

COMPUTATIONAL CHALLENGES IN DATA PRIVACY

FAIRNESS IN RISK ASSESSMENT

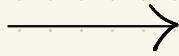
JULY 06, 2025

|

ROHIT VAISH



A collection
of numbers



MERGE SORT



Sorted list
of numbers

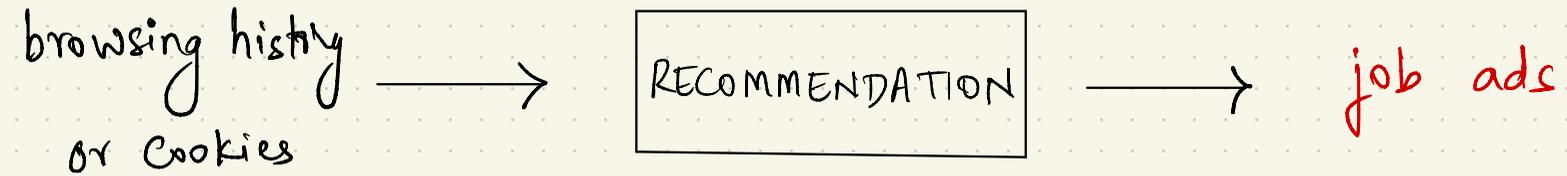
a graph/network
with
edge weights



DIJKSTRA

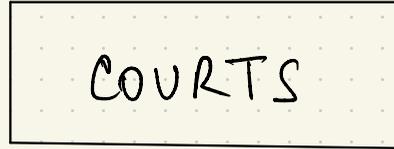
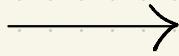


shortest path
between vertices

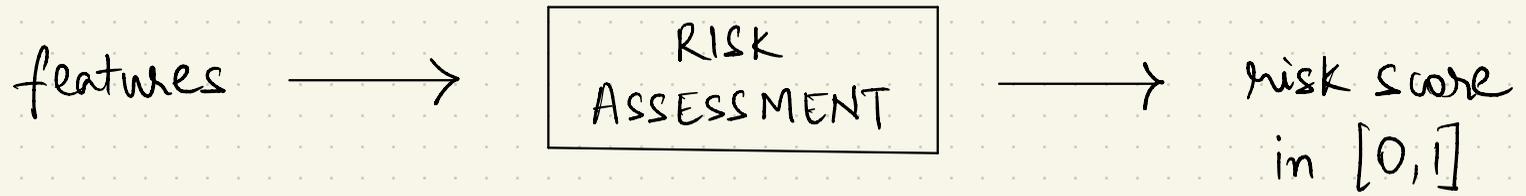




Criminal
records



bail decision



features



risk score
in $[0,1]$

weather data

medical data

chance of rain

probability of disease



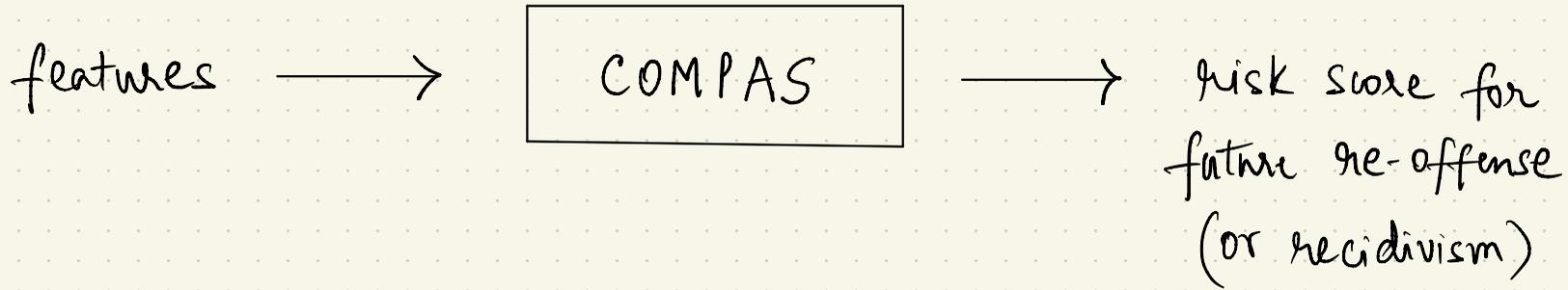
Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

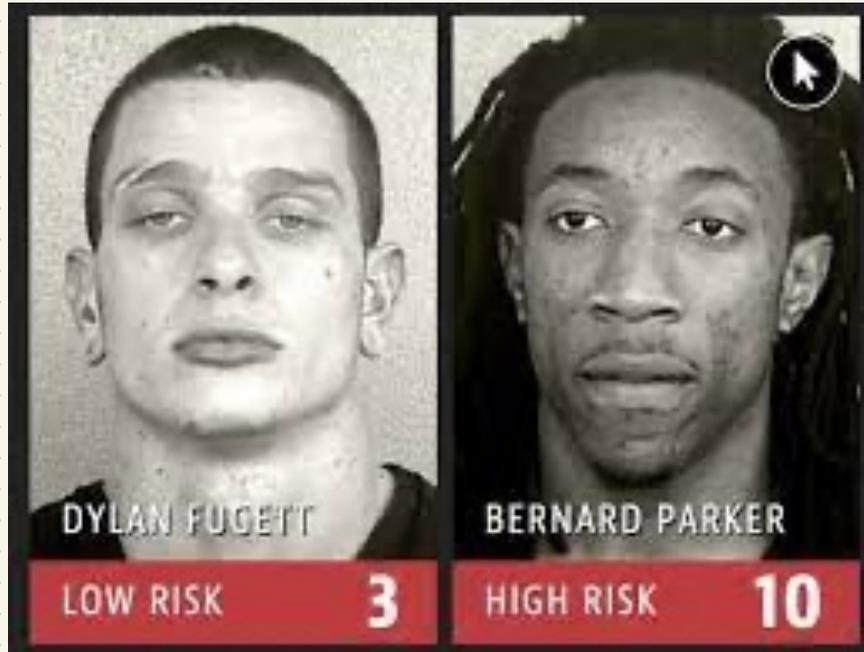
by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

May 23, 2016

Correctional Offender Management Profiling for Alternative Sanctions



ProPublica's Concern



ProPublica's Concern

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

ProPublica's Concern

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

African-American defendants more likely to be incorrectly labeled high-risk

White defendants more likely to be incorrectly labeled low-risk

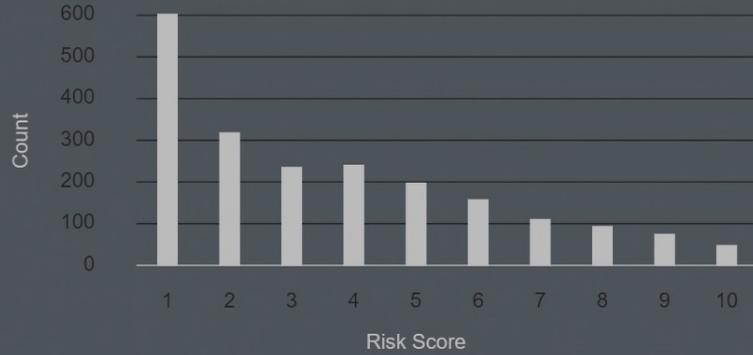
ProPublica's Concern

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

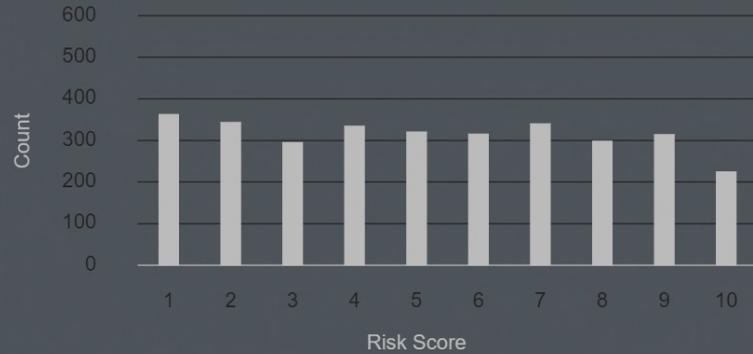
African-American defendants more likely to be incorrectly labeled high-risk

White defendants more likely to be incorrectly labeled low-risk

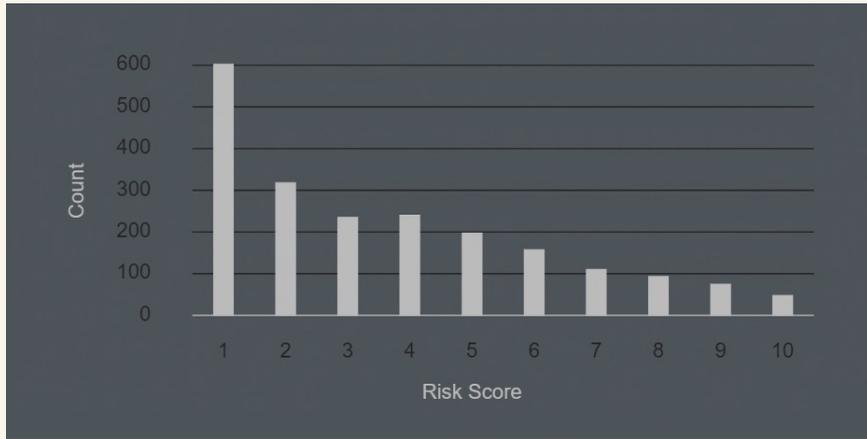
Imbalanced error rates



Black or White?



Black or White?



White



Black

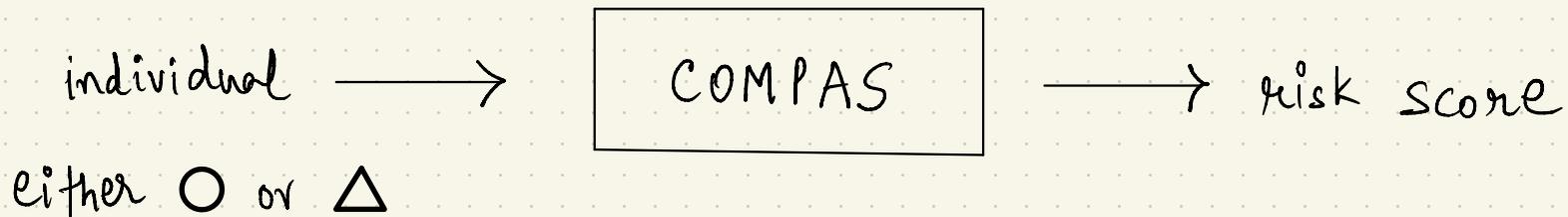
BALANCED ERROR RATES

BALANCED ERROR RATES

Two types/groups of individuals : \circ and \triangle

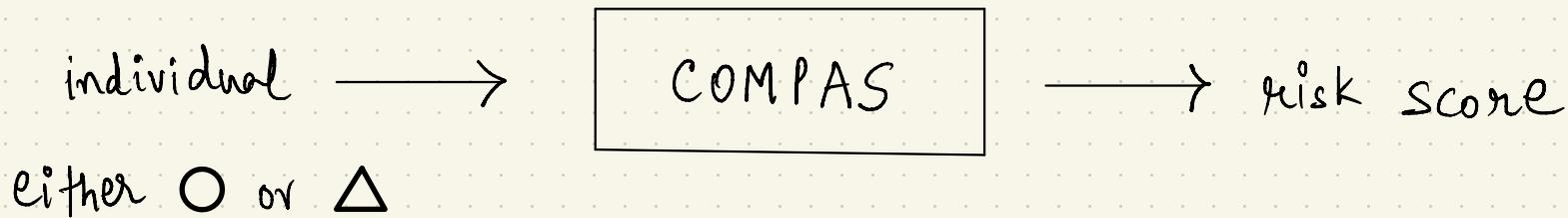
BALANCED ERROR RATES

Two types/groups of individuals : \bigcirc and \triangle



BALANCED ERROR RATES

Two types/groups of individuals : \bigcirc and \triangle



Actual outcomes

\bigcirc : \bigcirc reoffends

\triangle : \triangle reoffends

\bigcirc : \bigcirc doesn't reoffend

\triangle : \triangle doesn't reoffend

BALANCED ERROR RATES

Balance for red class :

$$\text{avg. score of } \bigcirc\text{s} = \text{avg. score of } \triangle\text{s}$$

BALANCED ERROR RATES

Balance for red class:

$$\text{avg. score of } \bigcirc\text{s} = \text{avg. score of } \triangle\text{s}$$

subject to *reaffording*, the score of a typical \bigcirc is same as
" " " " \triangle .

BALANCED ERROR RATES

Balance for red class :

$$\text{avg. score of } \bigcirc\text{s} = \text{avg. score of } \triangle\text{s}$$

subject to *reaffording*, the score of a typical \bigcirc is same as
" " " " \triangle .

Balance for green class :

$$\text{avg. score of } \bigcirc\text{s} = \text{avg. score of } \triangle\text{s}$$

BALANCED ERROR RATES

Balance for red class :

$$\text{avg. score of } \bigcirc\text{s} = \text{avg. score of } \triangle\text{s}$$

subject to reoffending, the score of a typical \bigcirc is same as
" " " " \triangle .

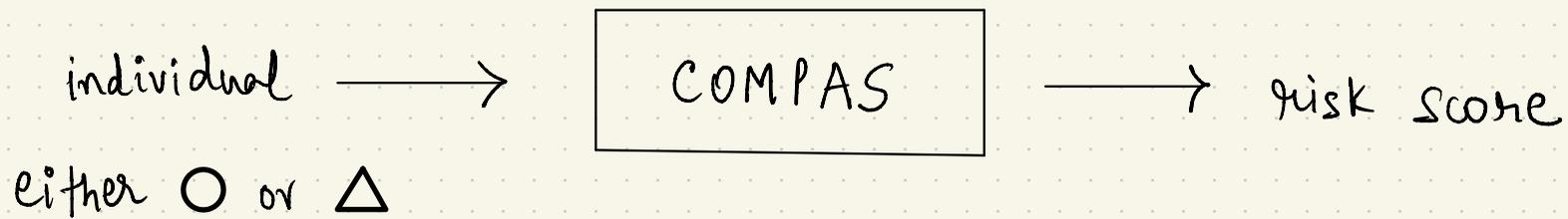
Balance for green class :

$$\text{avg. score of } \bigcirc\text{s} = \text{avg. score of } \triangle\text{s}$$

subject to not reoffending, the score of a typical \bigcirc is same as
" " " " \triangle .

COMPAS' Rebuttal

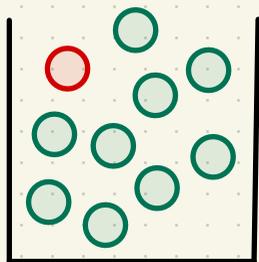
Calibration



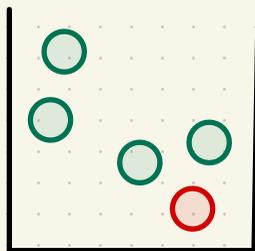


Can divide same-score individuals
into bins

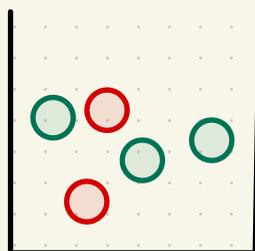
CALIBRATION



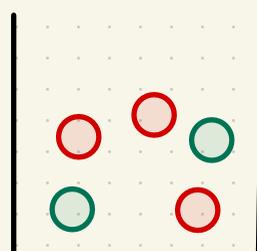
0.1



0.2

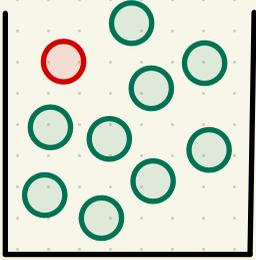


0.4

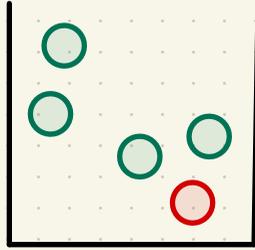


0.6

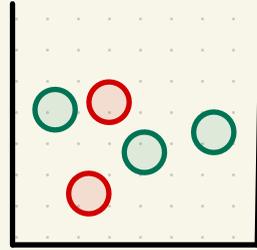
CALIBRATION



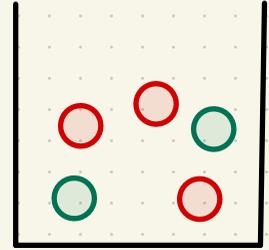
0.1



0.2



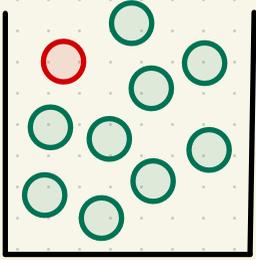
0.4



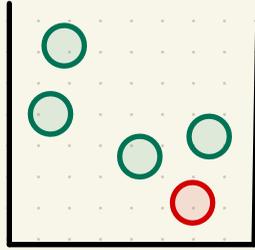
0.6

0.2 means 20% regardless of shape

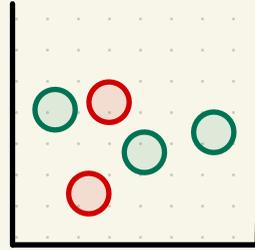
CALIBRATION



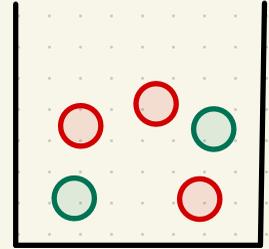
0.1



0.2



0.4

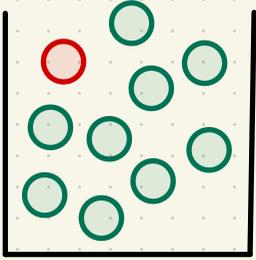


0.6

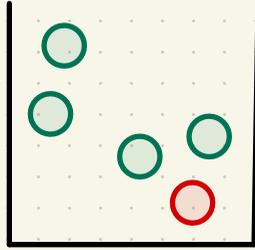
0.2 means 20% regardless of shape

(same holds for  )

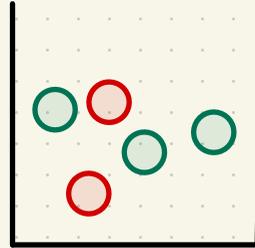
CALIBRATION



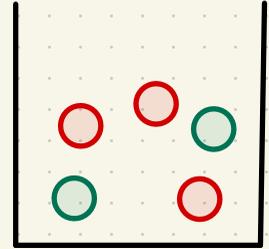
0.1



0.2



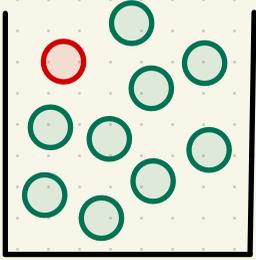
0.4



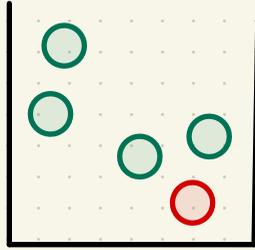
0.6

Why is calibration *meaningful*?

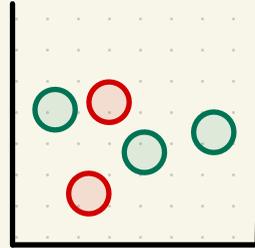
CALIBRATION



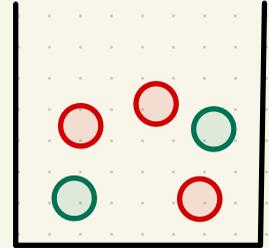
0.1



0.2



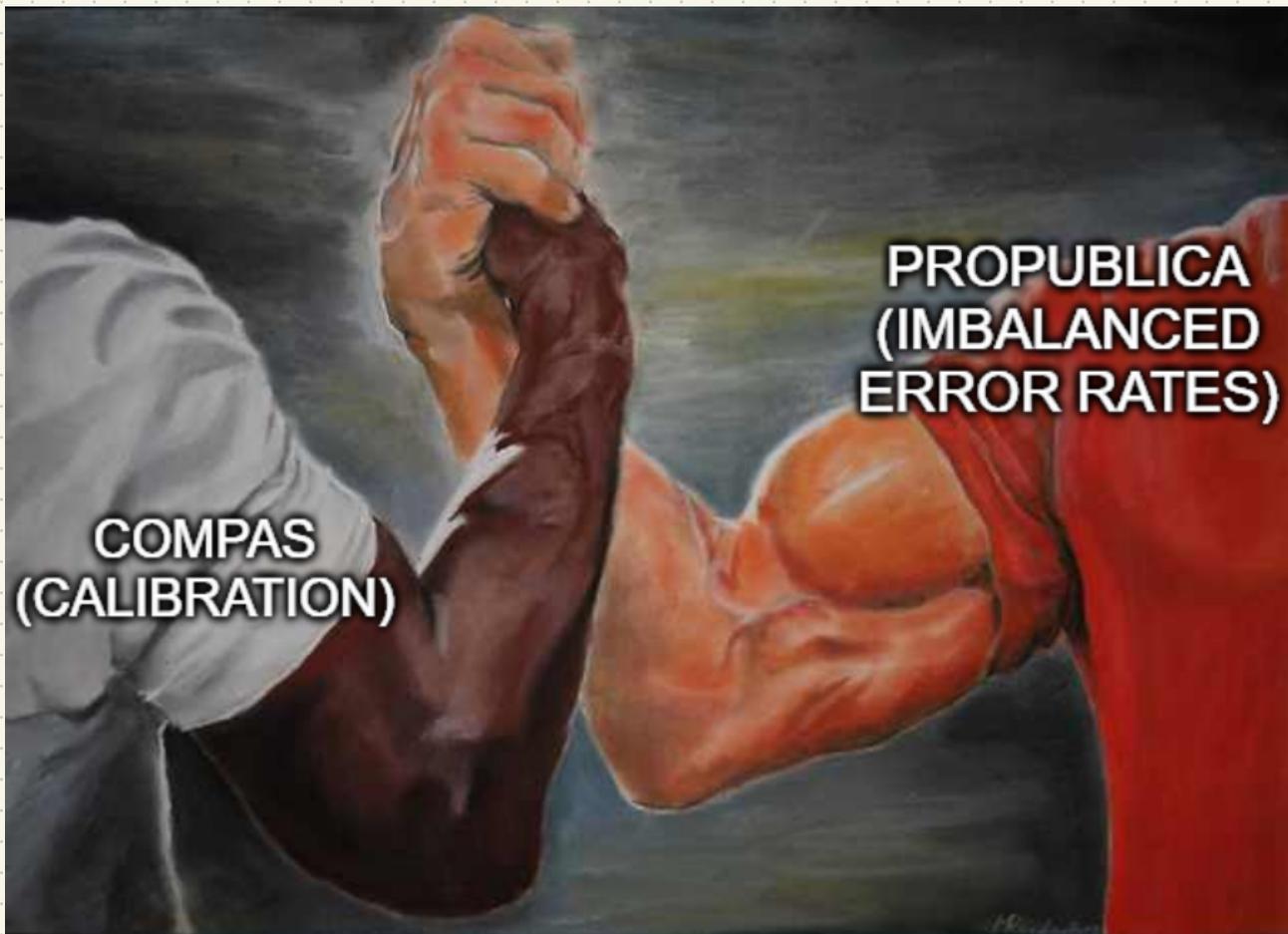
0.4



0.6

Why is calibration *meaningful*?

Treat same-score individuals similarly regardless of shape.



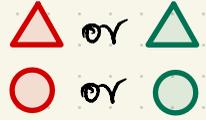
Can we have the best of both worlds?

Balanced error rates + calibration

BEST OF BOTH WORLDS?

BEST OF BOTH WORLDS?

Perfect prediction : Can correctly predict *only*
from the features.



BEST OF BOTH WORLDS?

Perfect prediction : Can correctly predict **only**  or 
 or 
from the features.

Calibration ?

BEST OF BOTH WORLDS?

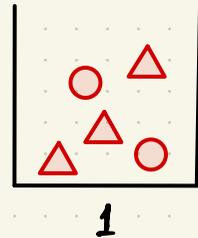
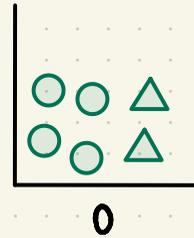
Perfect prediction : Can correctly predict **only**  or 
 or 
from the features.

Calibration? Yes!

BEST OF BOTH WORLDS?

Perfect prediction : Can correctly predict **only**  or 
 or 
from the features.

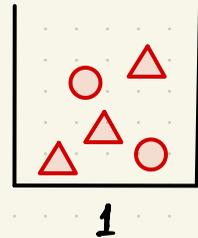
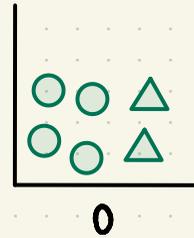
Calibration ? **Yes!**



BEST OF BOTH WORLDS?

Perfect prediction : Can correctly predict **only**  or  or  or 
from the features.

Calibration? **Yes!**

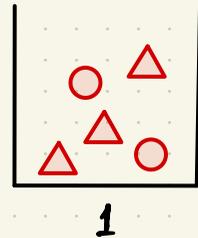
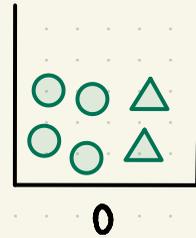


Balanced error rates?

BEST OF BOTH WORLDS?

Perfect prediction : Can correctly predict **only**  or 
 or 
from the features.

Calibration? **Yes!**

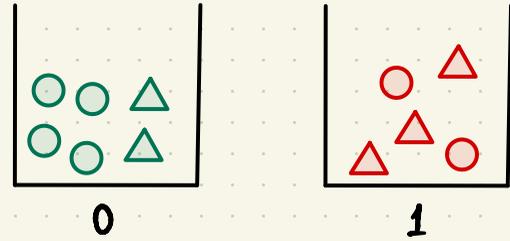


Balanced error rates? **Yes!**

BEST OF BOTH WORLDS?

Perfect prediction : Can correctly predict **only**  or 
 or 
from the features.

Calibration? **Yes!**



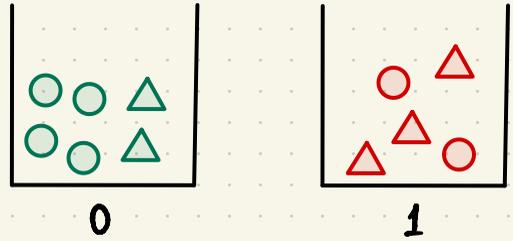
Balanced error rates? **Yes!**

avg score of s = avg score of 
" " s = " " 

BEST OF BOTH WORLDS?

Perfect [✓] prediction : Can correctly predict *only*  or 
 or 
from the features.

Calibration? Yes!



Balanced error rates? Yes!

avg score of s = avg score of 
" " s = " " 

BEST OF BOTH WORLDS?

Equal base rates: when $\frac{\# \triangle_s}{\# \Delta_s} = \frac{\# \circ_s}{\# \bigcirc_s} = p \text{ (say)}$

BEST OF BOTH WORLDS?

Equal base rates: when $\frac{\# \triangle_s}{\# \Delta_s} = \frac{\# \circ_s}{\# \bigcirc_s} = p \text{ (say)}$

Calibration?

BEST OF BOTH WORLDS?

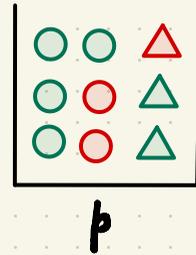
Equal base rates: when $\frac{\# \triangle_s}{\# \Delta_s} = \frac{\# \circ_s}{\# \bigcirc_s} = p \text{ (say)}$

Calibration? Yes!

BEST OF BOTH WORLDS?

Equal base rates: when $\frac{\# \triangle_s}{\# \Delta_s} = \frac{\# \circ_s}{\# \circ_s} = p$ (say)

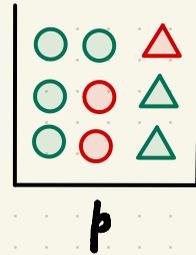
Calibration? Yes!



BEST OF BOTH WORLDS?

Equal base rates: when $\frac{\# \triangle_s}{\# \Delta_s} = \frac{\# \circ_s}{\# \circ_s} = p$ (say)

Calibration? Yes!

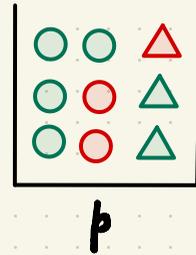


Balanced error rates?

BEST OF BOTH WORLDS?

Equal base rates: when $\frac{\# \triangle_s}{\# \Delta_s} = \frac{\# \circ_s}{\# \circ_s} = p$ (say)

Calibration? Yes!

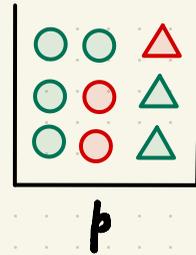


Balanced error rates? Yes!

BEST OF BOTH WORLDS?

Equal base rates: when $\frac{\# \triangle_s}{\# \Delta_s} = \frac{\# \circ_s}{\# \bigcirc_s} = p$ (say)

Calibration? Yes!

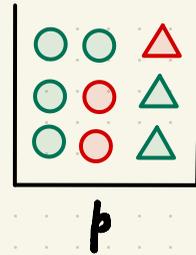


Balanced error rates? Yes! avg score of $\triangle_s =$ avg score of $\circ_s = p$
" " $\triangle_s =$ " " $\bigcirc_s = p$

BEST OF BOTH WORLDS?

Equal base rates:  when $\frac{\# \triangle_s}{\# \Delta_s} = \frac{\# \circ_s}{\# \bigcirc_s} = p$ (say)

Calibration? Yes!



Balanced error rates? Yes! avg score of $\triangle_s =$ avg score of $\circ_s = p$
" " $\triangle_s =$ " " $\bigcirc_s = p$

BEST OF BOTH WORLDS?

Calibration

Balanced error rates

Perfect prediction

✓

✓

Equal base rates

✓

✓

BEST OF BOTH WORLDS?

	Calibration	Balanced error rates
Perfect prediction	✓	✓
Equal base rates	✓	✓

Any others?

BEST OF BOTH WORLDS?

	Calibration	Balanced error rates
Perfect prediction	✓	✓
Equal base rates	✓	✓

Any others?

NO!

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Theorem : Any risk assessment method that satisfies calibration and balanced error rates must either be perfect or have equal base rates.

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Theorem : Any risk assessment method that satisfies calibration and balanced error rates must either be perfect or have equal base rates.

- * Not about computational power
- * Impossibility of "balancing averages" by assigning scores

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

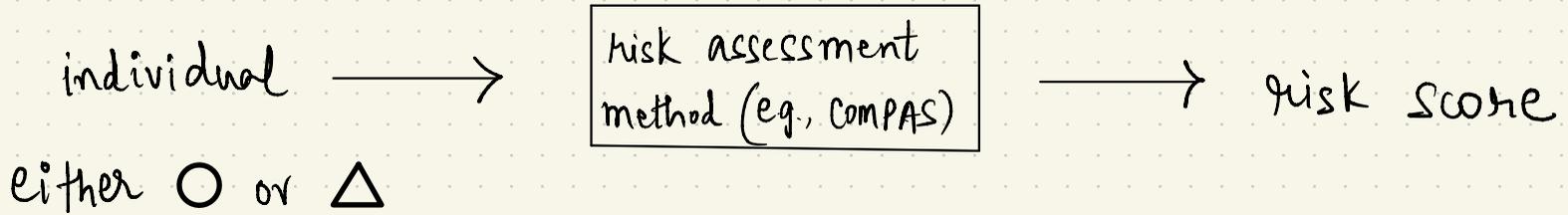
Proof sketch:



[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:



[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:



$$n_{\circ} = \# \circ + \# \ominus$$

$$n_{\triangle} = \# \triangle + \# \triangle$$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch: By calibration:

Bin with score b_1 :

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch: By calibration:

Bin with score b_1 : $b_1 \times [\# \Delta s \text{ in bin } b_1] = \# \Delta s \text{ in bin } b_1$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch: By calibration:

Bin with score b_1 : $b_1 \times [\# \Delta s \text{ in bin } b_1] = \# \Delta s \text{ in bin } b_1$

Bin with score b_2 : $b_2 \times [\# \Delta s \text{ in bin } b_2] = \# \Delta s \text{ in bin } b_2$

⋮

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch: By calibration:

Bin with score b_1 : $b_1 \times [\# \Delta_s \text{ in bin } b_1] = \# \Delta_s \text{ in bin } b_1$

Bin with score b_2 : $b_2 \times [\# \Delta_s \text{ in bin } b_2] = \# \Delta_s \text{ in bin } b_2$

\vdots

Summing over bins Total score of all $\Delta_s = \# \Delta_s$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch: By calibration:

$$\text{Bin with score } b_1: \quad b_1 \times [\# \Delta_s \text{ in bin } b_1] = \# \Delta_s \text{ in bin } b_1$$

$$\text{Bin with score } b_2: \quad b_2 \times [\# \Delta_s \text{ in bin } b_2] = \# \Delta_s \text{ in bin } b_2$$

\vdots

$$\text{Summing over bins} \quad \text{Total score of all } \Delta_s = \# \Delta_s$$

$$\text{Similarly,} \quad \text{Total score of all } 0_s = \# 0_s$$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch: By calibration:

Bin with score b_1 : $b_1 \times [\# \Delta_s \text{ in bin } b_1] = \# \Delta_s \text{ in bin } b_1$

Bin with score b_2 : $b_2 \times [\# \Delta_s \text{ in bin } b_2] = \# \Delta_s \text{ in bin } b_2$

⋮

Summing over bins Total score of all $\Delta_s = \# \Delta_s$

Similarly, Total score of all $\circ_s = \# \circ_s$



[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:

r = avg score of all red shapes

g = " " " green "

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch: r = avg score of all red shapes
 g = " " " green "

By balanced error rates :

$$\text{avg score of } \triangle_s = \text{avg score of } \circ_s = r$$

$$\text{avg score of } \triangle_s = \text{avg score of } \circ_s = g$$

[Kleinberg, Mullainathan, Raghavan ITCs 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:

Total score of all Δ_s = Total score of Δ_s + Total score of Δ_s

[Kleinberg, Mullainathan, Raghavan ITCs 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:

$$\begin{aligned} \text{Total score of all } \Delta_s &= \text{Total score of } \Delta_s + \text{Total score of } \triangle_s \\ &= \# \Delta_s \cdot p + \# \triangle_s \cdot q \end{aligned}$$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:

$$\begin{aligned} \text{Total score of all } \Delta_s &= \text{Total score of } \Delta_s + \text{Total score of } \triangleleft_s \\ &= \# \Delta_s \cdot p + \# \triangleleft_s \cdot q \\ &= \# \Delta_s \cdot p + (n_\Delta - \# \Delta_s) \cdot q \end{aligned}$$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:

$$\begin{aligned} \text{Total score of all } \Delta_s &= \text{Total score of } \Delta_s + \text{Total score of } \triangle_s \\ &= \# \Delta_s \cdot p + \# \triangle_s \cdot q \\ &= \# \Delta_s \cdot p + (n_{\Delta} - \# \Delta_s) \cdot q \\ &= \# \Delta_s \quad \text{by } \text{☺} \end{aligned}$$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:

$$\# \Delta_s = \# \Delta_s \cdot p + (n_{\Delta} - \# \Delta_s) \cdot q$$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch:

$$\# \Delta_s = \# \Delta_s \cdot r + (n_\Delta - \# \Delta_s) \cdot g$$

$$\Rightarrow g = \frac{\# \Delta_s (1 - r)}{n_\Delta - \# \Delta_s}$$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

Proof sketch: $\# \Delta_s = \# \Delta_s \cdot r + (n_\Delta - \# \Delta_s) \cdot g$

$$\Rightarrow g = \frac{\# \Delta_s (1 - r)}{n_\Delta - \# \Delta_s}$$

Similarly,

$$g = \frac{\# \circ_s (1 - r)}{n_\circ - \# \circ_s}$$

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

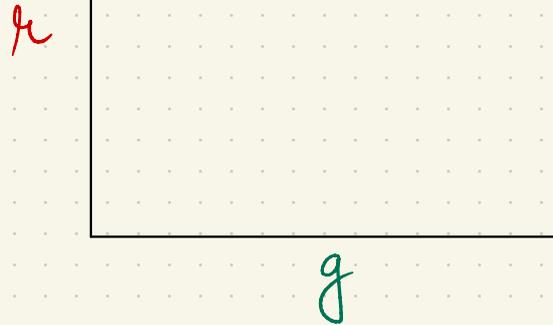
Proof sketch:

$$\# \Delta_s = \# \Delta_s \cdot r + (n_{\Delta} - \# \Delta_s) \cdot g$$

$$\Rightarrow g = \frac{\# \Delta_s (1 - r)}{n_{\Delta} - \# \Delta_s}$$

Similarly,

$$g = \frac{\# \circ_s (1 - r)}{n_{\circ} - \# \circ_s}$$



[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

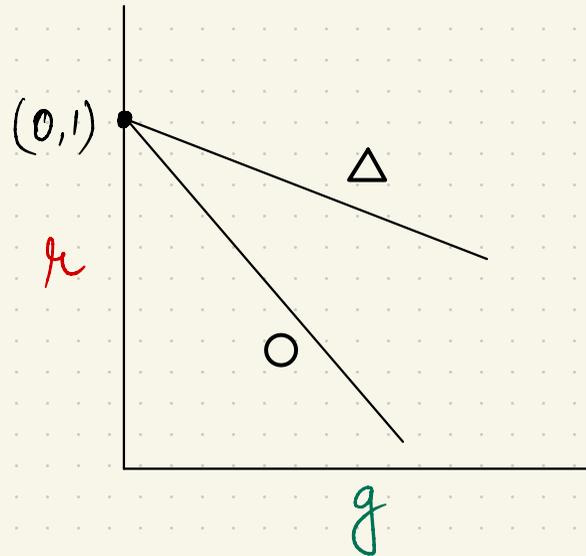
Proof sketch:

$$\# \Delta_s = \# \Delta_s \cdot r + (n_{\Delta} - \# \Delta_s) \cdot g$$

$$\Rightarrow g = \frac{\# \Delta_s (1 - r)}{n_{\Delta} - \# \Delta_s}$$

Similarly,

$$g = \frac{\# \circ_s (1 - r)}{n_{\circ} - \# \circ_s}$$



[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

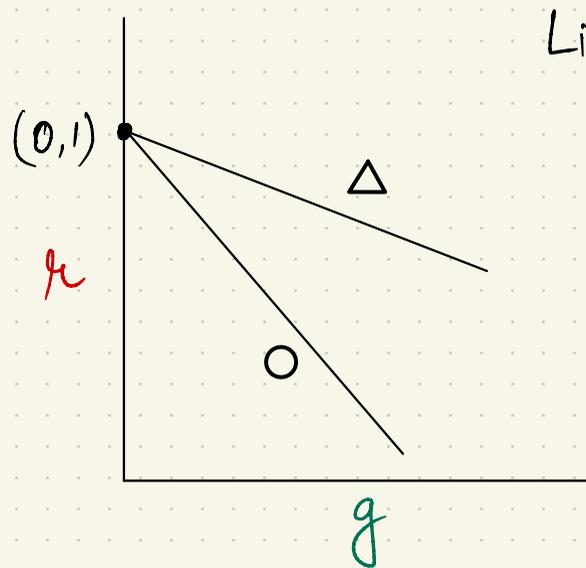
Proof sketch:

$$\# \Delta_s = \# \Delta_s \cdot r + (n_{\Delta} - \# \Delta_s) \cdot g$$

$$\Rightarrow g = \frac{\# \Delta_s (1 - r)}{n_{\Delta} - \# \Delta_s}$$

Similarly,

$$g = \frac{\# \circ_s (1 - r)}{n_{\circ} - \# \circ_s}$$



[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

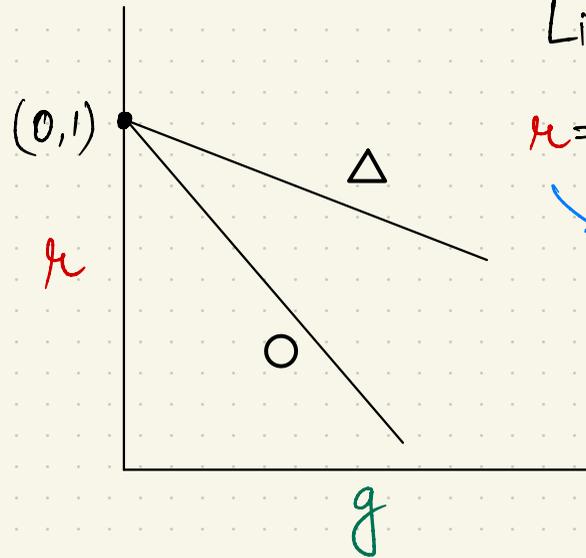
Proof sketch:

$$\# \Delta_s = \# \Delta_s \cdot r + (n_{\Delta} - \# \Delta_s) \cdot g$$

$$\Rightarrow g = \frac{\# \Delta_s (1 - r)}{n_{\Delta} - \# \Delta_s}$$

Similarly,

$$g = \frac{\# \circ_s (1 - r)}{n_{\circ} - \# \circ_s}$$



Lines intersect when

$$r = 1, g = 0$$

\rightarrow perfect prediction

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

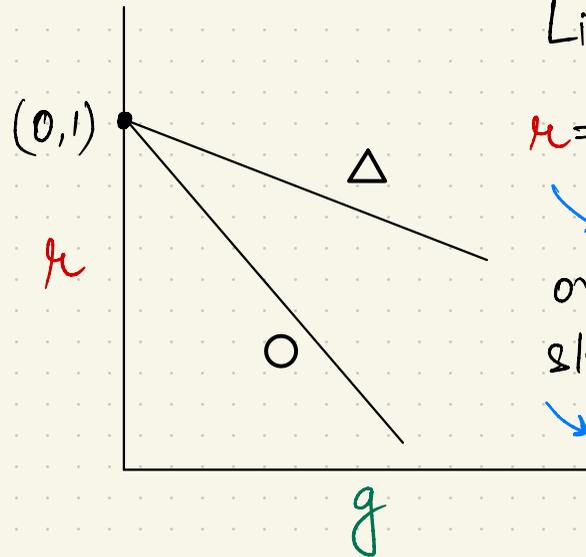
Proof sketch:

$$\# \Delta_s = \# \Delta_s \cdot r + (n_\Delta - \# \Delta_s) \cdot g$$

$$\Rightarrow g = \frac{\# \Delta_s (1 - r)}{n_\Delta - \# \Delta_s}$$

Similarly,

$$g = \frac{\# \circ_s (1 - r)}{n_\circ - \# \circ_s}$$



Lines intersect when

$$r = 1, g = 0$$

\rightarrow perfect prediction

or

slopes are equal

\rightarrow equal base rates

[Kleinberg, Mullainathan, Raghavan ITCS 2017]

Thm: Calibration and balanced error rates \Rightarrow perfect prediction or equal base rates

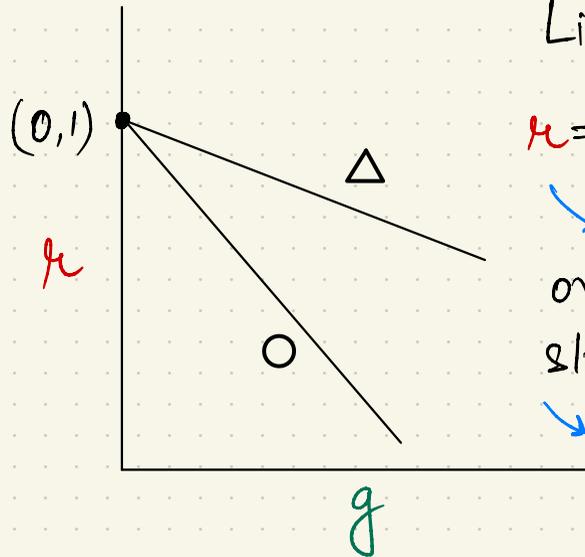
Proof sketch:

$$\# \Delta_s = \# \Delta_s \cdot r + (n_{\Delta} - \# \Delta_s) \cdot g$$

$$\Rightarrow g = \frac{\# \Delta_s (1 - r)}{n_{\Delta} - \# \Delta_s}$$

Similarly,

$$g = \frac{\# \circ_s (1 - r)}{n_{\circ} - \# \circ_s}$$



Lines intersect when

$$r = 1, g = 0$$

\rightarrow perfect prediction

or

slopes are equal

\rightarrow equal base rates



EXTENSIONS

* Approximate calibration and approx. balanced error rates also in conflict.

EXTENSIONS

* Approximate calibration and approx. balanced error rates also in conflict.

* Binary classification [Chouldechova 2017; Corbett-Davies, Pierson, Feller, Goel 2016]

EXTENSIONS

- * **Approximate** calibration and **approx.** balanced error rates also in conflict.
- * **Binary** classification [Chouldechova 2017; Corbett-Davies, Pierson, Feller, Goel 2016]
- * Workarounds
 - balanced error rates (w/o calibration) [Hardt, Price, Srebro 2016]
 - calibration and one-sided error only

FIND OUT MORE AT

"Fairness and Machine Learning: Limitations and Opportunities"

Book by Solon Barocas, Moritz Hardt and Arvind Narayanan

<https://fairmlbook.org>