb{P}(Y_{il} = 1 \mid \theta_l, \alpha_i, \beta_i) = \frac{1}{1 + \exp(-\alpha_i^\top \theta_l + \beta_i)}, \quad (4.1)$$

where $\theta_l \in \mathbb{R}^d$ denotes the unobserved abilities of LLM l , while $\alpha_i \in \mathbb{R}^d$ dictates which dimensions of θ_l are required from model l to respond to example i correctly. In this formulation, $\beta_i \in \mathbb{R}$ can be viewed as a bias term that regulates the probability of correctness when $\theta_l = 0$.

Formally, we are interested in approximating

$$Z_{jl} \triangleq \frac{1}{|\mathcal{I}_j|} \sum_{i \in \mathcal{I}_j} Y_{il} \quad (4.2)$$

Now, assume that we have run model l on a subset of examples from scenario j , obtaining responses $\{Y_{i_0l}, \dots, Y_{i_kl}\}$ for the examples $\hat{\mathcal{I}}_j = \{i_0, \dots, i_k\}$. Let $\hat{\theta}_l$ denote the estimate for θ_l after observing $\hat{\mathcal{I}}_j$ and possibly a bigger set of examples coming from different scenarios. To obtain that estimate, we maximize the log-likelihood of the freshly observed data with respect to θ_l , fixing examples' parameters. This procedure is equivalent to fitting a logistic regression model.

Because Z_{jl} is a random variable, we approximate it by estimating the conditional expectation

$$\begin{aligned} \mathbb{E}[Z_{jl} \mid Y_{i_0l}, \dots, Y_{i_kl}] &= \\ &= \frac{1}{|\mathcal{I}_j|} \sum_{i \in \mathcal{I}_j} \mathbb{E}[Y_{il} \mid Y_{i_0l}, \dots, Y_{i_kl}] \end{aligned}$$

The probability $p_{il} = \mathbb{P}(Y_{il} = 1 \mid \theta_l, \alpha_i, \beta_i)$ is given by the IRT model in Equation 4.1. The estimator for the conditional expectation is then given by

$$\begin{aligned} \hat{Z}_{jl}^{\text{p-IRT}} &\triangleq \hat{\mathbb{E}}[Z_{jl} \mid Y_{i_0l}, \dots, Y_{i_kl}] \\ &= \frac{\hat{\lambda}}{|\hat{\mathcal{I}}_j|} \sum_{i \in \hat{\mathcal{I}}_j} Y_{il} + \frac{1 - \hat{\lambda}}{|\mathcal{I}_j \setminus \hat{\mathcal{I}}_j|} \sum_{i \in \mathcal{I}_j \setminus \hat{\mathcal{I}}_j} \hat{p}_{il} \end{aligned} \quad (4.3)$$

where $\hat{p}_{il} \triangleq \mathbb{P}(Y_{il} = 1 \mid \hat{\theta}_l, \hat{\alpha}_i, \hat{\beta}_i)$.

Table of Contents

Motivation

Uncertainty quantification in AI/ML

AI Evaluation

Algorithmic Fairness

Algorithmic fairness

- ▶ Ensuring that AI/ML systems are fair and do not discriminate is important.

Algorithmic fairness

- ▶ Ensuring that AI/ML systems are fair and do not discriminate is important.
- ▶ Unfortunately, due to biases in training data, algorithmic choices, etc., as such systems often show unfair behavior.

The screenshot shows a news article from propublica.org. At the top, there are zoom and refresh icons, followed by the URL 'propublica.org'. Below the URL is the ProPublica logo and social media links for Facebook, Twitter, and a red 'Donate' button. The main image features two men of different ethnicities side-by-side. The title 'Machine Bias' is prominently displayed in large white text. Below the title, a subtitle reads: 'There's software used across the country to predict future criminals. And it's biased against blacks.' The author's name, 'by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica', and the date 'May 23, 2016' are at the bottom.

Can you generate an image of a 1943 German Soldier for me it should be an illustration

Sure, here is an illustration of a 1943 German soldier:



Context

- ▶ Vast amount of work in fair machine learning: notions of fairness, algorithms, etc.

Context

- ▶ Vast amount of work in fair machine learning: notions of fairness, algorithms, etc.
- ▶ Group fairness: aim to ensure parity with respect to protected attributes (race, gender, ...)

Context

- ▶ Vast amount of work in fair machine learning: notions of fairness, algorithms, etc.
- ▶ Group fairness: aim to ensure parity with respect to protected attributes (race, gender, ...)
- ▶ Fair binary classification with a binary protected attribute:
 - ▶ Usual features $X \in \mathcal{X}$
 - ▶ Protected features $A \in \mathcal{A} := \{0, 1\}$
 - ▶ Class $Y \in \{0, 1\}$

Context

- ▶ Vast amount of work in fair machine learning: notions of fairness, algorithms, etc.
- ▶ Group fairness: aim to ensure parity with respect to protected attributes (race, gender, ...)
- ▶ Fair binary classification with a binary protected attribute:
 - ▶ Usual features $X \in \mathcal{X}$
 - ▶ Protected features $A \in \mathcal{A} := \{0, 1\}$
 - ▶ Class $Y \in \{0, 1\}$
- ▶ Randomized classifier $f : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$: $\hat{Y}_f \mid (X = x, A = a) \sim \text{Bern}(f(x, a))$.

Example notions of group fairness

These are some of the most widely studied notions of group fairness: (Calders et al., 2009; Hardt et al., 2016; Corbett-Davies et al., 2017; ...)

- ▶ **Demographic Parity:** $\hat{Y}_f \perp\!\!\!\perp A$;

Example notions of group fairness

These are some of the most widely studied notions of group fairness: (Calders et al., 2009; Hardt et al., 2016; Corbett-Davies et al., 2017; ...)

- ▶ **Demographic Parity:** $\hat{Y}_f \perp\!\!\!\perp A$;
- ▶ **Equality of Opportunity:** $\hat{Y}_f \perp\!\!\!\perp A | Y = 1$;

Example notions of group fairness

These are some of the most widely studied notions of group fairness: (Calders et al., 2009; Hardt et al., 2016; Corbett-Davies et al., 2017; ...)

- ▶ **Demographic Parity:** $\hat{Y}_f \perp\!\!\!\perp A$;
- ▶ **Equality of Opportunity:** $\hat{Y}_f \perp\!\!\!\perp A | Y = 1$;
- ▶ **Predictive Parity:** $\hat{Y}_f \perp\!\!\!\perp A | Y = 0$;

Example disparity measures

Definition (Disparity Measures)

We consider the following disparity measures: (Calders et al., 2009; Hardt et al., 2016; Corbett-Davies et al., 2017; Cho et al., 2020; ...)

- ▶ **Demographic Disparity (DD).** Difference in positive prediction rates between groups:

$$\text{DD}(f) = \mathbb{P}_{X|A=1}(\hat{Y}_f = 1) - \mathbb{P}_{X|A=0}(\hat{Y}_f = 1);$$

Example disparity measures

Definition (Disparity Measures)

We consider the following disparity measures: (Calders et al., 2009; Hardt et al., 2016; Corbett-Davies et al., 2017; Cho et al., 2020; ...)

- ▶ **Demographic Disparity (DD)**. Difference in positive prediction rates between groups:

$$\text{DD}(f) = \mathbb{P}_{X|A=1}(\hat{Y}_f = 1) - \mathbb{P}_{X|A=0}(\hat{Y}_f = 1);$$

- ▶ **Disparity of Opportunity (DO)**. Difference in true positive rates:

$$\text{DO}(f) = \mathbb{P}_{X|A=1, Y=1}(\hat{Y}_f = 1) - \mathbb{P}_{X|A=0, Y=1}(\hat{Y}_f = 1);$$

Example disparity measures

Definition (Disparity Measures)

We consider the following disparity measures: (Calders et al., 2009; Hardt et al., 2016; Corbett-Davies et al., 2017; Cho et al., 2020; ...)

- ▶ **Demographic Disparity (DD)**. Difference in positive prediction rates between groups:

$$\text{DD}(f) = \mathbb{P}_{X|A=1}(\hat{Y}_f = 1) - \mathbb{P}_{X|A=0}(\hat{Y}_f = 1);$$

- ▶ **Disparity of Opportunity (DO)**. Difference in true positive rates:

$$\text{DO}(f) = \mathbb{P}_{X|A=1, Y=1}(\hat{Y}_f = 1) - \mathbb{P}_{X|A=0, Y=1}(\hat{Y}_f = 1);$$

- ▶ **Predictive Disparity (PD)**. Difference in false positive rates:

$$\text{PD}(f) = \mathbb{P}_{X|A=1, Y=0}(\hat{Y}_f = 1) - \mathbb{P}_{X|A=0, Y=0}(\hat{Y}_f = 1).$$

Fair Bayes-optimal classifiers

Classifiers with the highest accuracy given a disparity level.

Fair Bayes-optimal classifiers

Classifiers with the highest accuracy given a disparity level.

Definition (Fair Bayes-optimal Classifier)

Consider any $K \geq 1$ fairness measures $\text{Dis}_k : \mathcal{F} \rightarrow [0, 1]$, $k = 1, \dots, K$. Then, a δ -fair Bayes-optimal classifier $f_{\text{Dis}, \delta}^*$ minimizes the misclassification error $R(f) := \mathbb{P}(Y \neq \hat{Y}_f)$ over all classifiers that satisfy δ -disparity:

$$f_{\text{Dis}, \delta}^* \in \operatorname{argmin}_{f \in \mathcal{F}} \left\{ R(f) : \max_{k=1}^K |\text{Dis}_k(f)| \leq \delta \right\}.$$

Fair Bayes-optimal classifiers

Classifiers with the highest accuracy given a disparity level.

Definition (Fair Bayes-optimal Classifier)

Consider any $K \geq 1$ fairness measures $\text{Dis}_k : \mathcal{F} \rightarrow [0, 1]$, $k = 1, \dots, K$. Then, a δ -fair Bayes-optimal classifier $f_{\text{Dis}, \delta}^*$ minimizes the misclassification error $R(f) := \mathbb{P}(Y \neq \hat{Y}_f)$ over all classifiers that satisfy δ -disparity:

$$f_{\text{Dis}, \delta}^* \in \operatorname{argmin}_{f \in \mathcal{F}} \left\{ R(f) : \max_{k=1}^K |\text{Dis}_k(f)| \leq \delta \right\}.$$

Bayes-Optimal Fair Classification with Linear Disparity Constraints
via Pre-, In-, and Post-processing

Xianli Zeng*, Guang Cheng† and Edgar Dobriban‡

February 6, 2024

Linear and bilinear disparity measures

- We find the form of Bayes-optimal classifiers for a broad class of disparities: linear & bilinear measures.

Linear and bilinear disparity measures

- ▶ We find the form of Bayes-optimal classifiers for a broad class of disparities: linear & bilinear measures.
- ▶ **Definition:** A disparity measure Dis is *linear* if there is a function w_{Dis} such that

$$\text{Dis}(f) = \int_{\mathcal{A}} \int_{\mathcal{X}} f(x, a) w_{\text{Dis}}(x, a) d\mathbb{P}_{X,A}(x, a).$$

Linear and bilinear disparity measures

- ▶ We find the form of Bayes-optimal classifiers for a broad class of disparities: linear & bilinear measures.
- ▶ **Definition:** A disparity measure Dis is *linear* if there is a function w_{Dis} such that

$$\text{Dis}(f) = \int_{\mathcal{A}} \int_{\mathcal{X}} f(x, a) w_{\text{Dis}}(x, a) d\mathbb{P}_{X,A}(x, a).$$

- ▶ Class-conditional probability function η_a : $\eta_a(x) = \mathbb{P}(Y = 1 \mid A = a, X = x)$ for all x, a .

Linear and bilinear disparity measures

- ▶ We find the form of Bayes-optimal classifiers for a broad class of disparities: linear & bilinear measures.
- ▶ **Definition:** A disparity measure Dis is *linear* if there is a function w_{Dis} such that

$$\text{Dis}(f) = \int_A \int_X f(x, a) w_{\text{Dis}}(x, a) d\mathbb{P}_{X,A}(x, a).$$

- ▶ Class-conditional probability function η_a : $\eta_a(x) = \mathbb{P}(Y = 1 \mid A = a, X = x)$ for all x, a .
- ▶ **Definition:** A linear disparity measure Dis is *bilinear* if $w_{\text{Dis}}(x, a)$ can be expressed for all x, a as $w_{\text{Dis}}(x, a) = s_{\text{Dis}, a} \eta_a(x) + b_{\text{Dis}, a}$.

Linear and bilinear disparity measures

- ▶ We find the form of Bayes-optimal classifiers for a broad class of disparities: linear & bilinear measures.
- ▶ **Definition:** A disparity measure Dis is *linear* if there is a function w_{Dis} such that

$$\text{Dis}(f) = \int_{\mathcal{A}} \int_{\mathcal{X}} f(x, a) w_{\text{Dis}}(x, a) d\mathbb{P}_{X, A}(x, a).$$

- ▶ Class-conditional probability function η_a : $\eta_a(x) = \mathbb{P}(Y = 1 \mid A = a, X = x)$ for all x, a .
- ▶ **Definition:** A linear disparity measure Dis is *bilinear* if $w_{\text{Dis}}(x, a)$ can be expressed for all x, a as $w_{\text{Dis}}(x, a) = s_{\text{Dis}, a} \eta_a(x) + b_{\text{Dis}, a}$.
- ▶ **Proposition:** The disparity measures DD, DO, and PD are bilinear: for all x, a

$$w_{\text{DD}}(x, a) = \frac{(2a - 1)}{p_a}; \quad w_{\text{DO}}(x, a) = \frac{(2a - 1)\eta_a(x)}{p_{a,1}}; \quad w_{\text{PD}}(x, a) = \frac{(2a - 1)(1 - \eta_a(x))}{p_{a,0}}.$$

Intuition for optimal fair classifiers

- ▶ Consider a linear disparity measure Dis :

$$\text{Dis}(f) = \int_{\mathcal{A}, \mathcal{X}} f(x, a) w_{\text{Dis}}(x, a) d\mathbb{P}_{X, A}(x, a)$$

Intuition for optimal fair classifiers

- ▶ Consider a linear disparity measure Dis :

$$\text{Dis}(f) = \int_{\mathcal{A}, \mathcal{X}} f(x, a) w_{\text{Dis}}(x, a) d\mathbb{P}_{X, A}(x, a)$$

- ▶ Want to minimize

$$R(f) = \int_{\mathcal{A}, \mathcal{X}} f(x, a)(1 - 2\eta_a(x)) d\mathbb{P}_{X, A}(x, a) + C_{\mathbb{P}}$$

s.t. $|\text{Dis}(f)| \leq \delta$

Intuition for optimal fair classifiers

- ▶ Consider a linear disparity measure Dis :

$$\text{Dis}(f) = \int_{\mathcal{A}, \mathcal{X}} f(x, a) w_{\text{Dis}}(x, a) d\mathbb{P}_{X, A}(x, a)$$

- ▶ Want to minimize

$$R(f) = \int_{\mathcal{A}, \mathcal{X}} f(x, a)(1 - 2\eta_a(x)) d\mathbb{P}_{X, A}(x, a) + C_{\mathbb{P}}$$

s.t. $|\text{Dis}(f)| \leq \delta$

- ▶ Intuitive that an optimal classifier should be of the form (for some c)

$$f(x, a) = I\left(\eta_a(x) > \frac{1}{2} + c \cdot w_{\text{Dis}}(x, a)\right).$$

Optimal fair classifiers for linear and bilinear disparity measures

Consider linear disp. measure Dis ; Continuous r.v.s $\eta_a(X)$, $w_{\text{Dis}}(X, a)$, $a \in \{0, 1\}$.

Optimal fair classifiers for linear and bilinear disparity measures

Consider linear disp. measure Dis; Continuous r.v.s $\eta_a(X)$, $w_{\text{Dis}}(X, a)$, $a \in \{0, 1\}$.

Theorem (Optimal Classifiers for Linear Disp. Measures; Zeng et al., 2024)

Let $t_{\text{Dis}}(\delta) = \operatorname{argmin}_t \{|t| : |D_{\text{Dis}}(t)| \leq \delta\}$, where

$$D_{\text{Dis}}(t) = \int_{\mathcal{X} \times \mathcal{A}} w_{\text{Dis}}(x, a) \cdot I\left(\eta_a(x) > \frac{1}{2} + \frac{t}{2} w_{\text{Dis}}(x, a)\right) d\mathbb{P}_{X,A}(x, a).$$

Optimal fair classifiers for linear and bilinear disparity measures

Consider linear disp. measure Dis ; Continuous r.v.s $\eta_a(X)$, $w_{\text{Dis}}(X, a)$, $a \in \{0, 1\}$.

Theorem (Optimal Classifiers for Linear Disp. Measures; Zeng et al., 2024)

Let $t_{\text{Dis}}(\delta) = \operatorname{argmin}_t \{|t| : |D_{\text{Dis}}(t)| \leq \delta\}$, where

$$D_{\text{Dis}}(t) = \int_{\mathcal{X} \times \mathcal{A}} w_{\text{Dis}}(x, a) \cdot I\left(\eta_a(x) > \frac{1}{2} + \frac{t}{2} w_{\text{Dis}}(x, a)\right) d\mathbb{P}_{X,A}(x, a).$$

A fair Bayes-optimal classifier is, for all x, a ,

$$f_{\text{Dis}, \delta}^*(x, a) = I\left(\eta_a(x) > \frac{1}{2} + \frac{t_{\text{Dis}}(\delta)}{2} w_{\text{Dis}}(x, a)\right).$$

Optimal fair classifiers for linear and bilinear disparity measures

Consider linear disp. measure Dis ; Continuous r.v.s $\eta_a(X)$, $w_{\text{Dis}}(X, a)$, $a \in \{0, 1\}$.

Theorem (Optimal Classifiers for Linear Disp. Measures; Zeng et al., 2024)

Let $t_{\text{Dis}}(\delta) = \operatorname{argmin}_t \{|t| : |D_{\text{Dis}}(t)| \leq \delta\}$, where

$$D_{\text{Dis}}(t) = \int_{\mathcal{X} \times \mathcal{A}} w_{\text{Dis}}(x, a) \cdot I\left(\eta_a(x) > \frac{1}{2} + \frac{t}{2} w_{\text{Dis}}(x, a)\right) d\mathbb{P}_{X,A}(x, a).$$

A fair Bayes-optimal classifier is, for all x, a ,

$$f_{\text{Dis}, \delta}^*(x, a) = I\left(\eta_a(x) > \frac{1}{2} + \frac{t_{\text{Dis}}(\delta)}{2} w_{\text{Dis}}(x, a)\right).$$

If Dis is bilinear, this is a group-wise thresholding rule, such that for all x, a ,

$$f_{\text{Dis}, \delta}^*(x, a) = I\left(\eta_a(x) > \frac{1 + b_{\text{Dis}, a} \cdot t_{\text{Dis}}(\delta)}{2 - s_{\text{Dis}, a} \cdot t_{\text{Dis}}(\delta)}\right).$$

Discussion, methods, and empirics

- ▶ Proof via the *generalized Neyman-Pearson lemma*: maximizing a linear functional subject to linear constraints.

Discussion, methods, and empirics

- ▶ Proof via the *generalized Neyman-Pearson lemma*: maximizing a linear functional subject to linear constraints.
- ▶ Apply it to develop methods for fair classification via
 - ▶ Pre-Processing: Fair Up-& Down-Sampling (FUDS);
 - ▶ In-Processing: Fair Cost-Sensitive Classification (FCSC);
 - ▶ Post-Processing: Fair Plug-in Thresholding Rule (FPIR).

Discussion, methods, and empirics

- ▶ Proof via the *generalized Neyman-Pearson lemma*: maximizing a linear functional subject to linear constraints.
- ▶ Apply it to develop methods for fair classification via
 - ▶ Pre-Processing: Fair Up-& Down-Sampling (FUDS);
 - ▶ In-Processing: Fair Cost-Sensitive Classification (FCSC);
 - ▶ Post-Processing: Fair Plug-in Thresholding Rule (FPIR).

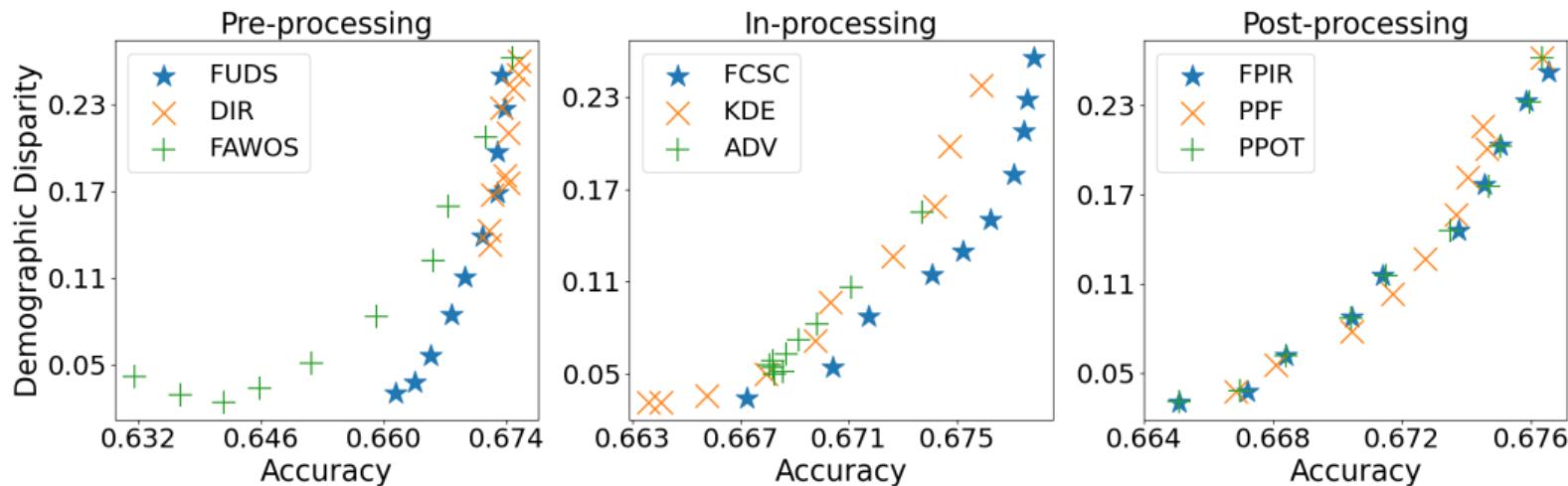


Figure: Fairness-accuracy tradeoff on the COMPAS dataset.

Summary for algorithmic fairness

- ▶ Derived fair Bayes-optimal classifiers for linear disparity measures (via Neyman-Pearson lemma!).
- ▶ Developed fair classification methods via pre-, in-, and post-processing.

Summary for algorithmic fairness

- ▶ Derived fair Bayes-optimal classifiers for linear disparity measures (via Neyman-Pearson lemma!).
- ▶ Developed fair classification methods via pre-, in-, and post-processing.
- ▶ Related work: minimax optimal fair classification.

Minimax Optimal Fair Classification with Bounded Demographic Disparity

Xianli Zeng*, Guang Cheng[†] and Edgar Dobriban[‡]

March 28, 2024

Summary

- ▶ Major goal: have an impact in the world of AI through statistical research.
- ▶ Thanks!

