# Can I Trust My Fairness Metric?
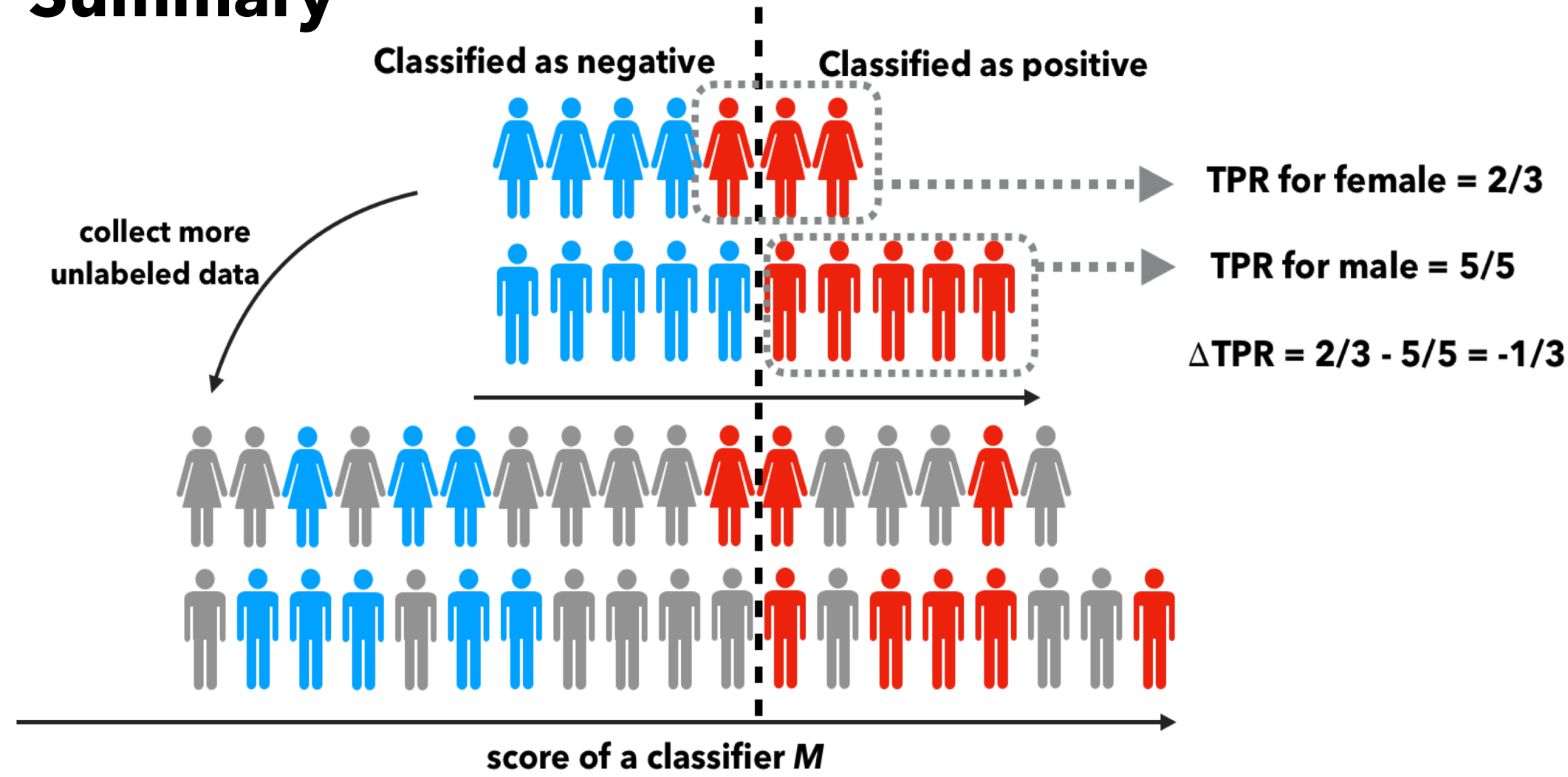# Assessing Fairness with Unlabeled Data and Bayesian Inference

Disi Ji[1], Padhraic Smyth[1], Mark Steyvers[2]
[1]Department of Computer Science,  [2]Department of Cognitive Sciences
University of California, Irvine

**UCIRVINE**

## Summary



Classified as negative | Classified as positive

TPR for female = 2/3

TPR for male = 5/5

ΔTPR = 2/3 - 5/5 = -1/3

collect more unlabeled data

score of a classifier $M$
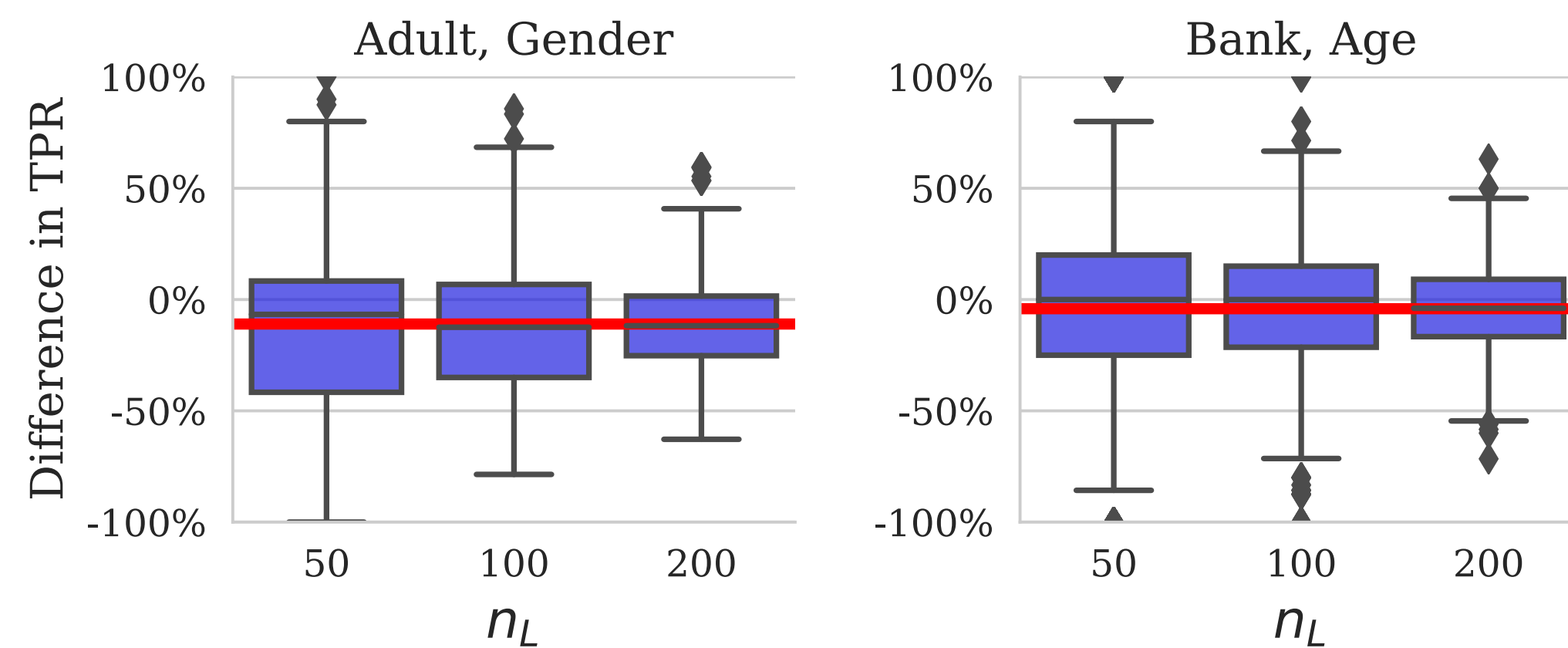
➤ **Equality of opportunity**:  equal TPR across different groups
➤ Due to small sample size, the estimated TPRs are **noisy**!

➤ Contribution:
➤ 1. **Quantify uncertainty** in fairness metrics using Bayesian methods
➤ 2. **Reduce uncertainty** of fairness by leveraging unlabeled data

## Frequentist-based Estimates Have High Variance
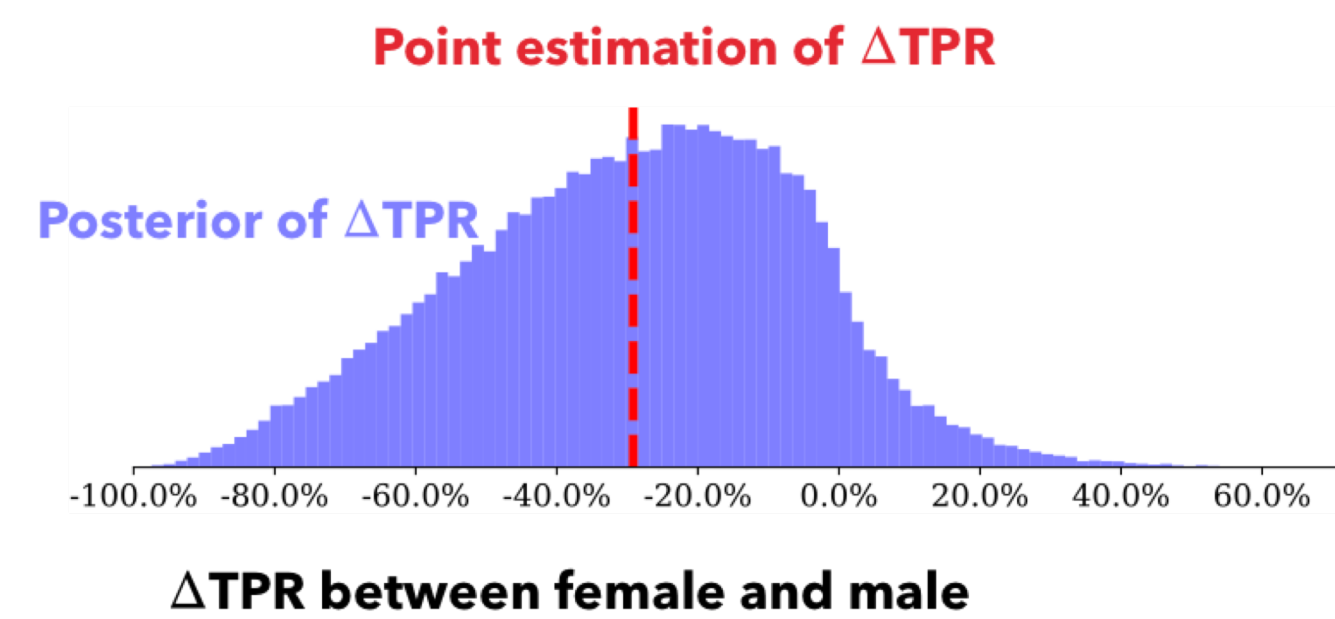


Adult, Gender

Bank, Age

Difference in TPR

$n_L$

➤ High variability for the estimated TPRs relative to the true TPRs (shown in red) as a function of the number of labeled examples.
➤ In many cases the estimates are two or three or more times larger than the true difference.
➤ A relatively large percentage of the estimates have the opposite sign of the true difference, potentially leading to mistaken conclusions
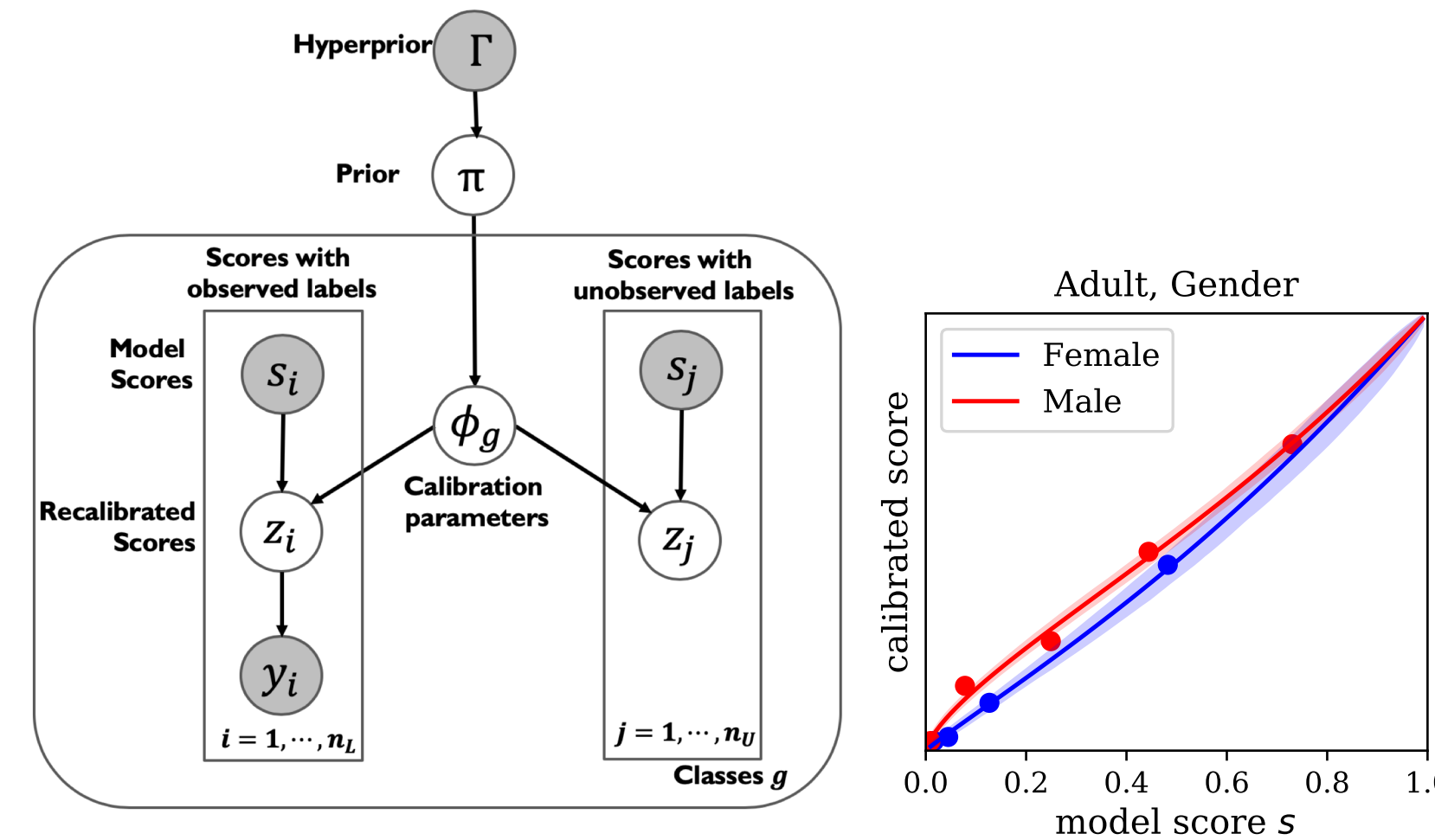
## Quantify Uncertainty of Fairness Assessment

For two groups $g = 0,1$, and $n_L$ labeled data $D_L$:
➤ **Groupwise performance metric** $\theta_g = P(\hat{y} = 1 | y = 1, g)$
$$\theta_g \sim Beta(\alpha_g, \beta_g)$$
➤ **Correctness** of the prediction model for $i$: $I_i = I(\hat{y}_i = y_i),\ 1 \leq i \leq n_L$:
$$I_i \sim Bernoulli(\theta_g)$$
➤ **Group fairness metric**: $\Delta = \theta_1 - \theta_0$
➤ Obtain posterior distribution $P(\Delta | D_L)$ via Monte Carlo samples



**Point estimation of ΔTPR**

**Posterior of ΔTPR**

-100.0% -80.0% -60.0% -40.0% -20.0% 0.0% 20.0% 40.0% 60.0%

**ΔTPR between female and male**

## Reduce Uncertainty with Unlabeled Data



Hyperprior Γ

Prior π

Scores with observed labels | Scores with unobserved labels

Model Scores $s_i$ | $s_j$

$\phi_g$ Calibration parameters

Recalibrated Scores $z_i$ | $z_j$

$y_i$

$i = 1, \cdots, n_L$ | $j = 1, \cdots, n_U$

Classes $g$

Adult, Gender
— Female
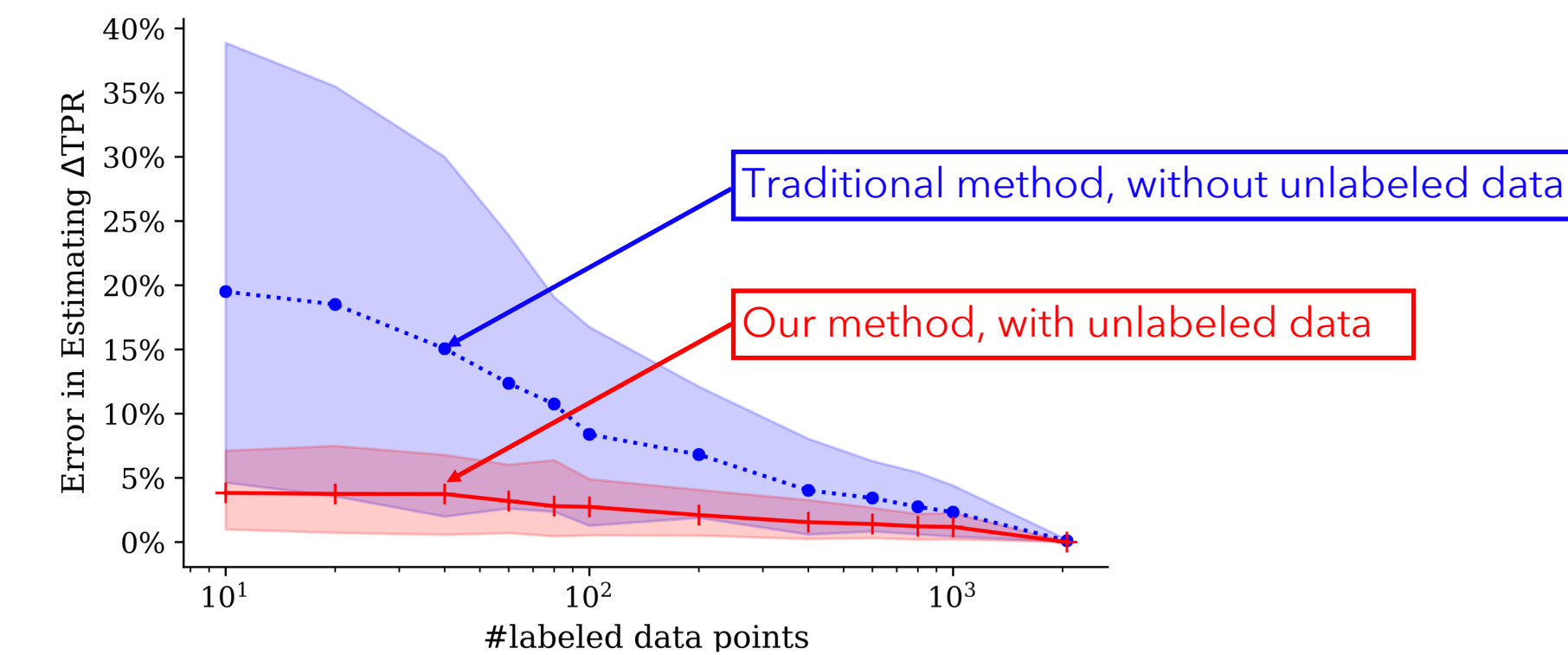— Male

calibrated score

model score $s$

We treat each $z_j, j = 1, \ldots, n_U$ as a latent variable per example. The high level steps of the approach are as follows:

1. Use the $n_L$ labeled examples to **estimate groupwise calibration functions with parameters $\varphi_g$**, that transform the (potentially) uncalibrated scores $s$ of the model to calibrated scores.  More specifically, we perform Bayesian inference to obtain posterior samples from $P(\varphi_g | D_L)$ for the groupwise calibration parameters $\varphi_g$.
2. **Obtain posterior samples of recalibrated scores** from $P\varphi_g(z_j | D_L, s_j)$ for each unlabeled example $j = 1, \ldots, n_U$, conditioned on posterior samples of the $\varphi_g$'s.
3. Use posterior samples from the $z_j$'s, combined with the labeled data, to **generate estimates of the groupwise metrics $\theta_g$ and the difference in metrics Δ.**
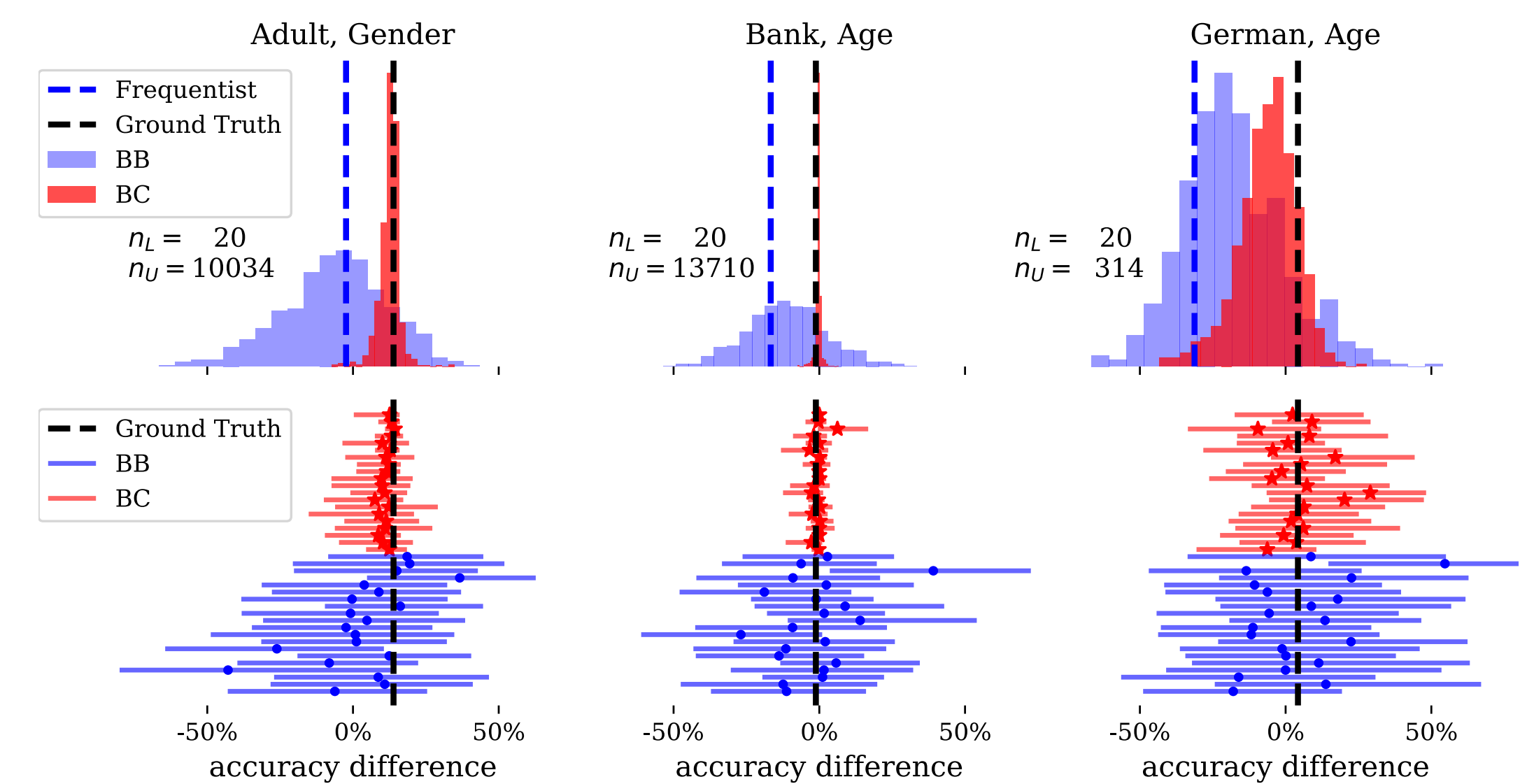
## Experimental Results

| Dataset | Test Size | $G$ | $P(g = 0)$ | $P(y = 1)$ |
|---|---|---|---|---|
| Adult | 10054 | gender, race | 0.68, 0.86 | 0.25 |
| Bank | 13730 | age | 0.45 | 0.11 |
| German | 334 | age, gender | 0.79, 0.37 | 0.17 |
| Compas-R | 2056 | gender, race | 0.7, 0.85 | 0.69 |
| Compas-VR | 1337 | gender, race | 0.8, 0.34 | 0.47 |
| Ricci | 40 | race | 0.65 | 0.50 |

**Example**: assess Δ TPR of COMPAS- Recividism, Race



Error in Estimating ΔTPR

Traditional method, without unlabeled data

Our method, with unlabeled data

#labeled data points

With **10** labeled data and **~2000** unlabeled data, error in estimating  TPR is **5%** for our method versus 20% with only labeled data

**Illustrative Results:** Posterior density (samples) and frequentist estimates (dotted vertical blue lines) for the difference in group accuracy with **20** random labeled examples for both the BB (beta-binomial) and BC (Bayesian calibration) methods



Adult, Gender | Bank, Age | German, Age

- - Frequentist
- - Ground Truth
BB
BC

$n_L = 20$
$n_U = 10034$

$n_L = 20$
$n_U = 13710$

$n_L = 20$
$n_U = 314$

- - Ground Truth
BB
BC

-50%  0%  50% | -50%  0%  50% | -50%  0%  50%

accuracy difference