# AUROC vs. AUPRC under Class Imbalance

**Matthew McDermott,** Haoran Zhang, Lasse Hansen, Giovanni Angelotti, Jack Gallifant

Berkowitz Postdoctoral Fellow

matthew_mcdermott@hms.harvard.edu
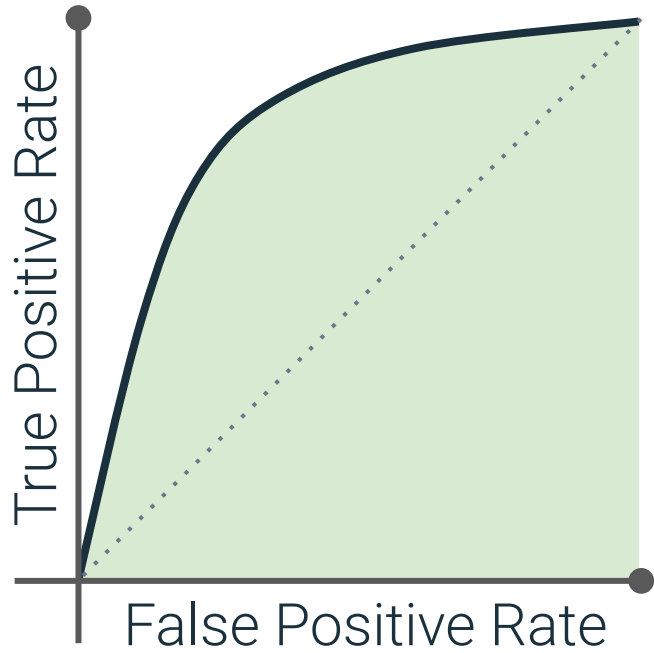
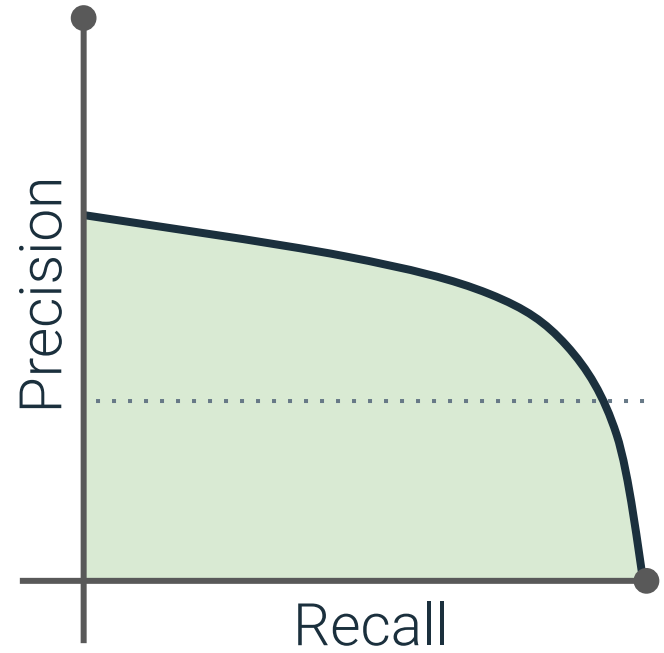HARVARD MEDICAL SCHOOL | BLAVATNIK INSTITUTE BIOMEDICAL INFORMATICS

COLUMBIA | COLUMBIA UNIVERSITY DEPARTMENT OF BIOMEDICAL INFORMATICS

AUROC

True Positive Rate

False Positive Rate

AUPRC

Precision

Recall

# AUPRC: Better under class imbalance?



**Daniel Rosenberg**
Jun 7, 2022 · 6 min read · ▶ Listen

## Unbalanced Data? Stop Using ROC-AUC and Use AUPRC Instead

Advantages of AUPRC when measuring performance in the presence of data imbalance — clearly explained

**Imbalanced data**

AUC-ROC is less sensitive to class imbalance than AUC-PR. In an imbalanced dataset, where one class is much more prevalent than the other, the ROC

**Tam D Tran-The** Follow
Nov 29, 2021 · 4 min read · ▶ Listen

## Precision-Recall Curve is More Informative than ROC in Imbalanced Data: Napkin Math & More

prevalent and there is low value in
ion-Recall curve is preferred over ROC

**Samuele Mazzanti** Follow
Apr 30, 2020 · 13 min read · ✦ Member-only · ▶ Listen

## Why You Should Stop Using the ROC Curve

The most popular metric may not be as meaningful as you think
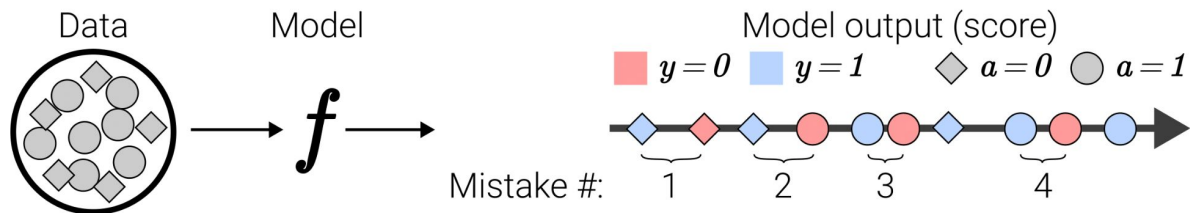
# A probabilistic view reveals not!

*Theorem* 1. Let $\mathcal{X}, \mathcal{Y} = 0, 1$ represent a paired feature and binary classification label space from which i.i.d. samples $(x, y) \in \mathcal{X} \times \mathcal{Y}$ are drawn via the joint distribution over the random variables $\mathsf{x}, \mathsf{y}$. Let $f : \mathcal{X} \to (0, 1)$ be a binary classification model outputting continuous probability scores over this space. Then,
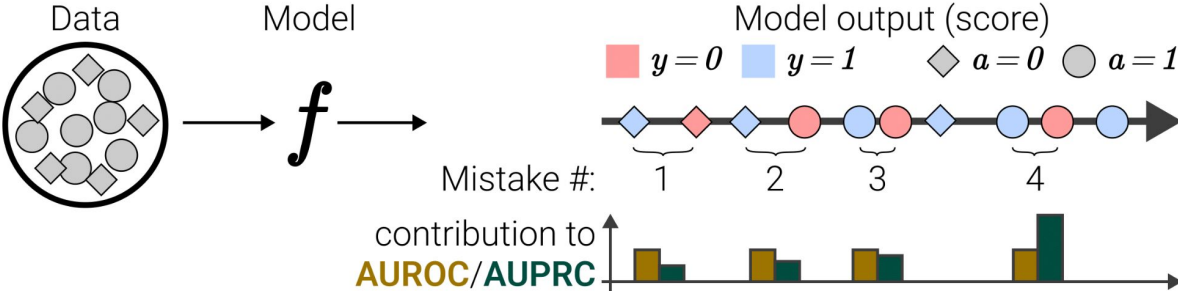
$$\text{AUROC}(f) = 1 - \mathbb{E}_{t \sim f(\mathsf{x})|\mathsf{y}=1} \left[ \text{FPR}(f, t) \right]$$

$$\text{AUPRC}(f) = 1 - P_{\mathsf{y}}(y = 0) \mathbb{E}_{t \sim f(\mathsf{x})|\mathsf{y}=1} \left[ \frac{\text{FPR}(f, t)}{P_{\mathsf{x}}(f(x) > t)} \right]$$
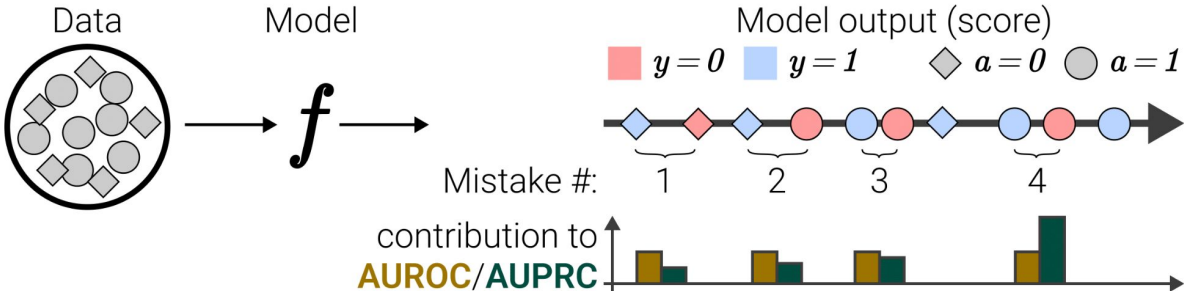
# Mistake order correction reveals AUPRC biases

# Mistake order correction reveals AUPRC biases
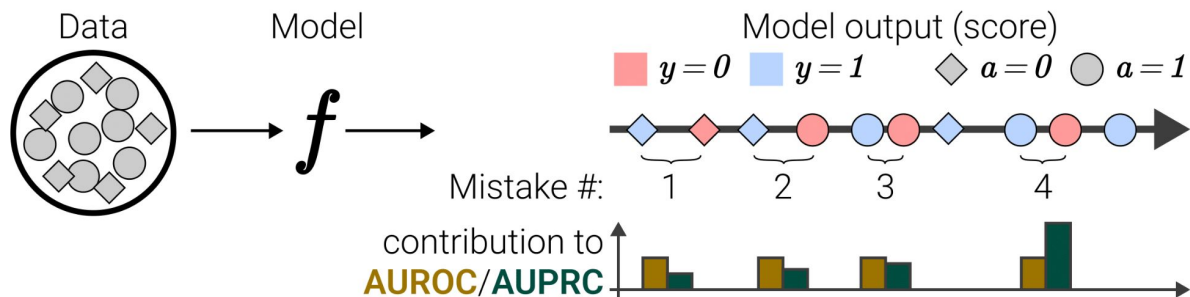
# Mistake order correction reveals AUPRC biases

# Mistake order correction reveals AUPRC biases

Data

Model

Model output (score)

🟥 $y = 0$   🟦 $y = 1$   ◇ $a = 0$   ⬤ $a = 1$

$f$

Mistake #:   1   2   3   4

contribution to
**AUROC**/**AUPRC**

**Small Molecule Screening**

$\rightarrow f \rightarrow$ Screen top-$k$ hits

Favor high score region to optimize top-$k$

# Mistake order correction reveals AUPRC biases

Data

Model

Model output (score)

$y = 0$    $y = 1$    $a = 0$    $a = 1$

$f$

Mistake #:   1    2    3    4

contribution to
**AUROC**/**AUPRC**

**Cancer Screening**

$\rightarrow f \rightarrow$ $\begin{cases} 1 : \text{Further tests} \\ 0 : \text{No follow up} \end{cases}$

Favor low-score region to catch all disease

# Mistake order correction reveals AUPRC biases
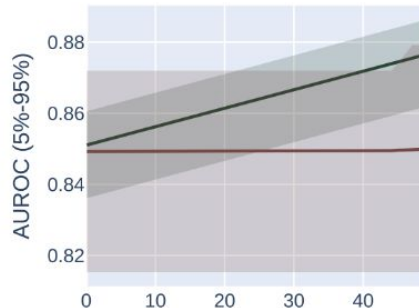
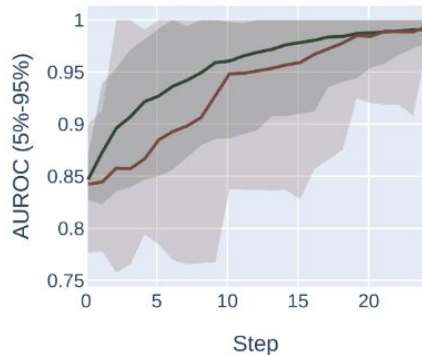# Synthetic experiments verify this bias



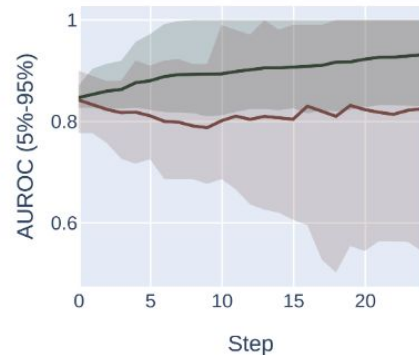**a)** Mistakes Fixed by AUROC

**b)** Mistakes Fixed by AUPRC

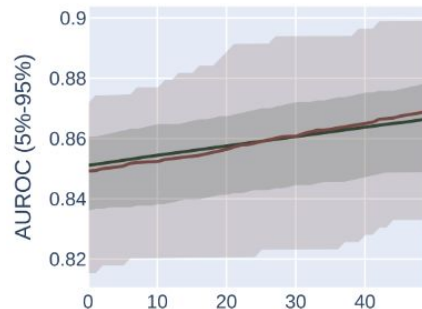**c)** Permutation optimization by AUROC
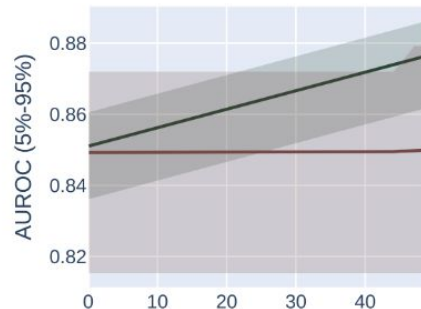
**d)** Permutation optimization by AUPRC

Low-prevalence Group    High-prevalence Group
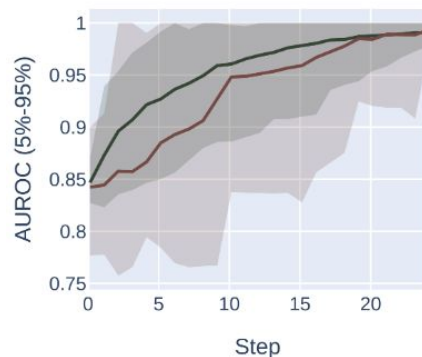
# Synthetic experiments verify this bias
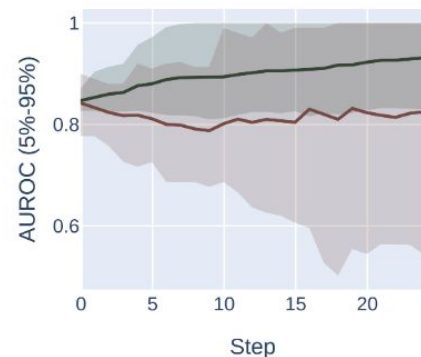


**a)** Mistakes Fixed by AUROC

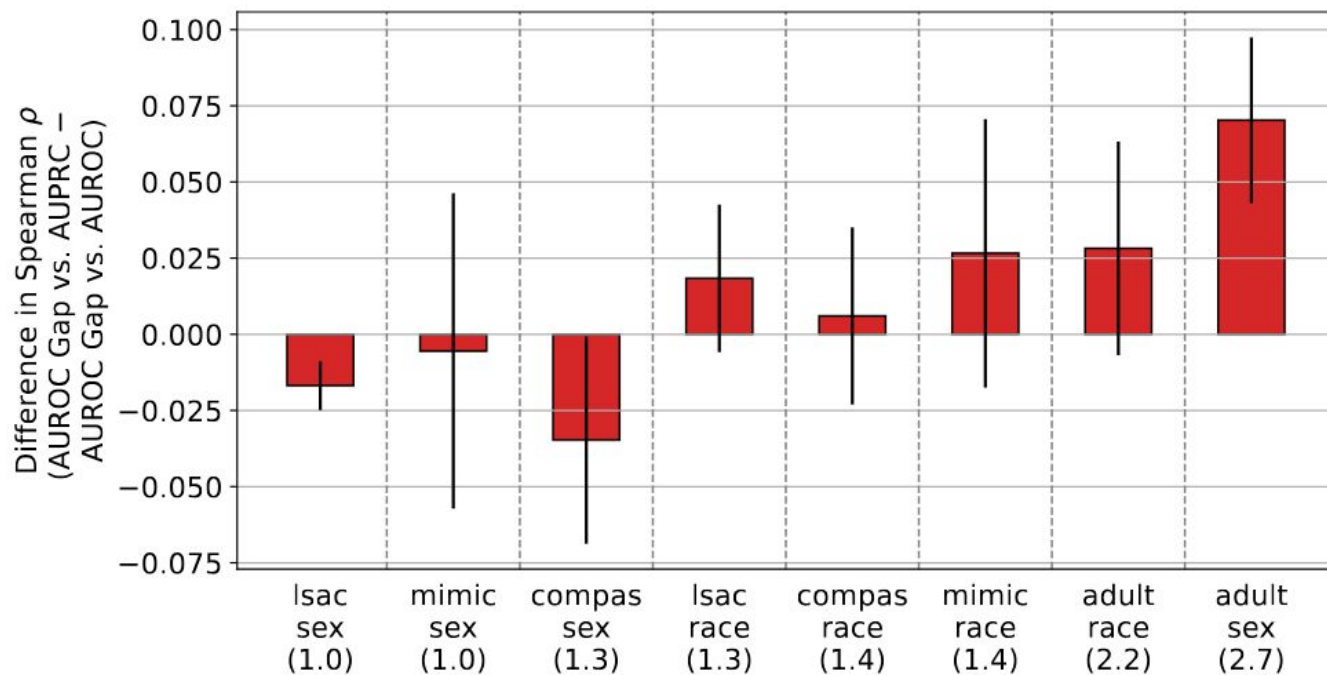**b)** Mistakes Fixed by AUPRC

**c)** Permutation optimization by AUROC

**d)** Permutation optimization by AUPRC

Low-prevalence Group — High-prevalence Group

Extent to which AUPRC favors high-prevalence group!
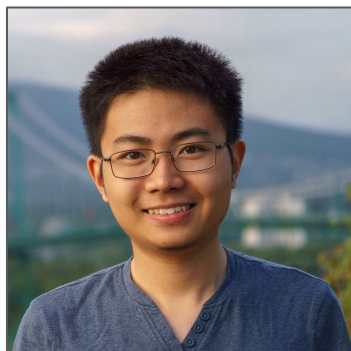
# This is not just synthetic -- hyperparameter tuning shows this effect!

# Acknowledgements



**Matthew McDermott**
HMS

Haoran Zhang
MIT

Lasse Hansen
Aarhus University
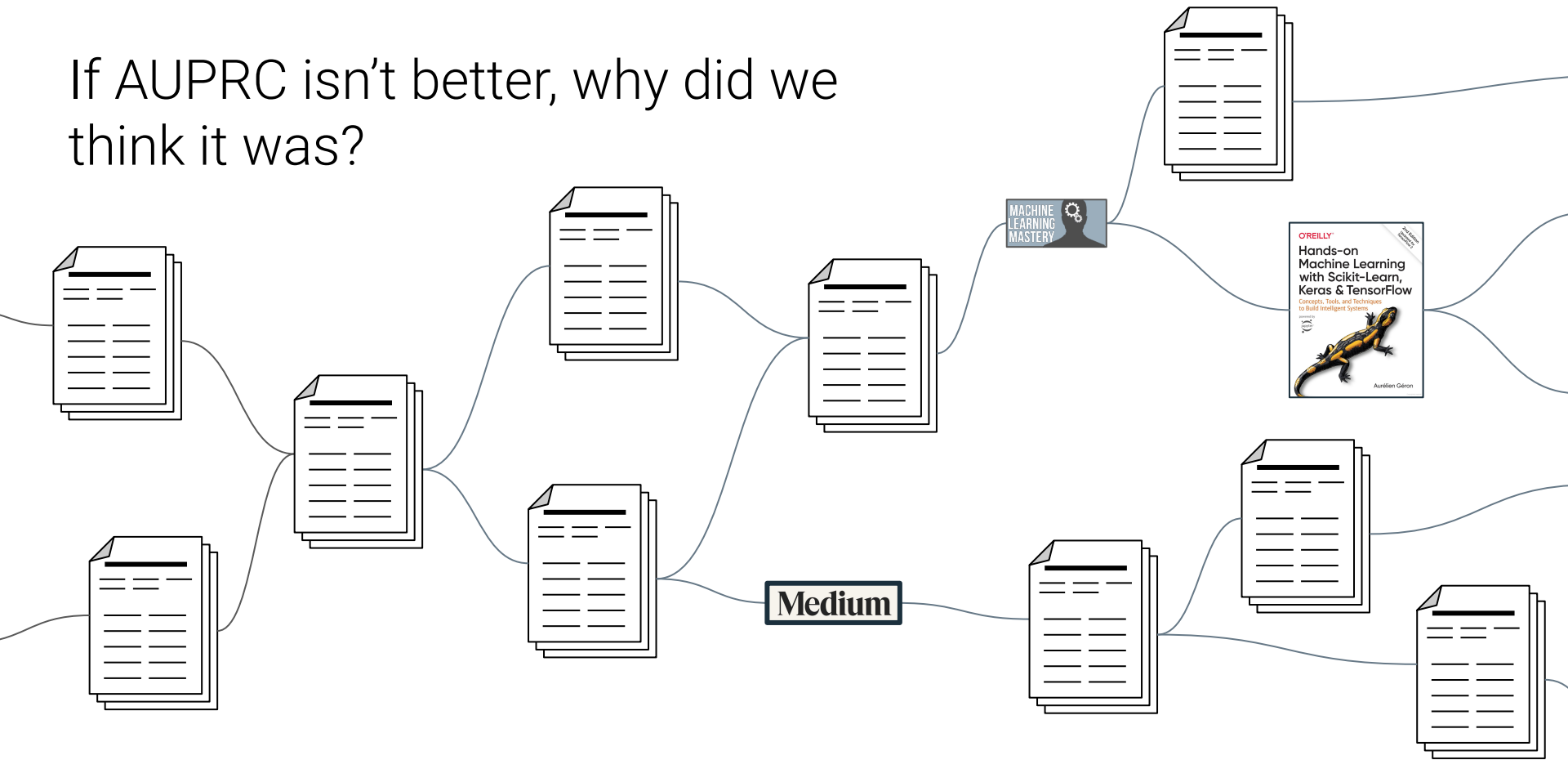
Giovanni Angelotti
IRCCS Humanitas

Jack Gallifant
MIT

If AUPRC isn't better, why did we think it was?

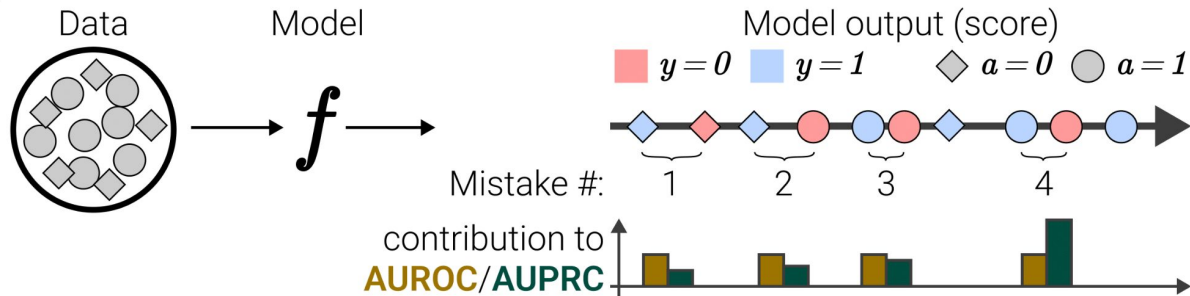If AUPRC isn't better, why did we think it was?

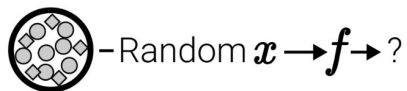# If AUPRC isn't better, why did we think it was?

1. Use of AUPRC justified by class imbalance in cases where other metrics are more appropriate are common.

2. Significant rates of mis-citation and misattribution of this claim.

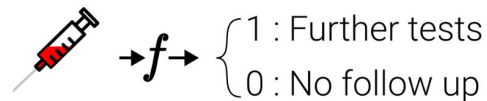3. Inaccurate and overly simplistic arguments are widespread.
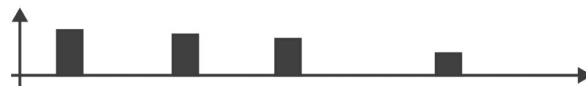
**a)**

Data     Model           Model output (score)

$y=0$    $y=1$    $a=0$    $a=1$

$f$

Mistake #:   1    2    3     4

contribution to
**AUROC**/**AUPRC**

**b) Methodology Comparison**      Favor mistakes equally

$-$Random $x \to f \to$ ?

**c) Cancer Screening**      Favor low-score region to catch all disease

$\to f \to$ $\begin{cases} 1 : \text{Further tests} \\ 0 : \text{No follow up} \end{cases}$

**d) Public Health Intervention**      Favor per-group regions differently

$\to f \to$ $\begin{cases} a=1 : \text{top-}k_1 \\ a=0 : \text{top-}k_0 \end{cases}$

**e) Small Molecule Screening**      Favor high score region to optimize top-$k$

$\to f \to$ Screen top-$k$ hits

19