

Theoretical Guidelines

A framework for Multi-A(rmed)/B(andid) testing with
online FDR control

Yanis Daci, Clément Hardy

February 2019

Contents

1	Introduction	3
2	Le méta-algorithme	4
3	Etude des différentes procédures	5
3.1	Online FDR control	5
3.2	Multi Armed Bandit	5
3.3	Hypothèse et p-valeur	6
3.4	Algorithme MAB modifié	7
4	MAB avec <i>online FDR control</i>	9
5	Implémentation et résultats	11
5.1	Best-arm identification	11
5.2	MAB-LORD: best-arm identification with online FDR control	12
6	Critiques et pistes de réflexions	15
6.1	Critiques	15
6.2	Pistes de reflexion	16
7	Conclusion	18

1 Introduction

Les fausses découvertes dans les papiers de recherche sont devenues un sujet de préoccupation. L'article "Why Most Published Research Findings are false" nous informe que certains domaines tels que la médecine ou la psychologie sont plus touchés que d'autres.

Ces fausses découvertes ont, selon cette étude, quelques points en communs : la plupart sont basées sur des observations expérimentales et ont une légitimité statistique faible. En effet, les expériences et les études faites pour obtenir les résultats n'ont pas été répliquées suffisamment de fois.

Elles peuvent être également dûes à d'autres problématiques. Statistiquement parlant, tester un grand nombre d'hypothèses mène toujours à en accepter quelques-unes ; un biais dans l'étude peut contribuer à en accepter une fautive. Ces problématiques peuvent être de différentes natures :

- des échantillons trop petits, de mauvaise qualité
- un trop grand nombre d'hypothèses testées
- conflit d'intérêt (une entreprise peut avoir un grand intérêt à accepter une étude)
- une trop faible rigueur dans les expériences
- ...

De même, des hypothèses correctes peuvent être réfutées. Ces erreurs sont communément appelées erreurs de type I et II.

Les fausses découvertes sont difficilement évitables et ne peuvent être que réduites. L'objectif est donc d'étudier le false discovery rate afin de le contrôler. Cela revient à chercher un seuil d'acceptation qui minimise les faux positifs tout en gardant un nombre de vrais positifs correcte.

Le false discovery rate n'est autre que l'espérance de faux positif sur le nombre total de positifs. Ainsi, en utilisant le tableau ci-dessous, le FDR est donné par : $FDR = E\left(\frac{A}{A+B}\right)$

	Hypothèse vraie	Hypothèse fautive
Test significatif	A	B
Test non significatif	C	D

Le FDR est un critère très utilisé lorsque de nombreuses hypothèses doivent être testées. Cependant, dans de nombreux domaines (marketing, publicité...) , les hypothèses ne sont pas disponibles en même temps et arrivent sous forme de stream. Dans ce cas, l'acceptation de la nouvelle hypothèse est déterminée sans avoir connaissance des prochaines. Cette dernière est donc basée uniquement sur les précédentes décisions. Il s'agit donc de contrôler le FDR pour des tests successifs, en ligne (online FDR control).

Statistiquement, une hypothèse ayant une p-value inférieure à une certaine valeur (couramment notée α) est considérée significative. Généralement, la valeur choisie est 0.05.

2 Le méta-algorithme

L'article propose un méta-algorithme afin d'éventuellement déterminer pour chaque itération un meilleur bras que celui de contrôle.

Il s'articule de la manière suivante :

1. L'utilisateur fixe un taux de contrôle du FDR α qu'il désire.
2. Pour chaque $j= 1,2,\dots$:
 - Un bras de contrôle et un certain nombre de bras alternatifs sont attribués à l'expérience j .
 - Une procédure *online-FDR* renvoie un α_j qui est fonction des précédentes valeurs $\{P^l\}_{l=1}^{j-1}$.
 - Une procédure MAB est exécutée et retourne un bras recommandé si celle ci se termine automatiquement. Elle prend en entrées le bras de contrôle et $K(j)$ bras alternatifs, un niveau de confiance α_j , et optionnellement un précision $\epsilon \geq 0$.
 - Au cours de la procédure MAB, une p-valeur toujours valide est construite continuellement pour chaque instant t en utilisant uniquement les échantillons collectés jusqu'à maintenant et qui proviennent de l'expérience j .
 - Quand la procédure MAB se termine à l'instant t , que ce soit automatiquement ou par un critère d'arrêt définie par l'utilisateur, si le bras avec la moyenne empirique la plus élevée n'est pas le bras de contrôle et $P_t^j \leq \alpha_j$, le bras de contrôle est alors rejeté au détriment de ce dernier.

A présent, nous allons étudier l'ensemble des procédures utilisées dans cet algorithme.

3 Etude des différentes procédures

3.1 Online FDR control

Dans le cas d'hypothèses arrivant successivement, l'une des méthodes pour contrôler le FDR est d'avoir une succession de niveau significatif $\alpha_1, \dots, \alpha_i$. Chaque hypothèse possède sa propre p-value (p_i). L'hypothèse H_i est alors rejetée si sa p-value est inférieure au niveau significatif α_i .

Nous pouvons noter ces rejets R_i :

$$R_i = \begin{cases} 1 & \text{si } p_i < \alpha_i \\ 0 & \text{sinon} \end{cases}$$

Cependant, nous souhaitons que le false discovery rate soit inférieur à un certain seuil fixé α . Ainsi, le seuil α_i doit dépendre des précédents tests. Il est donc résultat d'une fonction impliquant les précédents $R_j, 0 \leq j < i$.

3.2 Multi Armed Bandit

Nous commencerons cette partie par une brève explication de l'A/B testing. Puis, nous étudierons le multi armed bandit afin de dégager son intérêt par rapport à l'A/B testing.

L'A/B testing est une méthode permettant de décider quelle variante d'un objet est la meilleure. Il peut être utilisé dans le cadre d'un site web, d'un traitement médical... Après un certain temps, les variantes sont mises en compétition pour déterminer la meilleure.

Prenons la création d'un vaccin comme exemple. Durant la phase de test, chaque variante du vaccin est injectée à la même quantité de patient. A la fin de cette phase, les résultats sont statistiquement étudiés et celle qui a obtenu les résultats les plus probants est conservée.

Le principal inconvénient d'une telle méthode est que la variante ayant obtenu les meilleurs résultats n'est pas forcément la meilleure. En effet, le hasard ou un biais a pu avoir une importance durant la période de test.

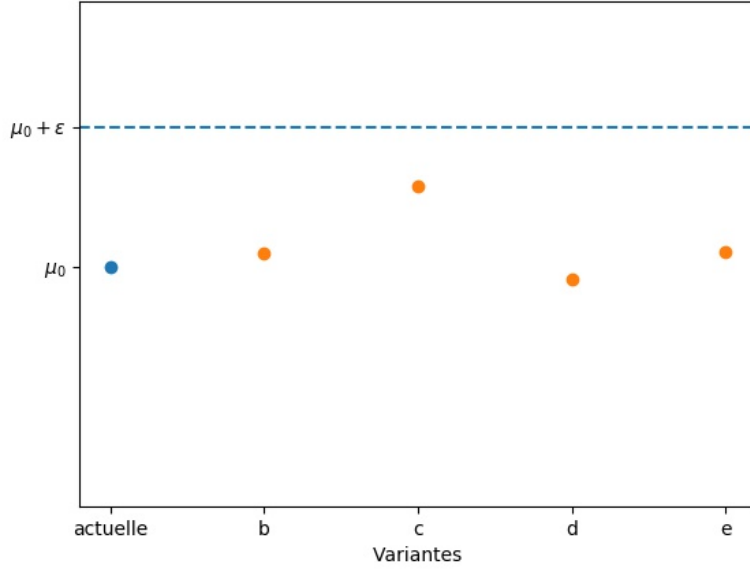
Le Multi Armed bandit ne présente pas cet inconvénient. En effet, à défaut d'avoir une période de test conduisant à la conservation brutale d'une unique variante, le multi armed bandit détermine progressivement la meilleure.

Dans le cas d'un site web, le trafic est de plus en plus redirigé vers la meilleure version sans pour autant supprimer les autres qui continuent d'attirer un petit panel.

Cependant, la mise en application d'une telle méthode en online FDR control peut poser problème. En effet, pour certains acteurs comme une entreprise ou un laboratoire, la mise en application de la variante considérée comme la meilleure peut engendrer différents coûts. (coût financier, humain...)

Ainsi, si le gain de la nouvelle variante est trop faible par rapport à l'ancienne, le gain total peut être négatif.

Illustrons ce fait par un exemple :



Sur le graphique ci dessus, bien que certaines variantes (en orange) soient supérieures à l'actuelle, nous pouvons remarquer qu'elles restent en dessous de la ligne formée par ϵ . Celle ci peut être considérée comme la ligne de rentabilité. Il n'y a donc pas de raisons valables pour l'acteur en question de changer de variante.

Il y a donc un intérêt à privilégier la variante actuelle par rapport aux autres. Mais cet aspect pose un problème pour déterminer une hypothèse de test ainsi que des p-valeurs. Dans ce contexte, les p-valeurs (du MAB avec online FDR control) ne sont plus indépendantes entre elles car le choix entre les variantes n'est plus uniforme.

L'article de recherche "A framework for Multi-A/B testing with online FDR control" tente d'y remédier.

3.3 Hypothèse et p-valeur

L'hypothèse exposée ci-dessus peut être formulée de la sorte:

$$H_o : \mu_0 + \epsilon \geq \mu_i, \forall i = 1, \dots, K$$

avec μ_0 la moyenne des récompenses de la variante actuelle et μ_i les moyennes des alternatives.

Soit ϵ la perte engendrée par la mise en place d'une alternative.

On cherche donc une variante ayant la meilleure espérance de gain tout en minimisant la perte sèche engendrée par sa mise en place.

Maintenant qu'une hypothèse de test a été déterminée, nous pouvons nous focaliser sur

les p-valeurs.

Tout d'abord, une p-valeur doit avoir une distribution uniforme sous l'hypothèse nulle. En effet si on note F la fonction de distribution, on a:

$$\begin{aligned} P_{\theta_0}(p_{value} < \alpha) &= P_{\theta_0}(F(T) < \alpha) \\ &= P_{\theta_0}(T < F^{-1}(\alpha)) \\ &= F(F^{-1}(T)) \\ &= \alpha \end{aligned}$$

Dans le cas du multi armed bandit, le nombre "d'échantillons" dans chaque variante n'est pas identique.

Ainsi, la moyenne estimée pour chaque variante est dépendante des autres. A titre d'exemple, le trafic d'un site web dépend des résultats précédents.

Pour cette raison, une suite de p-valeur doit être définie.

Elles sont définies par:

$$P_t = \min_s \min_{i=1, \dots, K} P_{i,s}$$

où $P_{i,s}$ peut être exprimée par:

$$P_{i,s} = \sup \{ \gamma \in [0, 1] \mid LCB_i(t) \leq UCB_0(t) + \epsilon \}$$

3.4 Algorithme MAB modifié

Les hypothèses nulles privilégient la variante actuelle par le biais de ϵ . L'algorithme LUCB ne peut donc pas fonctionner. De légères modifications sont nécessaires.

Commençons par présenter cet algorithme:

Tout d'abord, il échantillonne l'historique de chaque variante.

Puis, il estime une condition d'arrêt. Si la condition est vérifiée, l'algorithme s'arrête et retourne la variante avec la plus grande moyenne.

Sinon, il se poursuit en "ajoutant" un échantillon dans une variante. La sélection de cette variante suit une règle précise.

Nous désirons déterminer la variante avec la plus grande moyenne. Aini, on définit $m = 1$. Notons $\hat{\mu}_i$ ($i = 1, \dots, K$, avec K le nombre de variante) l'estimateur de la moyenne de la i ième variante.

Pour chaque moyenne de ces variantes, nous pouvons calculer un intervalle de confiance.

Notons I_i^- (resp I_i^+) la borne inférieur (resp borne supérieur) de la i ième moyenne.

Notons a la variante ayant la plus faible estimation de la moyenne, b celle ayant l'estimation la plus élevée et I_a^+ la borne supérieure de a , I_b^- la borne inférieure de b .

L'algorithme s'arrête lorsque $(\mu_a - I_a^+) - (\mu_b - I_b^-) < \epsilon$.

C'est à dire, lorsque la borne supérieure de la plus petite moyenne estimée et la borne inférieure de la plus grande s'entrecroisent sur une portion inférieure à ϵ .

Ici, la différence vient du fait que l'hypothèse de base est

$$H_o : \mu_0 + \epsilon \geq \mu_i, \forall i = 1, \dots, K$$

et non

$$H_o : \mu_0 \geq \mu_i, \forall i = 1, \dots, K$$

Cette hypothèse modifie le critère d'arrêt comme on peut l'observer dans l'étape 2 (a). En effet, les conditions d'arrêt imposent que la variante choisie ait une moyenne supérieure à celle courante (via $LCB_{h_t}(t) > UCB_0(t) + \epsilon$) et que la différence entre la borne supérieure la plus élevée et la borne inférieure de la moyenne la plus élevée soit inférieure à ϵ .

4 MAB avec *online FDR control*

Maintenant que nous avons présenté l'algorithme MAB modifié, nous allons nous intéresser à la manière dont le control online du FDR est intégré.

Pour cela, nous devons définir le false ratio discovery. Précédemment, nous avons mentionné que le false discovery rate pouvait s'exprimer comme l'espérance du nombre de fausse découverte divisé par le nombre total d'hypothèse rejeté.

Dans l'article, la variable R_i est utilisée. Nous l'avons défini dans la partie 1.1. En effet, celle ci indique si l'hypothèse nulle a été rejeté. Nous pouvons alors l'utiliser pour exprimer le false discovery rate.

Le nombre de fausse découverte peut s'exprimer en sommant les variables R_i pour lesquelles l'hypothèse H_0 est vérifiée. Ainsi, le false discovery rate s'écrit:

$$FDR = E \left(\frac{\sum_{j \in H_0} R_j}{\sum_{i=1} R_i} \right)$$

Maintenant, intéressons nous à l'algorithme utilisé. Pour cela, commençons par présenter brièvement "generalized alpha investing".

Le generalized alpha investing est utilisé pour générer une séquence de variable binaire à partir d'une séquence de p-valeur. Chaque variable vaut 1 si l'hypothèse nulle (H_j) est rejetée, 0 sinon.

C'est à dire:

$$R_j = \begin{cases} 1 & \text{si } p_j \leq \alpha_j \\ 0 & \text{sinon} \end{cases}$$

Cette variable indique donc la décision du test j .

Or, α_j est une fonction de $(R_1, R_2, \dots, R_{j-1})$. Pour connaître sa valeur, nous avons besoin de mettre à jour plusieurs autres fonctions.

Pour commencer, la fonction de potentiel W qui décroît dans le cas d'une acceptation et croît sinon. Dans le cas du "generalized alpha investing" la mise à jour se fait comme suit:

$$\begin{aligned} W(0) &= w_0 \geq 0 \\ W(j) &= W(j-1) - \varphi(R_1^{j-1}) + R_j \psi(R_1^{j-1}) \end{aligned}$$

avec ψ_j et φ des fonctions de (R_{j-1}) .

Mais dans le cas du "MAB with online FDR control", la mise à jour ne suit pas la même méthode. En effet, l'augmentation de la fonction potentiel est de $\alpha - W(0)$ en cas de rejet; où α est le niveau d'erreur fixé.

De plus, on utilise cette fois ci une séquence décroissante monotone $\gamma_i \geq 0$ tel que $\sum_{i=1}^{\infty} \gamma_i = 1$. Elle sert à connaitre le "reste" du niveau d'erreur.

Ainsi, la mise a jour de W s'écrit:

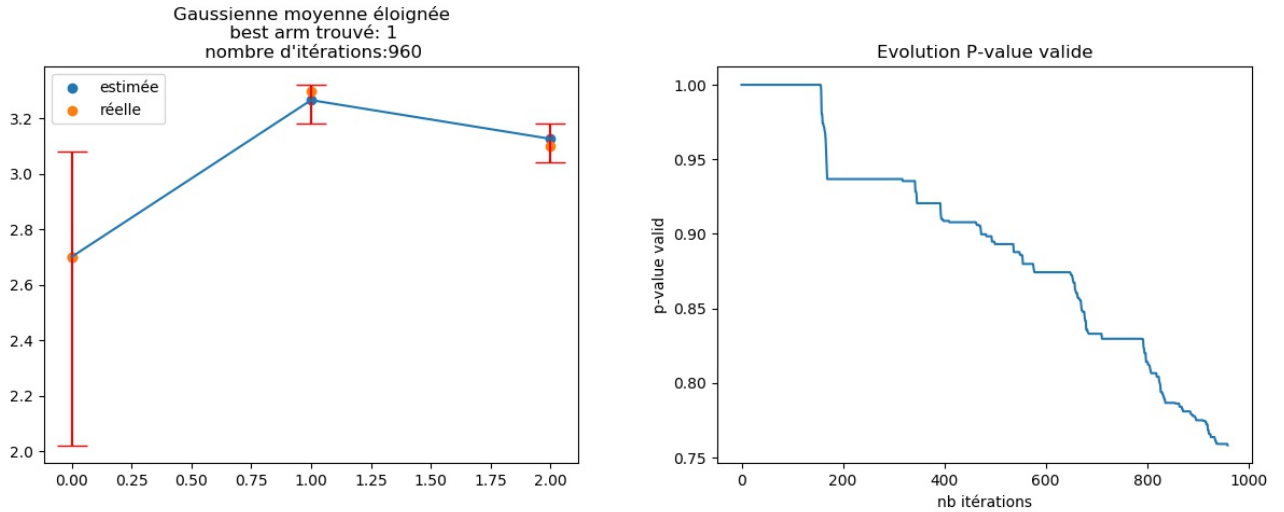
$$W(j+1) = W(j) - \alpha_j + R_j(\alpha - W(0))$$

5 Implémentation et résultats

Les résultats ci dessous ont été obtenu après exécution du fichier **BestArm_MABLord.py** contenu dans le fichier zip.

5.1 Best-arm identification

Les simulations ont été réalisées en utilisant des gaussiennes avec des moyennes différentes pour chacune d'elles mais une variance identique. Pour commencer, regardons les résultats fournis par l'algorithme dans le cas où ϵ vaut 0.



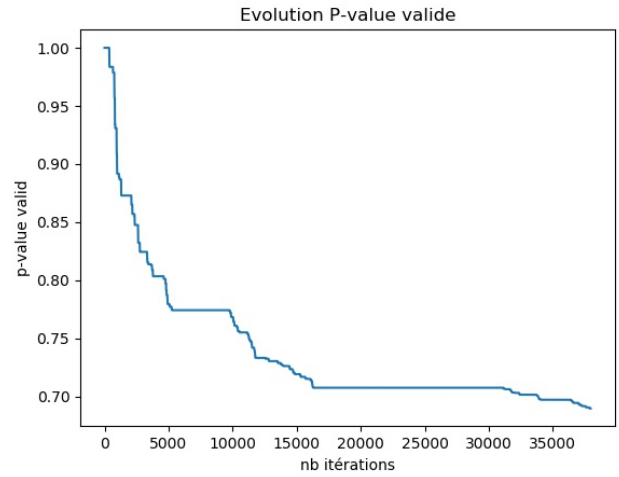
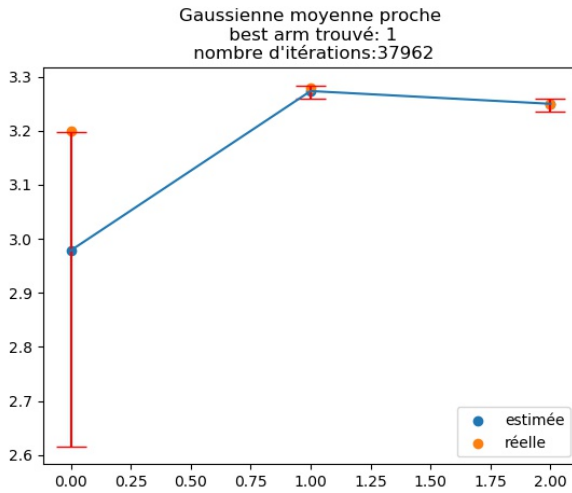
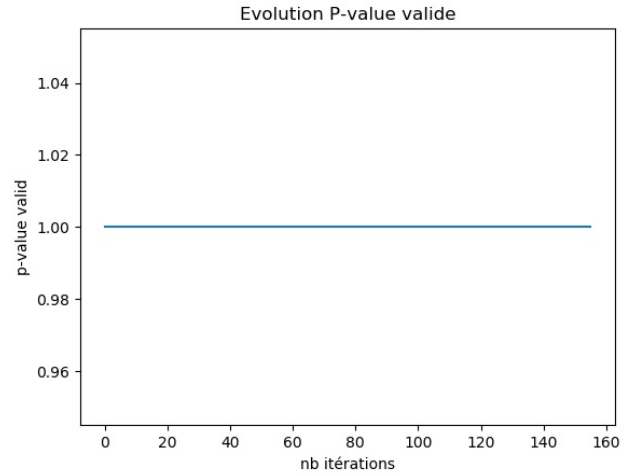
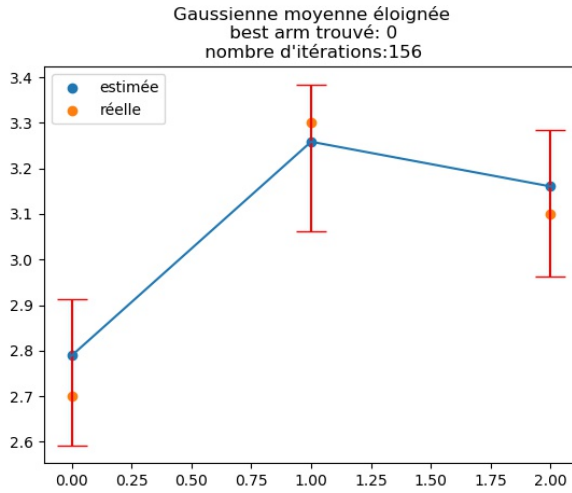
Tout d'abord, nous pouvons remarquer que la meilleure variante trouvée est la première ce qui est cohérent par rapport aux moyennes prises. De même, nous pouvons remarquer que l'intervalle de confiance (à 95%) de la variante 0 ne chevauche plus celui de la meilleure variante. La p-value valide est quant à elle décroissante au fur et à mesure des itérations. En effet, la p-value valide est le min des p-value. Ainsi, comme le nombre de p-value croît, il est logique que la p-value valide décroisse.

Pour ϵ valant 0.8, nous obtenons les résultats suivant:

Cette fois ci, le coût à payer pour changer de variante (ϵ) est supérieur au gain possible. L'algorithme ne se trompe pas et annonce la variante courante comme meilleure variante. Un fait intéressant à observer est que la p-value valide est constante, elle vaut 1. Les p-value qui sont le sup d'un ensemble vide sont fixées à 1. Ainsi, la p-value valide vaut elle aussi 1.

Les résultats suivant sont obtenus en modifiant les moyennes des gaussiennes et en les prenant très proches les unes des autres.

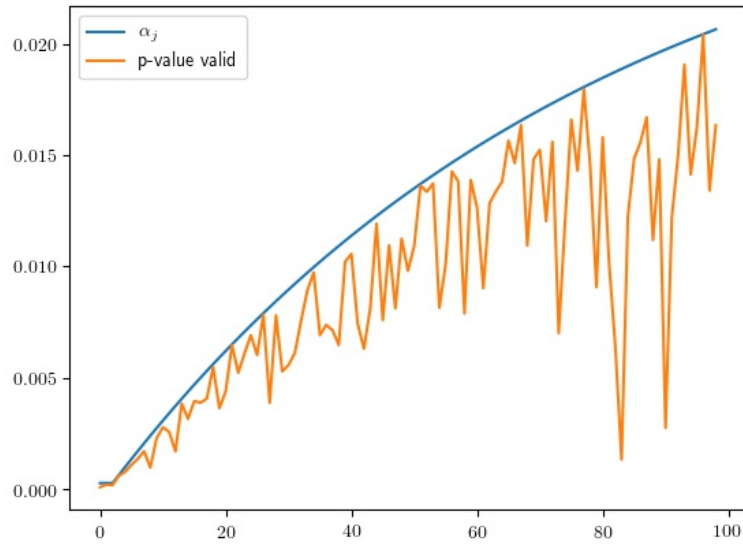
Remarquons comme attendu que l'algorithme a besoin de beaucoup plus de pull sur certaines variantes pour trouver l'optimale. Notons que l'estimation de la moyenne de la variante courante n'est pas très bonne. Cela est dû au nombre trop faible de pull initial qui n'a probablement pas permis d'avoir un meilleur résultat. De plus, comme la moyenne a été estimée à la baisse, l'algorithme n'a pas trouvé nécessaire de continuer à pull sur cette



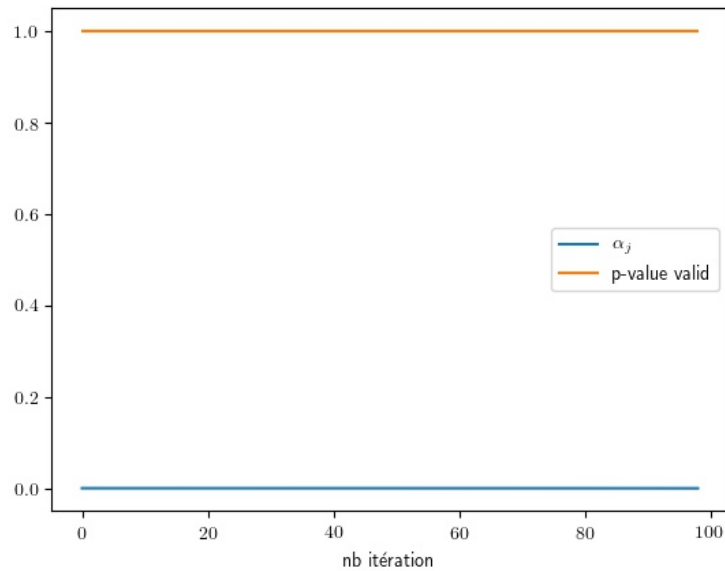
variante (pensant que sa moyenne était beaucoup plus basse) et c'est ainsi concentré sur les deux autres.

5.2 MAB-LORD: best-arm identification with online FDR control

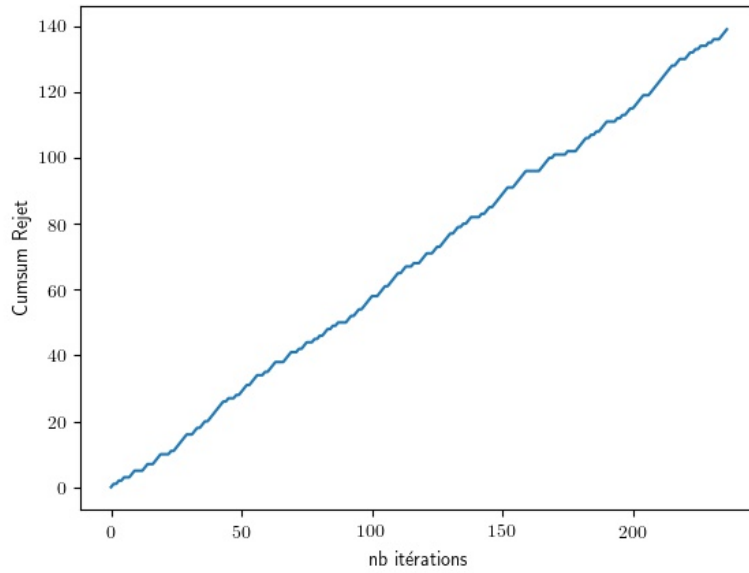
Dans le cas où une variante possède une moyenne plus élevée que la variante courante, l'hypothèse nulle doit être rejetée. Nous pouvons l'observer sur le graphique ci dessous: Les p-values restent en dessous des α_j .



A l'inverse, lorsque la variante courante est la meilleure, l'hypothèse est toujours acceptée comme visible sur le graphique suivante. La p-value vaut toujours 1 et est au dessus des α_j .



Si nous prenons le cas où les moyennes des variantes sont relativement proches, nous pouvons remarquer que tour à tour l'algorithme rejette l'hypothèse ou la conserve. Ce phénomène est observable sur le graphique suivant étant donné que la progression de la somme des R_i est certes à peu près linéaire mais inférieure à la droite d'équation $f(x) = x$.



Le temps d'arrêt T_j semble également avoir une grande importance dans ce cas. Un temps d'arrêt trop court menera à l'acceptation systématique de l'hypothèse même si cela est faux. Il doit donc être suffisamment élevé.

6 Critiques et pistes de réflexions

6.1 Critiques

L'article de recherche "A framework for Multi-A(rmed)/B(andid) testing with online FDR control" propose des arguments très intéressants mais comporte tout de même certains points critiquables.

Tout d'abord, une grande partie de l'analyse est basée sur le FDR et son optimisation. Cependant, cette variable ne prend en considération uniquement le nombre d'hypothèses vraies ou fausses qui ont été acceptées. Nous pouvons déjà dégager certaines limites quand à l'utilité de cette variable pour déterminer l'exactitude d'une découverte.

Le premier constat est qu'elle ne prend aucunement en considération les bonnes hypothèses qui ont été rejetées durant l'étude. Si ce n'était qu'une hypothèse majeure et véridique a été rejetée au cours de l'analyse, nous pouvons d'ores et déjà remettre en cause sa véridicité. Ce premier facteur met en lumière la nécessité de créer un système de "poids" d'hypothèses qui n'a pas été implémenté tout au long de l'article.

En effet, nous pouvons prendre comme exemple le nombre de fausses hypothèses acceptées. Dans le FDR, cette argument a autant d'importance que le nombre de vrai hypothèses acceptées. Cependant, la présence d'une seule fausse hypothèse dans une découverte suffit à la décrédibiliser totalement. A l'inverse, il est tout à fait normal de compter des hypothèses correctes lors d'une étude. Le FDR ne prend pas du tout ce facteur en considération et attribue la même importance à ces arguments.

Une solution serait donc de créer une nouvelle variable qui décrédibiliserait complètement une découverte avec ce qu'une seule fausse hypothèse sans oublier de pénaliser le rejet d'une vraie.

Puis, tout au long de l'article, la finalité est de déterminer si une découverte est véridique. Le résultat est donc binaire : soit une découverte est considérée comme vraie et est donc conservée ou alors elle est rejetée. Mais ce n'est pas aussi simple. En effet, une découverte peut être jugée biaisée sur certains critères mais ne pas l'être selon d'autre. Comme démontré dans le précédent paragraphe, le choix de la valeur d'évaluation importe énormément. Il faudrait donc évaluer à quel degrés une découverte est catégorisée d'inexacte et conserver celles qui ont le plus faible afin de les évaluer par d'autres méthodes. Bien entendu, une découverte qui utilise une ou des hypothèses absurdes auront un fort degrés de falsification et ne fera pas l'objet d'autres analyses.

Ensuite, dans aucune partie de l'article n'est introduit la notion du poids ou de l'importance d'une possible découverte. Il semble primordial de prendre en considération cet élément car il peut totalement altérer l'approche pour juger celle-ci.

Une découverte ayant un gros poids peut être composée de "sous découvertes". Prenons l'exemple d'un laboratoire qui cherche un remède contre le cancer. Dans ce cas précis, une

découverte sera jugée uniquement sur sa finalité et donc sur la production d'un possible remède. Tout autre résultat sera jugé de façon négative.

Cependant, au vue de l'importance et de la spécificité de l'objectif, on ne peut adopter la même approche que pour une découverte moins impactante. Lorsque le poids de la découverte est énorme, des critères plus spécifique doivent être adoptés.

Dans ces conditions, des "sous découvertes" pourraient être mises en valeur. Si un laboratoire n'aboutit pas à la création d'un remède, le travail effectué peut tout de même amener une avancée et ne peut pas être classé en tant que fausse découverte. Ce facteur n'est pas pris en consideration dans l'article. En cas de poids important, il faudrait effectuer une analyse plus précise sur chaque hypothèse pour eventuellement en conserver certaines.

Enfin, chaque découverte a un objectif unique avec son propre secteur d'application. Ainsi, il semble difficile de mettre en place une méthode générale tout comme le propose l'article. Par exemple, si une équipe travaille sur la détection de cancer chez des patients, il est moins grave de prédire qu'un patient est malade alors qu'il ne l'est pas que d'assurer qu'une personne est en bonne santé alors que c'est tout le contraire.

Chaque problème dispose donc de ces propres caracteristiques. Les méthodes de l'article ne les prennent pas en compte et peuvent donc être soumises à interrogation.

6.2 Pistes de reflexion

Dans cette sous partie, nous énumérerons les differentes interrogations rencontrées pendant la lecture de l'article et nous essayerons d'apporter des premiers éléments de réponse.

Tout d'abord, nous nous sommes demandés si certaines hypothèses ou découvertes pouvaient être rejetées du fait de l'impossibilité technique à les vérifier ou à les produire. Si le niveau technologique actuel ne permet pas de réaliser une expérience, doit on forcement la rejeter ? Demain, il sera peut être possible de la pratiquer. Avoir déclaré la découverte fausse apparait alors comme une erreur.

Il faudrait alors bien évaluer son potentiel pour déterminer si cette dernière n'a pas été jugée précipitamment.

Il faudrait alors opter pour un algorithme qui stockerait les differentes variantes et qui les mettrait à jour en fonction des prochaines.

Enfin, notre dernière interrogation portait sur l'impact d'un tel algorithme, même parfait, sur le jugement final de son utilisateur. En effet, même en admettant qu'il sache déterminer sans erreur la véracité d'une découverte, il est difficile à croire qu'un utilisateur accepte une réponse allant à l'encontre de ses croyances.

De nombreux exemples à travers l'histoire peuvent démontrer ce fait. Nicolas Copernic, malgré des preuves scientifiques tengibles, n'a pas réussi à convaincre de son vivant la théorie de l'héliocentrisme. Ce n'est qu'approximativement deux cent quatre vingt années après sa

mort que sa théorie fût acceptée.

Cette réticence s'explique par le fait que cette découverte allait à l'encontre des croyances et théories de l'ensemble de la société et des chercheurs de l'époque.

Même en fournissant des preuves tangibles, une théorie sera difficile à accepter si elle va à l'encontre des croyances et connaissances de la majorité. Ainsi, le résultat d'un algorithme pourrait très facilement être remis en cause et, finalement, ne pèsera pas ou peu sur l'appréciation du lecteur.

7 Conclusion

Les fausses découvertes sont omniprésentes dans de nombreux domaines.

Ainsi, l'article "A framework for Multi-A(rmed)/B(andid) testing with online FDR control" propose une procédure basée sur le contrôle du FDR afin de les déterminer et de les contourner. Pour ce faire, une méthode itérative est utilisée qui détermine α_j grâce à la "online-FDR procedure". Puis, elle implémente une "MAB procedure" retournant un bras recommandé (recommended arm) et construisant une p-valeur toujours valide. Enfin, si le bras avec la moyenne empirique la plus élevée n'est pas le bras de contrôle et que $P_t^j \leq \alpha_j$, le bras de contrôle est rejeté en faveur de ce dernier.

Malgré des résultats statistiquement concluants, cette méthode comporte certaines limites. En effet, elle se base sur le contrôle du FDR qui ne permet pas de gérer les différences pouvant être rencontrées dues aux spécificités des domaines d'application. Un laboratoire cherchant un remède contre le cancer n'aura pas les mêmes priorités qu'une entreprise cherchant à optimiser son chiffre d'affaire. Ce facteur n'est pas pris en compte tout au long de l'article. De plus, en fonction du poids d'une découverte, celle-ci peut être composée de "sous-découvertes" qui devraient être chacune analysées séparément. Si elle est jugée faussée, certaines de ces composantes ne le sont pas forcément. La sauvegarde de certaines hypothèses serait donc judicieuse.

Enfin, même dans le cas d'un algorithme parfait, son utilité peut être remise en cause par certaines personnes. Au final, le jugement d'un résultat se base surtout sur les connaissances et les croyances du lecteur. Un algorithme pourrait aider à la décision de l'utilisateur mais ne peut remplacer le jugement personnel qu'il porte sur la découverte. Ce dernier conditionne très souvent la décision finale.

References

- [1] Fanny, Yang Aaditya, Ramdasy Kevin, Jamieson Martin, J. Wainwrighty. *A framework for Multi-A(rmed)/B(audit) testing with online FDR control*. 18 Nov 2017.
- [2] A. Javanmard and A. Montanari, *Online rules for control of false discovery rate and false discovery exceedance,*” The Annals of Statistics, 2017.