

Improving fairness in predictions and decision making

Christophe Giraud

Laboratoire de Mathématiques d'Orsay,
Université Paris Saclay

IHP, Paris 2025

Content of the talk

Objectives

- to highlight fairness issues in data science
- to overview some statistical approaches
- to provide examples of statistical contributions

Plan of the talk

- A (biased!) introduction to algorithmic fairness
- Glimpse at two contributions in online learning

Disclaimer: I do not consider myself as an expert of this topic

Fairness in Machine Learning: a major societal concern

Machine Learning is ubiquitous in daily life



PRODUCTS ▾

CUSTOMERS ▾

PRICING

RESOURCES ▾

REQUEST A DEMO

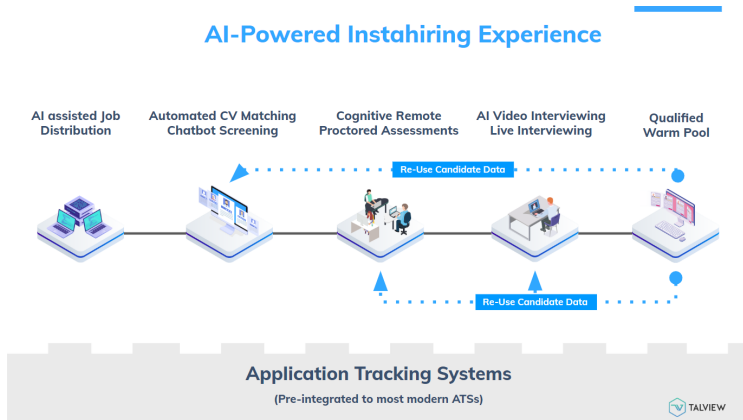


Talent Assessment | 16 Min Read

How AI-based HR Chatbots are Simplifying Pre-screening

Fairness in Machine Learning: a major societal concern

Machine Learning is ubiquitous in daily life



Fairness in Machine Learning: a major societal concern

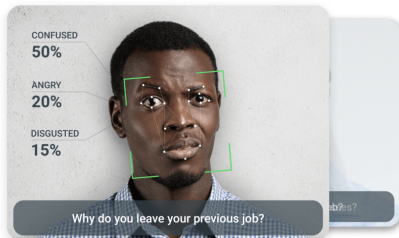
Machine Learning is ubiquitous in daily life

Emotion Analysis



LATEST AI/ML EMOTION RECOGNITION
TECHNOLOGY FOR VIDEO

- ✓ Get more data from recorded interviews
- ✓ Reduce personal bias
- ✓ Make data-driven hiring decisions
- ✓ Speed up recruiting process



Source [easyhire.me](https://www.easyhire.me)

Fairness in Machine Learning: a major societal concern

Machine Learning is ubiquitous in daily life

05-17-19

Schools are using software to help pick who gets in. What could go wrong?

Admissions officers are increasingly turning to automation and AI with the hope of streamlining the application process and leveling the playing field.

Fairness in Machine Learning: a major societal concern

Machine Learning is ubiquitous in daily life

SCIENCE ADVANCES | RESEARCH ARTICLE

RESEARCH METHODS

The accuracy, fairness, and limits of predicting recidivism

Julia Dressel and Hany Farid*

Algorithms for predicting recidivism are commonly used to assess a criminal defendant's likelihood of committing a crime. These predictions are used in pretrial, parole, and sentencing decisions. Proponents of these systems argue that big data and advanced machine learning make these analyses more accurate and less biased than humans. We show, however, that the widely used commercial risk assessment software COMPAS is no more accurate or fair than predictions made by people with little or no criminal justice expertise. In addition, despite COMPAS's collection of 137 features, the same accuracy can be achieved with a simple linear predictor with only two features.

Copyright © 2018
The Authors, some
rights reserved;
exclusive licensee
American Association
for the Advancement
of Science. No claim to
original U.S. Government
Works. Distributed
under a Creative
Commons Attribution
NonCommercial
License 4.0 (CC BY-NC).

Fairness in Machine Learning: a major societal concern

Machine Learning is ubiquitous in daily life, and it is used for sensitive decisions such as

- admission in university,
- bank loan,
- job recruitment,
- justice decision,
- medical diagnosis,
- ...

Promises of ML in decision-making

ML can improve decision-making

ML can be more accurate, more objective and more fair than humans, since algorithms can

- incorporate more data, and more factors in a complex analysis,
- and are not subject to personal biases, tiredness, emotional factors, etc

5 Benefits of Recruiting Automation



Boosts recruiter efficiency



Enhances hiring team communication



Improves candidate experiences



Fast-forwards candidate screening



Nurtures Candidate Engagement

Actuality of ML decision-making

Discriminations also happen in ML prediction

Many ML systems have been shown to produce unfair outcomes.

Some famous past examples:

- **Hiring AI** from Amazon was discriminating against female candidate on some jobs
- **Google Ad** was proposing higher-paying executive jobs more likely to men than women
- **COMPAS** was falsely predicting recidivism twice more likely for African-American than for Caucasian-American.

COMPAS recidivism algorithm in action

Source: ProPublica

<p>VERNON PRATER</p> <p>Prior Offenses 2 armed robberies, 1 attempted armed robbery</p> <p>Subsequent Offenses 1 grand theft</p> <p>LOW RISK 3</p>	<p>BRISHA BORDEN</p> <p>Prior Offenses 4 juvenile misdemeanors</p> <p>Subsequent Offenses None</p> <p>HIGH RISK 8</p>
--	---

<p>DYLAN FUGETT</p> <p>LOW RISK 3</p>	<p>BERNARD PARKER</p> <p>HIGH RISK 10</p>
---	---

<p>JAMES RIVELLI</p> <p>LOW RISK 3</p>	<p>ROBERT CANNON</p> <p>MEDIUM RISK 6</p>
--	---

<p>JAMES RIVELLI</p> <p>Prior Offenses 1 domestic violence aggravated assault, 1 grand theft, 1 petty theft, 1 drug trafficking</p> <p>Subsequent Offenses 1 grand theft</p> <p>LOW RISK 3</p>	<p>ROBERT CANNON</p> <p>Prior Offense 1 petty theft</p> <p>Subsequent Offenses None</p> <p>MEDIUM RISK 6</p>
--	--

Where does the unfairness come from?

Main potential causes of unfairness in data science

- [intentional discrimination]
- **historical biases in learning datasets**
- inadvertent bias in evaluations (biased proxy)
- inadvertent bias from data sampling: learning dataset not representative of the target population
- **inadvertent bias from algorithm objectives: focus on the benefit for majority group**

Fairness in Machine Learning: a major societal concern

Societal concern

- Standard use of ML can lead to unfair and discriminating decisions,
 - Machine Learning is ubiquitous in many sensitive decisions: bank loan, admission to university, job recruitment, crime recidivism prediction, etc
-
- Fairness in decision-making is an important topic;
 - The statistical community has an important role to play for providing
 - ▶ conceptual ideas
 - ▶ competitive algorithms with provable performances
 - ▶ theoretical insights
 - ▶ **education of the next generation of data scientists**
 - In collaboration with experts from human science and policy makers.

EU regulation for AI

Regulation on AI now include some fairness requirement



EUROPEAN COMMISSION

Brussels, 21.4.2021

COM(2021) 206 final

2021/0106(COD)

Proposal for a

REGULATION OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

LAYING DOWN HARMONISED RULES ON ARTIFICIAL INTELLIGENCE (ARTIFICIAL INTELLIGENCE ACT) AND AMENDING CERTAIN UNION LEGISLATIVE ACTS

EU regulation for AI

Regulation on AI now include some fairness requirement

PROHIBITED ARTIFICIAL INTELLIGENCE PRACTICES

Article 5

1. **The following artificial intelligence practices shall be prohibited:**
 - (a) the placing on the market, putting into service or use of an AI system that deploys subliminal techniques beyond a person's consciousness in order to materially distort a person's behaviour in a manner that causes or is likely to cause that person or another person physical or psychological harm;
 - (b) the placing on the market, putting into service or use of an AI system that exploits any of the vulnerabilities of a specific group of persons due to their age, physical or mental disability, in order to materially distort the behaviour of a person pertaining to that group in a manner that causes or is likely to cause that person or another person physical or psychological harm;
 - (c) the placing on the market, putting into service or use of AI systems by public authorities or on their behalf for the evaluation or classification of the trustworthiness of natural persons over a certain period of time based on their social behaviour or known or predicted personal or personality characteristics, with the social score leading to either or both of the following:
 - (i) detrimental or unfavourable treatment of certain natural persons or whole groups thereof in social contexts which are unrelated to the contexts in which the data was originally generated or collected;
 - (ii) detrimental or unfavourable treatment of certain natural persons or whole groups thereof that is unjustified or disproportionate to their social behaviour or its gravity;

EU regulation for AI

Regulation on AI now include some fairness requirement

→ Jean-Michel Loubes' talk

French anti-discrimination law (LOI no 2008-496)

Direct discrimination

Constitue une **discrimination directe** la situation dans laquelle, **sur le fondement** de son origine, de son sexe, de sa situation de famille, [...] une prétendue race ou une religion déterminée, **une personne est traitée de manière moins favorable qu'une autre ne l'est, ne l'a été ou ne l'aura été dans une situation comparable.**

Indirect discrimination

Constitue une **discrimination indirecte** une disposition, **un critère ou une pratique neutre en apparence**, mais susceptible d'**entraîner**, pour l'un des motifs mentionnés au premier alinéa, **un désavantage particulier** pour des personnes par rapport à d'autres personnes, **à moins** que cette disposition, ce critère ou cette pratique **ne soit objectivement justifiée** par un but légitime [...]

1- A (small) tour in algorithmic fairness

Algorithmic fairness

3 main statistical points of view for improving fairness

- 1 **Individual fairness** aims to treat similar people similarly (individual notions);
- 2 **Group fairness** seeks to comply to fairness criteria at the sub-population level (statistical notions);
- 3 **Causal fairness** tries to identify causes of unfairness in order to act on them (causal notions).

Learning framework

Notation

- Outcome $Y \in \mathcal{Y}$
- Covariate $X \in \mathcal{X}$
- Sensitive attribute $S \in \mathcal{S}$ (observed or not)
- Predictor: $f : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{Y}$ (possibly $f : \mathcal{X} \rightarrow \mathcal{Y}$)
- Prediction: $F = f(X, S)$

Individual fairness

Individual fairness: Lipschitz mapping

Lipschitz mapping: *treating similar people similarly*

For f a randomized predictor and two distances d on \mathcal{X} and D on the probability distribution on \mathcal{Y}

$$D(\text{law}(f(x)), \text{law}(f(x'))) \leq d(x, x').$$

This definition encodes the notion *treating similar people similarly*.

Caveat: the design of the distance d is critical and can lead to unfair decisions.

Ex: a job may require self-confidence and hard-concentration skills. Your hiring system may be Lipschitz, but strongly favor self-confidence compared to hard-concentration skills, leading to discrimination.

Group fairness

Group fairness

General principle

To comply to some fairness criteria at the sub-population level (statistical notions)

Ex: we want to treat equally women and men.

Caveat: group fairness focuses at the group level, so (alone) it does not enforce fairness at the individual level

Group fairness: no Disparate Treatment

no Disparate Treatment

The predictor f complies to no disparate treatment, if it does not use the sensitive attribute S

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

It protects against *pure intentional discrimination* (negative or positive!)

Caveats:

- 1 when X and S are correlated, it does not protect against unfairness or discrimination
- 2 difficult to protect against indirect discrimination when S is ignored

Group fairness: meritocratic fairness

Equalized Odds

$$F \perp\!\!\!\perp S \mid Y$$

Ex: (binary classification)

$$\mathbb{P}[F = 1 | S = 1, Y] \cong \mathbb{P}[F = 1 | S = 0, Y]$$

- Equalized Odds encodes a notion related to meritocracy
- There are many variants

Caveat: strongly subject to biases in learning datasets

Group fairness: Demographic parity

Demographic parity

$$F \perp\!\!\!\perp S$$

Ex: (binary classification)

$$\mathbb{P}[F = 1|S = 1] \cong \mathbb{P}[F = 1|S = 0]$$

Demographic parity promotes *diversity* and can be related to affirmative action policies.

Caveat: the outcome is not taken into account

A quick remark: fairness is only part of the game

Demographic Parity (DP) alone is meaningless: it is very simple to comply to DP, just provide a prediction F at random, independent of (X, S) . 😞

⇒ DP must be coupled with **risk minimisation**

Ex: (binary classification) an ideal classifier should be

$$F_{DP}^* \in \underset{F: \mathbb{P}[F=1|S=1]=\mathbb{P}[F=1|S=0]}{\operatorname{argmin}} \mathbb{P}[F \neq Y]$$

Remark (continued)

Ex: Assume you want to hire 3% of the students studying maths. You want to hire the **best ones**, while complying to **DP** in terms of gender.

Then, you will hire the top 3% female students and the top 3% male students.

⇒ Risk requirement enforces some individual fairness: we observe that your recruitment is not only fair between groups (in terms of DP), but also within a group (as far it is homogeneous), since the best students of each group are hired.

Remark: it sounds like **DP+risk minimisation** may have something to do with quantile adjustment...

→ Christophe Denis' talk

Demographic Parity (DP) and no Disparate Treatment (nDT)

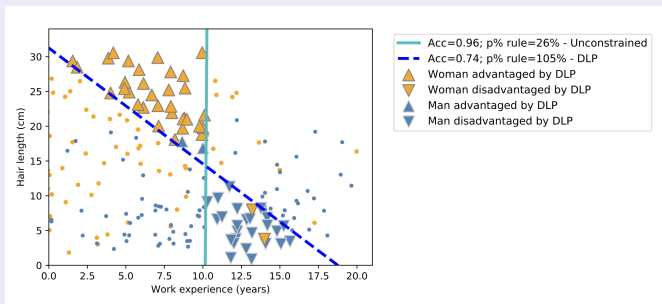
Limitation to previous example: Disparate Treatment is prohibited by many regulations.

Demographic Parity (DP) with no Disparate Treatment (nDT)?

- It is possible to enforce DP together with nDT, by using (more or less implicitly) proxies for estimating S .
- Hence DP can be enforced even when DT is prohibited.

Is it desirable to combine DP with nDT?

Undesirable side effects on a toy example of hiring data



source: Z.C. Lipton, A. Chouldechova, J. McAuley, NeurIPS 2018

- An unconstrained classifier (vertical line) hires candidates based on work experience, yielding higher hiring rates for men than for women.
- A nDT classifier (dashed diagonal) achieves DP by differentiating based on an irrelevant attribute (hair length).

The nDT hurts some short-haired women, flipping their decisions to reject, and helps some long-haired men. 😞

Current regulation can hurt fairness

→ nDT can be detrimental when trying to enforce some group-fairness properties like DP: re-ordering within groups can occur 😞

Yet, Disparate Treatment is prohibited by many regulations preventing from a safe application of DP (or other criteria)

Group fairness: Equal Opportunity

Equal Opportunity

Denote by \mathcal{Y}_+ the set of positive outcomes

$$F \perp\!\!\!\perp S \mid Y \in \mathcal{Y}_+$$

Ex: (binary classification) with $\mathcal{Y}_+ = \{1\}$

$$\mathbb{P}[F = 1|S = 1, Y = 1] \cong \mathbb{P}[F = 1|S = 0, Y = 1]$$

Equal Opportunity requires *equal True Positive rates* across groups: *successful people* should be given the *same chance* in all groups.

Group fairness: Predictive parity

Predictive parity (test fairness)

$$Y \in \mathcal{Y}_+ \perp\!\!\!\perp S \mid F \in \mathcal{Y}_+$$

Ex: (binary classification) with $\mathcal{Y}_+ = \{1\}$

$$\mathbb{P}[Y = 1 | S = 1, F = 1] \cong \mathbb{P}[Y = 1 | S = 0, F = 1]$$

Predictive parity asks for *equal fraction of correct positive prediction* across groups.

Somewhat related to group-wise calibration (below)

This criterion can be evaluated, even in partial monitoring scheme where we observe the outcome Y only when $F = 1$.

Performance fairness: Group-wise calibration

Group-wise calibration

$$\mathbb{E}[Y|S, F] \cong F$$

Ex: (binary classification) for a score $F \in [0, 1]$

$$\mathbb{P}[Y = 1|S = 1, F] \cong \mathbb{P}[Y = 1|S = 0, F] \cong F$$

The prediction are *calibrated for each group*.

Performance fairness: Equal group-wise risk

Equal group-wise risk

For a loss function ℓ

$$\mathbb{E}[\ell(Y, F)|S] \cong \mathbb{E}[\ell(Y, F)]$$

Equal risk for each group.

Performance fairness: Group-wise no regret

Group-wise no regret

$$\max_{s \in \mathcal{S}} \left\{ \mathbb{E} [\ell(Y, F) | S = s] - \min_{f_s} \mathbb{E} [\ell(Y, f_s(X)) | S = s] \right\} = o(1)$$

Each group enjoys a *no regret* prediction.

And many other criteria...

A large zoology

Demographic parity	$F \perp\!\!\!\perp S$
Equalized odds	$F \perp\!\!\!\perp S Y$
Equal opportunity	$F \perp\!\!\!\perp S Y \in \mathcal{Y}_+$
Predictive parity	$Y \in \mathcal{Y}_+ \perp\!\!\!\perp S F \in \mathcal{Y}_+$
Group-wise calibration	$\mathbb{E}[Y S, F] \cong F$
Equal group-wise risk	$\mathbb{E}[\ell(Y, F) S] \cong \mathbb{E}[\ell(Y, F)]$
...	...

with some incompatible notions!

The famous COMPAS case

The Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) is a software which aims to predict recidivism risk.

ProPublica compared COMPAS predictions across ethnicity groups in the USA. It exhibits a large violation of the Equalized Odds criteria.

The COMPAS developers argue yet that COMPAS (almost) complies to Predictive parity.

Chouldechova (2017) and Kleinberg *et al.* (2017) show that it is impossible to comply simultaneously to Equalized Odds and Predictive parity, unless $Y \perp\!\!\!\perp S$.

Finding a balance between different notions

Relaxing fairness criteria

- Fairness criteria are imperfect mathematical transposition of qualitative ideas;
- Evaluations of fairness criteria are subjected to uncertainties;
- Some fairness criteria are incompatible;

so, it is wise to

- introduce some quantitative measures of violation of the fairness criteria;
- seek for a good trade-off between different fairness criteria and regret (Pareto frontier).

Learning with fairness constraints

Typical approach:

- 1 Introduce some quantitative fairness constraints $C_1(F), C_2(F), \dots$

Ex: $C_1(F) = |\mathbb{P}[F = 1|S = 1] - \mathbb{P}[F = 1|S = 0]|$

- 2 Minimize the risk $R(F)$ under these constraints

$$F^* \in \underset{F \in \mathcal{F}: C_1(F) \leq \delta_1, \dots, C_K(F) \leq \delta_K}{\operatorname{argmin}} R(F)$$

As usual, the risk $R(F)$, but also the constraints $C_k(F)$ cannot be directly computed

→ two main approaches: In-Processing and Post-Processing

Learning with fairness constraints: In-Processing

In-Processing

A typical In-Processing approach is to

- 1 Take empirical versions $\hat{R}(F)$ and $\hat{C}_k(F)$ of the risk and of the constraints
($\hat{R}(F)$ may include some regularisation terms)
- 2 Minimise the empirical version of the initial problem

$$\hat{F} \in \underset{F \in \mathcal{F}: \hat{C}_1(F) \leq \delta_1, \dots, \hat{C}_K(F) \leq \delta_K}{\operatorname{argmin}} \hat{R}(F)$$

→ some examples in online learning below

Learning with fairness constraints: Post-Processing

Post-Processing

- 1 Express the optimal fair predictor

$$F^* \in \underset{F \in \mathcal{F} : C_1(F) \leq \delta_1, \dots, C_K(F) \leq \delta_K}{\operatorname{argmin}} R(F)$$

in terms of the risk-optimal predictor $F^\dagger \in \operatorname{argmin}_{F \in \mathcal{F}} R(F)$:

→ very often, F^* can be obtained via a (quite) simple transformation $F^* = \mathcal{T}(F^\dagger)$

- 2 Compute a predictor \hat{F}^\dagger with a small risk $R(\hat{F}^\dagger)$ with classical technics
- 3 Plug-in: compute the estimator by applying the transformation of the first step $\hat{F} = \mathcal{T}(\hat{F}^\dagger)$

→ Christophe Denis' talk

Another approach: Pre-Processing

Pre-Processing: learning fair features

The rationale is

- 1 to find a mapping ϕ such that

$$\phi(X, S) \perp\!\!\!\perp S$$

- 2 to apply classical ML algorithms with features $\phi(X, S)$ instead of X

Example: learning from data the mapping

$$\phi \in \underset{\phi(X, S) \perp\!\!\!\perp S}{\operatorname{argmin}} \mathbb{E} [d(X, \phi(X, S))]$$

Causal fairness

Causal fairness: identifying causes of unfairness

Causal fairness aims to identify sources of unfairness.

Causal fairness

Typical approach

- The relations between attributes (X, S) and their influence on outcome Y is modeled by structural equations
- These structural equations capture the influence of sensitive attributes
- The objective is to remove all discriminatory influences

Pros: causal fairness is an individual fairness notion 😊

Caveat: the notions of causal fairness heavily rely on the causal model. The accuracy of this model is critical. 😞

→ alternative notions based on optimal transport in Jean-Michel Loubes and Fanny Jourdan's talks 😊

Causal fairness: Counterfactual fairness

Counterfactual fairness

A predictor F is fair if, for any individual (x, s) , the distribution of F is unchanged, had the **same individual** been of type s' .

More formally,

$$\mathbb{P}[F_{S \leftarrow s}(U) = y | X = x, S = s] = \mathbb{P}[F_{S \leftarrow s'}(U) = y | X = x, S = s],$$

where the distribution of $F_{S \leftarrow s'}(U) = y$ given $(X = x, S = s)$ is obtained by computing F with the intervention $do(S = s')$ and with the latent variables U distributed according to the conditional distribution of U given $(X = x, S = s)$.

An example of counterfactually fair predictor is when, in the causal graph, it does not depend on a descendant of the sensitive attribute.

Causal fairness: Counterfactual fairness

Counterfactual fairness

A predictor F is fair if, for any individual (x, s) , the distribution of F is unchanged, had the **same individual** been of type s' .

More formally,

$$\mathbb{P}[F_{S \leftarrow s}(U) = y | X = x, S = s] = \mathbb{P}[F_{S \leftarrow s'}(U) = y | X = x, S = s],$$

where the distribution of $F_{S \leftarrow s'}(U) = y$ given $(X = x, S = s)$ is obtained by computing F with the intervention $do(S = s')$ and with the latent variables U distributed according to the conditional distribution of U given $(X = x, S = s)$.

An example of counterfactually fair predictor is when, in the causal graph, it does not depend on a descendant of the sensitive attribute.

Causal fairness: No unresolved discrimination

Resolving attribute

A resolving attribute is an attribute that is influenced by the sensitive attribute in a non-discriminatory manner.

Ex: body strength for hiring piano movers

No unresolved discrimination

A prediction has no unresolved discrimination if there exists no path from the sensitive attribute to the prediction, except via a resolving variable.

Summary

Three main statistical notions of algorithmic fairness

- 1 **Individual fairness** aims to treat similar people similarly (individual notions);
- 2 **Group fairness** seeks to comply to fairness criteria at the sub-population level (statistical notions);
- 3 **Causal fairness** tries to identify causes of unfairness in order to act on them (causal notions).

Three main algorithmic approaches

- **Pre-processing:** aim to remove biases from data
- **In-processing:** produce prediction by minimizing empirical risk under empirical fairness constraints
- **Post-processing:** take predictions from standard predictors as input, and adjust them to comply to fairness requirements

2- Improving fairness in online learning?

Collaborators and references



A unified approach to fair online learning via Blackwell approachability

E. Chzhen, C. Giraud, G. Stoltz; NeurIPS 2021 (spotlight).

Small Total-Cost Constraints in CBwK, with Application to Fairness

E. Chzhen, C. Giraud, Z. Li, G. Stoltz; NeurIPS 2023

Parameter-free projected gradient descent

E. Chzhen, C. Giraud, G. Stoltz; arXiv:2305.19605

Collaborators and references



The price of unfairness in linear bandits with biased feedback

S. Gaucher, A. Carpentier, C. Giraud; NeurIPS 2022.

Contextual online setting

Covariate and sensitive attribute

Each request is characterized by a covariate $x \in \mathcal{X}$ (observed) and a sensitive attribute $s \in \{-1, +1\}$ (observed or not).

Informal description of a typical setting

At each epoch $t = 1, 2, \dots$

- The Learner observes a context (x_t, s_t) or x_t only
- The Learner performs an action (or prediction) a_t
- The Learner observes a feedback y_t and suffers a regret r_t (stochastic or adversarial)

Goal of the learner

To minimize the cumulative regret $\sum_t r_t$, while complying to some fairness criteria (and possibly some other constraints).

Instantiating fairness constraints in online learning

Fairness cost

Fairness criteria can be encoded as vector valued cost constraints.

Example: demographic parity

The empirical demographic parity criteria (for $a_t \in \{0, 1\}$)

$$\left| \frac{1}{p_1 T} \sum_{t \leq T; s_t=1} a_t - \frac{1}{p_{-1} T} \sum_{t \leq T; s_t=-1} a_t \right| = \tilde{O}(T^{-1/2})$$

can be encoded as

$$\sum_{t \leq T} c_t = \tilde{O}(\sqrt{T}) \quad \text{with} \quad c_t := \begin{bmatrix} a_t s_t / p_{s_t} \\ -a_t s_t / p_{s_t} \end{bmatrix}.$$

Our contributions

Informal objective

$$\min_{\sum_{t \leq T} c_t \leq \tilde{O}(\sqrt{T})} \sum_{t \leq T} r_t.$$

Two points of view

1 In adversarial setting:

- ▶ the fair learning problem can be formulated as a contextual approachability problem,
- ▶ Blackwell theory can be adapted to handle this setting.

2 In stochastic bandit setting:

- ▶ the fairness objective falls into the Contextual Bandit with Knapsack (CBwK) framework,
- ▶ the theory for CBwK must be improved to handle $\tilde{O}(\sqrt{T})$ constraints (and signed cost).

Adversarial Setting :

a Contextual Blackwell Approachability Perspective

A unified approach to fair online learning via Blackwell approachability.

E. Chzhen, C. Giraud, G. Stoltz; NeurIPS 2021 (spotlight).

Online learning setting: formal description

We model our fair online learning problem as a contextual learning game between the Learner and Nature.

Stochastic attributes (context)

At each time t , the attributes (x_t, s_t) are sampled according to \mathbf{Q} , independently from the past.

Nature (un)awareness

Let G denotes Nature (un)awareness mapping

- Nature *awareness* $G(x, s) = (x, s)$,
- Nature *unawareness*: $G(x, s) = x$.

Nature is an adverse player

At each time t , Nature observes $G(x_t, s_t)$ and outputs an adversarial feedback y_t .

Fair online learning as a contextual approachability problem

Encoding the objectives of the learner

We can encode the learning objectives (vanishing-regret, demographic parity, etc) via

- a vector-valued payoff function $\mathbf{m}(a_t, y_t, x_t, s_t)$
- and a target set \mathcal{C} .

The learning objective is to comply to $\frac{1}{T} \sum_{t=1}^T \mathbf{m}(a_t, y_t, x_t, s_t) \rightarrow \mathcal{C}$.

Examples of targets (to be combined)

Criterion	Vector payoff function \mathbf{m}	Closed convex target set \mathcal{C}
Demographic parity	$\mathbf{m}_{\text{DP}}(a, s) = \left(\frac{a}{p-1} \mathbf{1}_{s=-1}, \frac{a}{p_1} \mathbf{1}_{s=1} \right)$	$\mathcal{C}_{\text{DP}} = \{(u, v) \in \mathbb{R}^2 : u - v \leq \delta\}$
No-regret	$\mathbf{m}_{\text{reg}}(a, y, x, s) = (f(a, y, x, s) - f(a', y, x, s))_{a' \in \mathcal{A}}$	$\mathcal{C}_{\text{reg}} = [0, +\infty)^N$
Group-calibration	$\mathbf{m}_{\text{gr-cal}}(a, y, s) = ((a' - y) \mathbf{1}_{s=s'} / \gamma_{s'})_{a' \in \mathcal{A}, s' \in \mathcal{S}}$	$\mathcal{C}_{\text{gr-cal}} = \{\mathbf{v} \in \mathbb{R}^{N \mathcal{S} } : \ \mathbf{v}\ _1 \leq \varepsilon\}$
Equalized payoffs	$\mathbf{m}_{\text{eq-pay}}(a, y, x, s) = \left(\frac{f(a, y, x, s')}{\gamma_{s'}} \mathbf{1}_{s=s'} \right)_{s' \in \mathcal{S}}$	$\mathcal{C}_{\text{eq-pay}} = \{(u, v) \in \mathbb{R}^2 : u - v \leq \varepsilon\}$

Online learning setting

Learning setting

For $t = 1, 2, \dots$

- 1 Simultaneously,
 - ▶ the Learner chooses $(\mathbf{p}_t^x)_{x \in \mathcal{X}}$ based on $(\mathbf{m}_T, x_T, s_T)_{T \leq t-1}$
 - ▶ Nature chooses $(\mathbf{q}_t^{G(x,s)})_{(x,s) \in \mathcal{X} \times \mathcal{S}}$ based on $(a_T, y_T, x_T, s_T)_{T \leq t-1}$
- 2 (x_t, s_t) are sampled according to \mathbf{Q} , independently from the past
- 3 Simultaneously
 - ▶ the Learner observes x_t , and picks an action $a_t \in \mathcal{A}$ according to $\mathbf{p}_t^{x_t}$
 - ▶ Nature observes $G(x_t, s_t)$, and picks $y_t \in \mathcal{Y}$ according to $\mathbf{q}_t^{G(x_t, s_t)}$
- 4 The Learner observes the payoff $\mathbf{m}_t = \mathbf{m}(a_t, y_t, x_t, s_t)$ and (x_t, s_t) , while Nature observes (a_t, y_t, x_t, s_t) .

Aim: The Learner wants to ensure that $\bar{\mathbf{m}}_T := \frac{1}{T} \sum_{t=1}^T \mathbf{m}_t \rightarrow \mathcal{C}$ a.s. for some target set \mathcal{C} .

Online learning setting

Learning setting

For $t = 1, 2, \dots$

- 1 Simultaneously,
 - ▶ the Learner chooses $(\mathbf{p}_t^x)_{x \in \mathcal{X}}$ based on $(\mathbf{m}_T, x_T, s_T)_{T \leq t-1}$
 - ▶ Nature chooses $(\mathbf{q}_t^{G(x,s)})_{(x,s) \in \mathcal{X} \times \mathcal{S}}$ based on $(a_T, y_T, x_T, s_T)_{T \leq t-1}$
- 2 (x_t, s_t) are sampled according to \mathbf{Q} , independently from the past
- 3 Simultaneously
 - ▶ the Learner observes x_t , and picks an action $a_t \in \mathcal{A}$ according to $\mathbf{p}_t^{x_t}$
 - ▶ Nature observes $G(x_t, s_t)$, and picks $y_t \in \mathcal{Y}$ according to $\mathbf{q}_t^{G(x_t, s_t)}$
- 4 The Learner observes the payoff $\mathbf{m}_t = \mathbf{m}(a_t, y_t, x_t, s_t)$ and (x_t, s_t) , while Nature observes (a_t, y_t, x_t, s_t) .

Aim: The Learner wants to ensure that $\bar{\mathbf{m}}_T := \frac{1}{T} \sum_{t=1}^T \mathbf{m}_t \rightarrow \mathcal{C}$ a.s. for some target set \mathcal{C} .

Online learning setting

Learning setting

For $t = 1, 2, \dots$

- 1 Simultaneously,
 - ▶ the Learner chooses $(\mathbf{p}_t^x)_{x \in \mathcal{X}}$ based on $(\mathbf{m}_\tau, x_\tau, s_\tau)_{\tau \leq t-1}$
 - ▶ Nature chooses $(\mathbf{q}_t^{G(x,s)})_{(x,s) \in \mathcal{X} \times \mathcal{S}}$ based on $(a_\tau, y_\tau, x_\tau, s_\tau)_{\tau \leq t-1}$
- 2 (x_t, s_t) are sampled according to \mathbf{Q} , independently from the past
- 3 Simultaneously
 - ▶ the Learner observes x_t , and picks an action $a_t \in \mathcal{A}$ according to \mathbf{p}_t^x
 - ▶ Nature observes $G(x_t, s_t)$, and picks $y_t \in \mathcal{Y}$ according to $\mathbf{q}_t^{G(x_t, s_t)}$
- 4 The Learner observes the payoff $\mathbf{m}_t = \mathbf{m}(a_t, y_t, x_t, s_t)$ and (x_t, s_t) , while Nature observes (a_t, y_t, x_t, s_t) .

Aim: The Learner wants to ensure that $\bar{\mathbf{m}}_T := \frac{1}{T} \sum_{t=1}^T \mathbf{m}_t \rightarrow \mathcal{C}$ a.s. for some target set \mathcal{C} .

Assumption: fast enough sequential estimation of \mathbf{Q}

The Player can build estimators $(\hat{\mathbf{Q}}_t)_{t \geq 1}$ of the unknown distribution \mathbf{Q} such that

$$\mathbb{E} \left[\text{TV}^2(\hat{\mathbf{Q}}_t, \mathbf{Q}) \right] \leq c (\log(t))^{-3} \quad \forall t \geq 2 \quad (1)$$

Theorem : Contextual Blackwell approachability

If $\mathcal{C} \subset \mathbb{R}^d$ is closed convex, \mathbf{m} is bounded, and (1) is satisfied, then

$\exists (\mathbf{p}_t^x)_{x \in \mathcal{X}, t \geq 1}$ such that $\forall (\mathbf{q}_t^{G(x,s)})_{(x,s) \in \mathcal{X} \times \{0,1\}, t \geq 1}$ we have $\bar{\mathbf{m}}_T \xrightarrow{a.s.} \mathcal{C}$

if and only if $\forall (\mathbf{q}^{G(x,s)})_{(x,s) \in \mathcal{X} \times \{0,1\}} \exists (\mathbf{p}^x)_{x \in \mathcal{X}}$ such that

$$\mathbf{m}(\mathbf{p}, \mathbf{q}, \mathbf{Q}) := \int_{\mathcal{X} \times \mathcal{S}} \mathbf{m}(\mathbf{p}^x, \mathbf{q}^{G(x,s)}, x, s) d\mathbf{Q}(x, s) \in \mathcal{C}$$

Contextual Blackwell strategy

Set $\mathbf{m}(\mathbf{p}, \mathbf{q}, \hat{\mathbf{Q}}_t) := \int \mathbf{m}(\mathbf{p}^x, \mathbf{q}^{G(x,s)}, x, s) d\hat{\mathbf{Q}}_t(x, s)$. At stage $t + 1$, choose

$$(\mathbf{p}_{t+1}^x)_{x \in \mathcal{X}} \in \operatorname{argmin}_{(\mathbf{p}^x)_x} \max_{(\mathbf{q}^{G(x,s)})_{x,s}} \langle \bar{\mathbf{m}}_t - \Pi_{\mathcal{C}} \bar{\mathbf{m}}_t, \mathbf{m}(\mathbf{p}, \mathbf{q}, \hat{\mathbf{Q}}_t) \rangle$$

Caveats

Caveat 1: the target set \mathcal{C} has to be known

The results can be extended (at the price of some technicalities) to the case where we only have a consistent super-estimate $\hat{\mathcal{C}}_t$ of \mathcal{C} .

Caveat 2: computational cost of projection

Computing the projection $\Pi_{\mathcal{C}}$ can be computationally expensive.

Caveat 3: pessimistic Pareto frontier and slow rates

- The adversarial setting leads to pessimistic Pareto frontier (trade-off) between the different criteria;
- The rates are governed by the estimation rate $\text{TV}(\hat{\mathbf{Q}}_t, \mathbf{Q})$, which is typically slow outside the finite case.

Stochastic setting:

A Contextual Bandit with Knapsack perspective

Small Total-Cost Constraints in CBwK, with Application to Fairness

E. Chzhen, C. Giraud, Z. Li, G. Stoltz; NeurIPS 2023

Stochastic setting

Learning problem

- The learner observes $\tilde{x}_t = (x_t, s_t) \stackrel{\text{i.i.d.}}{\sim} \mathbf{Q}$
- The learner chooses a policy $\pi_t : \tilde{\mathcal{X}} \rightarrow \mathcal{P}(\mathcal{A})$, and picks an action $a_t \sim \pi_t(\tilde{x}_t)$,
- The learner receives a feedback y_t and a fairness cost c_t such that

$$\mathbb{E}[y_t | \mathcal{F}_t] = f(\tilde{x}_t, a_t) \quad \text{and} \quad \mathbb{E}[c_t | \mathcal{F}_t] = c(\tilde{x}_t, a_t).$$

- The learner suffers a regret $r_t = \text{OPT} - y_t$ (described below).

Example: Demographic Parity

$$c(\tilde{x}_t, a_t) = \begin{bmatrix} a_t s_t / p_{s_t} \\ -a_t s_t / p_{s_t} \end{bmatrix}.$$

Optimal policy and regret

Optimal static policy and regret

The optimal static feedback is

$$\text{OPT}(\mathbf{Q}, f, c) := \max_{\pi : \mathbb{E}_{\mathbf{Q}}[\sum_{a \in \mathcal{A}} c(\tilde{X}, a)\pi_a(\tilde{X})] \leq \delta_T} \mathbb{E}_{\mathbf{Q}} \left[\sum_{a \in \mathcal{A}} f(\tilde{X}, a)\pi_a(\tilde{X}) \right]$$

and the regret is

$$r_t = \text{OPT}(\mathbf{Q}, f, c) - y_t.$$

Learning Objective

Minimize the cumulative regret $\sum_{t \leq T} r_t$ while complying to the fairness

constraint $\frac{1}{T} \sum_{t \leq T} c_t \leq \delta_T$ (w.h.p.).

CBwK problem

- 1 we recognize a Contextual Bandit with Knapsack (CBwK) problem
- 2 but state of the art theory can only handle $\delta_T = T^{-1/4}$ (or \mathcal{X} finite), which is too large for fairness constraints, where we typically wish to have $\delta_T = \tilde{O}(T^{-1/2})$

Learning assumption

Assumption: UCB and LCB

We can built UCB and LCB such that with probability $\geq 1 - \delta$

$$\hat{f}_t^{\text{UCB}}(.,.) \approx f(.,.) + \tilde{O}_\delta(1/\sqrt{t})$$

$$\hat{c}_t^{\text{LCB}}(.,.) \approx c(.,.) + \tilde{O}_\delta(1/\sqrt{t})$$

Examples

Linear or logistic model : when

$$f(x, a) = \eta(\varphi(x, a)^T \theta_a) \quad \text{and} \quad c(x, a) = \eta(\psi(x, a)^T \beta_a),$$

with $\eta(u) = u$ or $\eta(u) = e^u / (1 + e^u)$, we can use variant of LinUCB or LogisticUCB1.

A first idea

Idea1: playing empirical optimal static policy

Choose a_t according to a policy $\hat{\pi}_t$ maximizing $\text{OPT}(\hat{\mathbf{Q}}_t, \hat{f}_t^{\text{UCB}}, \hat{c}_t^{\text{LCB}})$.

Issues

- 1 The analysis of

$$\text{OPT}(\mathbf{Q}, f, c) - \text{OPT}(\hat{\mathbf{Q}}_t, \hat{f}_t^{\text{UCB}}, \hat{c}_t^{\text{LCB}})$$

produces some $\text{TV}(\hat{\mathbf{Q}}_t, \mathbf{Q})$ terms, leading to slow rates / large fairness violation.

- 2 Solving $\text{OPT}(\hat{\mathbf{Q}}_t, \hat{f}_t^{\text{UCB}}, \hat{c}_t^{\text{LCB}})$ is computationally expensive

Lagrangian version

Lagrangian formulation

$$\begin{aligned}\text{OPT}(\mathbf{Q}, f, c) &= \max_{\pi : \mathbb{E}_{\mathbf{Q}}[\sum_{a \in \mathcal{A}} c(\tilde{X}, a) \pi_a(\tilde{X})] \leq \delta_T} \mathbb{E}_{\mathbf{Q}} \left[\sum_{a \in \mathcal{A}} f(\tilde{X}, a) \pi_a(\tilde{X}) \right] \\ &= \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{Q}} \left[\sum_{a \in \mathcal{A}} \pi_a(\tilde{X}) \left(f(\tilde{X}, a) - \langle \lambda, c(\tilde{X}, a) - \delta_T \rangle \right) \right] \\ \text{strong duality} \rightarrow &= \min_{\lambda \geq 0} \max_{\pi} \mathbb{E}_{\mathbf{Q}} \left[\sum_{a \in \mathcal{A}} \pi_a(\tilde{X}) \left(f(\tilde{X}, a) - \langle \lambda, c(\tilde{X}, a) - \delta_T \rangle \right) \right] \\ &= \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{Q}} \left[\max_{a \in \mathcal{A}} \left\{ f(\tilde{X}, a) - \langle \lambda, c(\tilde{X}, a) - \delta_T \rangle \right\} \right]\end{aligned}$$

Two immediate benefits

- 1 for a fixed λ the problem is separable, and \mathbf{Q} can be forgotten;
- 2 we only need to learn the optimal $\lambda^* \in \mathbb{R}^d \implies$ parametric rates. 😊

High-level algorithm: Primal-dual descent-ascent

Iterate

- **full optimisation on primal variable:** pick

$$a_t \in \operatorname{argmax}_{a \in \mathcal{A}} \left\{ \hat{f}_t(\tilde{x}_t, a) - \langle \lambda_{t-1}, \hat{c}_t(\tilde{x}_t, a) - \delta_T \rangle \right\}$$

- **projected subgradient step on dual variable:** update

$$\lambda_t = (\lambda_{t-1} + \gamma (\hat{c}_t(\tilde{x}_t, a_t) - \delta_T))_+$$

Issues

- 1 **Benign issue:** we must replace δ_T by $\delta'_T = \delta_T - \tilde{O}(1/\sqrt{T})$ to prevent from violation of the fairness criteria due to random fluctuations
- 2 **Major issue:** to satisfy the constraints, we need to set $\gamma \approx \|\lambda^*\|/\sqrt{T}$



Choice of step size γ

Bounds (informal)

For a fixed step size $\gamma > 0$, we have w.h.p.

- $\|\text{constraint violation}\| = \tilde{O}\left(\sqrt{T} + \frac{1 \vee \|\lambda^*\|}{\gamma}\right)$
- $\text{Regret} = \tilde{O}\left(\gamma T + \|\lambda^*\| \sqrt{T}\right)$.

So best γ is $\gamma^* = (1 \vee \|\lambda^*\|)/\sqrt{T}$:

- $\|\text{constraint violation}\| = \tilde{O}\left(\sqrt{T}\right)$
- $\text{Regret} = \tilde{O}\left((1 \vee \|\lambda^*\|)\sqrt{T}\right)$.

Mispecified γ

If we simply set $\gamma = 1/\sqrt{T}$, then we have

$$\|\text{constraint violation}\| = \tilde{O}\left((1 \vee \|\lambda^*\|)\sqrt{T}\right) \quad \text{☹️}$$

Tuning γ

Good old doubling trick

- Start from $\gamma = 1/\sqrt{T}$
- Tracking the constraint violation at each epoch t , we can detect from the bound

$$\|\text{constraint violation}\| = \tilde{O}\left(\sqrt{t} + \frac{1 \vee \|\lambda_*\|}{\gamma}\right)$$

if our current choice of γ is too small

- If so, double γ .

Adaptive algorithm

Adaptive version

Iterate: for $t \geq 1$

- Pick $a_t \in \operatorname{argmax}_{a \in \mathcal{A}} \left\{ \hat{f}_t(\tilde{x}_t, a) - \langle \lambda_{t-1}, \hat{c}_t(\tilde{x}_t, a) - \delta'_T \rangle \right\}$
- Update $\lambda_t = \left(\lambda_{t-1} + \frac{2^k}{\sqrt{T}} (\hat{c}_t(\tilde{x}_t, a_t) - \delta'_T) \right)_+$

Until $\left\| \left(\sum_{\tau=T_k}^t c_\tau - (t - T_k + 1) \delta'_T \right)_+ \right\| > \tilde{O}(\sqrt{T})$

Then: increase k by one, set $T_k = t + 1$, and **iterate** again.

Theory

Regret bound

For $\delta_T \geq \tilde{O}(T^{-1/2})$, the above algorithm fulfills with probability at least $1 - \delta$

$$\sum_{t \leq T} r_t \leq \tilde{O}_\delta \left((1 \vee \|\lambda^*\|) \sqrt{T} \right) \quad \frac{1}{T} \sum_{t \leq T} c_t \leq \delta_T.$$

Suitable for fairness constraints 😊

Optimality?

A proof scheme suggests that this regret is optimal.

Concluding remark

An important fairness issue in decision making, not addressed in this presentation, is the problem of **misalignment between evaluations and objectives**.

This question does not fall into the framework described in this presentation, but it is an important question in order to improve overall fairness in decision making.