

**Bienvenue à toutes et à tous,
nous allons commencer dans un instant.**



Contributions à l'apprentissage semi-supervisé : **équité** et **étiquetage** dans les problèmes à classes multiples

François HU

07 février 2023

Remerciements



Caroline Hillairet,
Professeure à l'ENSAE



Romuald Elie,
Professeur à l'Université
Gustave Eiffel



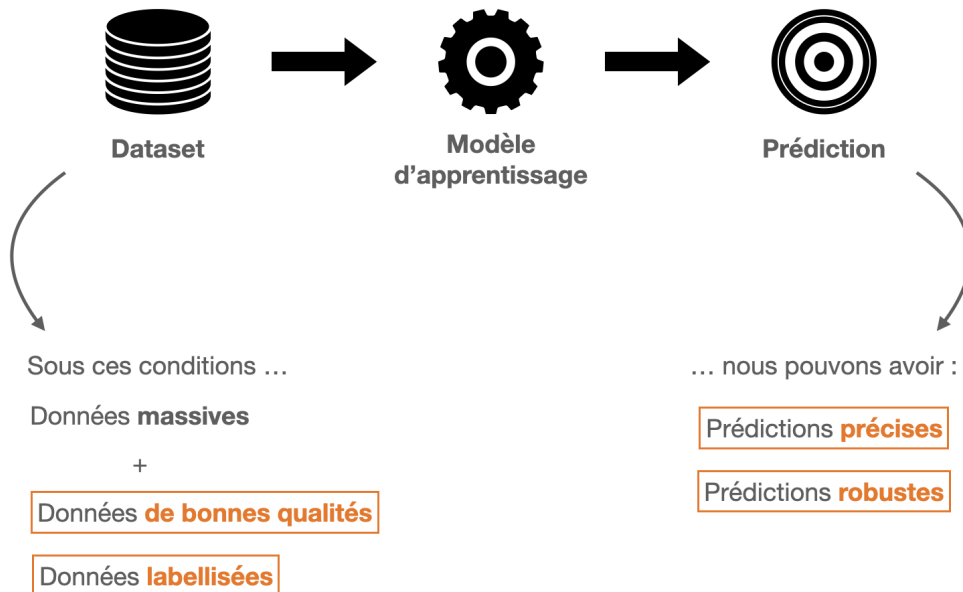
Marc Juillard,
Directeur du Data Hub de
Société Générale Assurances

Sommaire

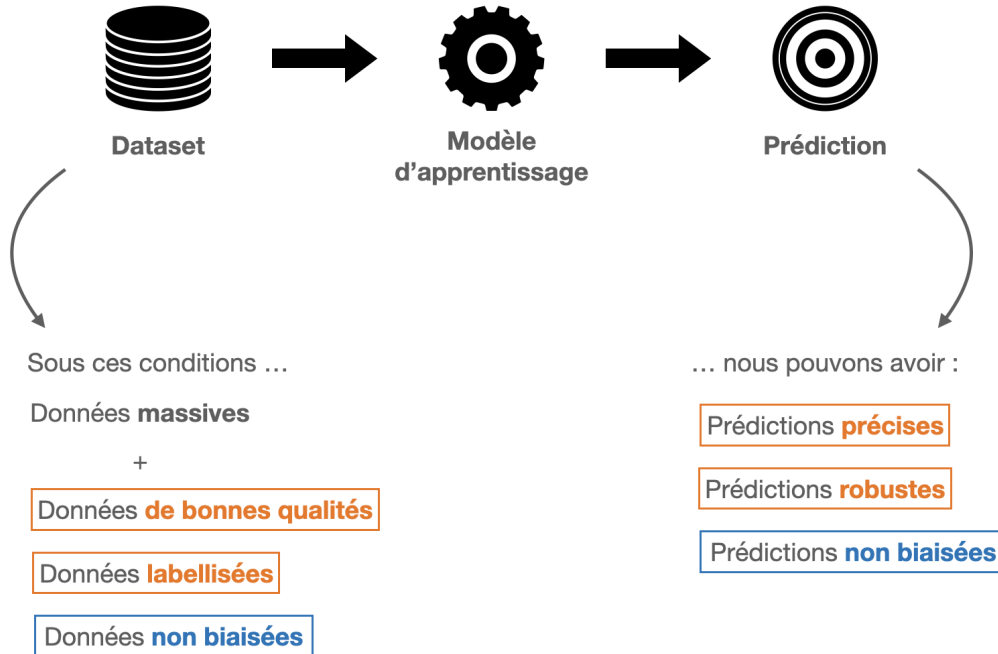
- Introduction
- **Défi 1** : étiquetage dynamique
- **Défi 2** : équité dans la classification multi-classes
- Conclusion + extensions

Mots-clés : apprentissage statistique, contrôle stochastique, optimisation.

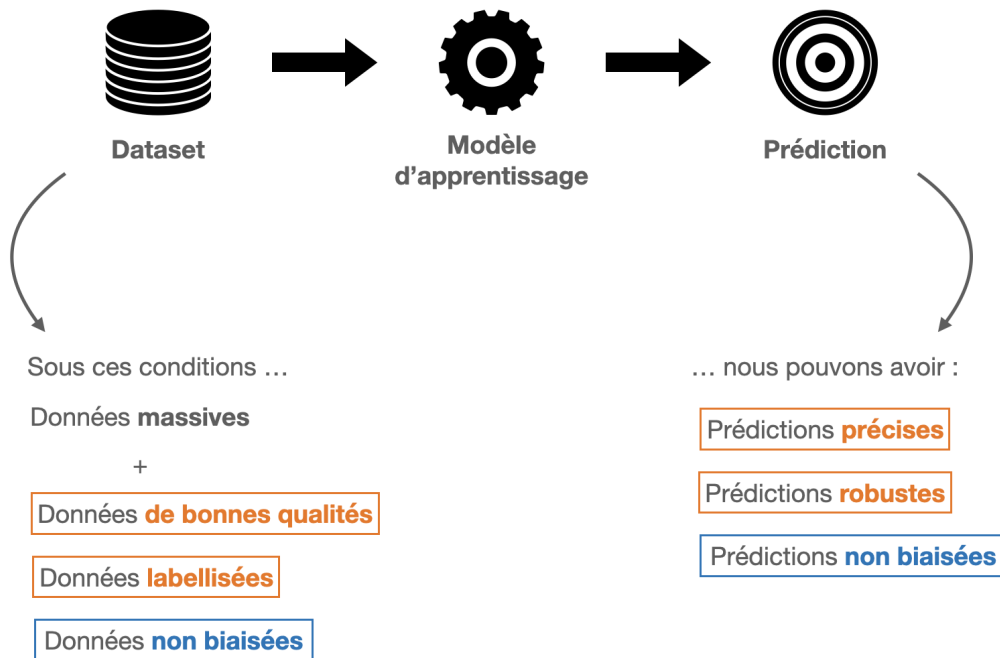
Contexte général : processus d'apprentissage statistique



Contexte général : processus d'apprentissage statistique



Contexte général : processus d'apprentissage statistique



Quelques exemples en assurance :

Catégorisation des photos de **voitures endommagées**



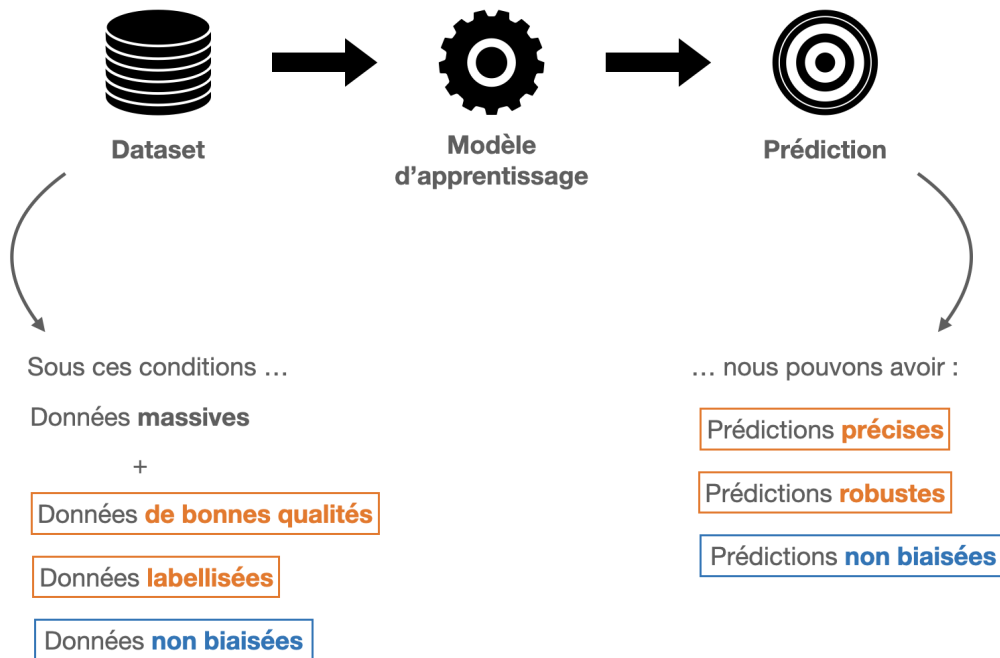
Analyse des **données télématiques**



Détection de la conformité **GDPR**



Contexte général : processus d'apprentissage statistique



Quelques exemples en assurance :

Catégorisation des photos de **voitures endommagées**



Problème d'**étiquetage** :
besoin d'assureurs experts

Analyse des **données télématiques**



Problème **confidentialité** :
dédire certaines var. sensibles

Détection de la conformité **GDPR**



Problèmes d'**étiquetage** et d'**équité** :
besoin de juristes

Défi 1 : problèmes **d'étiquetage** **Défi 2** : problème **d'équité**

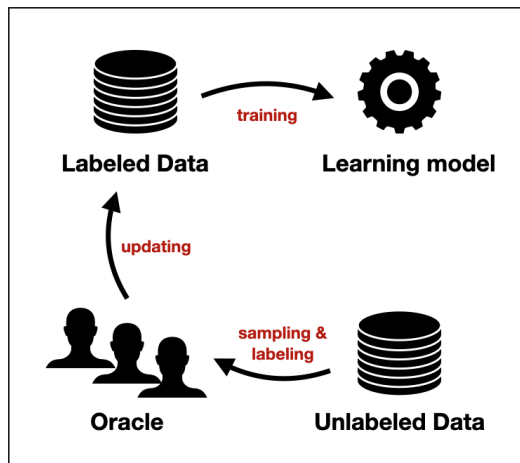
Défis et contributions

➔ **Défi 1** : entraîner un modèle ML avec un budget d'étiquetage limité

Défi 2 : garantir l'équité algorithmique dans les problèmes multi-classes

Défi 1 : quelques idées et leurs limites

Étiquetage en parallèle

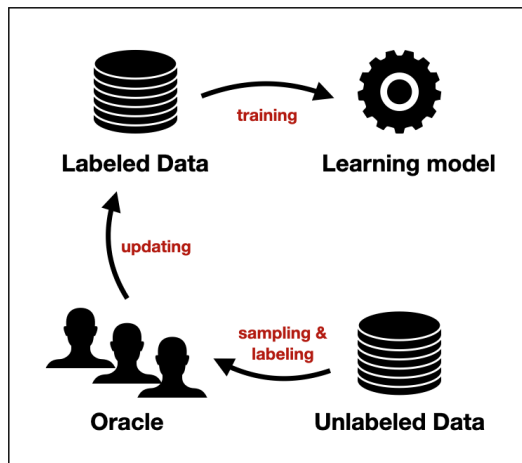


✘ coûteux et long

✔ génère des labels

Défi 1 : quelques idées et leurs limites

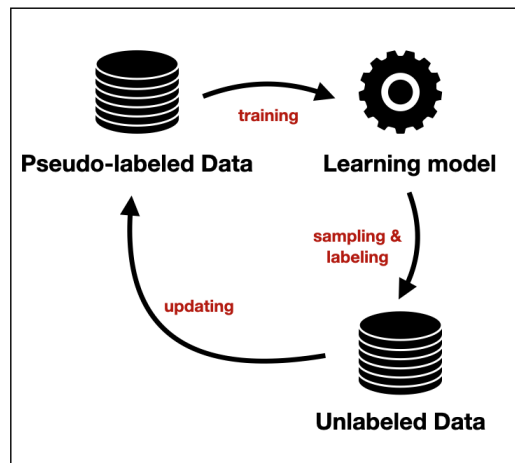
Étiquetage en parallèle



✗ coûteux et long

✓ génère des labels

Apprentissage semi-supervisé

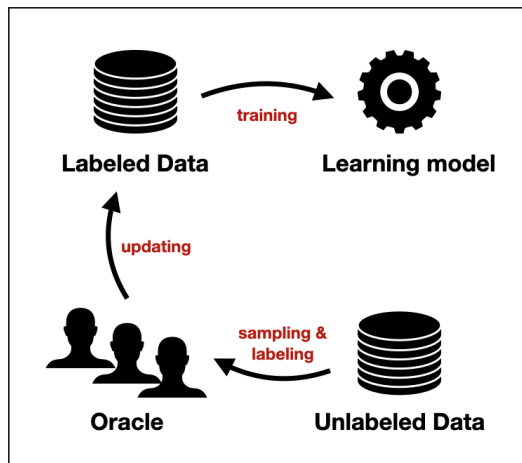


✓ moins coûteux, moins long

✗ génère des **pseudo**-labels

Défi 1 : Active Learning (AL)

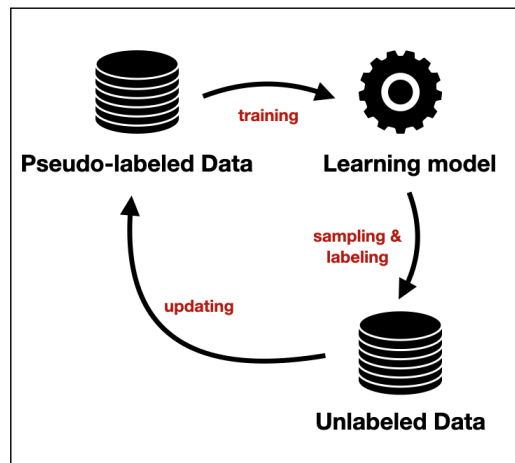
Étiquetage en parallèle



✗ coûteux et long

✓ génère des labels

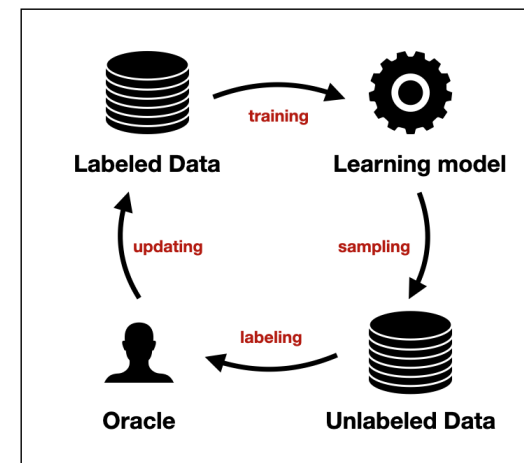
Apprentissage semi-supervisé



✓ moins coûteux, moins long

✗ génère des **pseudo**-labels

Apprentissage actif (*Active Learning*)



✓ moins coûteux, moins long

✓ génère des labels

Exemple dans le cas multi-classes : entropie de Shannon

Algorithm 1 Outline of active learning (AL) process

Input: h base estimator, train-set $\mathcal{D}^{(train)}$, pool-set $\mathcal{D}_X^{(pool)}$

Step 1. Fit h on the train-set $\mathcal{D}^{(train)}$

Step 2. Given a **score** $l(x, h)$, we sample:

$$x^* = \operatorname{argmax}_{x \in \mathcal{D}_X^{(pool)}} \{l(x, h)\}$$

Example: Entropy-based [Sha48]

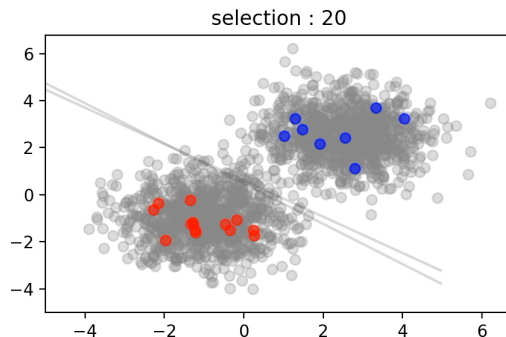
$$l(x, h) = - \sum_{k=1}^K \mathbb{P}(h(x) = k|x) \log \mathbb{P}(h(x) = k|x)$$

Step 3. If y^* is its label then we update:

$$\mathcal{D}^{(train)} = \mathcal{D}^{(train)} \cup \{(x^*, y^*)\}$$

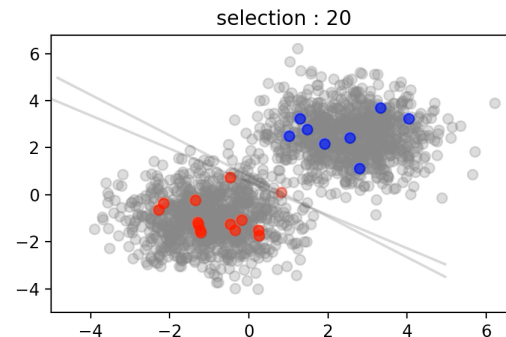
$$\mathcal{D}_X^{(pool)} = \mathcal{D}_X^{(pool)} - \{x^*\}$$

Step 4. Return to **step 1** until convergence.



Apprentissage
passif

(Passive learning)



Apprentissage
actif

(Active learning)

Exemple : entropie de Shannon

Algorithm 1 Outline of active learning (AL) process

Input: h base estimator, train-set $\mathcal{D}^{(train)}$, pool-set $\mathcal{D}_X^{(pool)}$

Step 1. Fit h on the train-set $\mathcal{D}^{(train)}$

Step 2. Given a score $l(x, h)$, we sample:

$$x^* = \operatorname{argmax}_{x \in \mathcal{D}_X^{(pool)}} \{l(x, h)\}$$

Example: Entropy-based [Sha48]

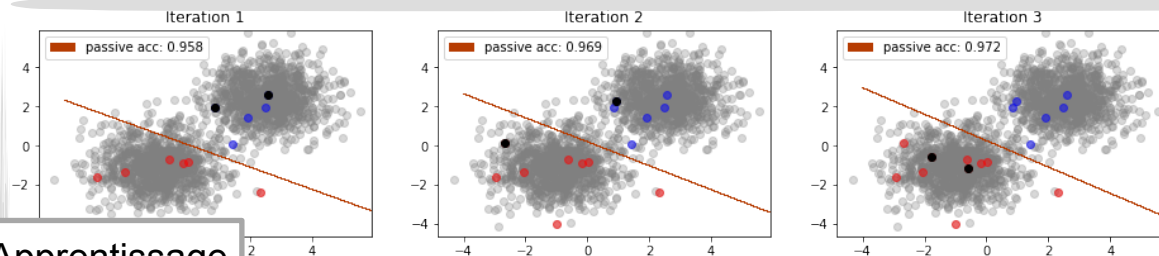
$$l(x, h) = - \sum_{k=1}^K \mathbb{P}(h(x) = k|x) \log \mathbb{P}(h(x) = k|x)$$

Step 3. If y^* is its label then we update:

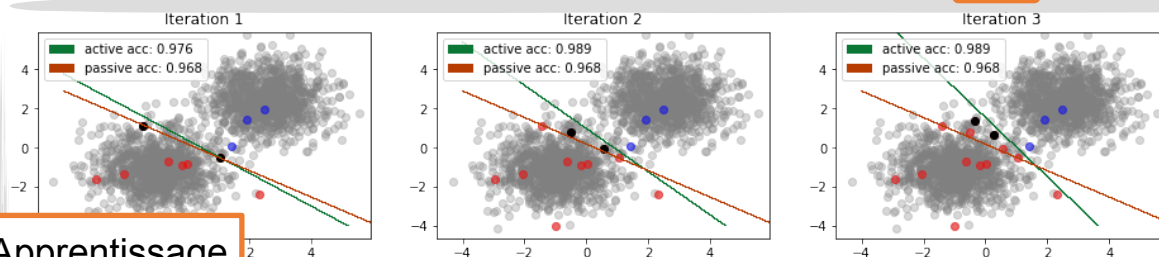
$$\mathcal{D}^{(train)} = \mathcal{D}^{(train)} \cup \{(x^*, y^*)\}$$

$$\mathcal{D}_X^{(pool)} = \mathcal{D}_X^{(pool)} - \{x^*\}$$

Step 4. Return to **step 1** until convergence.



Apprentissage passif



Apprentissage actif

Batch Mode Active Learning (BMAL)

Dataset : NPS de Sogécap (scores attribués par les clients sur un produit)



- **Verbatims** (\mathcal{X}) explication du score par le client, encodé par doc2vec
- **Analyse des sentiments** $\mathcal{Y} = \{ \text{score} \leq T, \text{score} > T \}$, T à déterminer

Processus d'AL : XGBoost + Entropie de Shannon

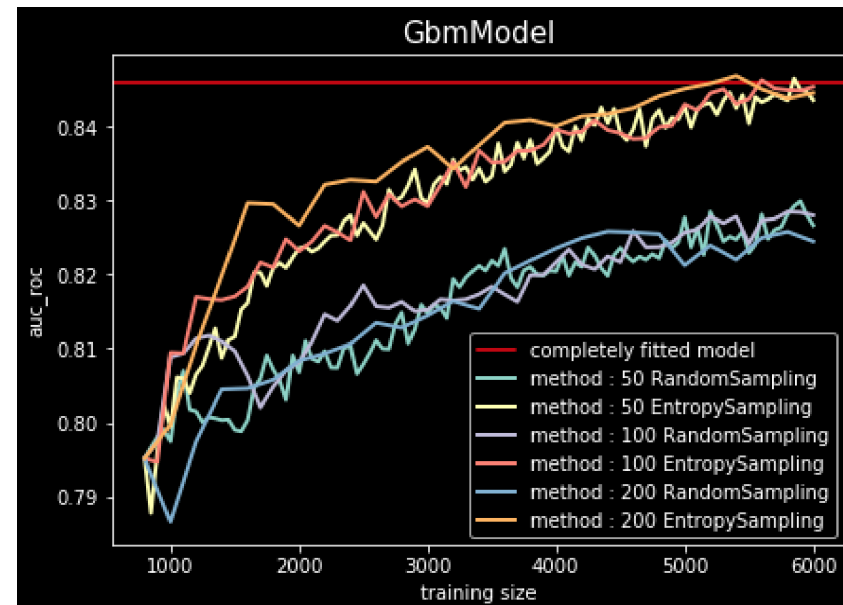


- **BMAL** : à chaque itération AL, échantillonner les b meilleures instances
- $b \in \{50, 100, 200\}$

Pourquoi utiliser **BMAL** en pratique ?



1. utile pour l'**étiquetage parallèle**
2. éviter le coût des **délais de ré-apprentissage**.



Batch Mode Active Learning (BMAL)

Dataset : NPS de Sogécap (scores attribués par les clients sur un produit)



- **Verbatims** (\mathcal{X}) explication du score par le client, encodé par doc2vec
- **Analyse des sentiments** $\mathcal{Y} = \{ \text{score} \leq T, \text{score} > T \}$, T à déterminer

Processus d'AL : XGBoost + Entropie de Shannon



- **BMAL** : à chaque itération AL, échantillonner les b meilleures instances
- $b \in \{50, 100, 200\}$

Pourquoi utiliser **BMAL** en pratique ?

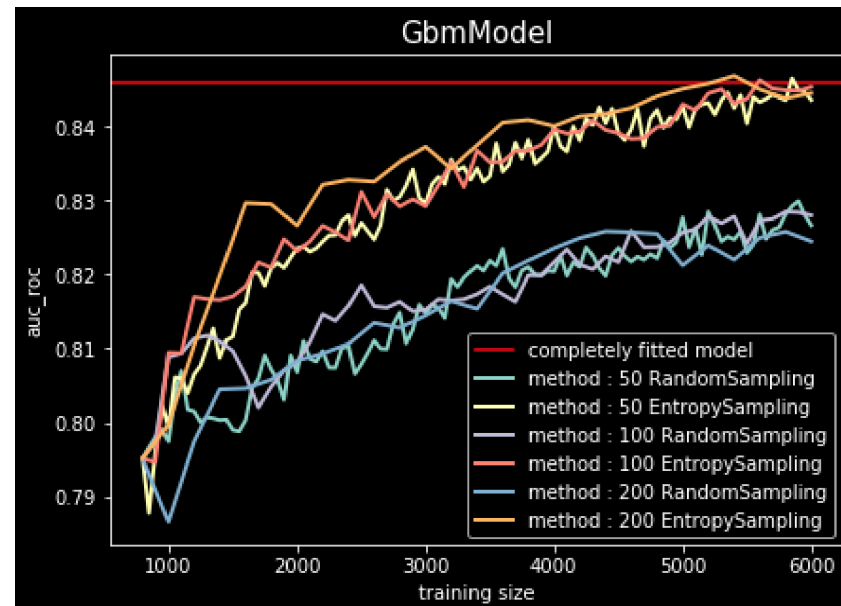


1. utile pour l'**étiquetage parallèle**
2. éviter le coût des **délais de ré-apprentissage**.

Objectif : trouver une séquence de tailles de lots $(b_i)_i$ avec un compromis entre



- **Maximiser** la performance du modèle
- **Réduire** le nombre d'itérations AL



Contribution : considérer BMAL comme un **Processus de Décision Markovien** (PDM)

Dynamic-Size Batch Mode Active Learning (DS-BMAL)

Dynamique des processus d'état (mouvement brownien W_t) :

$$\begin{cases} \text{(performance/qualité)} & dQ_t = \mu(B_t, b_t) \cdot Q_t(1 - Q_t) \cdot dt + \sigma(B_t, b_t) \cdot Q_t(1 - Q_t) \cdot dW_t \\ \text{(taille du train-set)} & dB_t = b_t \cdot dt \end{cases}$$

Problème d'optimisation

U fonction d'utilité qui modélise l'aversion au risque du modèle étiquetage-performance τ le temps d'arrêt

$$V_0 = \sup_b \mathbb{E} \left[U(Q_\tau) - \int_0^\tau C(b_s) ds \right] \quad \text{(Fonction de valeur)}$$

C coût supposé convexe p.r. au batch b

Paramètres :

- μ et σ sont déduits par analyse numérique :

$$\mu = \mu(B, b) \propto \frac{b}{B} \text{ et } \sigma = \sigma(B, b) \propto \frac{b}{B}$$

- C et U sont des fonctions de puissance standard :

$$C(b) \propto b^2 \text{ et } U(Q) = Q^p, p \in (0,1).$$

Dynamic-Size Batch Mode Active Learning (DS-BMAL)

Dynamique des processus d'état (mouvement brownien W_t):

- (performance/qualité) $\begin{cases} dQ_t = \mu(B_t, b_t) \cdot Q_t(1 - Q_t) \cdot dt + \sigma(B_t, b_t) \cdot Q_t(1 - Q_t) \cdot dW_t \\ dB_t = b_t \cdot dt \end{cases}$
- (taille du train-set)

Problème d'optimisation

U fonction d'utilité qui modélise l'aversion au risque du modèle étiquetage-performance τ le temps d'arrêt

$$V_0 = \sup_b \mathbb{E} \left[U(Q_\tau) - \int_0^\tau C(b_s) ds \right] \quad (\text{Fonction de valeur})$$

C coût supposé convexe p.r. au batch b

Paramètres :

- μ et σ sont déduits par analyse numérique :
- $\mu = \mu(B, b) \propto \frac{b}{B}$ et $\sigma = \sigma(B, b) \propto \frac{b}{B}$
- C et U sont des fonctions de puissance standard :

$$C(b) \propto b^2 \text{ et } U(Q) = Q^p, p \in (0,1).$$

Principe de la programmation dynamique

$$V_t := v(Q_t, B_t) = \sup_{b_s, s \in [t, \tau]} \mathbb{E} \left[v(Q_{t+h}, B_{t+h}) - \int_t^{t+h} C(b_s) ds \mid \mathcal{F}_t^W \right]$$

Hamilton Jacobi Bellman (HJB) equation

Equation de HJB dans $[0,1] \times [0, B_{MAX}]$ nous donne :

$$\sup_{b \geq 0} \left\{ \underbrace{\mu Q(1-Q) \frac{\partial V}{\partial Q}(Q, B) + b \frac{\partial V}{\partial B}(Q, B) + \frac{1}{2} \sigma^2 Q^2(1-Q)^2 \frac{\partial^2 V}{\partial Q^2}(Q, B) - C(b)}_{= A(Q, B, b, V)} \right\} = 0$$

Avec les **conditions aux bords** :

$$v(0^+, B) = U(0)$$

$$v(1^-, B) = U(1)$$

$$v(Q, B_{MAX}) = U(Q) \quad \text{pour } Q \in (0,1)$$

Dynamic-Size Batch Mode Active Learning (DS-BMAL)

Dynamique des processus d'état (mouvement brownien W_t):

- (performance/qualité) $\begin{cases} dQ_t = \mu(B_t, b_t) \cdot Q_t(1 - Q_t) \cdot dt + \sigma(B_t, b_t) \cdot Q_t(1 - Q_t) \cdot dW_t \\ dB_t = b_t \cdot dt \end{cases}$
- (taille du train-set)

Problème d'optimisation

U fonction d'utilité qui modélise l'aversion au risque du modèle étiquetage-performance τ le temps d'arrêt

$$V_0 = \sup_b \mathbb{E} \left[U(Q_\tau) - \int_0^\tau C(b_s) ds \right] \quad (\text{Fonction de valeur})$$

C coût supposé convexe p.r. au batch b

Paramètres :

- μ et σ sont déduits par analyse numérique :
- $\mu = \mu(B, b) \propto \frac{b}{B}$ et $\sigma = \sigma(B, b) \propto \frac{b}{B}$
- C et U sont des fonctions de puissance standard :

$$C(b) \propto b^2 \text{ et } U(Q) = Q^p, p \in (0,1).$$

Principe de la programmation dynamique

$$V_t := v(Q_t, B_t) = \sup_{b_s, s \in [t, \tau]} \mathbb{E} \left[v(Q_{t+h}, B_{t+h}) - \int_t^{t+h} C(b_s) ds \mid \mathcal{F}_t^W \right]$$

Hamilton Jacobi Bellman (HJB) equation

Equation de HJB dans $[0,1] \times [0, B_{MAX}]$ nous donne :

$$\sup_{b \geq 0} \left\{ \underbrace{\mu Q(1-Q) \frac{\partial V}{\partial Q}(Q, B) + b \frac{\partial V}{\partial B}(Q, B) + \frac{1}{2} \sigma^2 Q^2(1-Q)^2 \frac{\partial^2 V}{\partial Q^2}(Q, B) - C(b)}_{= A(Q, B, b, V)} \right\} = 0$$

Avec les **conditions aux bords** : **Résolution numérique pour trouver b^* optimal:**

$$v(0^+, B) = U(0)$$

$$v(1^-, B) = U(1)$$

$$v(Q, B_{MAX}) = U(Q) \text{ pour } Q \in (0,1)$$

- Différences finies + Algorithme de Howard

Résultats empiriques

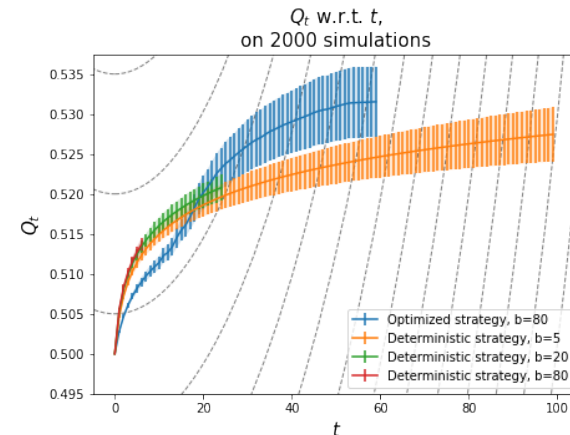
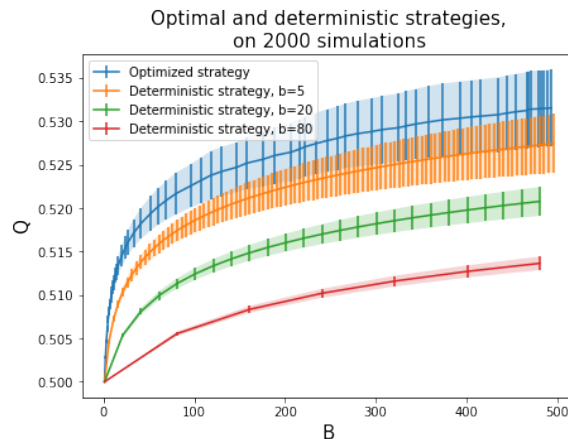
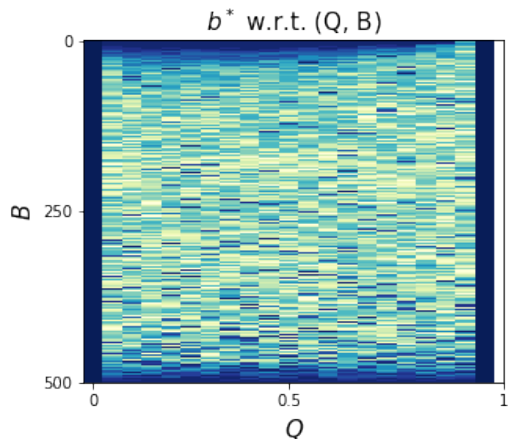


Figure a - Heatmap de b^* p.r. au processus (B, Q)

Figure b et c - stratégie **optimale** vs **déterministe** avec $b \in \{5, 20, 80\}$ et $(B_0, Q_0) = (0, 0.5)$.



La stratégie optimisée (DS-BMAL) montre de meilleurs résultats que les stratégies déterministes en termes de :

- **qualité du modèle** en étiquetant moins au début du processus
- **réduction des délais de réapprentissage.**

En effet, cette stratégie dynamique permet de réduire considérablement le nombre d'itérations.

Défis et contributions

Défi 1 : entraîner un modèle ML avec un budget d'étiquetage limité

Contributions : Des analyses numériques présentent DS-BMAL une stratégie optimale d'étiquetage qui réduit considérablement le nombre d'itérations d'AL tout en gardant une bonne performance du modèle. Cela permettrait d'améliorer les conditions d'étiquetage pour les experts humains.

➔ **Défi 2** : garantir l'équité algorithmique dans les problèmes multi-classes

Défi 2 : équité dans la classification multi-classes



Données : (variables, variable sensible, label-) $\sim \mathbb{P}$ sur $\mathcal{X} \times \mathcal{S} \times [K]$. Considérons $\mathcal{S} = \{+1, -1\}$, **distribution** de S : $(\pi_s)_{s \in \mathcal{S}}$

$\underbrace{\hspace{10em}}_X \quad \underbrace{\hspace{10em}}_S \quad \underbrace{\hspace{10em}}_Y$

Défi 2 : équité dans la classification multi-classes



Données : (variables, variable sensible, label-) $\sim \mathbb{P}$ sur $\mathcal{X} \times \mathcal{S} \times [K]$. Considérons $\mathcal{S} = \{+1, -1\}$, **distribution** de $S : (\pi_s)_{s \in \mathcal{S}}$

Parité démographique (Demographic Parity or DP) : classifier $g \in \mathcal{H}$ est

- **Exactement équitable** si

$$\mathbb{P}(g(X, S) = k | S = 1) = \mathbb{P}(g(X, S) = k | S = -1) \quad \forall k \in [K].$$

- **Approximativement** (ou ε -) **équitable** si pour $\varepsilon \geq 0$,

$$\left| \mathbb{P}(g(X, S) = k | S = 1) - \mathbb{P}(g(X, S) = k | S = -1) \right| \leq \varepsilon \quad \forall k \in [K].$$

Exemple : entretien d'embauche



Défi 2 : équité dans la classification multi-classes



Données : (variables, variable sensible, label-) $\sim \mathbb{P}$ sur $\mathcal{X} \times \mathcal{S} \times [K]$. Considérons $\mathcal{S} = \{+1, -1\}$, **distribution** de $S : (\pi_s)_{s \in \mathcal{S}}$

Parité démographique (Demographic Parity or DP) : classifier $g \in \mathcal{H}$ est

- **Exactement équitable** si

$$\mathbb{P}(g(X, S) = k | S = 1) = \mathbb{P}(g(X, S) = k | S = -1) \quad \forall k \in [K].$$

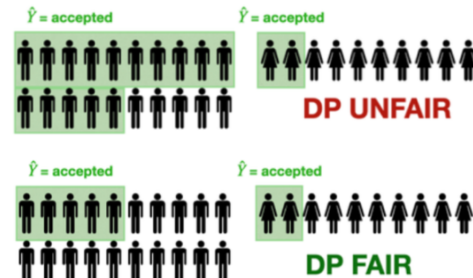
- **Approximativement** (ou ϵ -) **équitable** si pour $\epsilon \geq 0$,



$$\left| \mathbb{P}(g(X, S) = k | S = 1) - \mathbb{P}(g(X, S) = k | S = -1) \right| \leq \epsilon \quad \forall k \in [K].$$

pour cette présentation !

Exemple : entretien d'embauche



Q : où réduire les biais ? **R** : réduire les biais ...

dans les données
(pre-processing)

lors de l'entraînement ML
(in-processing)

après l'entraînement ML
(post-processing)



Dataset



Modèle
d'apprentissage



Prédiction

Défi 2 : équité dans la classification multi-classes



Données : (variables, variable sensible, label-) $\sim \mathbb{P}$ sur $\mathcal{X} \times \mathcal{S} \times [K]$. Considérons $\mathcal{S} = \{+1, -1\}$, **distribution** de $S : (\pi_s)_{s \in \mathcal{S}}$

Parité démographique (Demographic Parity or DP) : classifier $g \in \mathcal{H}$ est

- **Exactement équitable** si

$$\mathbb{P}(g(X, S) = k | S = 1) = \mathbb{P}(g(X, S) = k | S = -1) \quad \forall k \in [K].$$

- **Approximativement** (ou ϵ -) **équitable** si pour $\epsilon \geq 0$,

$$\left| \mathbb{P}(g(X, S) = k | S = 1) - \mathbb{P}(g(X, S) = k | S = -1) \right| \leq \epsilon \quad \forall k \in [K].$$

pour cette présentation !

Exemple : entretien d'embauche



Q : où réduire les biais ? **R** : réduire les biais ...

Contribution :

- La plupart des travaux sur l'équité algorithmique : tâches **binaires** ou de **régression** !
- Cette thèse propose de renforcer l'équité algorithmique en **multi-classes** avec une stratégie **post-processing**



Dataset



Modèle
d'apprentissage



Prédiction

après l'entraînement ML
(post-processing)

Équité **approximative** : notations et objectif

Risque (de mauvaise classification) pour un classifieur $g : \mathcal{X} \times \{-1, 1\} \rightarrow [K]$

$$\mathcal{R}(g) := \mathbb{P}(g(X, S) \neq Y).$$

Scores : pour $k \in [K]$, on note

$$p_k(X, S) := \mathbb{P}(Y = k | X, S).$$

Classifieur de Bayes minimise le risque : $\forall (x, s) \in \mathcal{X} \times \mathcal{S}$

$$g^*(x, s) \in \arg \max_k p_k(x, s).$$

Équité **approximative** : notations et objectif

Risque (de mauvaise classification) pour un classifieur $g : \mathcal{X} \times \{-1, 1\} \rightarrow [K]$

$$\mathcal{R}(g) := \mathbb{P}(g(X, S) \neq Y).$$

Scores : pour $k \in [K]$, on note

$$p_k(X, S) := \mathbb{P}(Y = k | X, S).$$

Classifieur de Bayes minimise le risque : $\forall (x, s) \in \mathcal{X} \times \mathcal{S}$

$$g^*(x, s) \in \arg \max_k p_k(x, s).$$

Objectif

- **Minimiser le risque** : $g^* \in \arg \min_g \mathcal{R}(g)$
- **Garantir l'équité** : pour $\varepsilon \geq 0$, $\max_{k \in [K]} \left| \mathbb{P}(g^*(X, S) = k | S = 1) - \mathbb{P}(g^*(X, S) = k | S = -1) \right| \leq \varepsilon$. (dénoté $g^* \in \mathcal{G}_{\varepsilon\text{-fair}}$)

Considérons son **Lagrangien** et introduisons pour $\lambda^{(1)} = (\lambda_1^{(1)}, \dots, \lambda_K^{(1)}) \in \mathbb{R}_+^K$ et $\lambda^{(2)} = (\lambda_1^{(2)}, \dots, \lambda_K^{(2)}) \in \mathbb{R}_+^K$,

$$\begin{aligned} \mathcal{R}_{\lambda^{(1)}, \lambda^{(2)}}(g) := & \mathcal{R}(g) + \sum_{k=1}^K \lambda_k^{(1)} [\mathbb{P}(g(X, S) = k | S = 1) - \mathbb{P}(g(X, S) = k | S = -1) - \varepsilon] && \text{(risque-équitable)} \\ & + \sum_{k=1}^K \lambda_k^{(2)} [\mathbb{P}(g(X, S) = k | S = -1) - \mathbb{P}(g(X, S) = k | S = 1) - \varepsilon]. \end{aligned}$$

Équité **approximative** : prédicteur équitable optimal

Hypothèse de continuité : $t \mapsto \mathbb{P}(p_k(X, S) - p_j(X, S) \leq t \mid S = s)$ considéré comme **continu** pour tout $k, j \in [K]$ et $s \in \mathcal{S}$.

- **Remarque** : en pratique cette hypothèse est satisfaite en donnant une **petite perturbation uniforme**.

Soit $H : \mathbb{R}_+^{2K} \rightarrow \mathbb{R}$ la fonction $H(\lambda^{(1)}, \lambda^{(2)}) = \sum_{s \in \mathcal{S}} \mathbb{E}_{X|S=s} \left[\max_k \left(\pi_s p_k(X, s) - s(\lambda_k^{(1)} - \lambda_k^{(2)}) \right) \right] + \varepsilon \sum_{k=1}^K (\lambda_k^{(1)} + \lambda_k^{(2)})$.

Proposition

- Sous l'hypothèse de continuité, nous définissons $\lambda^{*(1)}, \lambda^{*(2)} \in \mathbb{R}_+^{2K}$ par

$$(\lambda^{*(1)}, \lambda^{*(2)}) \in \arg \min_{(\lambda^{(1)}, \lambda^{(2)}) \in \mathbb{R}_+^{2K}} H(\lambda^{(1)}, \lambda^{(2)}).$$

Alors, $g_{\varepsilon\text{-fair}}^* \in \arg \min_{g \in \mathcal{G}_{\varepsilon\text{-fair}}} \mathcal{R}(g)$ ssi $g_{\varepsilon\text{-fair}}^* \in \arg \min_{g \in \mathcal{G}} \mathcal{R}_{\lambda^{*(1)}, \lambda^{*(2)}}(g)$.

- De plus $\forall (x, s) \in \mathcal{X} \times \mathcal{S}$,

$$g_{\varepsilon\text{-fair}}^*(x, s) = \arg \max_{k \in [K]} \left(\pi_s p_k(x, s) - s(\lambda_k^{*(1)} - \lambda_k^{*(2)}) \right).$$

Équité **approximative** : prédicteur équitable optimal + estimateur

Hypothèse de continuité : $t \mapsto \mathbb{P}(p_k(X, S) - p_j(X, S) \leq t | S = s)$ considéré comme **continu** pour tout $k, j \in [K]$ et $s \in \mathcal{S}$.

• **Remarque** : en pratique cette hypothèse est satisfaite en donnant une **petite perturbation uniforme**.

Soit $H : \mathbb{R}_+^{2K} \rightarrow \mathbb{R}$ la fonction $H(\lambda^{(1)}, \lambda^{(2)}) = \sum_{s \in \mathcal{S}} \mathbb{E}_{X|S=s} \left[\max_k \left(\pi_s p_k(X, s) - s(\lambda_k^{(1)} - \lambda_k^{(2)}) \right) \right] + \varepsilon \sum_{k=1}^K (\lambda_k^{(1)} + \lambda_k^{(2)})$.

Proposition

• Sous l'hypothèse de continuité, nous définissons $\lambda^{*(1)}, \lambda^{*(2)} \in \mathbb{R}_+^{2K}$ par

$$(\lambda^{*(1)}, \lambda^{*(2)}) \in \arg \min_{(\lambda^{(1)}, \lambda^{(2)}) \in \mathbb{R}_+^{2K}} H(\lambda^{(1)}, \lambda^{(2)}).$$

Alors, $g_{\varepsilon\text{-fair}}^* \in \arg \min_{g \in \mathcal{G}_{\varepsilon\text{-fair}}} \mathcal{R}(g)$ ssi $g_{\varepsilon\text{-fair}}^* \in \arg \min_{g \in \mathcal{G}} \mathcal{R}_{\lambda^{*(1)}, \lambda^{*(2)}}(g)$.

• De plus $\forall (x, s) \in \mathcal{X} \times \mathcal{S}$,

$$g_{\varepsilon\text{-fair}}^*(x, s) = \arg \max_{k \in [K]} \left(\pi_s p_k(x, s) - s(\lambda_k^{*(1)} - \lambda_k^{*(2)}) \right).$$

Approche semi-supervisée :



• **Données étiquetées** : $\mathcal{D}_n = (X_i, S_i, Y_i)_{i=1, \dots, n}$
Entraîner les estimateurs $(\hat{p}_k)_k$ (ex. RF, SVM, ...)

Randomization trick $\bar{p}_k(X, S, \zeta_k) = \hat{p}_k(X, S) + \zeta_k$



• **Données non-étiquetées** : $X_1^s, \dots, X_{N_s}^s \stackrel{\text{iid}}{\sim} \mathbb{P}_{X|S=s}$

Fréquences empiriques $(\hat{\pi}_s)_s$ comme estimateurs de $(\pi_s)_s$

Estimateur post-processing :

Si $(\hat{\lambda}^{(1)}, \hat{\lambda}^{(2)}) \in \arg \min_{(\lambda^{(1)}, \lambda^{(2)}) \in \mathbb{R}_+^{2K}} \hat{H}(\lambda^{(1)}, \lambda^{(2)})$

$$\hat{g}_\varepsilon(x, s) = \arg \max_{k \in [K]} \left(\hat{\pi}_s \bar{p}_k(x, s, \zeta_k) - s(\hat{\lambda}_k^{(1)} - \hat{\lambda}_k^{(2)}) \right)$$

Équité **approximative** : garanties théoriques

Mesure d'iniquité (Unfairness measure). $\mathcal{U}(g) := \max_{k \in [K]} \left| \mathbb{P}(g(X, S) = k | S = 1) - \mathbb{P}(g(X, S) = k | S = -1) \right|$

Théorème

Garantie universelle d'équité. Il existe une constante $C > 0$ dépendant uniquement de K et de $\min_{s \in S} \pi_s$, tel que, pour tout estimateur \hat{p}_k

$$|\mathbb{E}[\mathcal{U}(\hat{g}_\varepsilon)] - \varepsilon| \leq \frac{C}{\sqrt{N}}$$

Consistance. Si l'estimateur est consistant avec la norme L_1 alors,

$$\mathbb{E}[\mathcal{R}(\hat{g}_\varepsilon)] \rightarrow \mathcal{R}(g_{\varepsilon\text{-fair}}^*) \text{ quand } n, N \rightarrow \infty$$



$n = \#$ données étiquetées



\hat{g}_ε asymptotiquement aussi performant que g^* en termes **d'équité** et de **précision**



$N = \#$ données non-étiquetées

Résultats sur des données synthétiques



Dataset :

données **synthétiques**, pour $k \in [K]$

$(X | Y = k) \sim$ mélange Gaussien

$(S | Y = k) \sim 2 \cdot \mathcal{B}(p) - 1$, si $k \leq \lfloor K/2 \rfloor$,

$(S | Y = k) \sim 2 \cdot \mathcal{B}(1-p) - 1$, si $k > \lfloor K/2 \rfloor$.

- p mesure le biais historique dans le dataset
- K représente le nombre de classes

$$\max_{k \in [K]} \left| \mathbb{P}(Y = k | S = 1) - \mathbb{P}(Y = k | S = -1) \right| = \left| \frac{2}{K} \cdot (2p - 1) \right|.$$

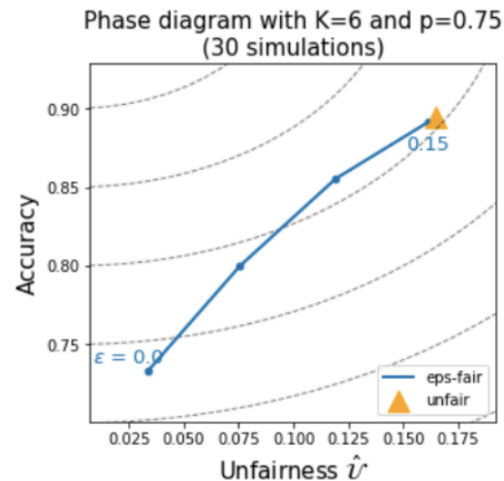
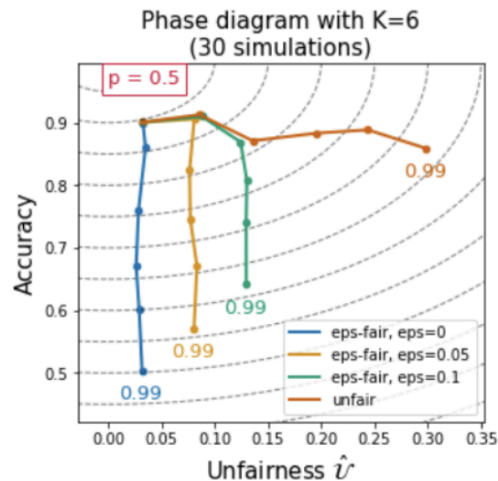


Modèle d'apprentissage :

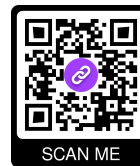
Random Forest (RF)



l'équité de g est mesurée par la version empirique de la mesure d'iniquité $\mathcal{U}(g)$



[Lien vers code/expérimentations \(Github\)](#)



Résultats sur des données réelles : cas binaire (1/2)



Dataset : DRUG, CRIME



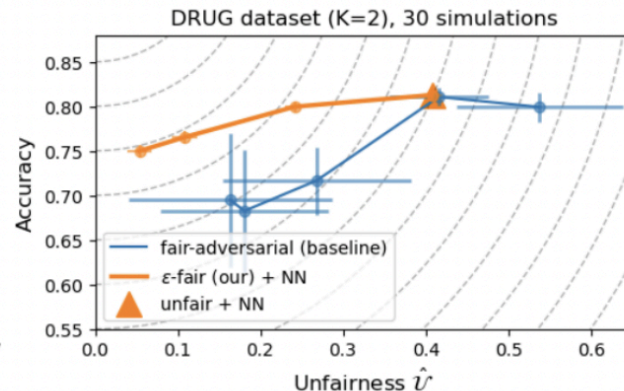
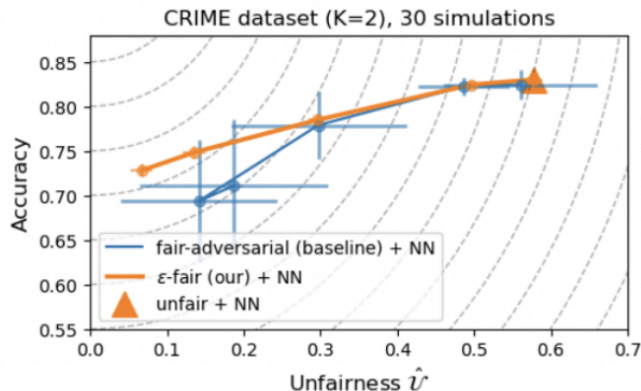
Modèle d'apprentissage :

Réseaux de Neurones (NN)



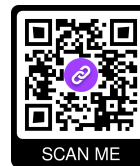
Benchmarks :

Fair-adversarial¹



[Lien vers code/expérimentations \(Github\)](#)

1: approche **in-processing** <https://aif360.readthedocs.io/en/stable/modules/generated/aif360.sklearn.inprocessing.AdversarialDebiasing.html>



Résultats sur des données réelles : cas binaire (2/2)



Dataset : DRUG, CRIME



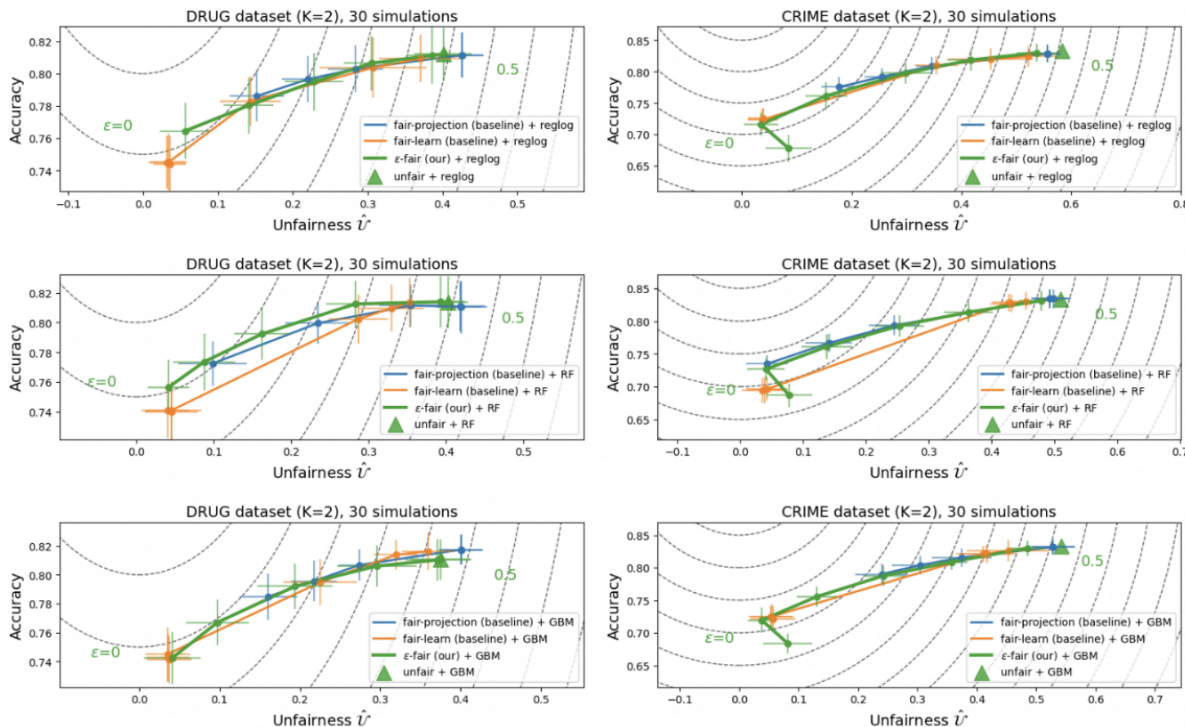
Modèle d'apprentissage :

Régression Logistique (**reglog**),
Random Forest (**RF**),
XGBoost (**GBM**)

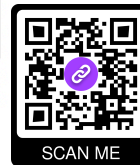


Benchmarks :

Fair-learn¹
Fair-projection²



1: approche **in-processing** pour binaire et régression <https://fairlearn.org/> (Agarwal et al., ICLR 2018)
2: approche **post-processing** pour multi-classes <https://github.com/HsiangHsu/Fair-Projection> (NeurIPS 2022)



Résultats sur des données réelles : cas multi-classes



Dataset : DRUG, CRIME

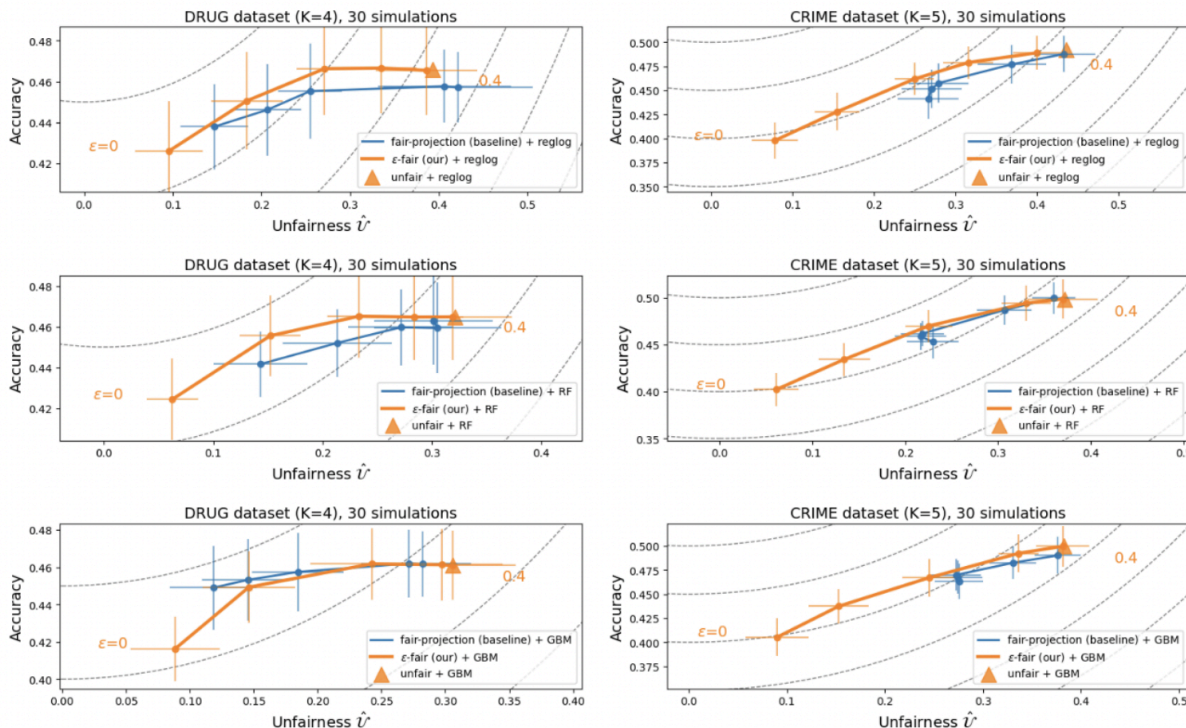


Modèle d'apprentissage :

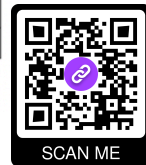
Régression Logistique (**reglog**),
Random Forest (**RF**),
XGBoost (**GBM**)



Benchmarks :
Fair-projection¹



1: approche **post-processing** pour multi-classes <https://github.com/HsiangHsu/Fair-Projection> (NeurIPS 2022)



Défis et contributions

Défi 1 : entraîner un modèle ML avec un budget d'étiquetage limité

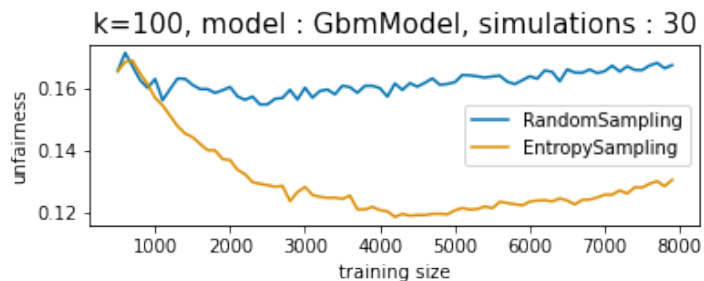
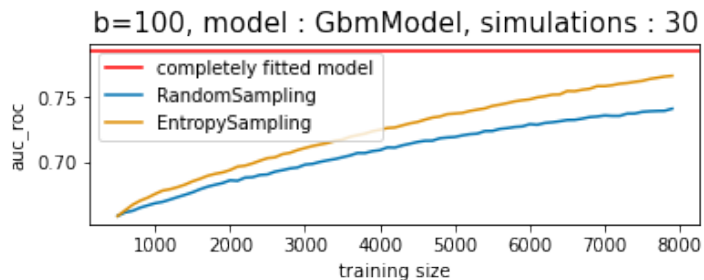
Contributions : Des analyses numériques présentent DS-BMAL une stratégie optimale d'étiquetage qui réduit considérablement le nombre d'itérations d'AL tout en gardant une bonne performance du modèle. Cela permettrait d'améliorer les conditions d'étiquetage pour les experts humains.

Défi 2 : garantir l'équité algorithmique dans les problèmes multi-classes

Contributions : Dans le cadre de la classification multi-classes, nous fournissons une règle de classification équitable optimale sous la contrainte DP. Nous traitons l'équité exacte et approximative et montrons que notre approche obtient des résultats remarquables sur divers ensembles de données synthétiques et réels. En particulier, notre algorithme est efficace pour faire respecter un niveau d'équité pré-spécifié.

Conclusion

Défi 1+2 : *fair active classifier* sur données d'assurance



Perspective

- Extension en **classification multi-labels** :

Au lieu de considérer **le risque de mauvaise classification**

$$R(g) = \mathbb{P}(g(X, S) \neq Y)$$

Nous considérons **le risque L_2**

$$R_2(g) = \mathbb{E} \left[\sum_{k=1}^K (1_{Y=k} - g_k(X, S))^2 \right]$$



Bientôt un **package** sur Python pour l'équité algorithmique en multi-classes

Références



INSTITUT
POLYTECHNIQUE
DE PARIS



ENSAE



IP PARIS

**Semi-supervised learning in insurance:
fairness and active learning**

Thèse de doctorat de l'Institut Polytechnique de Paris
préparée à l'École nationale de la statistique et de l'administration économique
École doctorale n° 574 École Doctorale de Mathématiques Hadamard (EDMH)
Spécialité de doctorat : Mathématiques appliquées

Thèse présentée et soutenue à Palaiseau, le 15 juin 2022, par

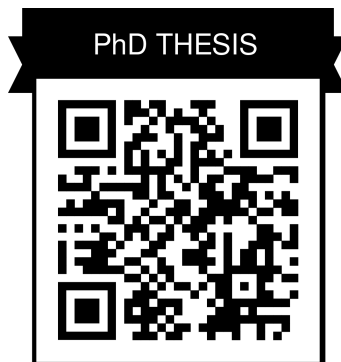
FRANÇOIS HU

Composition du Jury :

Christian-Yann ROBERT Professeur, École nationale de la statistique et de l'administration économique (CREST)	Président
Christophe DUTANG Maître de conférences, Université Paris Dauphine (GERSMADIS)	Rapporteur
Olivier WINTENBERGER Professeur, Université Pierre et Marie Curie (LPSM)	Rapporteur
Arthur CHAMPENTIER Professeur, Université du Québec à Montréal (QUANTACT)	Examineur
Stéphane LOISEL Professeur, Université Lyon 1 (SAP)	Examineur
Caroline HILLAIRET Professeure, École nationale de la statistique et de l'administration économique (CREST)	Directrice de thèse
Romuald ELIE Professeur, Université Gustave Eiffel (LAMA)	Directeur de thèse (Invité)
Marc JULLIARD Directeur du DataLab, Société Générale Assurances	Invité

Thèse de doctorat

NNT : 2022IPPAG002



Merci pour votre attention

