



Fairness-Aware Hyperparameter Optimization

André Cruz

Supervisors:

FEUP - MIEIC - 2019/2020

Pedro Saleiro, Feedzai
Carlos Soares, FEUP

Motivation

Objective

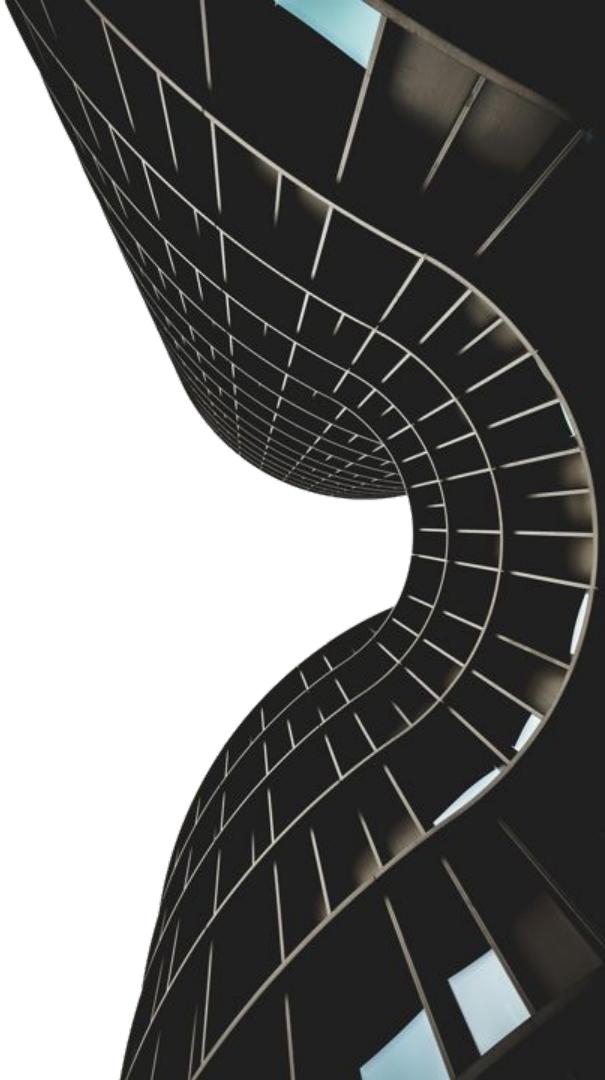
Related Work

Method: Fairband

Experimental Setup

Results & Discussion

Conclusion



Motivation

- ML has an increasing number of applications.
- Widespread reports of algorithmic bias.



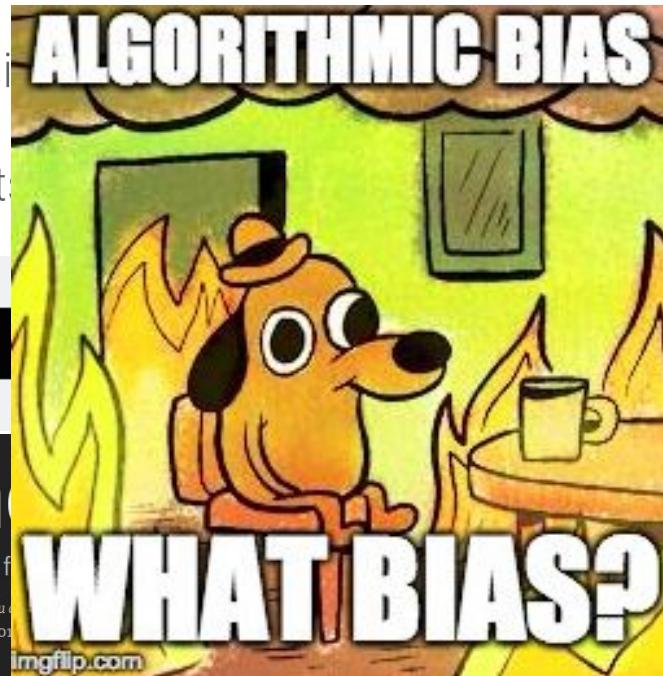
Machine Learning

There's software used across the country to predict if you're likely to commit a crime.

by Julia Angwin, Jeff Larson, Surya Mattu, and Laelaps

May 23, 2016

imgflip.com



es.

ems.

made courts more fair.

I Lending Evolves, and Blacks Face Trouble Getting Mortgages
The New York Times

Amazon scraps secret AI recruiting tool that showed bias against women

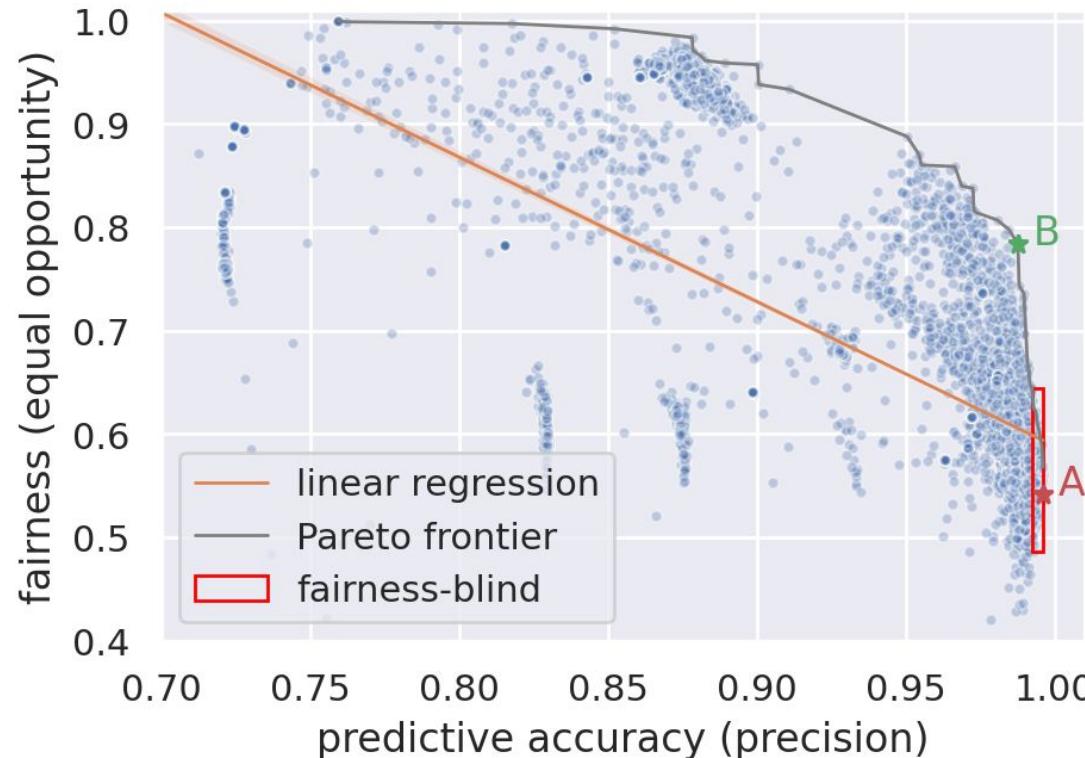
W I R E D

The Apple Card Didn't 'See' Gender—and That's the Problem

Problem Statement

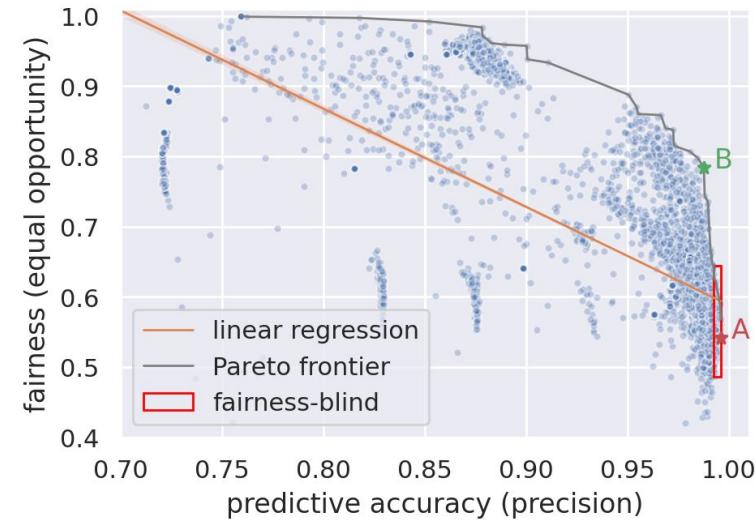
- Current landscape of fairness-aware ML lacks:
 - Practical methodologies
 - Tools for real-world practitioners
- Current bias mitigation methods are:
 - Model- and metric-dependent
 - Added complexity in ML pipelines

Fairness vs Accuracy



Fairness vs Accuracy

1. By optimizing solely for performance, we are unknowingly targeting unfair models.
2. Substantially fairer models can be reached with small decreases in performance.



Motivation

Objective

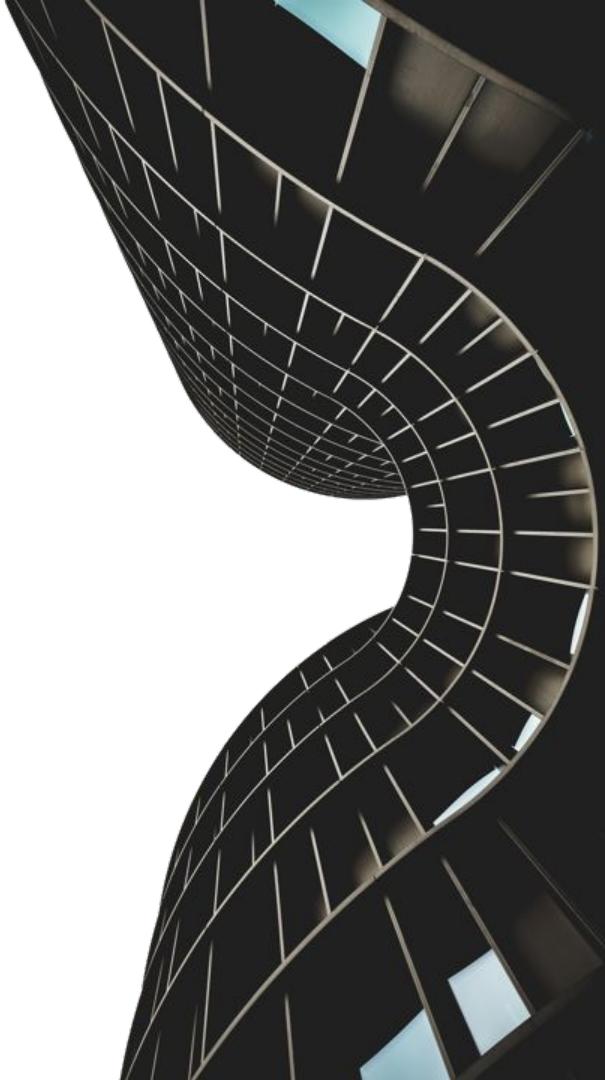
Related Work

Method: Fairband

Experimental Setup

Results & Discussion

Conclusion



Goals

- Enable model development that achieves good fairness-utility trade-offs.
- Method should be model- and metric-agnostic, and easily introduced into current ML pipelines.
- Out-of-the-box model selection, or targeting specific trade-offs.
- Study the fairness-utility trade-off on a real-world setting.



How ?

- Integrating fairness-awareness into hyperparameter optimization (HO).
 - Guide search towards fairer regions of hyperparameter space.
- As black-box optimization, HO is both model- and metric-agnostic.
- HO is already part of most real-world ML pipelines.
 - Potential to introduce fairness at no extra cost or effort.

What are Hyperparameters?

- Model hyperparameters
 - Number of estimators, learning rate, number of layers, ...
- Model type
 - LightGBM, Random Forest, Neural Network, ...
- Sampling hyperparameters
 - Use undersampling? Use oversampling? Balance prevalence across groups?



Any decision available to the Data Scientist

Motivation

Objective

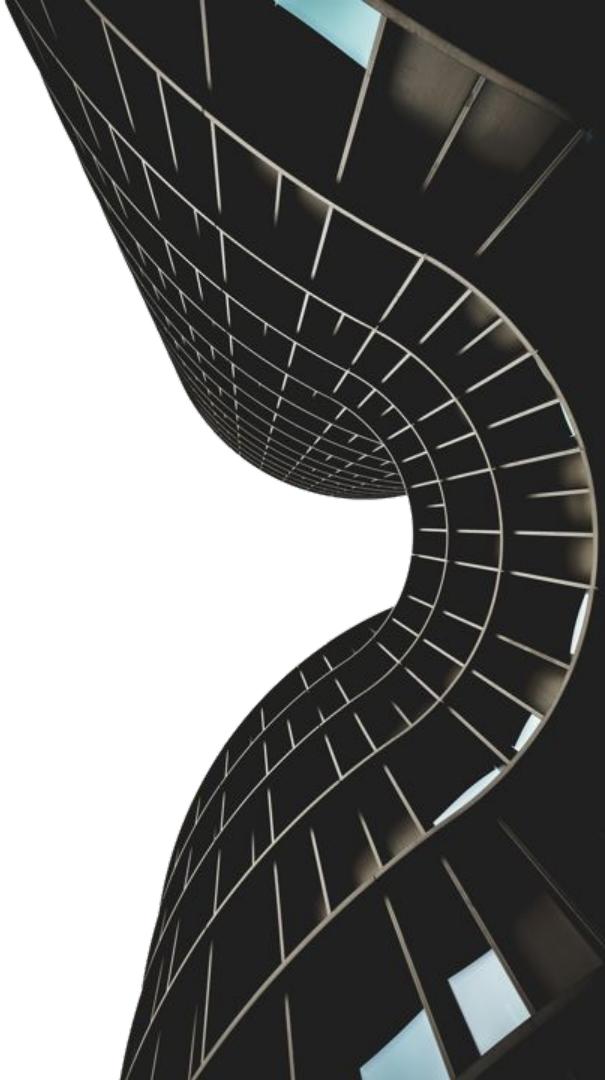
Related Work

Method: Fairband

Experimental Setup

Results & Discussion

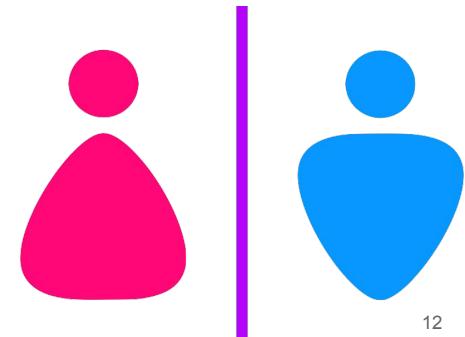
Conclusion



Algorithmic Bias

Disparate error rates among individuals from different sub-groups.

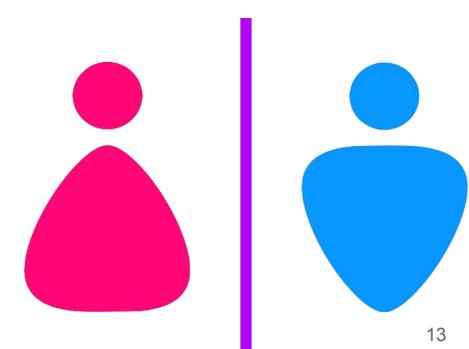
- An ML-powered social security program should successfully find individuals in need of assistance (*true positives*) in equal ratios among different ethnicities or genders.
 - Equal *true positive rates* (TPR).
- Bias reduction techniques attempt to balance these error rates.

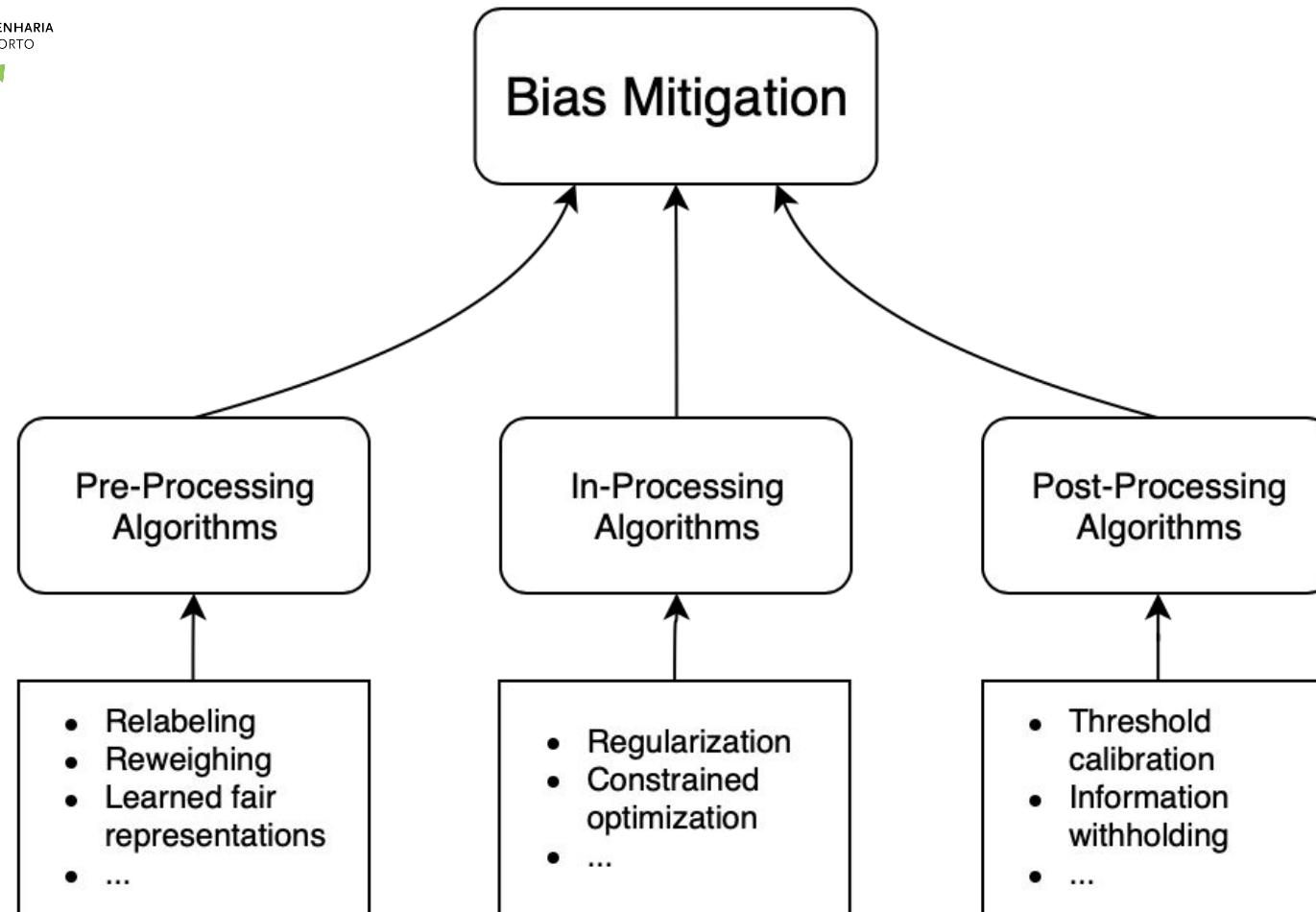


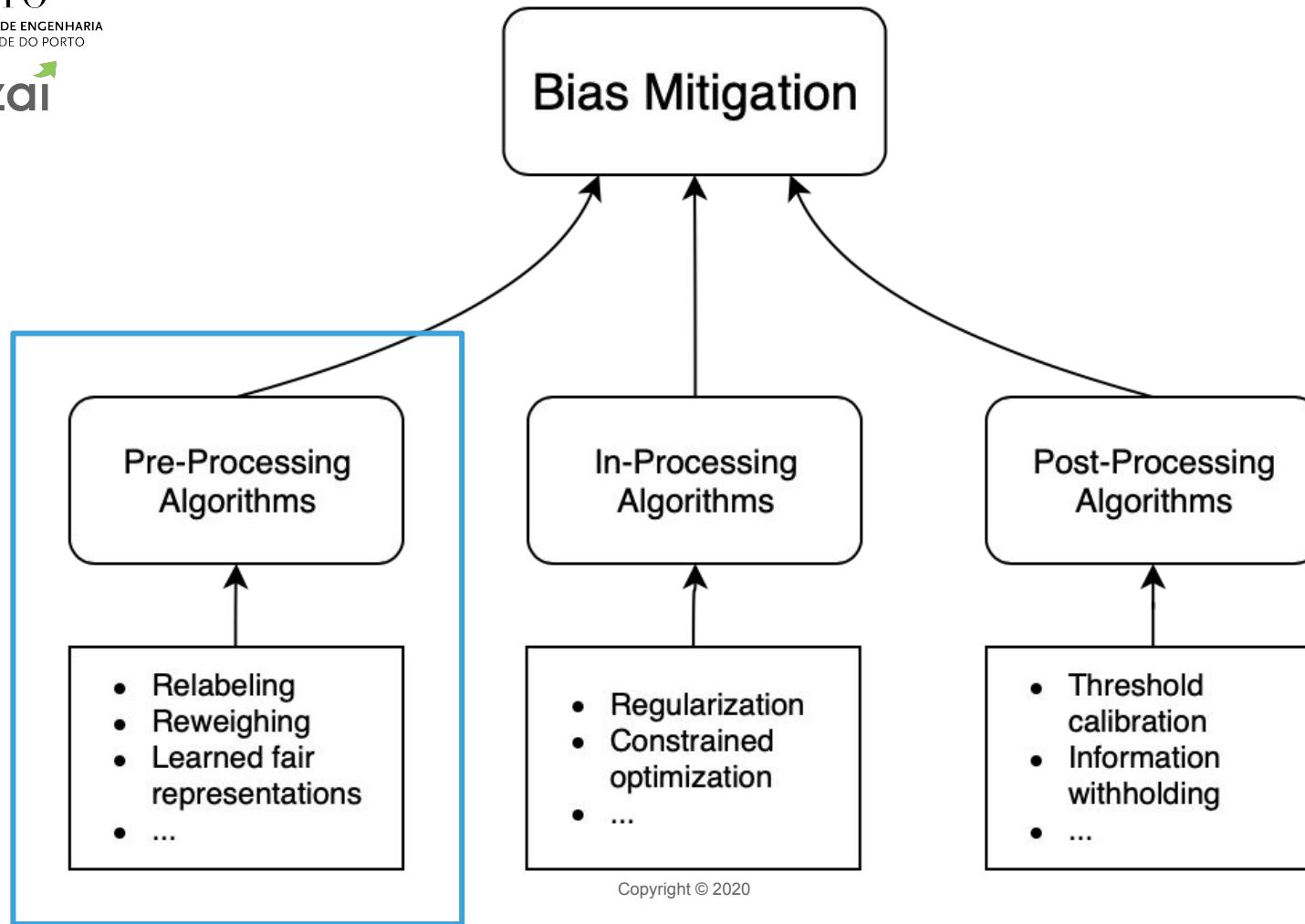
Fairness Metrics

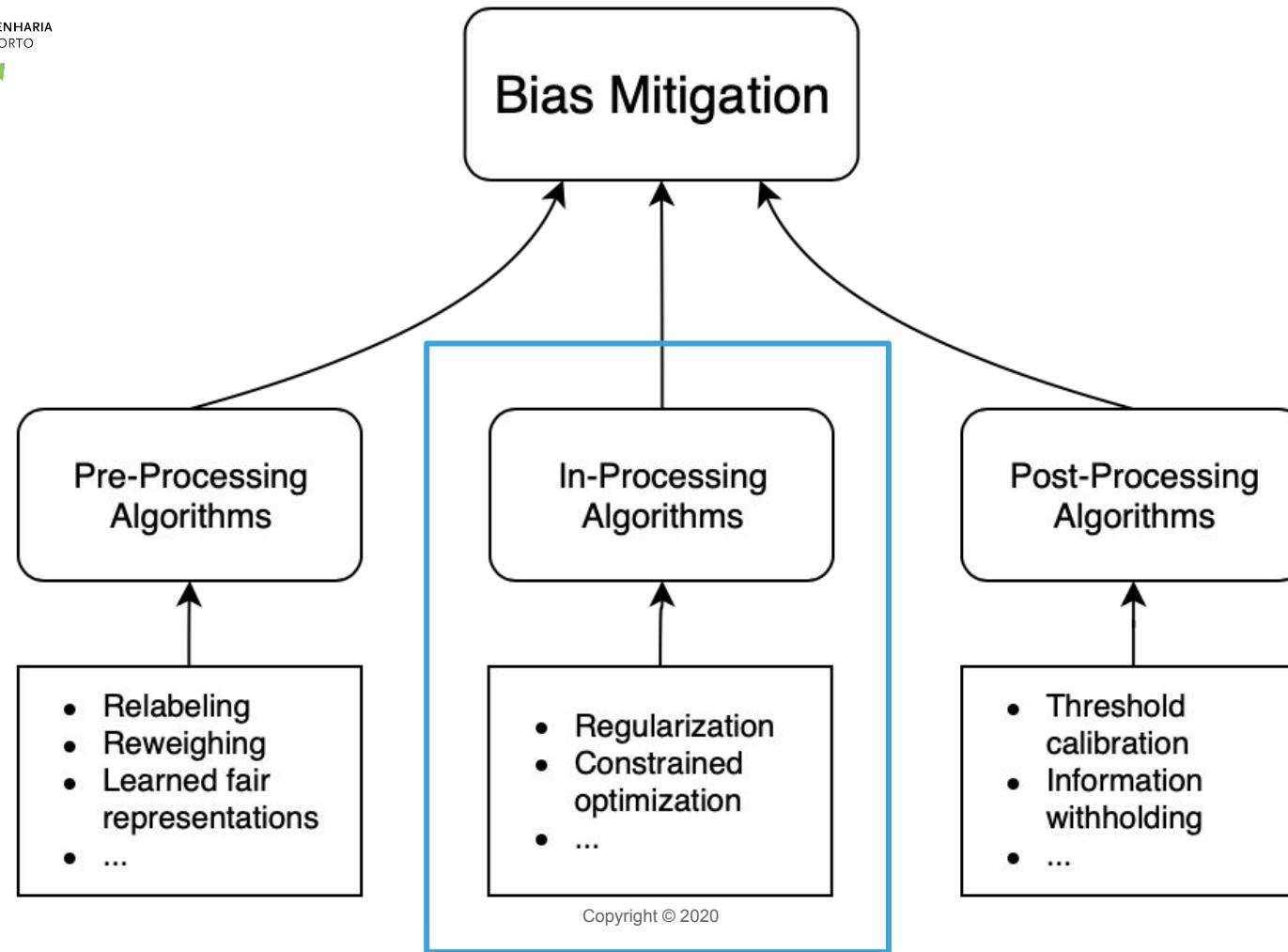
Group Fairness

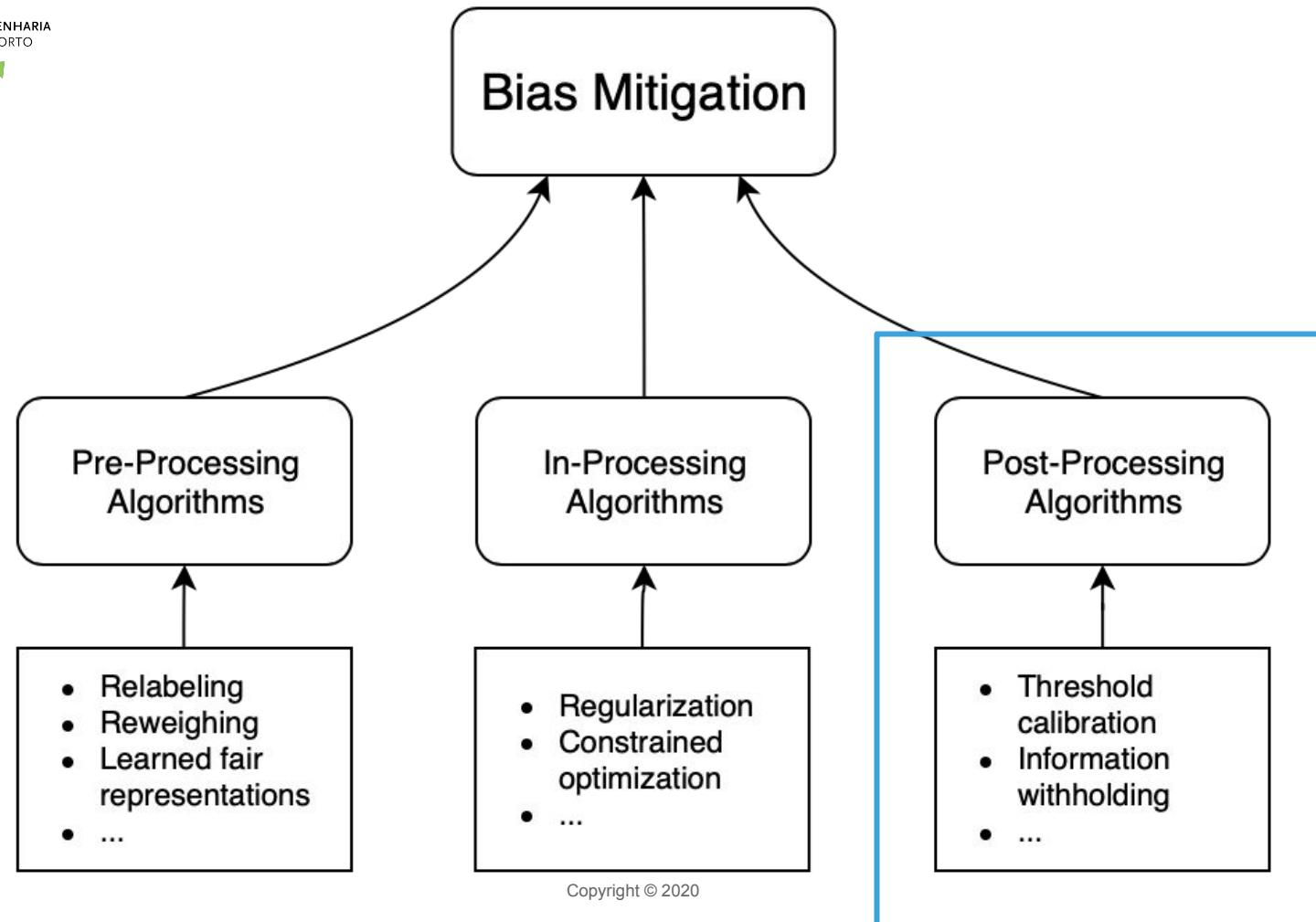
- Equalized Odds
 - Equal TPR and equal FPR
 - $P(\hat{Y}=1 | A=0, Y=y) = P(\hat{Y}=1 | A=1, Y=y)$, $y \in \{0,1\}$
- Equal Opportunity
 - Equal true positive rates
 - $P(\hat{Y}=1 | A=0, Y=1) = P(\hat{Y}=1 | A=1, Y=1)$
- Aggregates of any group-based metric.
 - FNR, FDR, FOR, TPR, TNR, PPR









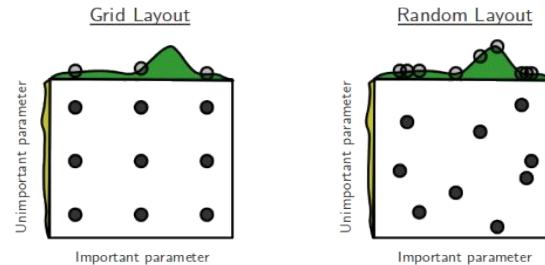


Current Methods

- *Pre-processing* methods cannot guarantee fairness in the end model.
 - As they act in the beginning of the ML pipeline.
- *In-processing* methods are metric-dependent and model-dependent.
 - Even non-existent for numerous algorithms (e.g., LightGBM).
- *Post-processing* methods are inherently sub-optimal.
 - Knowingly training a biased model and afterwards correct its outputs.
 - By acting in the end of the ML pipeline, we can only lose information.
- *All methods* invariably introduce complexity to real-world ML pipelines.



Random Search



Grid Search

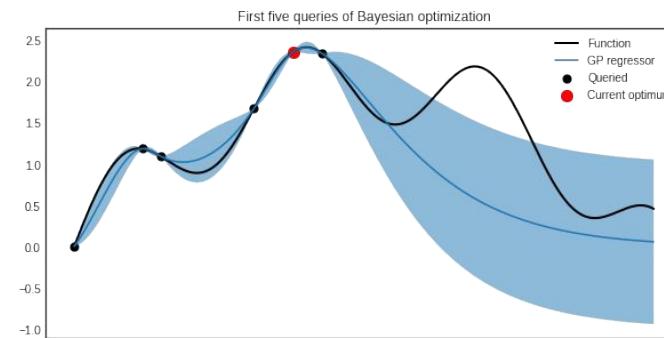
Bayesian Optimization

Hyperparameter Optimization

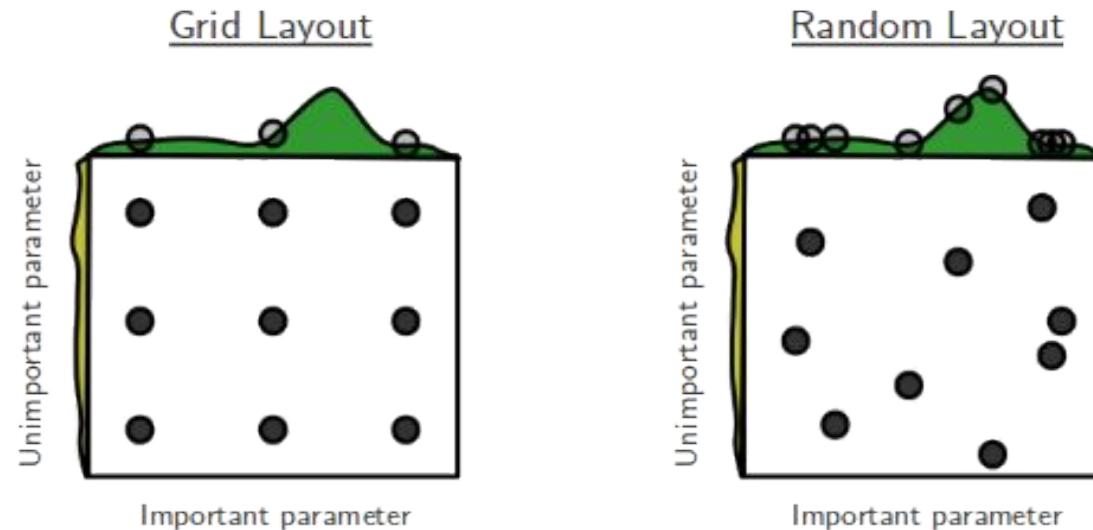
Successive Halving

Hyperband

...

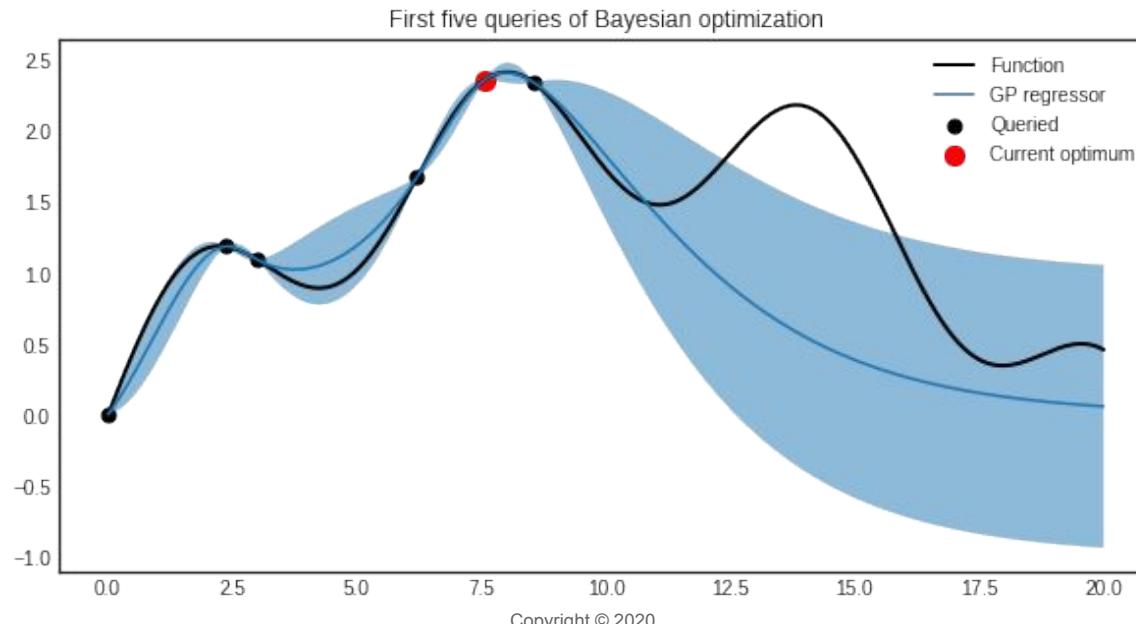


Random and Grid Search



Bayesian Optimization

A framework for model-based optimization of black-box functions



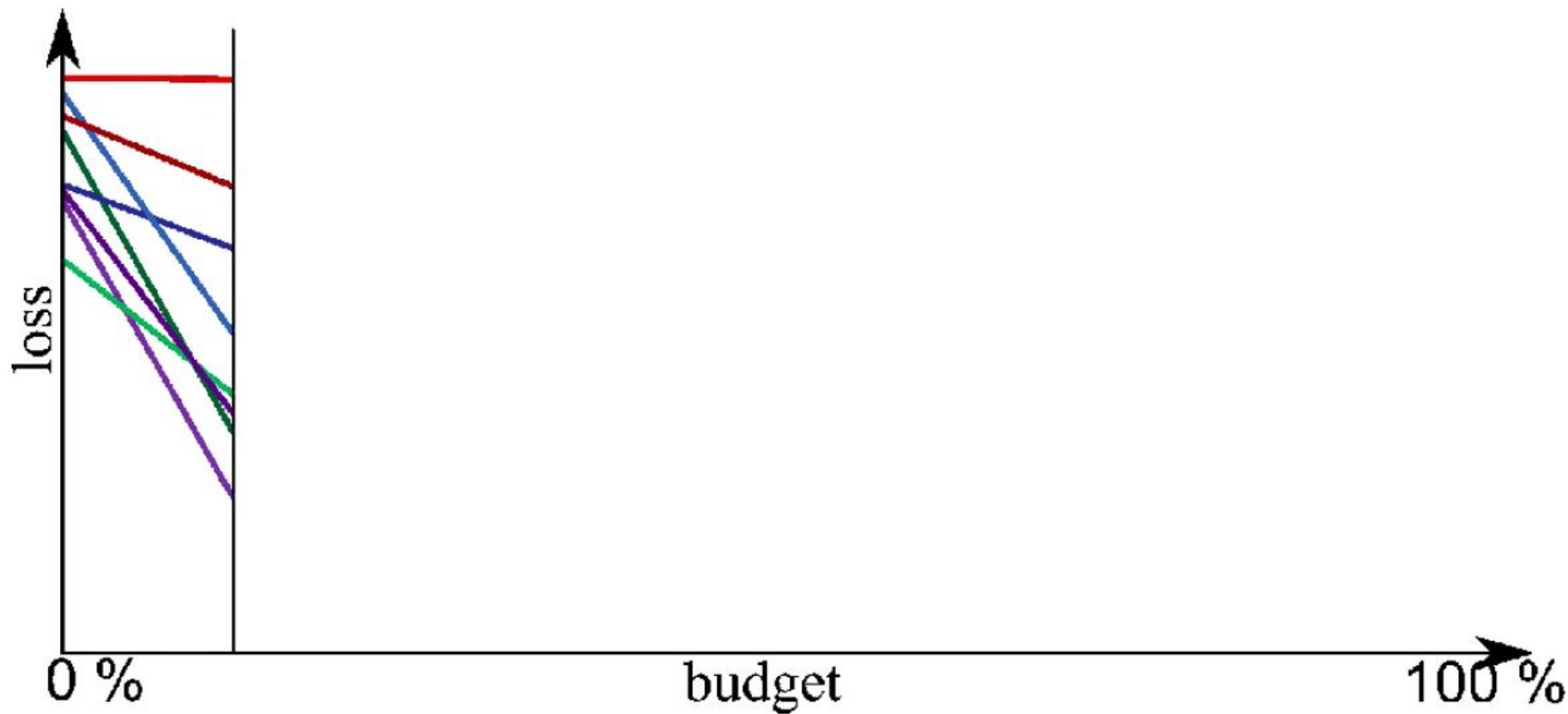
Successive Halving

- Multi-armed bandit setting:

A fixed amount of resources must be allocated between competing choices in a way that maximizes expected gain.

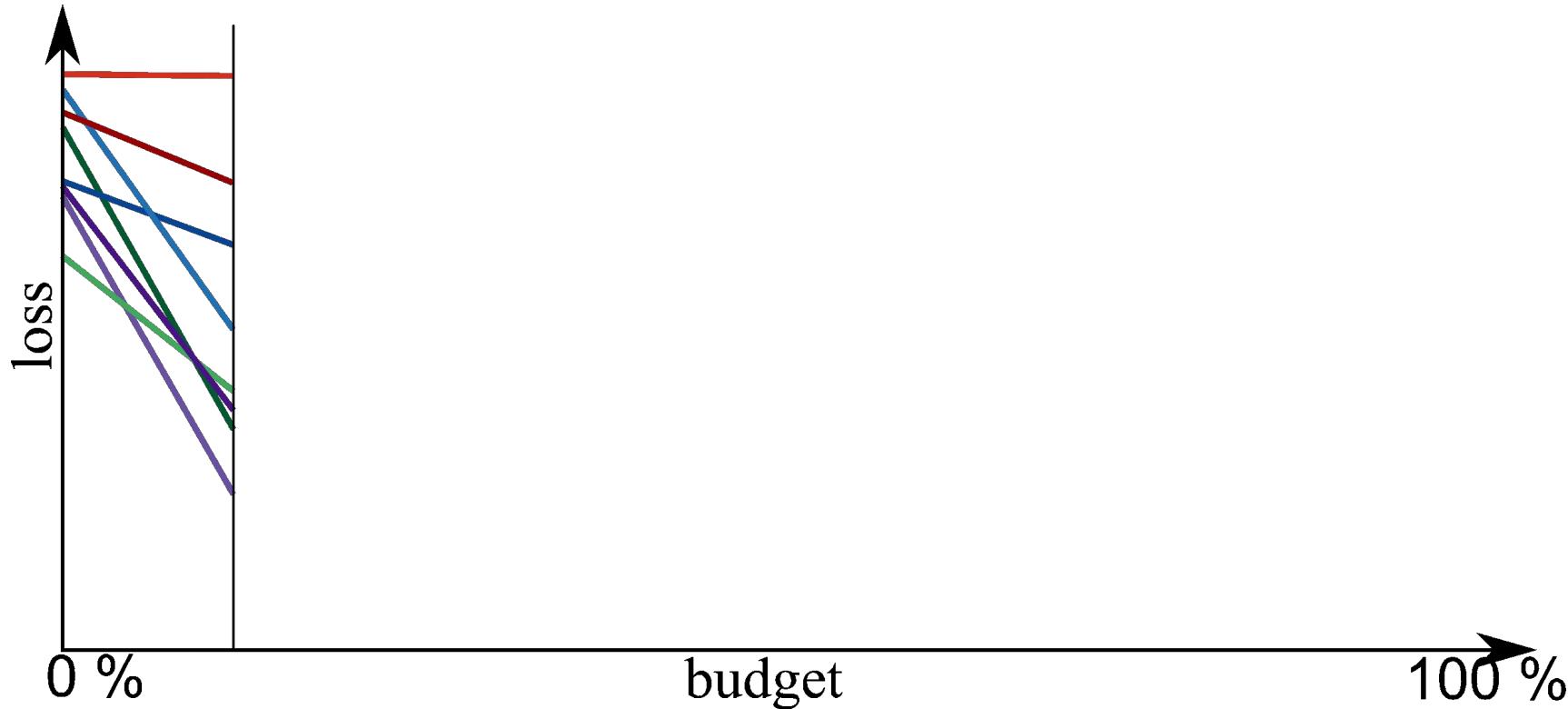


Successive Halving



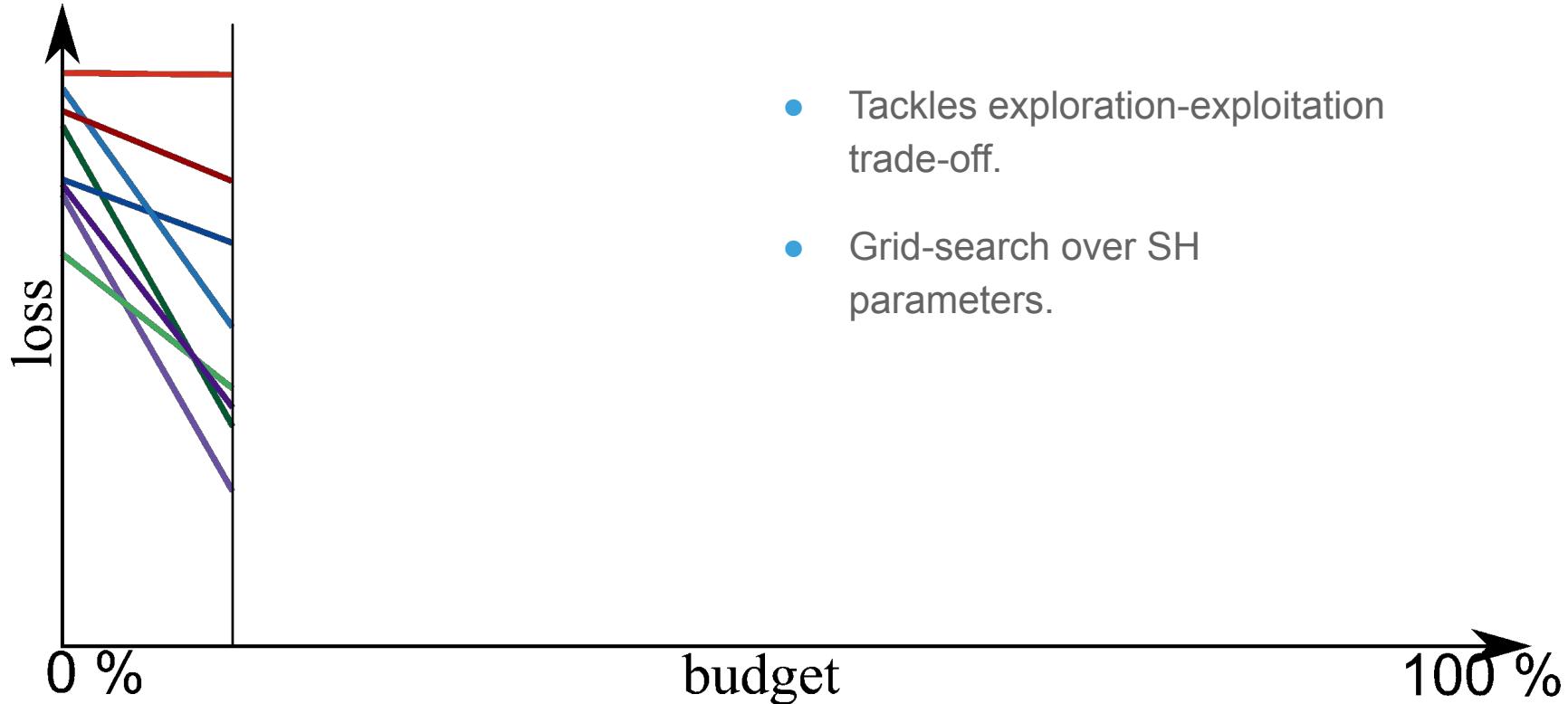
[Jamieson and Talwalkar. Non-stochastic best arm identification and hyper-parameter optimization. ICML 2016]

Successive Halving



[Jamieson and Talwalkar. Non-stochastic best arm identification and hyper-parameter optimization. ICML 2016]

Hyperband: SH \times n times



Motivation

Objective

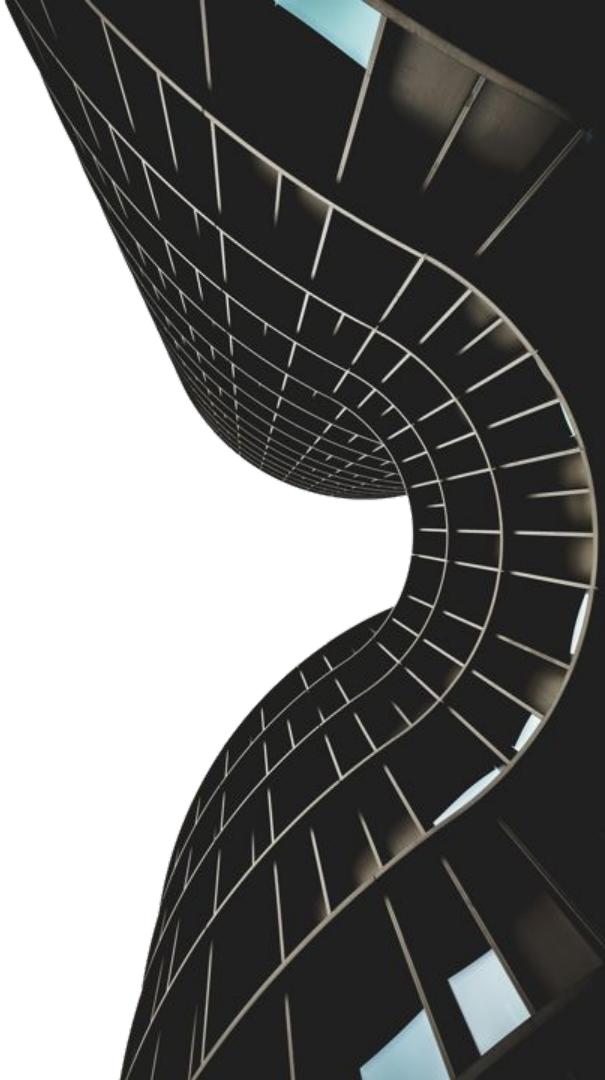
Related Work

Method: Fairband

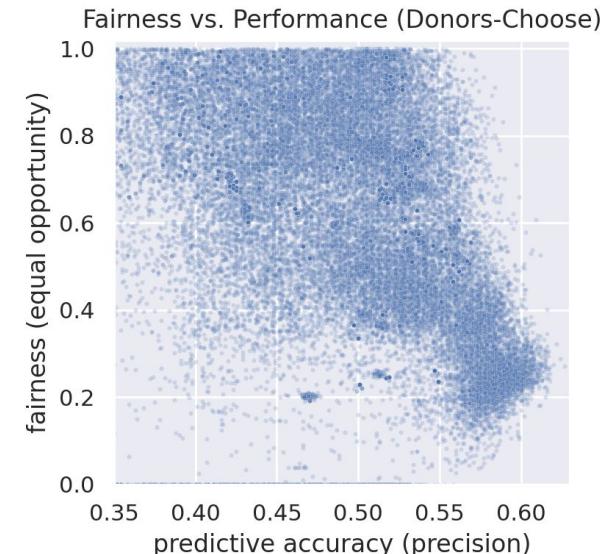
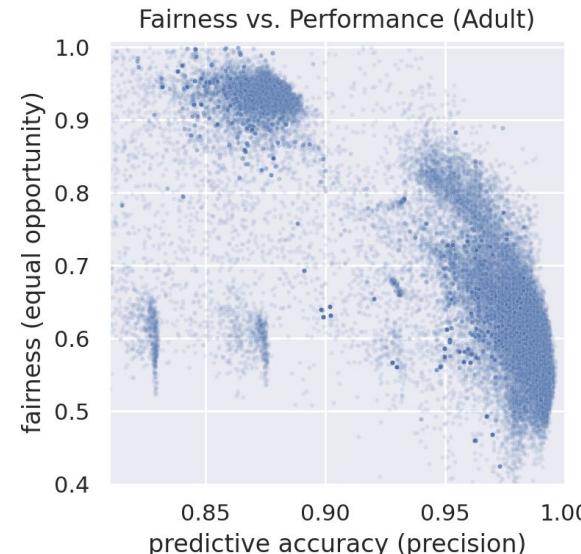
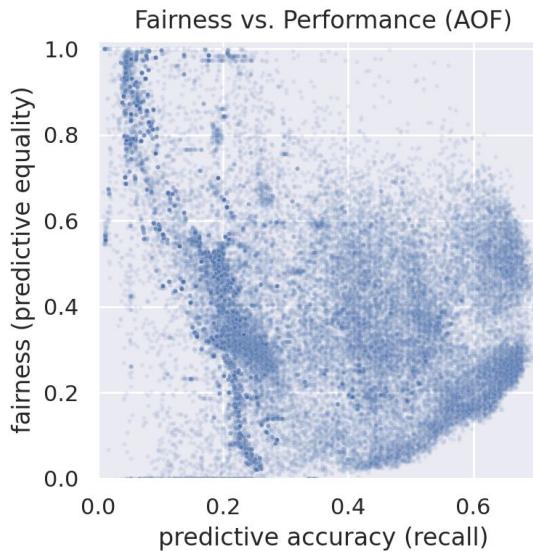
Experimental Setup

Results & Discussion

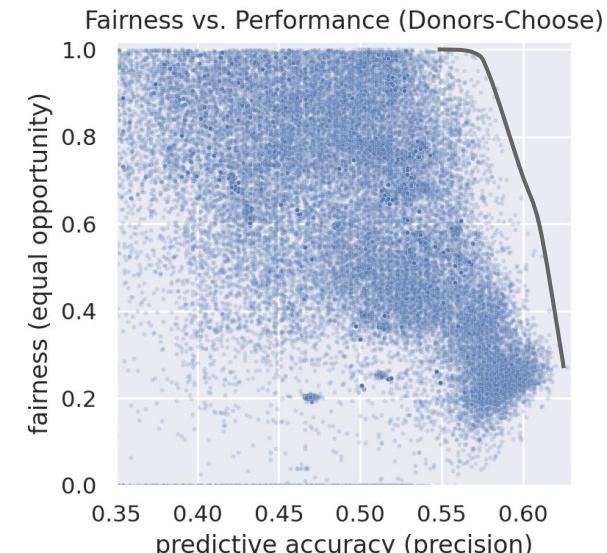
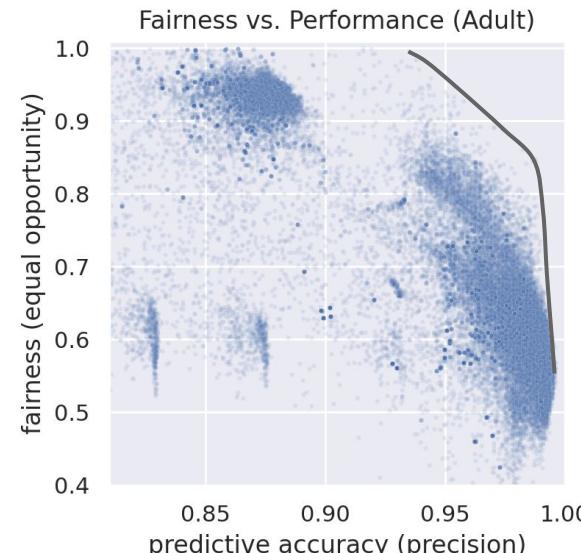
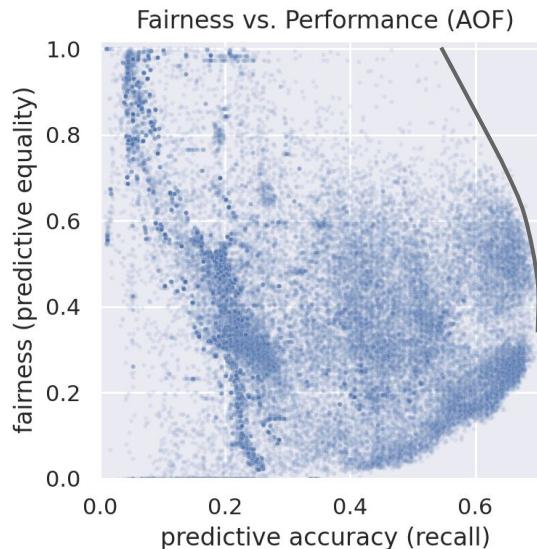
Conclusion



Fairness vs Accuracy



Fairness vs Accuracy



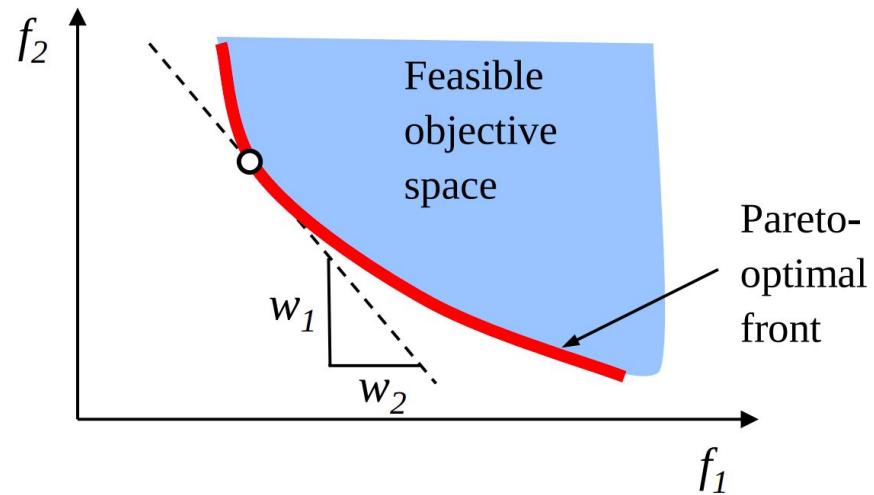
Scalarization

Convex case

Optimize as:

$$\arg \max_{\theta} w_1 \cdot f_1(\theta) + w_2 \cdot f_2(\theta)$$

- Can reach all optimal solutions by varying the weights.



* illustration for a minimization setting

Fairness as a goal

- Hypothesis:
 - If model a represents a better fairness-utility trade-off than model b with a short training budget, then this distinction is likely to be maintained with higher training budget.
- Thus: select models based on both fairness and performance.
- How: scalarize both objective functions.
 - Weighted by α parameter.

$$o = \alpha \cdot \text{performance} + (1 - \alpha) \cdot \text{fairness}$$

Fairband

Fairness-Aware Hyperparameter Optimization

- Built on top of resource-aware hyperparameter optimization methods.
 - Successive Halving and Hyperband.
- Model-agnostic, metric-agnostic, and efficient resource-usage.
- Highly exploratory.
 - Samples 6 times more configurations than RS or BO on equal budget.

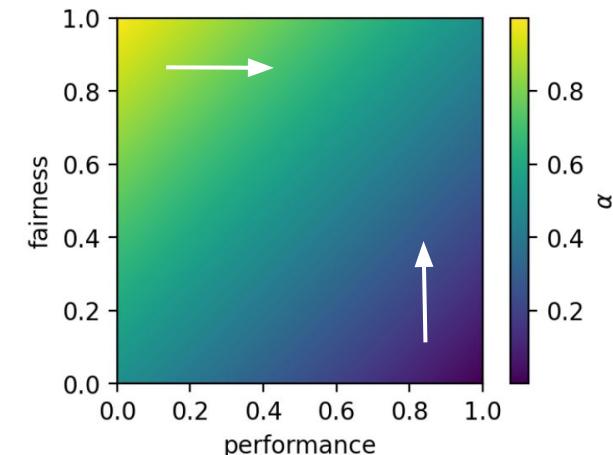


Dynamic α

- Selection of α requires specific task and domain knowledge.
- Guide α as search progresses based on metrics seen so far:

$$\alpha = 0.5 \cdot (\bar{f} - \bar{p}) + 0.5$$

- If models seen so far have high fairness but low performance: *promote performance*.
- Vice versa, promote fairness.



Fairband Algorithm

Algorithm 1

```

Input: maximum budget per configuration  $R$ ,
         $\eta$  (default  $\eta = 3$ ),
         $\alpha$  (default  $\alpha = \text{auto}$ )
1:  $s_{max} \leftarrow \lfloor \log_\eta(R) \rfloor$                                  $\triangleright$  define number of brackets
2:  $B \leftarrow (s_{max} + 1) \cdot R$                                       $\triangleright$  compute budget per bracket
3: for  $s \in \{s_{max}, s_{max} - 1, \dots, 0\}$  do            $\triangleright$  iterate through SH brackets, as per Li et al. (2016)
4:    $n \leftarrow \lceil \frac{B}{R} \cdot \frac{\eta^s}{s+1} \rceil$ ,  $r \leftarrow R \cdot \eta^{-s}$            $\triangleright$  choice of  $n$  versus  $B/n$  trade-off
5:    $T \leftarrow \text{get\_hyperparameter\_configurations}(n)$ 
6:   for  $i \in \{0, \dots, s\}$  do
7:      $n_i \leftarrow \lfloor n \cdot \eta^{-i} \rfloor$                                           $\triangleright$  run Successive Halving
8:      $r_i \leftarrow r \cdot \eta^i$                                                $\triangleright$  train  $n_i$  configurations
9:      $M \leftarrow \{\text{train\_with\_budget}(\lambda, r_i) : \lambda \in T\}$            $\triangleright$   $r_i$  training budget per config.
10:     $A \leftarrow \{\text{evaluate\_accuracy}(m_\lambda) : m_\lambda \in M\}$ 
11:     $F \leftarrow \{\text{evaluate\_fairness}(m_\lambda) : m_\lambda \in M\}$ 
12:    if  $\alpha = \text{auto}$  then                                          $\triangleright$  compute dynamic  $\alpha$  if applicable
13:       $\bar{f} \leftarrow \text{sum}(F)/|F|$                                           $\triangleright$  average fairness
14:       $\bar{a} \leftarrow \text{sum}(A)/|A|$                                           $\triangleright$  average accuracy
15:       $\alpha \leftarrow 0.5 \cdot (\bar{f} - \bar{a}) + 0.5$ 
16:     $O \leftarrow \{\alpha \cdot A[m_\lambda] + (1 - \alpha) \cdot F[m_\lambda] : m_\lambda \in M\}$   $\triangleright$  compute objective metric,  $o$ 
17:     $I \leftarrow \text{argsort}(O)$                                           $\triangleright$  sorted in descending order
18:     $k \leftarrow \lfloor n_i / \eta \rfloor$                                           $\triangleright$  number of configurations to keep
19:     $T \leftarrow T[I[0 : k]]$                                           $\triangleright$  select top  $k$  configurations
20: return  $\lambda^*$ , configuration with maximal intermediate goal seen so far

```

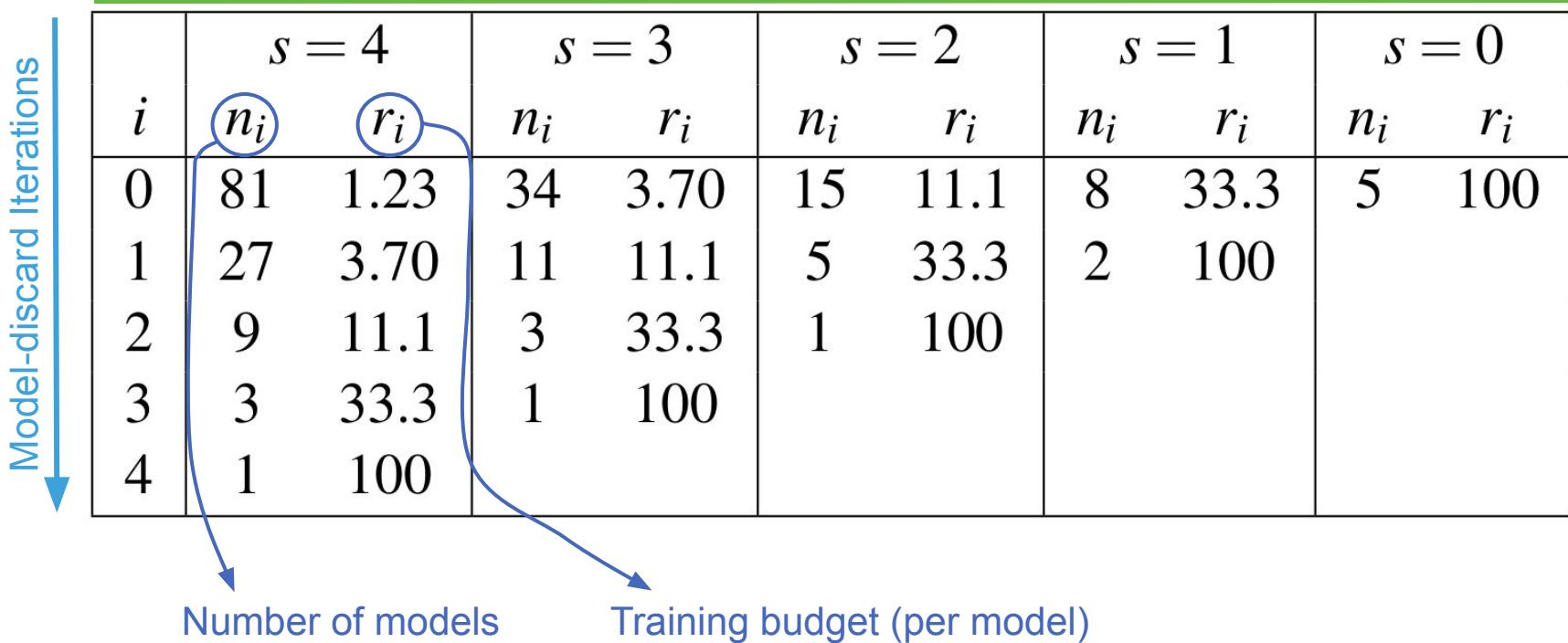
Hyperband

SH

Fairband

With $\eta = 3$ and $R = 100$

Optimization Brackets



Fairband

With $\eta = 3$ and $R = 100$

Select top-k based on α

i	$s = 4$		$s = 3$		$s = 2$		$s = 1$		$s = 0$	
	n_i	r_i								
0	81	1.23	34	3.70	15	11.1	8	33.3	5	100
1	27	3.70	11	11.1	5	33.3	2	100		
2	9	11.1	3	33.3	1	100				
3	3	33.3	1	100						
4	1	100								

121

46

21

10

5

= 206

Motivation

Objective

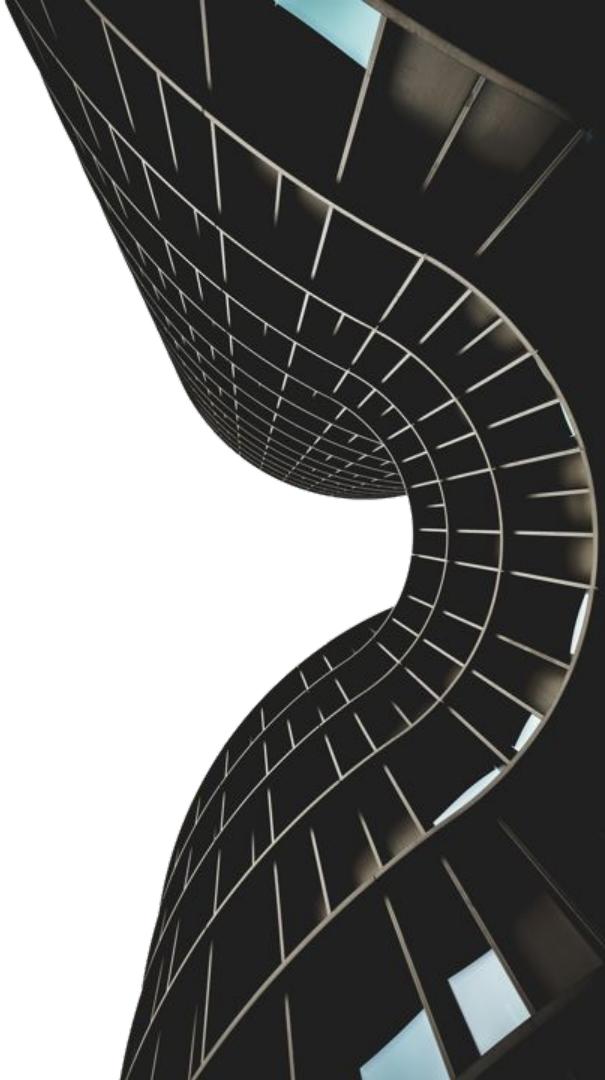
Related Work

Method: Fairband

Experimental Setup

Results & Discussion

Conclusion



Datasets

- Three benchmark datasets from the fairness literature
 - Donors Choose, Adult Census, COMPAS
- A large-scale real-world dataset on Account Opening Fraud
 - In-house *Feedzai* case study
 - Dubbed AOF

Donors Choose

- Identify projects (for K-12 schools) at risk of getting underfunded, to provide tailored interventions.
- Positive prediction: aid for that specific project.
- Setting: assistive.
- Fairness: balanced TPR.
- Sensitive attribute: poverty level.



Adult Census

- Data from the 1994 US Census (e.g., age, gender, ethnicity, occupation).
- Identify low-income individuals at need of assistance from a hypothetical social security program.
- Positive prediction: aid for that specific individual.
- Setting: assistive.
- Fairness: balanced TPR.
- Sensitive attribute: gender.



- Predict recidivism based on the person's criminal history.
- Positive prediction: high-risk score.
 - *E.g.*, defendant **not allowed to be released** on bail.
- Setting: punitive.
- Fairness: balanced FPR.
- Sensitive attribute: race.



AOF

Account Opening Fraud

- Online bank account opening fraud.
- Extremely imbalanced (~1% positive samples).
- Positive prediction: **not allowed** to open a bank account.
- Setting: punitive.
- Fairness: balanced FPR.
- Sensitive attribute: age.



Datasets

Dataset	Setting	Acc. Metric	Fairness Metric	Target Threshold	Sensitive Attribute
Donors Choose	assistive	precision	equal opportunity	1000 PP	poverty level
Adult	assistive	precision	equal opportunity	50% TPR	gender
COMPAS	punitive	precision	predictive equality	2% FPR	race
AOF	punitive	recall	predictive equality	5% FPR	age

Search Space

- Model type
 - LightGBM, Random Forest, Decision Tree, Logistic Regression, Neural Network
- Model hyperparameters
 - Number of estimators, maximum depth, number of neurons, number of layers, etc.
- Undersampling (only on AOF dataset)
 - 5%, 10%, or 20% positive samples

Motivation

Objective

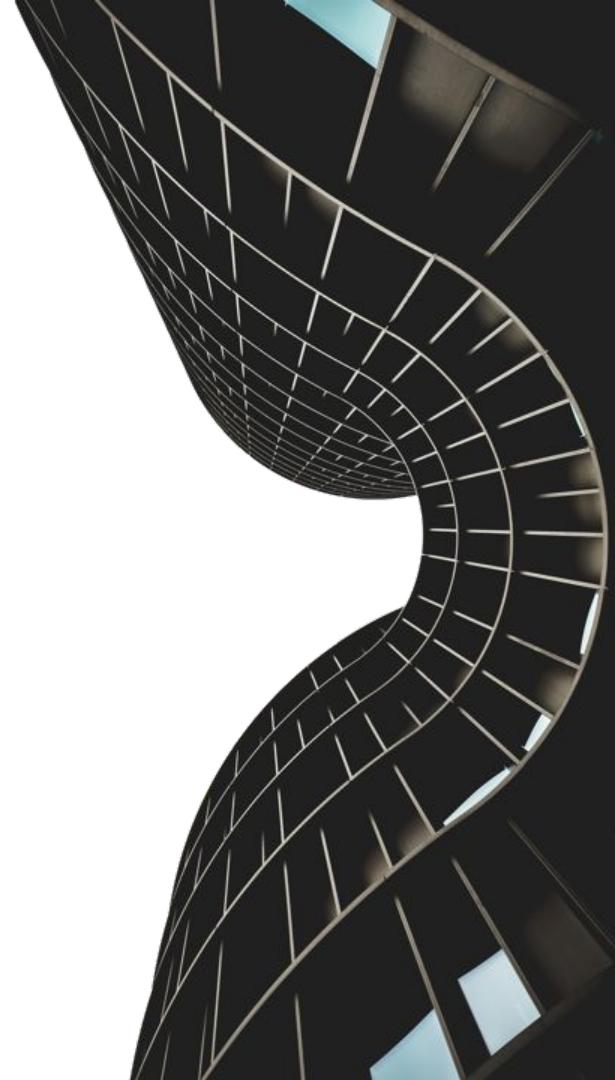
Related Work

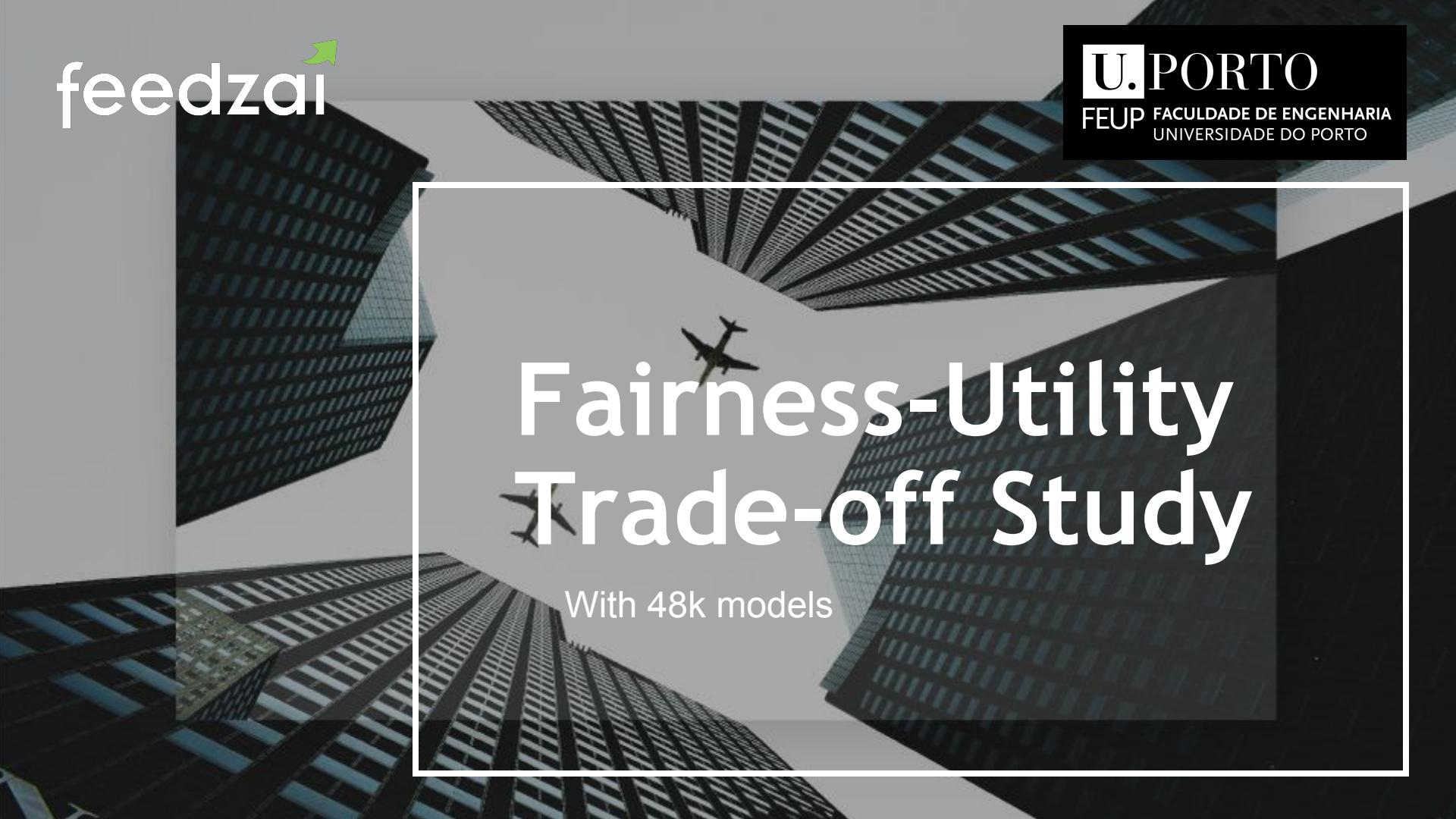
Method: Fairband

Experimental Setup

Results & Discussion

Conclusion

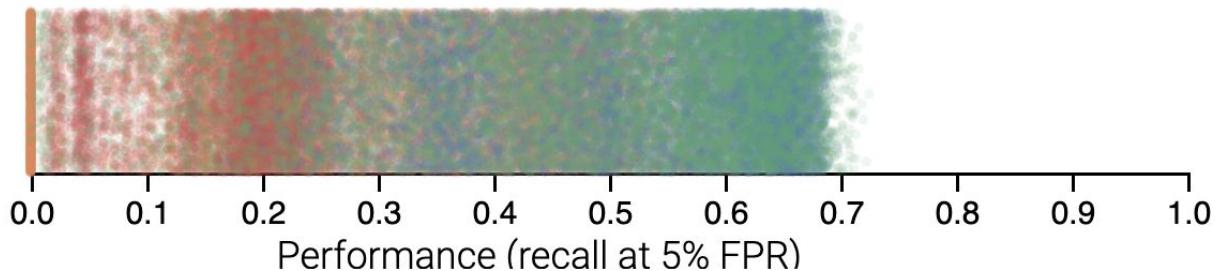


A black and white photograph of a city skyline with many skyscrapers. An airplane is flying across the sky between the buildings.

Fairness-Utility Trade-off Study

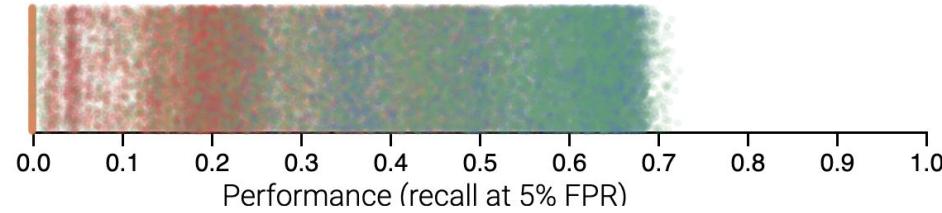
With 48k models

Model Performance

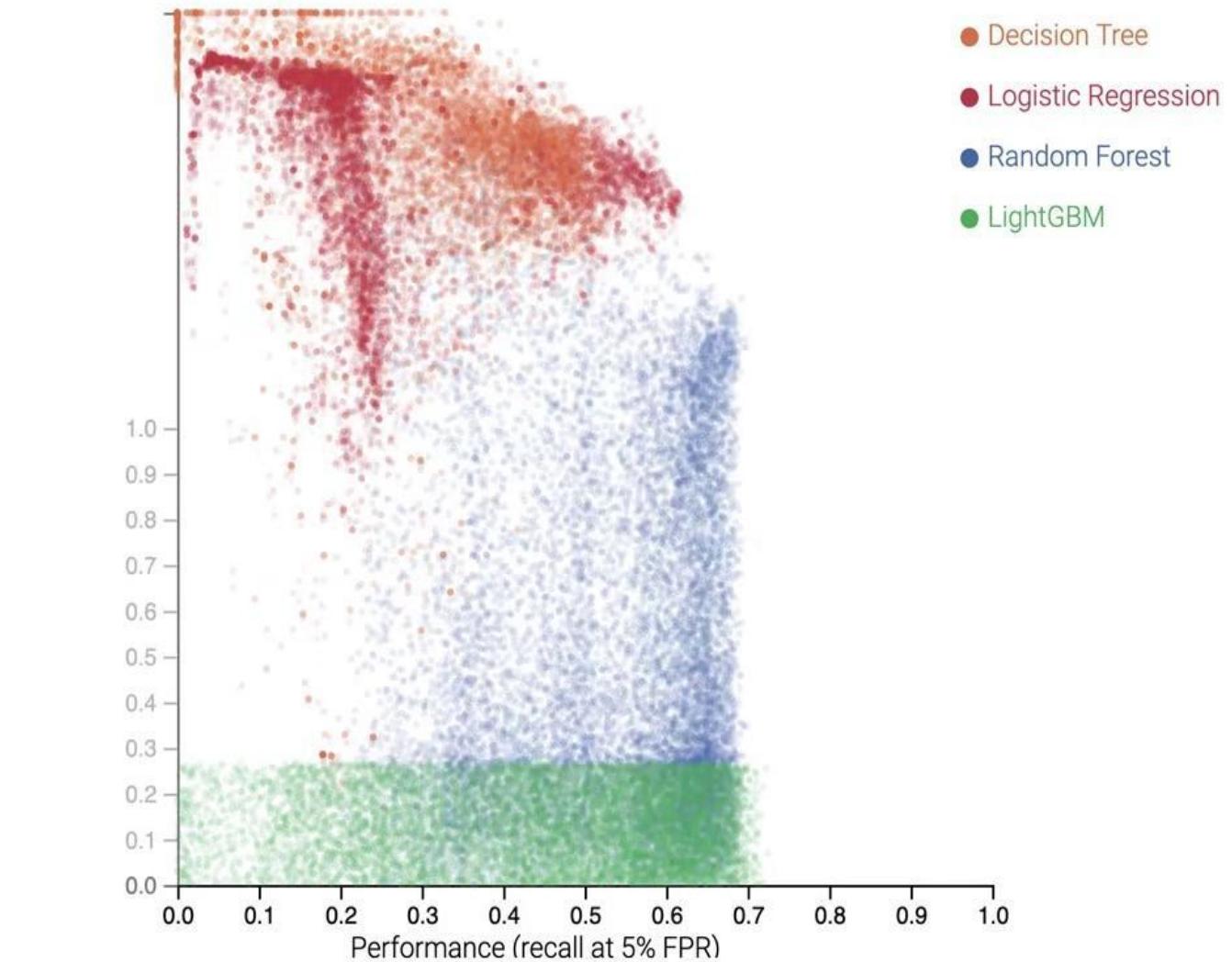


Model Performance

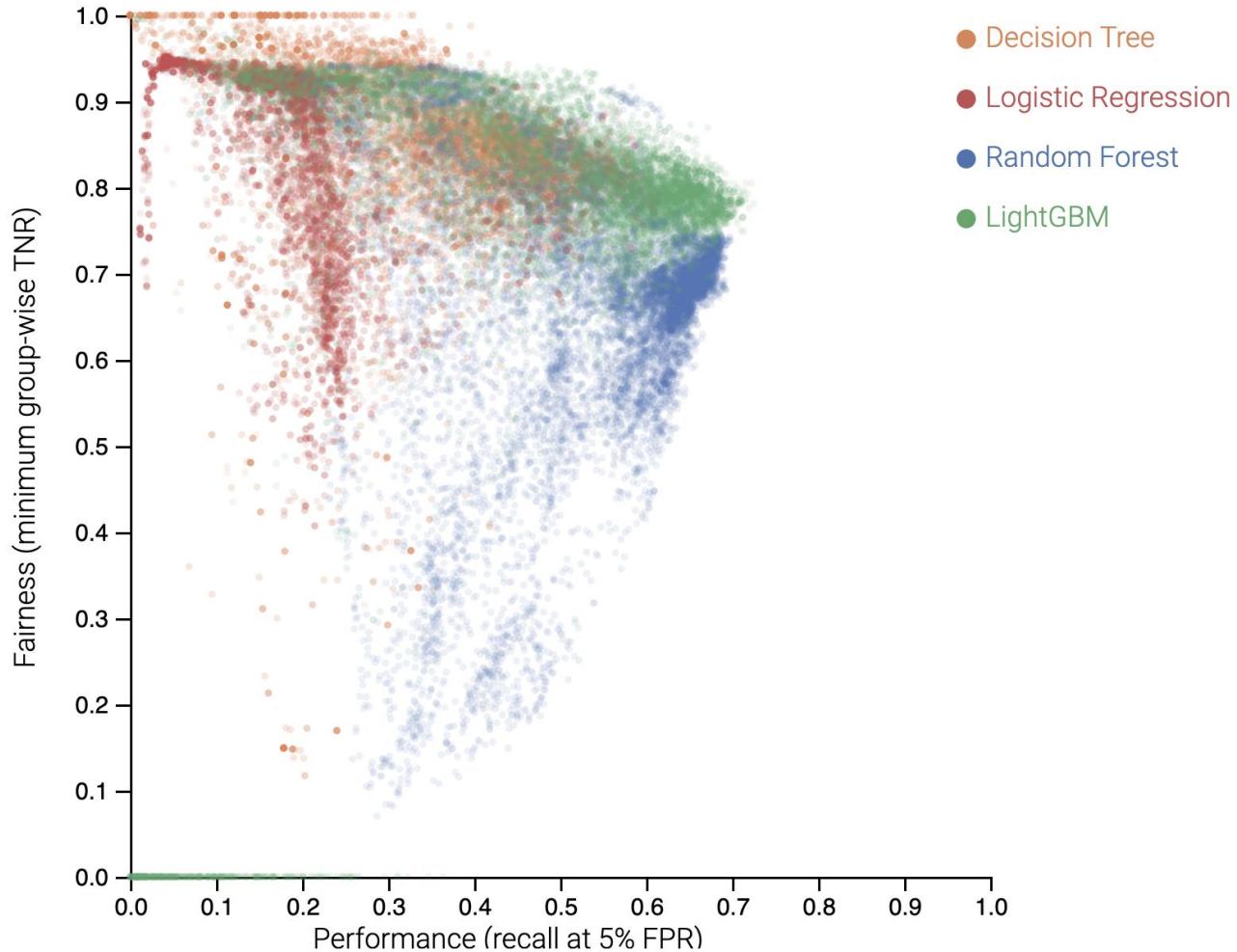
- Decision Tree
- Logistic Regression
- Random Forest
- LightGBM



Fairness vs. Performance



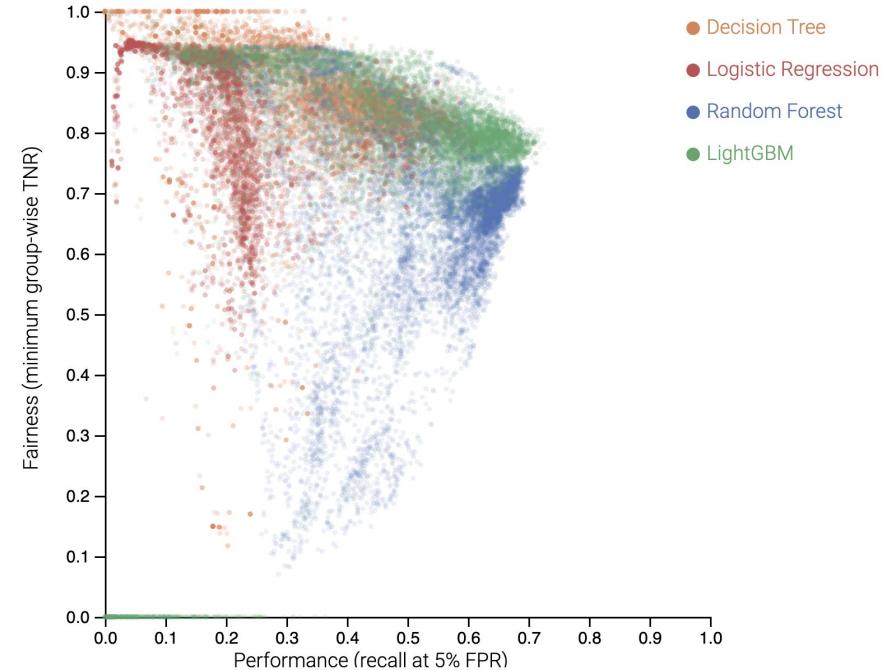
Fairness vs. Performance



Fairness-Utility Trade-off

It is possible to achieve substantially different fairness for the same performance.

Now it's just a question of guiding the search...



Fairband Results



Selected Models

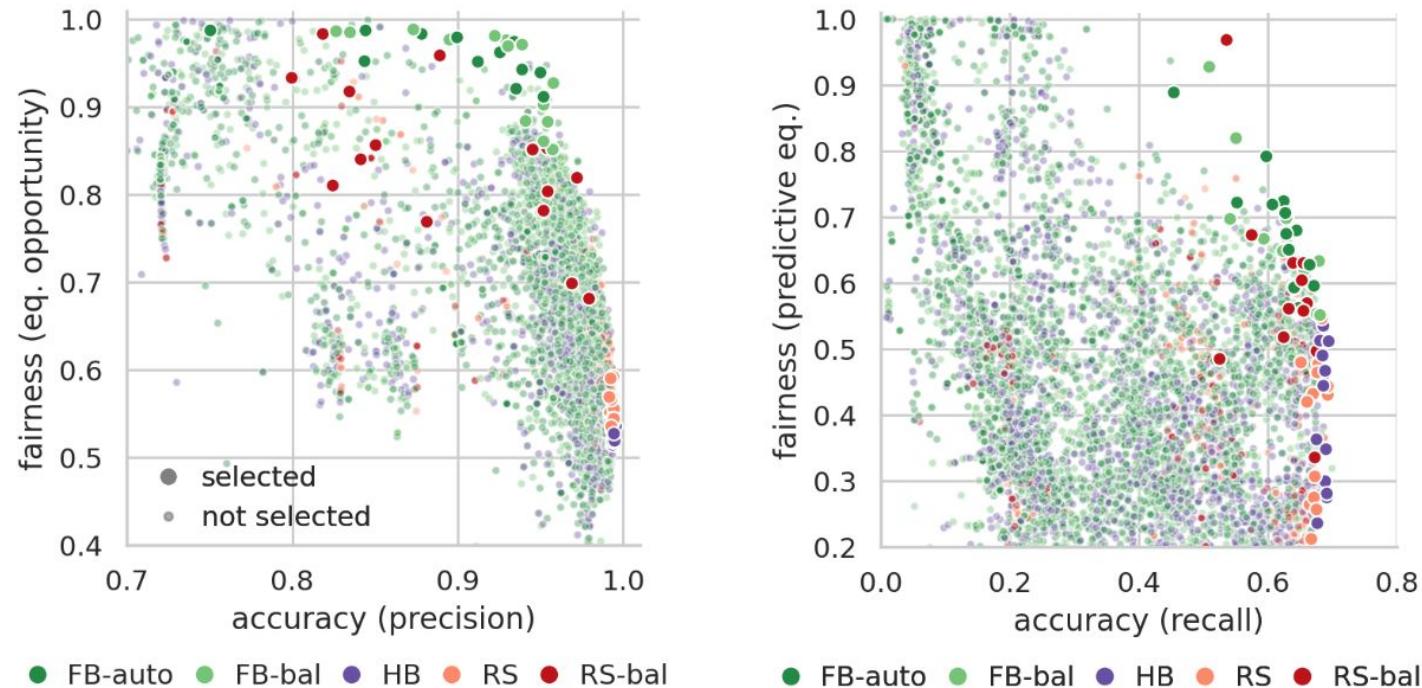


Figure 3: Fairness and predictive accuracy of selected models by hyperparameter optimization algorithm (Adult dataset on left plot, AOF on right plot).

Search Strategy

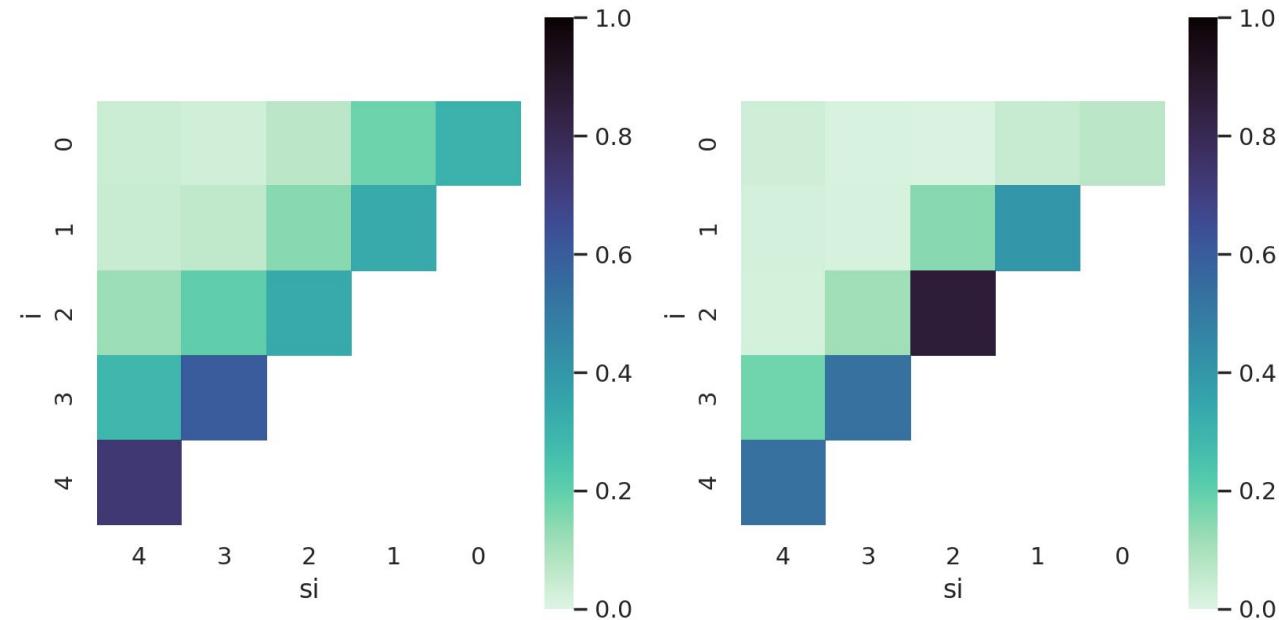


Figure 2: Average density of Pareto optimal models per FB-auto iteration (Adult dataset on left plot, AOF on right plot). Refer to Table 1 for information on the configurations at each iteration.

Fairband vs. State-of-the-Art

Averaged over 10 runs.

Statistical significance with
KS-test.

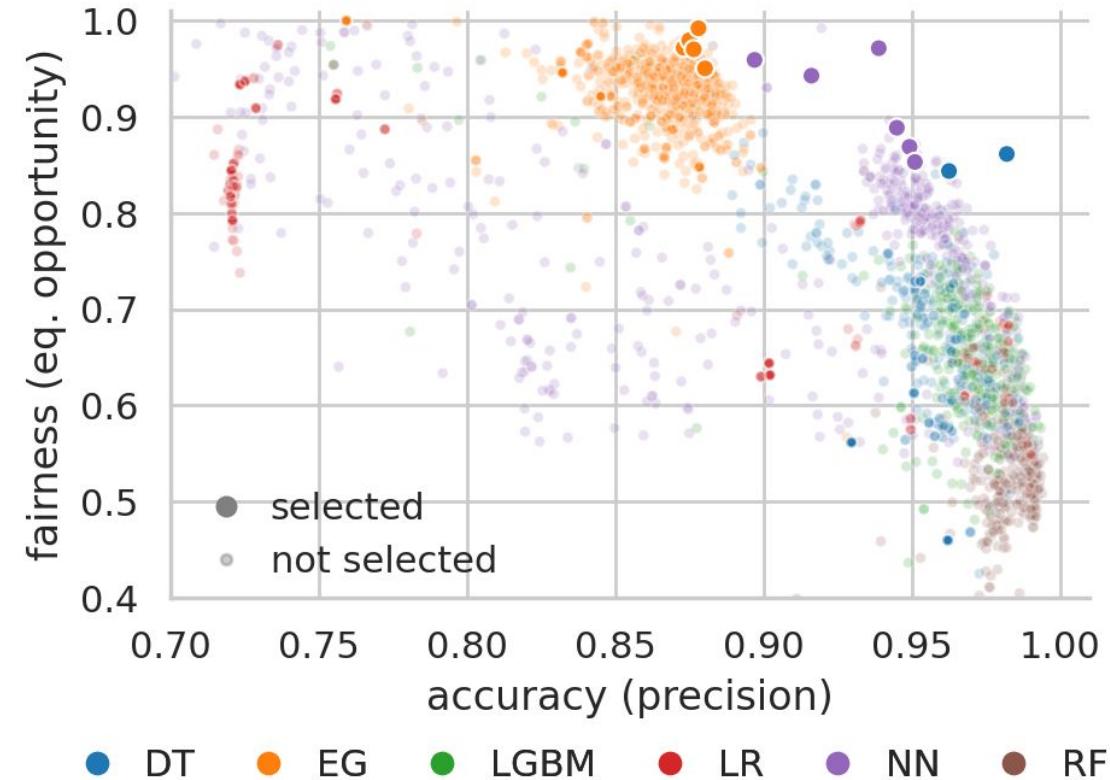
Algo.	Validation		Test	
	Predictive Acc.	Fairness	Predictive Acc.	Fairness
Donors Choose				
FB-auto	53.8 ^{▲♦}	97.9 ^{▲♦}	50.0 ^{▲♦}	87.3^{▲♦}
FB-bal	52.8 ^{▲♦}	98.6^{▲♦}	49.9 ^{▲♦}	86.1 ^{▲♦}
RS-bal	52.3 ^{▲♦}	95.2 ^{▲♦}	50.4 ^{▲♦}	82.2 ^{▲♦}
RS	59.9	26.8	53.1	32.8
HB	60.6	28.8	53.4	34.7
Adult				
FB-auto	90.6 ^{▲♦}	95.6^{▲♦}	90.1 ^{▲♦}	93.9^{▲♦}
FB-bal	91.9 ^{▲♦}	94.2 ^{▲♦}	79.7 ^{▲♦}	79.0 ^{▲♦}
RS-bal	89.8 ^{▲♦}	83.7 ^{▲♦}	90.5 ^{▲♦}	83.4 ^{▲♦}
RS	99.3	55.1	99.4	55.1
HB	99.4	53.2	99.4	53.3
COMPAS				
FB-auto	83.9 ^{▲◊}	97.0^{▲♦}	79.2 [▲]	41.6 [▲]
FB-bal	84.2 ^{▲◊}	96.0 ^{▲♦}	79.4 [▲]	42.7[▲]
RS-bal	80.8 ^{▲♦}	77.2 ^{▲♦}	74.8 ^{▲◊}	40.7 ^{▲♦}
RS	86.8	29.0	79.7	25.8
HB	88.6	17.7	82.6	24.1
AOF				
FB-auto	61.0 ^{▲♦}	69.7^{▲♦}	62.6 ^{▲♦}	74.9^{▲♦}
FB-bal	61.0 ^{▲♦}	68.8 ^{▲♦}	62.1 ^{▲♦}	72.9 ^{▲♦}
RS-bal	62.3 ^{▲♦}	61.5 ^{▲♦}	63.8 [▲]	66.6 ^{▲♦}
RS	67.3	37.8	67.4	40.4
HB	68.7	40.7	69.0	43.9

Fairband vs. Hyperband

Dataset	Abs. Difference (pp)		Rel. Difference (%)	
	Predictive Acc.	Fairness	Predictive Acc.	Fairness
Donors Choose	-3.4	+52.6	-6.4	+152
Adult	-9.3	+40.6	-9.4	+76.2
COMPAS	-3.4	+17.5	-4.1	+72.6
AOF	-6.4	+31.0	-9.3	+70.6
<i>Average</i>	<i>-5.6</i>	<i>+35.4</i>	<i>-7.3</i>	<i>+92.9</i>

Comparison of using Fairband (FB-auto) versus Hyperband (HB), on test results.

Optimizing bias
reduction
hyperparameters



Impact in Practice

Donors Choose

- Hyperband:
 - TPR low-poverty school: 6.25%
 - TPR high-poverty school: 1.79%
 - Global Precision: 60.6%
- Fairband:
 - TPR low-poverty school: 3.53%
 - TPR high-poverty school: 3.45%
 - Global Precision: 53.8%
- **93% higher TPR** for K12 schools with high poverty-level.

Impact in Practice

Adult

- Hyperband:
 - TPR males: 37.5%
 - TPR females: 69.8%
 - Global Precision: 99.4%
- Fairband:
 - TPR males: 49.2%
 - TPR females: 51.3%
 - Global Precision: 90.6%
- **32% higher TPR** for male low-income individuals.

Impact in Practice

COMPAS

- Hyperband:
 - FPR Caucasian: 0.8%
 - FPR Non-Caucasian: 2.7%
 - Global Precision: 88.6%
- Fairband:
 - FPR Caucasian: 1.9%
 - FPR Non-Caucasian: 1.9%
 - Global Precision: 83.9%
- **30% lower FPR** for non-caucasian individuals.

Impact in Practice

Bank Account Opening Fraud

- Hyperband:
 - FPR lower-aged applicants: 2.5%
 - FPR higher-aged applicants: 6.1%
 - Global Recall: 61.0%
- Fairband:
 - FPR lower-aged applicants: 3.8%
 - FPR higher-aged applicants: 5.4%
 - Global Recall: 68.7%
- **13% lower FPR** for higher-aged individuals.

Motivation

Objective

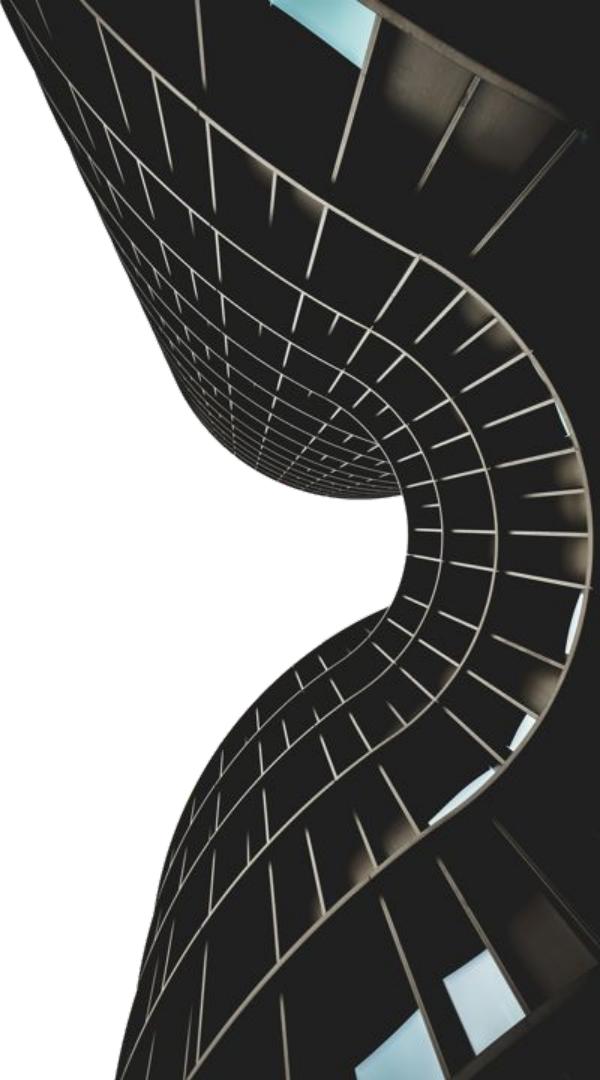
Related Work

Method: Fairband

Experimental Setup

Results & Discussion

Conclusion



Contributions

1. By optimizing solely for performance you are unknowingly targeting unfair models;
2. Fairband, an efficient method for multi-objective hyperparameter optimization.
 - o Model-agnostic, metric-agnostic, and maximizes fairness among multiple sub-groups.
3. A dynamic method of assigning the fairness-utility weight, α .
 - o Parameter-free Fairband.
4. State-of-the-art results: significantly improved fairness at a small performance cost, and no extra budget.
5. Hyperparameter optimization is an effective way to navigate the fairness-utility trade-off.

Thank You

Fairness-Aware Hyperparameter Optimization

André Cruz

Pedro Saleiro
Catarina Belém
Carlos Soares
Pedro Bizarro

Preprint at

<https://arxiv.org/abs/2010.03665>

Plots and data at

<https://github.com/feedzai/fair-automl>

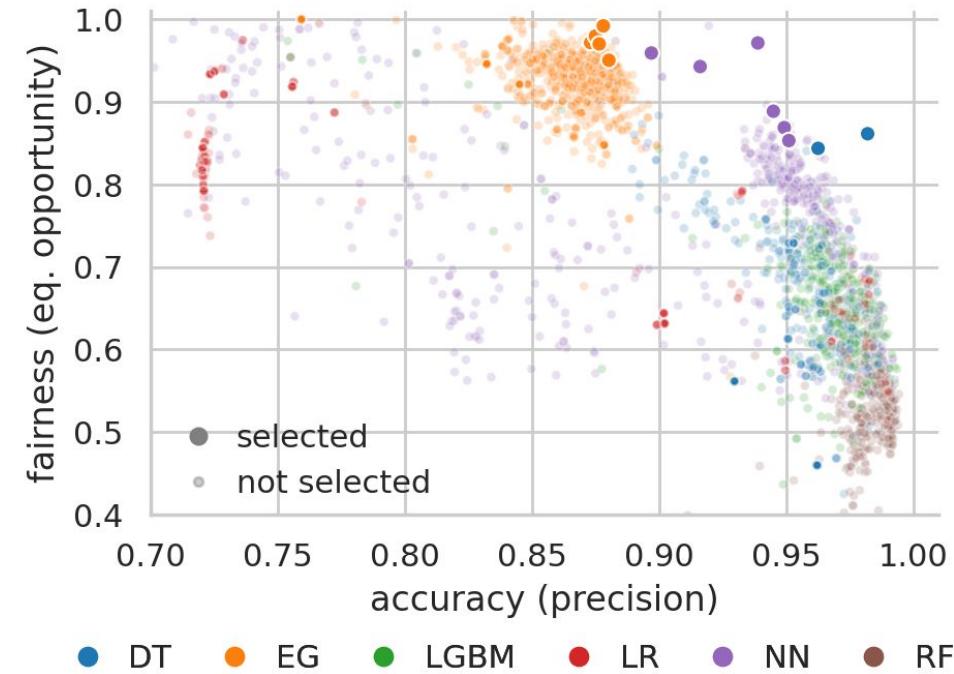
MSc Thesis at

<https://repositorio-aberto.up.pt/bitstream/10216/128959/2/414778.pdf>

Fairness-Aware Hyperparameter Optimization

André Cruz

Supervision:
Carlos Soares, FEUP.
Pedro Saleiro, Feedzai.



Annex

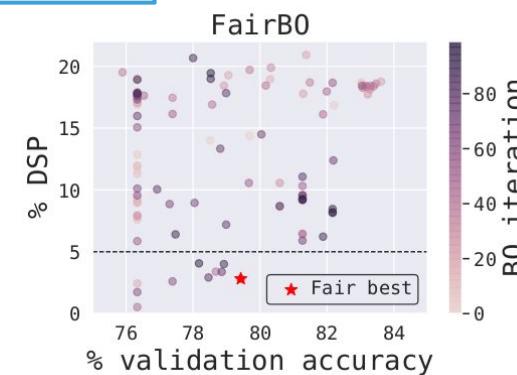
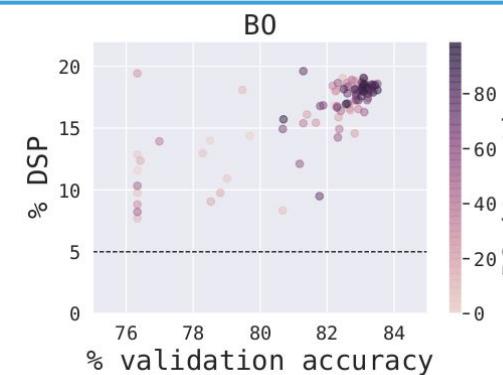
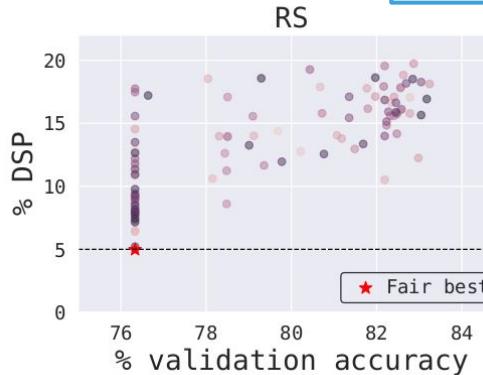
A. Related Work

Fair Bayesian Optimization

Perrone et al., arXiv, 9th June 2020

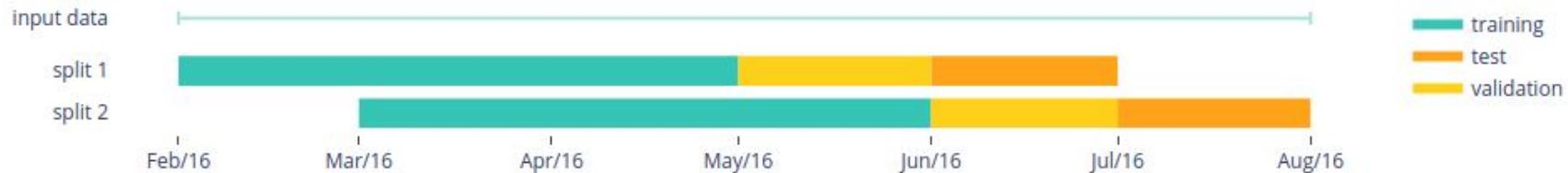
- Bayesian optimization approach
 - Blind to resource usage
 - Constraint may not be possible to fulfill
 - Tested on toy datasets, with non-standard metrics

$$cEI(\mathbf{x}) = \underline{EI(\mathbf{x})} P(c(\mathbf{x}) \leq \epsilon)$$



B. Experimental Setup

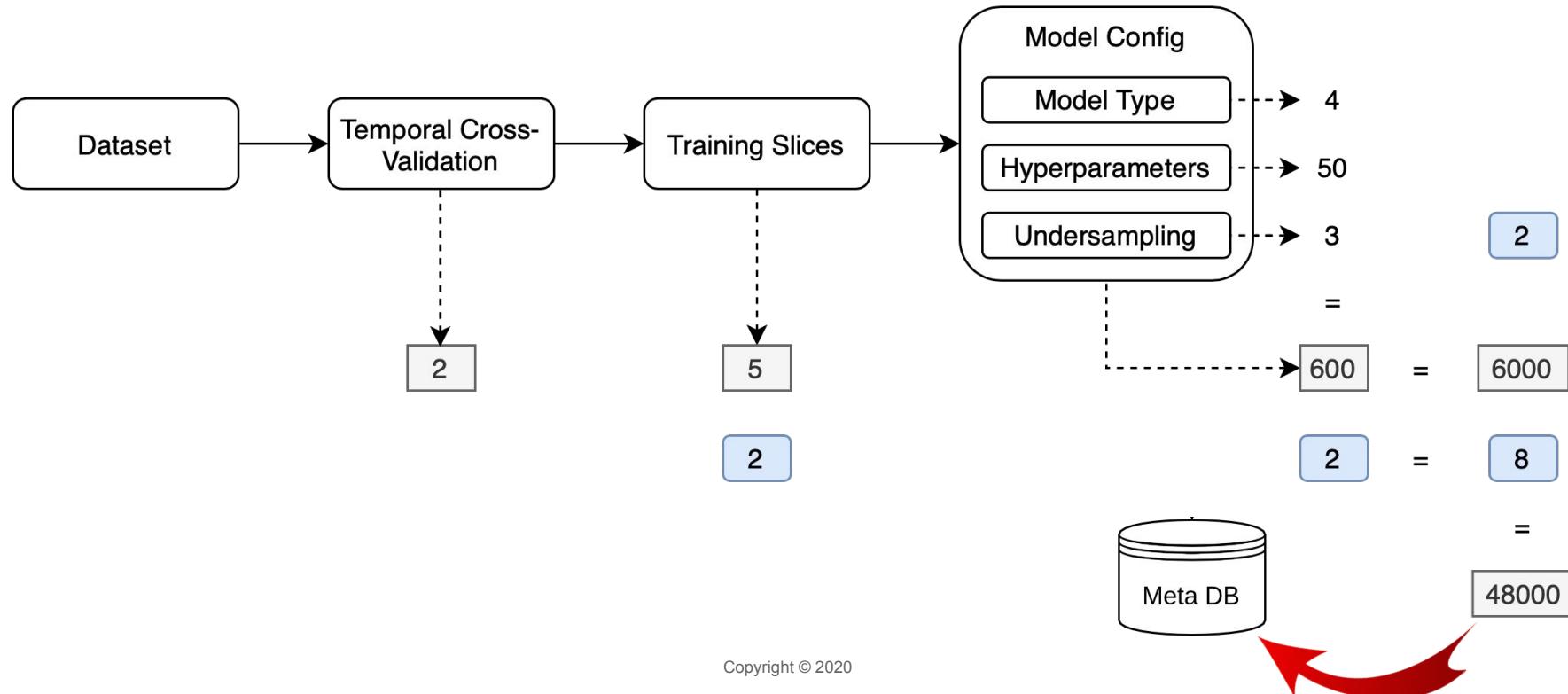
Temporal Cross-Validation



Data



Evaluation Grid



Evaluation Grid

