

Canadian Psychology / Psychologie canadienne

Pourquoi tout ce bruit autour de la valeur p? Quelques pistes de compréhension pour le non-expert

--Manuscript Draft--

Manuscript Number:	CAP-2022-0089R3
Full Title:	Pourquoi tout ce bruit autour de la valeur p? Quelques pistes de compréhension pour le non-expert
Abstract:	La valeur p (p-value) fait l'objet de vifs débats dans la communauté scientifique et influence l'évaluation que les intervenants font des recherches dans leur domaine de pratique. Certains souhaitent la proscrire, alors que d'autres veulent continuer à l'utiliser. Cet article présente des éléments de réflexion concernant les écarts entre sa réelle signification et l'usage commun qui en est fait. Il discute des définitions ou interprétations erronées qui lui sont associées, notamment par rapport aux seuils de la valeur p. Enfin, il présente des alternatives ou compléments à cette statistique. Son objectif est d'identifier et de proposer aux chercheurs en sciences sociales, humaines et de la santé les meilleures pratiques de l'analyse des données quantitatives et de développer en même temps une réflexion critique à l'égard de leurs résultats basés sur cette statistique, laquelle continue à attiser bien des passions.
Article Type:	Masked Article / Article anonyme
Keywords:	Valeur p; test d'hypothèse; précision; facteur de Bayes; comparaison de modèles
Corresponding Author:	Sébastien Béland, Ph.D. Université de Montréal: Universite de Montreal CANADA
Corresponding Author E-Mail:	sebastien.beland@umontreal.ca
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	Université de Montréal: Universite de Montreal
Other Authors:	Michael Cantinotti
	Denis Cousinea
	Marco Barroca-Paccard
	Marie-Aude Boislard
	Christian Bourassa
	Julien Bureau
	Pier-Olivier Caron
	Lucile Chanquoy
	Christophe Chénier
	Jean-François Daoust
	Éric Dionne
	Stéphanie Forté
	Éric Frenette
	Stéphanie Girard
	Vincent Grenon
	Bradley Harding
	Éric Lacourse
	Daniel Lalande

	Florent Michelot
	Quoc Dinh Nguyen
	Jean-Sébastien Renaud
	Jérôme St-Amand
	Floris Van Vugt
Author Comments:	Nous souhaitons soulever la grande qualité du travail de l'évaluateur 3
Corresponding Author's Secondary Institution:	
First Author:	Sébastien Béland, Ph.D.
Order of Authors Secondary Information:	
Manuscript Region of Origin:	CANADA
Opposed Reviewers:	
Order of Authors:	Sébastien Béland, Ph.D.
	Michael Cantinotti
	Denis Cousinea
	Marco Barroca-Paccard
	Marie-Aude Boislard
	Christian Bourassa
	Julien Bureau
	Pier-Olivier Caron
	Lucile Chanquoy
	Christophe Chénier
	Jean-François Daoust
	Éric Dionne
	Stéphanie Forté
	Éric Frenette
	Stéphanie Girard
	Vincent Grenon
	Bradley Harding
	Éric Lacourse
	Daniel Lalande
	Florent Michelot
	Quoc Dinh Nguyen
	Jean-Sébastien Renaud
	Jérôme St-Amand
	Floris Van Vugt

Pourquoi tout ce bruit autour de la valeur p? Quelques pistes de compréhension pour le non-expert

https://osf.io/prab2/?view_only=db563a9efa754b7d8c86defba5fc0416

- 1. Sébastien Béland, Université de Montréal, sebastien.beland@umontreal.ca:
auteur de correspondance**
2. Michael Cantinotti, Université du Québec à Trois-Rivières, Michael.Cantinotti@uqtr.ca
3. Denis Cousineau, Université d'Ottawa, Denis.Cousineau@uottawa.ca
4. Marco Barroca-Paccard, HEP Vaud, marco.barroca-paccard@hepl.ch
5. Marie-Aude Boislard, Université du Québec à Montréal, boislard-pepin.marie-aude@uqam.ca
6. Christian Bourassa, Université de Montréal, christian.bourassa.1@umontreal.ca
7. Julien Bureau, Université Laval, julien.bureau@fse.ulaval.ca
8. Pier-Olivier Caron, TÉLUQ, Pier-Olivier.Caron@teluq.ca
9. Lucile Chanquoy, Université Nice Sophia Antipolis, lucile.chanquoy@univ-cotedazur.fr
10. Christophe Chénier, Université de Montréal, christophe.chenier@umontreal.ca
11. Jean-François Daoust, Université de Sherbrooke, JF.Daoust@USherbrooke.ca
12. Éric Dionne, Université d'Ottawa, Eric.Dionne@uOttawa.ca
13. Stéphanie Forté, Université de Montréal, steph.forte1@gmail.com
14. Éric Frenette, Université Laval, Eric.Frenette@fse.ulaval.ca
15. Stéphanie Girard, Université du Québec à Trois-Rivières, Stephanie.Girard3@uqtr.ca
16. Vincent Grenon, Université de Sherbrooke, Vincent.Grenon@USherbrooke.ca
17. Bradley Harding, Université de Moncton, bradley.harding@umoncton.ca
18. Éric Lacourse, Université de Montréal, eric.lacourse@umontreal.ca
19. Daniel Lalande, Université du Québec à Chicoutimi, daniel_Lalande@uqac.ca
20. Florent Michelot, Université de Moncton, florent.michelot@umoncton.ca
21. Quoc Dinh Nguyen, Université de Montréal, nguyen.quoc.d@gmail.com
22. Jean-Sébastien Renaud, Université Laval, jean-sebastien.renaud@fmed.ulaval.ca
23. Jérôme St-Amand, Université du Québec en Outaouais, jerome.st-amand@uqo.ca
24. Floris Van Vugt, Université de Montréal, floris.van.vugt@umontreal.ca

Notes :

Texte révisé à la suite des remarques de l'évaluateur 3

Pourquoi tout ce bruit autour de la valeur p?
Quelques pistes de compréhension pour le non-expert

03/06/2024

Résumé

La valeur p (p -value) fait l'objet de vifs débats dans la communauté scientifique et influence l'évaluation que les intervenants font des recherches dans leur domaine de pratique. Certains souhaitent la proscrire, alors que d'autres veulent continuer à l'utiliser. Cet article présente des éléments de réflexion concernant les écarts entre sa réelle signification et l'usage commun qui en est fait. Il discute des définitions ou interprétations erronées qui lui sont associées, notamment par rapport aux seuils de la valeur p . Enfin, il présente des alternatives ou compléments à cette statistique. Son objectif est d'identifier et de proposer aux chercheurs en sciences sociales, humaines et de la santé les meilleures pratiques de l'analyse des données quantitatives et de développer en même temps une réflexion critique à l'égard de leurs résultats basés sur cette statistique, laquelle continue à attiser bien des passions.

Mots-clés

Valeur p ; test d'hypothèse ; précision ; facteur de Bayes ; comparaison de modèles

Déclaration d'importance publique

Un grand nombre d'études en psychologie rapportent des valeurs p . Malheureusement, cette statistique est souvent mal comprise et mal interprétée. Ce manuscrit vise à éclairer le non-expert sur la nature de la valeur p en plus de présenter quelques alternatives et stratégies pour mieux l'utiliser dans les écrits scientifiques.

Introduction

La valeur p (p -value) occupe une position importante dans les représentations entourant la recherche quantitative, au point où elle fait l'objet de nombreux commentaires humoristiques. Une recherche réalisée au début du mois de décembre 2023, avec l'expression « p value » et « même », génère ainsi plus d'un million de résultats sur Google.com. À titre d'illustration, nous avons observé un photomontage de ce type sur la porte d'un local à l'Université de Montréal (Figure 1).

[Insérer la Figure 1, ici]

Cela n'est pas étonnant, car des études ont montré qu'une valeur p inférieure à un seuil de 0,05 augmente les chances qu'un manuscrit soit évalué positivement par les pairs et publié (Easterbrook et al., 1991 ; Koletsi et al., 2009). D'autres études montrent aussi qu'il recevra davantage d'attention et qu'il sera plus souvent cité par la communauté scientifique (Greenberg, 2009). Ce seuil a reçu un nom, le *seuil alpha* (α , « the first source of error » selon Neyman et Pearson, 1928, p. 177), et les résultats inférieurs à ce seuil sont considérés par un grand nombre de chercheurs comme étant *statistiquement significatifs* alors que, dans le même temps, les limites qui sont associées à ce concept ont amené certains spécialistes à se positionner à l'encontre de son utilisation (Wasserstein, Schirm, et Lazar, 2019).

À titre d'exemple, Ioannidis (2019) a examiné 62 résumés d'articles dans le domaine de la santé. Sur un total de 141 valeurs p répertoriées, 86% étaient inférieures à 0,05. Devrait-on en conclure que les valeurs supérieures à 0,05 refléteraient des résultats moins dignes de mention que leurs contreparties présentant des valeurs p inférieures à

1 0,05? C'est ce que pense Melton (1965), éditeur du *Journal of Experimental Psychology*,
2 quand il a introduit la politique consistant à refuser les articles dont les valeurs p ne se
3 logeaient pas sous ce seuil (Loftus, 1993 ; Lykken, 1991).

4 Les étudiants sont au courant de cet enjeu. Plusieurs d'entre eux ont sans doute
5 aussi lu les appels à la prudence concernant la valeur p , par exemple, ceux d'Amrhein et
6 al. (2019) dans *Nature*. De nombreux articles provenant de toutes les disciplines
7 scientifiques et ayant comme objectif de clarifier cette statistique sont parus (Dixon,
8 2003 ; Goodman, 2008 ; Nahm, 2017 ; Hubbard et Lindsay, 2008 ; Lakens, 2021) dont
9 celui émanant de l'*American Statistical Association* (Wasserstein et Lazar, 2016) est
10 probablement l'un des plus... significatifs.

11 Le contenu de ce qui suit n'est pas nouveau et s'adresse aux chercheurs qui ne
12 sont pas des experts en méthodes quantitatives. Bien qu'il existe de nombreux écrits
13 critiques sur la valeur p , un grand nombre de chercheurs et d'intervenants continuent
14 d'utiliser mécaniquement cette statistique comme nous pouvons le constater dans les
15 écrits scientifiques et professionnels. L'objectif de ce texte consiste à clarifier ce qu'est la
16 valeur p et son seuil d'interprétation alpha en présentant leurs forces et limites principales
17 et quelques alternatives ou compléments intéressants. Nous nous insérons donc dans la
18 mouvance faisant la promotion de meilleures pratiques dans l'analyse de données
19 quantitatives en sciences sociales (voir Finkel et al., 2015, pour une discussion dans le
20 contexte de la psychologie). L'idée est d'éviter d'utiliser ce que Sijtsma (2016) ainsi que
21 Flake et Fried (2020) nomment les *pratiques de recherche discutables* (« *questionable
22 research practices* »). L'article conclut sur des pistes à envisager en vue d'améliorer les
23 pratiques en méthodologie de recherche et en interprétation des résultats de recherche.

Qu'est-ce que la valeur p ?

Le test d'hypothèse enseigné dans le paradigme statistique dominant (souvent identifié par l'acronyme anglophone *NHST* pour « *null hypothesis significance testing* ») utilise la valeur p pour évaluer si l'hypothèse nulle (symbolisée par H_0) peut être rejetée ou non (Denis, 2004). Pour prendre cette décision, quatre étapes sont habituellement franchies.

Premièrement, H_0 doit être définie. Traditionnellement, elle caractérise une situation de *statu quo*, selon laquelle il ne se produit rien de systématique (par exemple, deux groupes expérimentaux ne diffèrent pas sur une variable dépendante ou deux variables ne sont pas corrélées entre elles). Cette hypothèse nulle décrivant un paramètre dans la population sert de référence pour tester à quel point les données observées concordent avec elle. Cela correspond, par exemple, à une différence de moyennes négligeable ou une corrélation presque nulle.

Afin d'illustrer ce test, imaginons que nous souhaitions comparer les résultats en mathématiques des élèves de deux écoles différentes en utilisant les scores obtenus à une même épreuve. Dans ce cas, H_0 caractérise une situation où les moyennes des deux populations, notées μ_1 et μ_2 , ne diffèrent pas, ce qui peut s'écrire $H_0 : \mu_1 = \mu_2$, ou encore $H_0 : \mu_1 - \mu_2 = 0$. On cherche alors à savoir si cette hypothèse reste plausible alors que les moyennes observées de deux écoles seront vraisemblablement inégales. Autrement dit, la procédure vise à estimer la probabilité que la différence observée entre les moyennes soit le fruit d'une fluctuation d'échantillonnage et à décider qu'en dessous d'un certain seuil de probabilité (par exemple 5%), cette probabilité est trop faible pour qu'il soit plausible de l'attribuer à cette seule fluctuation.

Ceci étant dit, est-ce réaliste, plausible ou utile de spécifier une H_0 qui correspond à une différence exactement nulle entre deux groupes (certains auteurs anglophones parlent de « *nil hypothesis* »)? En effet, il est utopique de prédire la survenue d'une différence absolument nulle entre des groupes (Lykken, 1968, 1991). Par ailleurs, même si l'effet s'avère infime, il y aura toujours une taille d'échantillon suffisamment grande pour rendre cet effet significatif et rejeter l' H_0 (par exemple, Mitroff et Biggs, 2014, collectent un million de nouvelles observations par jour!). Pour cette raison Craver (1978) considère que « le test de signification statistique met en place un homme de paille, soit l'hypothèse nulle, et tente de le renverser [*knock down*] » (p. 381, traduction libre).

Deuxièmement, il est nécessaire de sélectionner un seuil de rejet nommé alpha (α). Ce seuil prend généralement la valeur de 0,05 en sciences sociales, mais parfois de 0,01 et, même, plus bas dans des domaines tels que la recherche en génétique et la physique ; nous y reviendrons plus loin. Puisque les statistiques inférentielles visent à évaluer la plausibilité d'un effet hypothétique dans la population, à partir d'échantillons observés, il est toujours possible de rejeter H_0 erronément à la suite de l'analyse des données. Ici, α représente un seuil définissant une probabilité d'erreur, l'erreur de type I, qui est établie *a priori* par le/la chercheur.se. C'est le seuil selon lequel H_0 n'est plus jugée plausible puisqu'il n'est pas raisonnable d'attribuer à la seule fluctuation d'échantillonnage la différence observée entre les échantillons. L'erreur de type I consiste donc à rejeter H_0 alors qu'il aurait fallu la conserver parce que la valeur p apparaît trop faible.

Troisièmement, il faut réaliser la collecte des données (la passation d'une épreuve en mathématiques dans notre exemple). Il faut pour cela idéalement choisir une méthode d'échantillonnage aléatoire, par exemple un échantillonnage par grappe (souvent utilisé dans les milieux scolaires où *des groupes d'étudiants* sont choisis aléatoirement, pas les élèves) ou encore un échantillonnage par strate (fréquent pour les sondages téléphoniques où l'on tente de préserver la représentativité des groupes d'âge, par exemple).

Quatrièmement, il faut choisir un test statistique pertinent. Dans notre cas, le test t pour échantillons indépendants est adéquat pour comparer les moyennes des scores provenant des élèves de deux groupes différents. La différence de moyennes observée est alors convertie sur une échelle standardisée qui correspond à une distribution théorique de référence nommée distribution t de Student. Celle-ci s'approche progressivement d'une distribution normale lorsque la taille de l'échantillon augmente.

Finalement, le test t permet de se prononcer sur le rejet ou non de H_0 . Pour un test unidirectionnel, le rejet de H_0 survient seulement si le t calculé dépasse une valeur critique qui dépend du seuil α prédéfini. Par exemple, avec un seuil α de 0,05 et un n de 5 000 pour chacun des deux groupes (donc, un total de 10 000 élèves), cela correspond à la valeur critique t de 1,65. Si nous obtenons une valeur t observée, par exemple, de 2,20, la procédure commande de rejeter H_0 : on peut considérer que les moyennes des deux groupes d'étudiants sont différentes. Par contre, un t observé de 1,60 amène à conclure au « non-rejet de H_0 ». Cette phrase sibylline a généré énormément de confusion dès les débuts (Fisher le reconnaît dès 1935, p. 16), au point que Howell (2010, p. 93) écrit « le problème de l'interprétation d'un non-rejet a tourmenté les étudiants de statistique depuis 75 ans » (traduction libre ; nous sommes maintenant presque rendus à

90 ans de tourments!). Il faut en retenir que l'absence de preuves n'est pas la preuve de l'absence.

Dans la pratique contemporaine où l'utilisation de logiciels est omniprésente, les chercheurs utilisent plutôt la valeur p pour décider du rejet ou non de H_0 . Cette valeur indique la probabilité d'obtenir une différence de moyennes égale ou plus extrême que celle observée lorsque H_0 est présumée vraie ; cette valeur p correspond à une valeur du t pour une taille d'échantillon donnée. Elle permet de jauger directement, sur une échelle de probabilité, la plausibilité des données observées si l'on postule que H_0 est vraie, le chercheur pouvant la confronter ou non à une probabilité-seuil s'il le juge opportun. Krueger et Heck (2019) insistent sur le fait que cette statistique doit être écrite $p(D|H_0)$, soit la probabilité des données qui sont observées (D) postulant (le signe « $|$ ») que l'hypothèse nulle H_0 est vraie dans la population d'où l'échantillon est extrait.

Qui es-tu, valeur p ?

Comme nous l'avons déjà mentionné, de nombreuses études ont discuté de l'idée que la valeur p est mal comprise (Cohen, 1994 ; Goodman, 2008 ; Greenland, et al., 2016 ; Haller et Krauss, 2002). Par exemple, plusieurs croient, à tort, que p représente la vraisemblance de H_0 étant donné les données observées (qu'on pourrait écrire $p(H_0 | D)$ dans la notation ci-dessus). Cette distinction est importante puisque la probabilité de A postulant B n'est pas égale à la probabilité de B postulant A. Ainsi, au Canada, si je vous présente une personne membre d'un ordre professionnel (A), la probabilité qu'elle soit psychologue en exercice (B) peut être faible, disons 4 sur 100, soit $p(B|A)$. Mais si je vous présente une psychologue en exercice, la probabilité qu'elle soit membre d'un ordre

professionnel, soit $p(A|B)$, est égale à 1 ou presque (si l'on exclut les personnes en situation d'infraction!).

Que sait-on, alors, sur cette statistique p ? Demidenko (2016) montre que, toutes choses étant par ailleurs égales, la valeur p diminue lorsque la taille de l'échantillon s'accroît, lorsque l'effet qui est mesuré est de plus grande taille ou les deux. Pour illustrer cette idée, utilisons la différence de moyennes d'anxiété suivante qui est établie à partir de deux populations d'élèves : $\mu_1 - \mu_2 = 2,5$ et les écarts-types dans les deux populations qui sont égales : $\sigma_1 = \sigma_2 = 10$. Nous allons progressivement accroître la taille de l'échantillon, comme indiqué dans la colonne de gauche du Tableau 1.

[Insérer le Tableau 1 ici]

Même lorsque les moyennes et les écarts-types des groupes restent inchangés, le test t mène à une autre interprétation lorsque la taille de l'échantillon augmente : nous ne rejetons pas H_0 , au seuil $\alpha = 0,05$, lorsque $n = 20$ et $n = 200$, alors que c'est le cas dès que n est égal ou supérieur à 250. Toutes choses étant égales par ailleurs, deux élèves de moins ($n = 248$) et nous aurions manqué cet effet. Cet exemple montre qu'adopter une posture catégorique, binaire, ne s'accorde pas avec la nature continue des phénomènes observés. Mentionnons également que la valeur p varie aussi en fonction d'autres variables telles que la taille d'effet¹, la variance et la forme de la distribution des données (Dahiru, 2008).

Une critique récurrente de la valeur p est que son interprétation dichotomique s'avère artificielle. Ainsi, pourquoi devrait-on se réjouir d'un $p = 0,049$ (ce qui mène à un rejet de H_0) et éprouver de la déception face à un $p = 0,051$? À l'évidence ces deux

¹ La taille de l'effet (*effect size*) mesure dans une population statistique la force de la relation entre deux ou plusieurs variables.

valeurs présentent des différences infimes. Une position rigide, tout-ou-rien, manque de nuances et pourrait nuire au développement des connaissances (McShane et Gal, 2016). Gelman et Stern (2006) rappelaient d'ailleurs qu'« une erreur statistique courante consiste à résumer les comparaisons par signification statistique, puis à établir une distinction nette entre les résultats significatifs et non significatifs » (p. 328, traduction libre) dans la discussion des résultats.

Pour reprendre un exemple en contexte scolaire, avec un échantillon total de 20 élèves, le chercheur mal avisé, décrit par Gelman et Stern (2006), écrira « après avoir collecté un échantillon par école, on trouve $p = 0,583$, ce qui montre une absence totale de différence entre les élèves des deux écoles ». Cette interprétation est erronée, car, d'une part, la valeur p ne permet pas de conclure à la véracité de H_0 (une absence de différence) et, d'autre part, la différence observée n'est pas totalement nulle. Le chercheur avisé fournira quant à lui plus de détails: « le premier échantillon présente une moyenne de 72,5 et le second, de 75,0, soit une différence de 2,5 points, avec un intervalle de confiance à 95% sur cette différence allant de -4,1 à 9,1 (écart-type regroupé $s_p = 10$). Cette différence est faible (d de Cohen de 0,25, intervalle de confiance à 95% allant de -0,63 à +1,13) et ne s'avère pas probante, $t(18) = 0,559$, $p = 0,583$. Obtenir ces résultats dans un contexte hypothétique où il n'y aurait pas de réelle différence entre les élèves des écoles est plausible ».

Enfin, mentionnons que Wasserstein et Lazar (2016) ont publié un article influent pour résumer la position de l'*American Statistical Association* sur ce qu'est la valeur p et ce qu'elle n'est pas. Nous reprenons intégralement ces six éléments : i) La valeur p indique à quel point les données sont compatibles avec un modèle statistique spécifique ;

ii) La valeur p ne mesure pas la probabilité qu'une hypothèse étudiée soit vraie ou la probabilité que des données soient produites seulement par la chance, soit une hypothèse non conditionnelle, $p(H_0)$; iii) Des conclusions scientifiques, d'affaires ou des décisions politiques ne devraient pas uniquement être basées sur la comparaison d'une valeur p avec un seuil prédéterminé ; iv) Une inférence adéquate requiert de rapporter toutes les informations sur le devis de recherche et la population sur laquelle l'on veut faire des inférences et de faire preuve de transparence sur les forces et les faiblesses de ce devis ; v) Une valeur p ne mesure pas la grandeur de l'effet ou l'importance dans un contexte appliqué d'un résultat ; et vi) La valeur p , à elle seule, ne présente pas un appui suffisant à l'égard d'un modèle ou d'une hypothèse. En statistique, une règle importante est que plus l'échantillon est grand, plus l'estimé sera précis. Or, en ce qui concerne la valeur p , plus l'échantillon est grand ($n = 500\,000$, voir plus), plus les résultats apparaîtront statistiquement significatifs pour un seuil donné.

Pourquoi un seuil de 0,05?

La valeur critique la plus communément utilisée est 0,05. Mais d'où vient ce seuil? Et pourquoi pas un autre seuil? Il semble que cette pratique se soit progressivement installée à la suite de la lecture d'un paragraphe qui provient de l'ouvrage phare du statisticien Fisher, *Statistical methods for research workers*, initialement publié en 1925. Nous reproduisons ici l'extrait traduit de ce paragraphe, où Fisher indiquait au sujet de l'interprétation du résultat à un test du khi carré:

« Ce qui nous intéresse, c'est savoir si la valeur p remet une hypothèse en cause. Si p se situe entre 0,1 et 0,9, il n'y a pas de raison de remettre en cause l'hypothèse nulle. Si la valeur p est inférieure à 0,02, il y a alors de fortes raisons de croire que l'hypothèse nulle ne parvient pas à expliquer ce qui est observé. Ainsi, nous nous

1 égarerons rarement si nous traçons un seuil à 0,05 » (Fisher, 1934, p. 82, traduction
2 libre).

3
4 Le seuil critique de 0,05 semble s'être tranquillement installé dans les us et coutumes de
5 la communauté scientifique en séparant cette idée du contexte de réflexion plus général
6 proposé par Fisher.

7 Certains se sont demandé s'il fallait utiliser un seuil plus restrictif tel que 0,005
8 (e.g., Benjamin et al., 2017). La réponse a été assez rapide (e.g., Trafimov et al., 2018 ;
9 Lakens, 2021): le problème émane de la façon dont la valeur p est conceptualisée plutôt
10 que du seuil en soi. Pour illustrer ce point, nous vous suggérons de faire le test de la tache
11 de café, à la Figure 2, avant de poursuivre votre lecture.

12 [Insérer la Figure 2, ici]

13 Dans le graphique de gauche, la valeur p est de 0,051. Selon le seuil alpha de
14 0,05, c'est un résultat non significatif, et Raphaëlle est déçue d'avoir obtenu ce résultat.
15 Dans le graphique de droite, la valeur p est de 0,049, et la différence est significative. Or,
16 bien malin est celui qui peut voir la différence entre les deux graphiques tellement ils sont
17 quasiment identiques (pour le graphique de gauche, les moyennes sont de 77,00 et 71,77,
18 pour un d de Cohen de 0,395, IC95% [-0,002, 0,789] ; pour le graphique de droite, les
19 moyennes sont de 77,56 et 71,21, pour un d de Cohen de 0,397, IC95% [0,002, 0,792]).
20 Comme le mentionnent les chercheurs étasuniens Rosnow et Rosenthal (1989) « Dieu
21 aime presque autant le seuil 0,06 qu'il aime le seuil 0,05 » (p. 1277, traduction libre).
22 Cette boutade montre qu'un seuil basé sur une valeur p rigide trace une distinction binaire
23 entre des résultats dont la différence est presque totalement indiscernable. Un seuil basé
24 sur une lecture critique de la valeur p permet de discriminer et de souligner un résultat
25 vraisemblablement exceptionnel, surtout lorsqu'elle est appuyée par une méta-analyse ou

des reproductions qui ajoutent un poids considérable à la valeur des inférences statistiques.

Fisher, qui est le premier à avoir systématisé le calcul de la valeur p dans ses écrits des années 1910-1930, arrondit cette valeur à des ordres de grandeur faciles à appréhender et permettant de produire des tableaux statistiques dans les publications: 1 sur 10 (soit 0,10), 1 sur 20 (soit 0,05) ou encore 1 sur 50 (soit 0,02). Pour Fisher, ce sont des raccourcis simples à comprendre pour indiquer à quel point une valeur p obtenue est faible ou non (Lehman, 2011). Par ailleurs, il n'utilise pas ce résultat de façon dichotomique, mais plutôt afin de graduer le résultat des analyses. Pour lui, la valeur p est un *indice* de la force des évidences contre H_0 (bien que formellement, ce soit la probabilité des données quand on postule H_0 ; Berger, 2003). Nuzzo (2014) explique qu'aux yeux de Fisher, une valeur p représente une façon informelle de juger de l'importance du résultat: est-ce qu'il mérite qu'on s'y intéresse un peu plus? Autrement dit, une petite valeur p serait un drapeau qui s'agite pour attirer l'attention des chercheurs et des utilisateurs des recherches et non une fin en soi.

Suivant cette logique, nous mentionnons Rafi et Greenland (2020), qui suggèrent de transformer la valeur p en « indice de surprenance des données » avec cette formule :

$$s = -\log_2 p \quad (1)$$

où s indique à quel point les données sont surprenantes dans le cadre de H_0 . Par exemple, lorsque s vaut 1, les données sont aussi surprenantes que d'obtenir pile en lançant une pièce de monnaie une fois (un cas pas du tout surprenant, d'où une valeur p de 0,50) ; lorsque s vaut 10, le résultat est aussi surprenant que d'obtenir 10 piles après 10 lancers (c.-à-d., très surprenant, et p vaut $0,50^{10}$, soit $\approx 0,001$). Pour un p de 0,05, l'indice de

surprenance de 4,3 correspond approximativement à la probabilité d'obtenir quatre piles sur quatre lancers (en effet, cette dernière probabilité est de 6,25%, soit proche de 5%). Si pareil résultat n'est pas exceptionnellement fréquent, il n'est pas non plus exceptionnellement rare. La valeur s ne remplacera pas la valeur p puisque les deux sont interchangeables, mais s a l'avantage de possiblement rendre une valeur p plus intuitive.

6

7 *Significatif... mais instable ?*

8 Imaginez la situation suivante: vous avez réalisé une étude où le résultat d'intérêt s'avère « significatif » (disons que $p = 0,04$) et votre réaction ressemble à celle illustrée à la Figure 1. Or, quelques minutes plus tard, un de vos participants vous contacte et, après réflexion, il préfère que vous n'utilisiez pas ses réponses dans votre recherche. Vous avez bien entendu obtenu un certificat d'éthique, et donc, la démarche est claire: vous effacez les données liées à ce participant et recrutez un nouveau participant pour le remplacer. La nouvelle valeur p sera-t-elle identique? Il y a très peu de chance que cela soit le cas, sauf si le participant de remplacement obtenait exactement le même score que le participant remplacé. Une simulation qui se veut une simple illustration sans prétention de généralisabilité montre qu'avec un seul changement dans un groupe de 50 personnes, la nouvelle valeur p obtenue pourra être aussi petite que 0,02 (Figure 1), mais aussi grande que 0,08 (Figure 3).

20

[Insérer la Figure 3, ici]

1 Murdoch et al. (2008) nous rappellent que la valeur p est de source aléatoire: elle dépend
 2 des données et si les données changent, la valeur p change elle aussi.²

3

4 **Quelques stratégies pour embrasser l'ère « post $p < 0,05$ »**

5 **En plus d'être mal comprise, la valeur p ne serait pas sans faille** (Goodman, 2008).

6 Ce constat avait déjà été établi dès 1933 par Neyman et Pearson. Certains suggèrent de
 7 rejeter cette statistique (Hurlbert et Lombardi, 2009, qui pointent entre autres la
 8 dichotomisation arbitraire des résultats en « significatifs et « non-significatifs » et la
 9 perception erronée que p concerne la probabilité relative $\frac{p(H_0)}{p(H_1)}$. Ce rejet de p a été adopté
 10 par le journal *Basic and Applied Social Psychology* où les valeurs p sont proscrites depuis
 11 2015 (Trafimow et Marks, 2015). La valeur p nous semble **toujours pertinente. Nous**
 12 **comprenons que, pour les chercheurs et les chercheuses pour qui la valeur p est habituelle**
 13 **dans l'usage des procédures statistiques, il puisse être difficile d'envisager de ne plus**
 14 **l'utiliser. Une alternative à la démarche radicale appliquée par le journal *Basic and***
 15 ***Applied Social Psychology* consisterait à inciter les chercheurs et les chercheuses à mieux**
 16 **comprendre ce qu'elle signifie et d'exiger que les résultats soient complémentés avec**
 17 **d'autres informations statistiques.** Dans ce qui suit, nous allons présenter quelques
 18 stratégies d'intérêt à ce sujet.

19

20 *Stratégies centrées sur la valeur p et le seuil α*

² Notons que cette critique s'applique à tous les types de mesure, pas seulement la valeur p .

Est-ce vraiment utile de formuler une H_0 où les moyennes sont exactement égales, comme on le voit dans un grand nombre d'études comparatives? En réalité, la population des élèves de deux écoles n'a jamais exactement la même moyenne à un examen ; la question est plutôt de savoir quelle sera l'ampleur des différences observées et avec quels phénomènes ces différences sont associées. Il est douteux qu'il existe, dans nos disciplines, une situation dans laquelle une H_0 de zéro serait plausible. Or, puisque la valeur p est une probabilité conditionnelle à H_0 (donc $p(D|H_0)$) et que H_0 est très peu plausible par principe, que vaut la valeur p ? Est-elle même interprétable? Plusieurs auteurs pensent que c'est une raison importante pour la laisser de côté, par exemple Hurlbert et Lombardi (2009), Cumming (2014), ainsi que McShane et al. (2019).

Cependant, les habitudes d'utilisation d'une statistique populaire sont souvent bien ancrées dans la pratique et, jusqu'à présent, les résistances individuelles et institutionnelles n'ont pas permis de concrétiser cet objectif.³ Goodman (2019) mentionnait à ce sujet:

« L'explication la plus simple relativement à l'utilisation de la valeur p n'est ni philosophique ni scientifique, mais bien sociologique : tout le monde l'utilise et tout le monde croit en sa valeur et c'est cette valeur qui détermine ce que constitue le savoir, et ce qui peut être avancé, publié, financé et promu. Peu importe la réelle signification de la valeur p , c'est elle qui fixe la valeur de ce qui est écrit » (p. 27, traduction libre).

Dans le contexte de la crise de reproductibilité (Pashler et Wagenmakers, 2012), certains auteurs suggèrent de réduire le seuil α . Par exemple, Benjamin et al. (2017),

³ Fait intéressant, le même genre de remarque est apparue concernant le coefficient alpha de Cronbach, dans la communauté en psychométrie. Par exemple, Sijtsma (2009) écrivait : « la seule et unique raison de rapporter l'alpha est que les meilleures revues ont tendance à accepter les articles qui utilisent des méthodes statistiques qui existent depuis longtemps, telles que l'alpha » (p. 118, traduction libre).

proposent de le réduire à 0,005 pour rejeter H_0 . L'impact de cette modification revient à baisser l'occurrence d'erreurs de type I (faux positifs), ce qui est l'effet recherché par ces auteurs. Néanmoins, ceci résulterait en une augmentation des erreurs de type II (faux négatifs), notée β . Ioannidis (2018) est toutefois très sceptique avec cette proposition. Nous estimons que ce type de solution mécanique comprend le même vice que l'heuristique habituelle: elle ne peut remplacer l'exercice de la pensée critique et mène à des pratiques fétichistes de la statistique plus qu'à un raisonnement nuancé et rigoureux concernant les données. Or, par l'effet de l'argument d'autorité (« c'est publié, donc cela doit être pertinent »), cette pratique influence également de nombreux utilisateurs des connaissances, allant des psychologues au grand public. De plus, utiliser un seuil plus restrictif tel que 0,005, toutes autres choses étant égales par ailleurs, se traduit par une perte de puissance statistique. Ceci peut impacter les études dont les tailles des échantillons sont plus faibles, par exemple dans le cas des populations difficiles à étudier ou pour des maladies ou troubles mentaux plus rares. Si les grands échantillons sont souhaitables, car ils génèrent une estimation plus précise des paramètres qui intéressent le chercheur, la qualité méthodologique d'une étude ne doit toutefois pas être réduite à ce seul enjeu.

Alternativement, plutôt que manipuler le seuil α , certains suggèrent de manipuler l'hypothèse⁴. Au lieu de tester une H_0 , pourquoi ne pas tester l'hypothèse de la présence d'un effet d'intérêt substantiel? Les tests de non-équivalence ont pour but de répondre à ce genre de question (Schuirmann, 1987). La difficulté de cette approche est cependant de définir ce que représente un « effet intéressant ». Cette définition doit être établie avant l'étude pour éviter d'être biaisée par l'observation des résultats, et doit revêtir un

⁴ Nous remercions un évaluateur anonyme de nous avoir rappelé cette alternative.

1 caractère consensuel. Au final, il apparaît toutefois que la valeur p de tels tests ne s'avère
 2 pas plus aisée d'interprétation.

3 D'autres auteurs déclarent qu'il faut devenir des utilisateurs de statistique avisés,
 4 tout spécialement concernant la valeur p (Wasserstein et Lazar, 2016). Ainsi, laisser les
 5 logiciels définir des valeurs arbitraires par défaut (par exemple, $\alpha = 0,05$, $\alpha = 0,01$, ou
 6 encore $\alpha = 0,005$) constitue un automatisme qui peut engendrer des conclusions
 7 erronées ou dénuées de mise en contexte.

8 Dans l'optique de rendre les chercheurs plus responsables, Lakens (2022) invite à
 9 garder à l'esprit quels sont les buts et les attentes des chercheurs (en particulier, quels
 10 sont les effets **auxquels** ils s'attendent, quels sont les effets qu'ils voudraient rejeter, quels
 11 sont les effets qu'ils trouveraient **intéressants** d'observer). Lakens (2022) suggère même
 12 qu'il serait acceptable d'accroître α si cela peut mieux servir les buts de l'étude et si ce
 13 choix est adéquatement argumenté.

14 Aguinis et al. (2010), de leur côté, suggèrent de contraindre le choix du α afin de
 15 mieux distinguer les erreurs de type I et les erreurs de type II acceptables dans un projet
 16 de recherche. Pour cela, ils proposent de calculer le niveau α avec une formule qu'ils
 17 nomment $\alpha_{\text{Désiré}}$, de la façon suivante:

$$\alpha_{\text{Désiré}} = \left[\frac{p(H_1)\beta}{1-p(H_1)} \right] \left(\frac{1}{\text{DRS}} \right) \quad (2)$$

19 où DRS est le « desired relative seriousness » **qui représente le danger de commettre une**
 20 **erreur de type I par rapport au danger de commettre une erreur de type II**, $p(H_1)$ est la
 21 probabilité que l'hypothèse alternative soit plausible (basé sur des connaissances
 22 obtenues *a priori*) et β est l'erreur de type II. La valeur critique $\alpha_{\text{Désiré}}$ est ainsi mieux
 23 justifiée et moins arbitraire qu'une valeur α « par défaut ». Supposons que nous fixons

1 $p(H_1) = 0,6$ en nous basant sur des études antérieures, avec $\beta = 0,15$ et $DRS = 2$, le
 2 $\alpha_{\text{Désiré}}$ serait alors égal à 0,11. Cela veut dire que 0,11 est l'erreur de type I jugée
 3 acceptable sélectionnée avant d'accomplir les analyses. En se référant aux résultats de la
 4 Figure 2, les deux tests rejetteraient H_0 . Plus le DRS sera élevé, plus l'accent sera mis sur
 5 l'importance d'éviter une erreur de type I.

6 Quant à eux, Blume et al. (2019) proposent l'idée d'une valeur p de deuxième
 7 génération, notée p_δ . Lorsque H_0 n'est pas contenue dans l'intervalle de confiance de
 8 95%, nous savons que la valeur p est plus petite que 0,05. Or, baser H_0 sur une seule
 9 valeur est artificiel et inintéressant en réalité. Voilà pourquoi ces auteurs utilisent un
 10 intervalle de valeurs défini par le chercheur (« zone d'indifférence » aux résultats), plutôt
 11 qu'une seule valeur, pour opérationnaliser H_0 . Cela permet de mieux formaliser
 12 l'information du contexte d'analyse en calculant la proportion de l'intervalle de confiance
 13 qui chevauche (ou pas) l'intervalle de valeurs de H_0 . Le lecteur curieux pourra en
 14 apprendre davantage sur cette approche en consultant l'interface visuelle suivante :
 15 <https://lucy.shinyapps.io/sgpvalue/>

16

17 *Stratégie centrée sur les tailles de l'effet et leur précision*

18 Des auteurs tel que Ranstam (2012) et Cumming (2014) proposent plus
 19 simplement de rapporter systématiquement les tailles d'effet (*effect size*) et leurs
 20 intervalles de confiance (IC). La taille de l'effet représente une estimation de la taille de
 21 la différence entre deux groupes, par exemple, qui n'est pas biaisée par la taille de
 22 l'échantillon comme la valeur p peut l'être. Comme le rappelle Glass, cité dans Sullivan
 23 et Feinn (2012, p. 279) : « la signification statistique est la partie la moins intéressante

1 des résultats d'une étude. C'est **donc elle** qui devrait être décrite et non sa simple
2 présence/non présence. Autrement dit, nous ne devrions pas nous demander si un
3 traitement génère un effet sur des individus, mais à quel point ce traitement présente un
4 effet sur ceux-ci » (traduction libre).

5 Les tailles d'effets sont de plusieurs sortes (voir le Tableau 2). Certaines sont
6 « brutes », c'est-à-dire tirées directement des données, comme la différence entre deux
7 moyennes. D'autres sont standardisées et sont plus abstraites, comme le d de Cohen, le η -
8 partiel et la corrélation, mais ont l'avantage d'être interprétables sur une échelle
9 universelle (XXXXX). Par exemple, la valeur du d de Cohen est de 0,395, IC95% [-
10 0,002, 0,790] pour la comparaison à gauche de la Figure 2 et de 0,397, IC95% [0,001,
11 0,792] pour la comparaison qui est à sa droite.

12 [Insérer le Tableau 2, ici]

13 Pour chaque statistique observée, il est important d'estimer sa précision en lui
14 assignant un intervalle de confiance (e.g., XXXXX). Rappelons que l'intervalle de
15 confiance représente une mesure de la précision d'une statistique (par exemple,
16 l'intervalle de confiance de la moyenne ou l'intervalle de confiance de la médiane). Si
17 une étude est conduite à plusieurs reprises et de la même façon, **nous pouvons estimer**
18 **pour chaque échantillon tiré au hasard**, une moyenne et un intervalle de confiance à 95%.
19 L'interprétation de l'intervalle de confiance à 95% indique que 95% de ces intervalles de
20 confiance contiennent la vraie moyenne de la population (le paramètre). Quand un
21 sondeur rapporte la moyenne des votes qui irait à un parti politique en vue d'une élection,
22 omettre l'intervalle de confiance revient à ignorer la précision du sondage (ou indiquer
23 subrepticement que l'échantillonnage était non probabiliste), ce dont le lecteur devrait

1 être informé directement. De plus, il est intéressant de constater que, contrairement à la
2 valeur p qui n'offre pratiquement aucune information utile dans une perspective de
3 réplication des résultats d'une étude, un intervalle de confiance à 95% pour une moyenne
4 représente une probabilité de 83% que la moyenne de l'étude répliquée dans des
5 conditions similaires se trouve entre le seuil inférieur et le seuil supérieur de cet intervalle
6 (Cumming, 2008). Enfin, l'intervalle de confiance permet de s'extraire d'une vision
7 dichotomique des résultats d'une procédure statistique. À ce sujet, l'intervalle de
8 confiance ne doit pas être perçu comme une simple duplication de la valeur p (dans le
9 sens où l'on peut se demander si l'intervalle inclut une différence de moyenne nulle),
10 mais bien comme un intervalle dont les bornes inférieures et supérieures apportent une
11 information pertinente pour développer une pensée méta-analytique. Par exemple, en
12 agréant les résultats de différentes études au sein d'un domaine de recherche, les
13 intervalles de confiance offrent une idée de l'ampleur des effets observés. Ces
14 informations permettent donc de réfléchir à la portée appliquée et théorique des résultats,
15 plus qu'en termes de présence ou d'absence d'effets « statistiquement significatifs »
16 (Dienes, 2008).

17 Les tailles d'effet les moins abstraites sont sans doute celles qu'il faut
18 recommander le plus, comme les tailles d'effets brutes. Baguley (2009) avance deux
19 éléments pour soutenir cette proposition. Tout d'abord, bien que la standardisation
20 permette la comparaison de mesures provenant d'échelles différentes, elle entraîne
21 plusieurs enjeux, entre autres un manque de robustesse (ex. : différentes manières de
22 calculer un même indice, manque de fidélité en raison des fluctuations de la variance
23 dans l'échantillon). Ensuite, l'utilisation d'indices faisant sens sur les échelles de mesure

d'origine facilite le jugement critique concernant l'interprétation et la portée des résultats, alors que la standardisation entraîne un niveau d'abstraction qui peut complexifier l'exercice d'un jugement critique sur les résultats. Cependant, Borenstein et coll. (2016, p. 25) rappellent à juste titre que les mesures doivent être toutes sur la même échelle, sinon, la standardisation est la seule façon de comparer ces mesures. *A contrario*, l'utilisation de tailles d'effets standardisées ouvre la porte à l'invocation rituelle de seuils arbitraires permettant de démarquer des tailles d'effet petites, moyennes et grandes pour le d de Cohen ou η^2 sans considérer si ces seuils sont porteurs d'un sens par rapport à un domaine de recherche spécifique ou par rapport à des résultats publiés antérieurement. Cela représente finalement un enjeu similaire à l'utilisation dichotomique de la valeur p . Afin que ces seuils ne soient pas arbitraires, Gignac et Szodorai (2016) réalisent une vaste méta étude afin d'identifier les quartiles des corrélations généralement rapportées en psychologies. Ces quartiles deviennent alors les frontières pour définir les tailles d'effets faibles et larges. Ils trouvent ,1 en deçà duquel une corrélation peut être considérée faible, et ,3 au-delà duquel une corrélation peut être considérée forte. Ce dernier nombre est bien inférieur au ,5 suggéré par Cohen (1992).

Stratégies centrées sur la comparaison de modèles

Le processus statistique expliqué jusqu'ici se base sur H_0 et son rejet (ou pas), en faveur de H_1 . Selon cette perspective, on ne peut pas soutenir H_0 avec des données, mais seulement la rejeter (ou ne pas être en mesure de la rejeter). Pour le dire autrement, on ne peut pas, avec les tests de H_0 , soutenir l'absence de différence entre deux conditions

1 pertinentes, même s'il est parfois scientifiquement pertinent de montrer que deux
2 interventions sont équivalentes.

3 En guise d'illustration, reprenons l'exemple des écoles et imaginons que la
4 seconde école a reçu un enseignement qui est présumé plus efficace que celui de la
5 première école (qui sert de point de référence). Après réflexion par rapport aux
6 connaissances sur les effets de taille minimalement pertinente associées à la mesure, il
7 apparaît que dans la seconde école, les scores devraient augmenter minimalement de 5
8 points en moyenne pour que le résultat soit jugé utile. Il serait alors intéressant de tester
9 cela en **formulant H_0** : $\mu_1 - \mu_2 = 5$. Sur le site OSF <https://osf.io/prab2/> se trouve un script
10 R qui montre comment il est aisé de tester une hypothèse nulle différente de zéro. Ceci
11 revient toutefois à soumettre à l'épreuve des données uniquement cette H_0 , ce qui peut
12 entraîner les inconvénients mentionnés au paragraphe précédent lorsqu'elle est rejetée.

13 Comme autre stratégie, nous pourrions aussi opposer H_0 : $\mu_1 - \mu_2 \leq 0$ à l'hypothèse
14 telle que H_1 : $\mu_2 - \mu_1 \geq 5$, c'est-à-dire que nous envisagerions dans ce cas que la seconde
15 classe aura 5 points de plus en moyenne que la première. Ainsi, en plus de tester des
16 hypothèses nulles non nulles, il est également possible de tester des hypothèses de non-
17 nullité (les tests d'équivalence).

18 Dans la suite de cette section, un duo d'hypothèses sert à représenter des modèles
19 permettant d'utiliser des approches basées sur le ratio de vraisemblance ou le facteur de
20 Bayes. Wasserstein et Lazar (2016) écrivent que certaines personnes préfèrent remplacer
21 les valeurs p par d'autres approches « comme des méthodes qui mettent **en avant**
22 l'estimation plutôt que les tests, tels les intervalles de confiance, les intervalles de
23 crédibilité ou les intervalles de prédiction, des méthodes bayésiennes ; des ratios de

1 vraisemblance ou des facteurs de Bayes, et d'autres approches comme la théorie de la
2 décision et des faux positifs » (p. 132, traduction libre).

3 Le ratio de vraisemblances RV constitue le rapport entre le maximum de
4 vraisemblance des données avec le modèle correspondant à l'hypothèse 1 donnée ci-
5 dessus et le maximum de vraisemblance des données avec le modèle correspondant à
6 l'hypothèse 0, que l'on exprime comme suit dans le cadre de la comparaison entre H_0 et
7 H_1 :

$$8 \quad RV = \frac{p_{MV}(D|H_1)}{p_{MV}(D|H_0)} = \frac{\text{Maximum de vraisemblance de } H_1}{\text{Maximum de vraisemblance de } H_0} \quad (3)$$

9 La vraisemblance des données indique à quel point des données telles que celles
10 observées auraient pu se produire lorsque le modèle est postulé. Sans surprise, un bon
11 modèle est celui qui suggère que les données observées s'avèrent probables. Par
12 extension, la vraisemblance des données devient la vraisemblance du modèle une fois les
13 données obtenues. Cette statistique présente l'avantage d'être intuitive à interpréter.
14 Ainsi, un $RV = 5$ indiquera que H_1 est cinq fois mieux soutenue par les données que H_0 .
15 À l'inverse, un RV de 0,2 (1 sur 5) indiquera que l'hypothèse nulle est cinq fois mieux
16 soutenue par les données que l'hypothèse alternative (ce ratio est donc réversible).
17 Évidemment, ce ratio ne fournit aucune information sur la comparaison avec n'importe
18 quel autre modèle qui n'a pas été estimé.

19 Le facteur de Bayes (FB)⁵ a aussi été proposé pour comparer des hypothèses. Il a
20 en outre l'avantage de contourner certaines limites de l'approche fréquentiste que Cohen

⁵ Nous ne prétendons pas offrir une quelconque introduction aux modèles bayésiens. Notre objectif vise plutôt à mentionner l'existence du FB comme alternative (ou complément) et de donner des éléments de compréhension minimaux à son sujet.

1 (1994) discute de manière particulièrement éloquente (c'est pourquoi nous laissons le
 2 passage en anglais) :

3 “What’s wrong with NHST? Well, among many other things, it does
 4 not tell us what we want to know, and we so much want to know what
 5 we want to know that, out of desperation, we nevertheless believe that it
 6 does! What we want to know is "Given these data, what is the
 7 probability that H_0 is true?" But as most of us know, what it tells us is
 8 "Given that H_0 is true, what is the probability of these (or more
 9 extreme) data?" These are not the same, as has been pointed out many
 10 times over the years...” (1994, p. 997)

12 En réalité, ce qui intéresse le chercheur est surtout de connaître la probabilité d’un
 13 modèle (ou d’une hypothèse) à partir des données qui ont été analysées. Et c’est une
 14 information qu’il est possible d’obtenir avec les statistiques bayésiennes, par exemple
 15 $p(H_0|D)$.

16 Mathématiquement, cette statistique prend la forme suivante :

$$17 \quad FB = \frac{p_{FB}(D|H_1)}{p_{FB}(D|H_0)} = \frac{\text{Vraisemblance marginale de } H_1}{\text{Vraisemblance marginale de } H_0}. \quad (4)$$

18 Le facteur de Bayes implique aussi un rapport de vraisemblance (similaire à celui de
 19 l’équation 3, mais dans certains cas seulement) auquel on ajoute l’information *a priori*
 20 que nous avons sur chacun des modèles, laquelle est *souvent* de l’information qui doit
 21 être formalisée (par exemple, des résultats provenant d’études précédentes ou certaines
 22 connaissances pertinentes du chercheur). Cette information est de nature subjective et le
 23 FB indique comment ce chercheur devrait réviser ces connaissances.

Des barèmes ont été proposés pour interpréter le facteur de Bayes, dont le plus connu a été fourni par Kass et Raftery (1995). Ici, nous allons plutôt rapporter celui de Kelter (2020) qui nous semble plus nuancé :

[Insérer la Figure 4, ici]

Par exemple, un FB égal à 20 montre que les données sont 20 fois mieux représentées par H_1 que par H_0 lorsque les distributions *a priori* sont pris en compte. Les données observées corroborent de manière probante l'hypothèse H_1 . Un FB de 0,05 offrirait, à l'opposé, un élément de preuve relativement fort en faveur de H_0 .

Lorsque les moyennes présentées à la Figure 2 sont analysées à l'aide du test t bayésien, nous observons un FB égal à 1,173, donc faiblement en faveur de H_1 . Ce résultat s'avère toutefois anecdotique selon le barème de la Figure 4. Un FB de 1,195, également en faveur de H_1 , est calculé à la droite de cette même figure. Encore une fois, ce résultat est considéré comme anecdotique. Pour ces analyses, la distribution Cauchy est utilisée en guise de distribution *a priori* (il serait possible d'envisager une autre loi de distribution, au besoin). Notez que le facteur de Bayes est de plus en plus intégré dans les logiciels, dont R, JASP (<https://jasp-stats.org/>) et SPSS, parmi d'autres.

Il est possible d'estimer un facteur de Bayes à partir d'une valeur p . Benjamin et Berger (2019) ainsi que Held et Ott (2018) proposaient d'utiliser une façon de calculer la limite supérieure du facteur de Bayes (notée « LFB »), ce qu'ils ont appelé le « Bayes factor bound ». Selon Benjamin et Berger (2019), le LFB :

« indique l'inférence la plus forte qui puisse être effectuée à partir des données. La rapporter pourrait informer les chercheurs des cas où une preuve apparemment forte (par exemple, une valeur p) n'est, en fait, pas très convaincante; ce qui éviterait à des chercheurs de conclure trop

rapidement, et à tort, sur la base d'une valeur $p \gg (p. 188, \text{traduction libre})$.

Le tableau qui suit présente des exemples de valeurs p et leur équivalent LFB .

[Insérer le Tableau 3, ici]

Par exemple, pour une valeur p de 0,05, LFB vaut 2,44, donc en faveur de H_1 et lorsque $p = 0,01$, LFB vaut 8,13, également en faveur de H_1 . Le premier reste anecdotique ; le second commence à être plus convaincant. Le LFB représente ainsi la valeur la plus élevée du FB qui soit en adéquation avec une valeur p observée. Selon cette stratégie, les valeurs p obtenues à la Figure 2, soit 0,051 et 0,049, génèrent une valeur $LFB = 2,44$ en faveur de H_1 .

Quelques pistes pour améliorer les pratiques d'analyse

Nous suggérons quatre pistes pour entrer dans l'ère « post $p < 0,05$ » qui nous semblent réalistes, afin d'aider les chercheurs à accomplir cette transition.

La simplicité des statistiques descriptives pour « sonder le terrain »

Un article scientifique raconte une histoire. Les statistiques descriptives et des graphiques pertinents (Tufte, 2001) constituent la première étape nécessaire pour illustrer l'effet qui nous intéresse. Ces informations amènent à réfléchir sur la crédibilité des résultats par rapport aux attentes du chercheur et des connaissances antérieures sur le sujet. Pour le dire autrement, les statistiques descriptives poussent à une réflexion conceptuelle plutôt qu'à l'application mécanique d'heuristiques ou de conventions.

1

2 *Dans la majorité des cas, la valeur p n'est pas suffisante*

3 Betensky (2019) rappelle que la valeur p nécessite un contexte plutôt qu'un seuil α ,
4 qui simplifie à outrance la prise de décision lors d'un test d'hypothèse. Ainsi, il faut
5 plutôt contextualiser la valeur p en fonction de l'échantillon et de la valeur de la taille de
6 l'effet minimalement pertinente (par exemple, le d de Cohen lors d'une comparaison
7 entre deux moyennes), soit celle qui a du sens pour le chercheur ou pour les personnes
8 visées par une intervention. Comme le titraient Gelman et Stern, dans leur article de
9 2006 : « The Difference Between “Significant” and “Not Significant” is not itself
10 statistically significant » (voir aussi le test de la tache de café, en guise d'illustration).
11 Nous mettons l'accent sur le fait que la valeur p ne s'avère pas, à elle seule, suffisante
12 pour conclure à la supériorité d'une hypothèse sur une autre.

13 Le chercheur doit avoir dans sa boîte à outils un faisceau de stratégies, dont plusieurs
14 sont présentées dans cet article. Ainsi, le rapport de vraisemblance et le facteur de Bayes
15 offrent une interprétation plus nuancée qu'un seul et unique seuil α . Mais il y a plus : il
16 est important de contextualiser la statistique utilisée par le chercheur dans une perspective
17 de décision pertinente plutôt que comme une conclusion scientifique isolée. Pour le dire
18 autrement, est-ce que ce qui a été trouvé est utile de façon appliquée? Et de plus, comme
19 « le doute est l'oreiller du savant » (Héger, 1895, p. 169), et ajouterions-nous, également
20 celui de tout bon lectorat scientifique, quels sont les aspects méthodologiques qui
21 viennent nuancer les résultats obtenus et qui doivent être transmis au lectorat pour lui
22 permettre de juger des résultats dans leur contexte ? Il faut accepter de se mettre « à nu »

1 méthodologiquement pour offrir une inspection détaillée de la démarche scientifique dont
2 les résultats découlent.

3

4 *Intégrer systématiquement des mesures de précision dans les écrits qui relèvent d'une*
5 *approche quantitative*

6 Une carence méthodologique revient souvent dans les sondages rapportés dans
7 certains journaux : l'absence d'intervalle de confiance fréquentiste ou d'intervalle de
8 crédibilité bayésien. Or, la mesure de la précision est fondamentale pour bien
9 contextualiser une statistique telle qu'une moyenne. De plus, cette information est
10 nécessaire lorsque vient le temps de comparer des hypothèses sur des paramètres
11 statistiques.

12

13 *Adopter un langage plus nuancé*

14 Plusieurs auteurs ont critiqué l'utilisation de la terminologie « statistiquement
15 significatif ». En effet, « $p < 0,05$ » a parfois été considéré à tort comme équivalent de
16 « Découverte scientifique ». Or, au moment où l'un d'entre nous écrit ces lignes, les
17 hirondelles volent une nouvelle fois plus bas qu'à l'habitude et l'orage gronde. Après de
18 multiples observations, cette relation apparaît statistiquement significative ($p < 0,001$),
19 mais faut-il pour autant en faire une théorie scientifique qui suggère que les hirondelles
20 volent plus bas pour se protéger de l'orage? La valeur p ne permet aucunement de
21 suggérer ceci. Un biologiste affirmera plutôt que c'est la pression d'une masse d'air qui

1 influence la position des insectes et donc la position du garde-manger des hirondelles...
2 (Cabaret, 2014).

3 Benjamin et Berger (2019) suggèrent d'utiliser le qualificatif « suggestif » plutôt que
4 « significatif » ; en français, il est possible d'utiliser « notable » pour un grand effet (dans
5 le sens *il vaut la peine de le noter*) et « cohérent » (c.-à-d., *cohérent avec cette*
6 *hypothèse*). Nous abondons dans le même sens en estimant qu'il faut adopter une écriture
7 prudente au moment de rapporter des résultats découlant de comparaisons d'hypothèses
8 ou de modèles. La « significativité » statistique peut constituer un guide approximatif sur
9 ce qu'il faut rapporter dans le texte, sur là où il faut mettre l'accent, mais n'est
10 certainement pas un guide pour développer une théorie.

11 Dans un contexte bayésien, lorsqu'un résultat incohérent apparaît (incompatible avec
12 une hypothèse), cette hypothèse perd en crédibilité au détriment d'une autre qui est plus
13 crédible. Contrairement à la perspective poppérienne qui présume que la « vérité »
14 scientifique peut à tout moment être réfutée, ici une seule réfutation (basée sur un critère
15 statistique) ne constitue pas un poids suffisant pour rejeter une théorie. Une réfutation
16 basée sur un critère statistique représente plutôt une invitation à identifier une explication
17 alternative.

18 Enfin, nous considérons que les analyses statistiques ne remplacent pas un bon devis
19 de recherche et une approche méthodologique rigoureuse. Ainsi, une avenue qui semble
20 prometteuse est le pré-enregistrement des devis de recherche, soit faire évaluer par les
21 pairs la question de recherche et la méthodologie avant de connaître les résultats d'une
22 étude. En effet, une bonne méthodologie peut se passer de statistiques, mais les
23 statistiques ne peuvent jamais se passer d'une bonne méthodologie.

1

2

Conclusion

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

Remerciement

18

19

20

Référence

21

22

23

XXXXX : 4 références sont cachées pour préserver l'anonymat des auteurs, en conformité avec les instructions de la revue.

- Amrhein, V., Greenland, S. et McShane, B. (2019). Scientists rise up against statistical significance. *Nature*, 567, 305-307. <https://doi.org/10.1038/d41586-019-00857-9>
- Baguley, T. (2009), Standardized or simple effect size: What should be reported? *British Journal of Psychology*, 100, 603-617. <https://doi.org/10.1348/000712608X377117>
- Benjamin, D. J., Berger, J. O., Johannesson, M. et 69 autres auteurs (2017). Redefine statistical significance. *Nature Human Behavior*, 2, 6-10. <https://doi.org/10.1038/s41562-017-0189-z>
- Berger, J. O. (2003). Could Fisher, Jeffreys and Neyman Have Agreed on Testing? *Statistical Science*, 18, 1-12. <https://doi.org/10.1214/ss/1056397485>
- Betensky, R. A. (2019) The p -Value Requires Context, Not a Threshold, *The American Statistician*, 73(sup1), 115-117. <https://doi.org/10.1080/00031305.2018.1529624>
- Blume, J. D., Greevy, R. A., Welty, V. F., Smith, J. R. et Dupont, W. D. (2019). An Introduction to Second-Generation p -Values, *The American Statistician*, 73(sup1), 157-167. <https://doi.org/10.1080/00031305.2018.1537893>
- Borenstein, M., Hedges, L. V., Higgins, J. P., & Rothstein, H. R. (2021). *Introduction to meta-analysis*. John Wiley & Sons.
- Cabaret, M. (2014). Pourquoi les hirondelles volent-elles bas avant le mauvais temps? [extrait audio]. Espace des sciences, Centre de culture scientifique technique et industrielle de Rennes. https://soundcloud.com/espacedessciences/pourquoi-les-hirondelles?utm_source=clipboard&utm_medium=text&utm_campaign=social_sharing
- Carver, R. P. (1978). The case against statistical significance testing. *Harvard Educational Review*, 48(3), 378-399. <https://doi.org/10.17763/haer.48.3.t490261645281841>
- Cohen, J. (1992). A power primer. *Psychological Bulletin*, 112(1), 155-159.
- Cohen, J. (1994). The earth is round ($p < .05$). *American Psychologist*, 49(12), 997-1003. <https://doi.org/10.1037/0003-066X.49.12.997>
- Cumming, G. (2014). The new statistics: Why and how. *Psychological Science*, 25(1), 7-29. <https://doi.org/10.1177/0956797613504966>
- Cumming, G. (2008). Replication and p intervals : p values predict the future only vaguely, but confidence intervals do much better. *Perspectives on Psychological Science*, 3(4), 286-300. <https://doi.org/10.1111/j.1745-6924.2008.00079.x>
- Dahiru T. (2008). P-Value, a true test of statistical significance? A cautionary note. *Annals of Ibadan postgraduate medicine*, 6(1), 21-26. <https://doi.org/10.4314/aipm.v6i1.64038>
- Demidenko, E. (2016). The p -Value You Can't Buy, *The American Statistician*, 70(1), 33-38. <https://doi.org/10.1080/00031305.2015.1069760>
- Denis, D. J. (2004) The modern hypothesis testing hybrid: R. A. Fisher's fading influence. *Journal de la Société Française de Statistique*, 145, 5-26.
- Dienes, Z. (2008). Understanding psychology as a science: An introduction to scientific and statistical inference. Palgrave Macmillan.
- Dixon, P. (2003). The p -value fallacy and how to avoid it. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale*, 57(3), 189-202. <https://doi.org/10.1037/h0087425>

- 1 Easterbrook, P. J., Berlin, J. A., Gopalan, R. et Matthews, D. R. (1991). Publication bias
2 in clinical research. *Lancet (London, England)*, 337(8746), 867-872.
3 [https://doi.org/10.1016/0140-6736\(91\)90201-y](https://doi.org/10.1016/0140-6736(91)90201-y)
- 4 Finkel, E. J., Eastwick, P. W. et Reis, H. T. (2015). Best research practices in
5 psychology: Illustrating epistemological and pragmatic considerations with the case
6 of relationship science. *Journal of Personality and Social Psychology*, 108(2), 275-
7 297. <https://doi.org/10.1037/pspi0000007>
- 8 Fisher, R.A. (1934). *Statistical methods for research workers* (5^e éd.). Oliver and Boyd:
9 Edinburgh.
- 10 Fisher, R. A. (1935). *The Design of Experiments*. Hafner Publishing Co.
- 11 Flake, J. K. et Fried, E. I. (2020). Measurement Schmeasurement: Questionable
12 Measurement Practices and How to Avoid Them. *Advances in Methods and*
13 *Practices in Psychological Science*, 3(4), 456-465.
14 <https://doi.org/10.1177/2515245920952393>
- 15 Gelman, A. et Stern, H. (2006). The Difference Between “Significant” and “Not
16 Significant” is not Itself Statistically Significant, *The American Statistician*, 60(4),
17 328-331. <https://doi.org/10.1198/000313006X152649>
- 18 Gignac, G. E., & Szodorai, E. T. (2016) Effect size guidelines for individual differences
19 researchers. *Personality and Individual Differences*, 102, 74-78.
20 <https://doi.org/10.1016/j.paid.2016.06.069>
- 21 Glass, G. V. (1976). Primary, secondary, and meta-analysis of research. *Educational*
22 *researcher*, 5(10), 3-8. <https://doi.org/10.2307/1174772>
- 23 Goodman S. N. (1999). Toward evidence-based medical statistics. 1: The P value fallacy.
24 *Annals of internal medicine*, 130(12), 995-1004. [https://doi.org/10.7326/0003-](https://doi.org/10.7326/0003-4819-130-12-199906150-00008)
25 [4819-130-12-199906150-00008](https://doi.org/10.7326/0003-4819-130-12-199906150-00008)
- 26 Goodman S. N. (2008). A dirty dozen: twelve p-value misconceptions. *Seminars in*
27 *hematology*, 45(3), 135-140. <https://doi.org/10.1053/j.seminhematol.2008.04.003>
- 28 Goodman, S. N. (2019). Why is Getting Rid of P-Values So Hard? Musings on Science
29 and Statistics, *The American Statistician*, 73(sup1), 26-30.
30 <https://doi.org/10.1080/00031305.2018.1558111>
- 31 Greenberg S. A. (2009). How citation distortions create unfounded authority: Analysis of
32 a citation network. *BMJ*, 339(b2680), 1-14. <https://doi.org/10.1136/bmj.b2680>
- 33 Greenland, S., Senn, S. J., Rothman, K. J., Carlin, J. B., Poole, C., Goodman, S. N., et
34 Altman, D. G. (2016) Statistical tests, P values, confidence intervals, and power: A
35 guide to misconceptions. *European Journal of Epidemiology*, 31, 337-350.
36 <https://doi.org/10.1007/s10654-016-0149-3>
- 37 Grove, J. (2015). *Social sciences and humanities faculties ‘to close’ in Japan after*
38 *ministerial intervention*. Times Higher Education.
39 [https://www.timeshighereducation.com/news/social-sciences-and-humanities-](https://www.timeshighereducation.com/news/social-sciences-and-humanities-faculties-close-japan-after-ministerial-intervention)
40 [faculties-close-japan-after-ministerial-intervention](https://www.timeshighereducation.com/news/social-sciences-and-humanities-faculties-close-japan-after-ministerial-intervention)
- 41 Haller, H., et Krauss, S. (2002). Misinterpretations of Significance: A Problem Students
42 Share with Their Teachers? *Methods of Psychological Research Online*, 7(1), 1-20.
- 43 Héger, P. (1895). De l’idéal. *Revue de Belgique*, 2(14), 148-169.
- 44 Held, L., et Ott, M. (2018). On P-Values and Bayes Factors, *Annual Review of Statistics*
45 *and Its Application*, 5, 393-419. [https://doi.org/10.1146/annurev-statistics-031017-](https://doi.org/10.1146/annurev-statistics-031017-100307)
46 [100307](https://doi.org/10.1146/annurev-statistics-031017-100307)

- Hubbard R, Lindsay RM. (2008). Why P Values Are Not a Useful Measure of Evidence in Statistical Significance Testing. *Theory & Psychology*, 18(1), 69-88.
<https://doi.org/10.1177/0959354307086923>
- Hung, H. M., O'Neill, R. T., Bauer, P., & Kohne, K. (1997). The Behavior of the P-Value When the Alternative Hypothesis Is True. *Biometrics*, 53, 11-22.
- Hurlbert, S. H. et Lombardi, C. M. (2009). Final collapse of the Neyman-Pearson decision theoretic framework and rise of the neoFisherian. *Annales Zoologici Fennici*, 46(5), 311-349. <https://doi.org/10.5735/086.046.0501>
- Howell, D. C. (2010). *Statistical Methods for Psychology* (7th edition). Wadsworth.
- Ioannidis J. (2018). The Proposal to Lower P Value Thresholds to .005. *JAMA*, 319(14), 1429-1430. <https://doi.org/10.1001/jama.2018.1536>
- Ioannidis J. (2019). Publishing research with P-values: Prescribe more stringent statistical significance or proscribe statistical significance?. *European heart journal*, 40(31), 2553-2554. <https://doi.org/10.1093/eurheartj/ehz555>
- Kass, R. E. et Raftery, A. E. (1995). Bayes Factors. *Journal of the American Statistical Association*, 90(430), 773-795. <https://doi.org/10.1080/01621459.1995.10476572>
- Kelter, R. (2020). Bayesian alternatives to null hypothesis significance testing in biomedical research: a non-technical introduction to Bayesian inference with JASP. *BMC Medical Research Methodology*, 20(142). <https://doi.org/10.1186/s12874-020-00980-6>
- Koletsis, D., Karagianni, A., Pandis, N., Makou, M., Polychronopoulou, A. et Eliades, T. (2009). Are studies reporting significant results more likely to be published?. *American journal of orthodontics and dentofacial orthopedics : official publication of the American Association of Orthodontists, its constituent societies, and the American Board of Orthodontics*, 136(5), 632.e1-632.e5.
<https://doi.org/10.1016/j.ajodo.2009.02.024>
- Lakens, D. (2021). The Practical Alternative to the p Value Is the Correctly Used p Value. *Perspectives on Psychological Science*, 16(3), 639-648.
<https://doi.org/10.1177/1745691620958012>
- Lakens, D. (2022). Sample Size Justification. Collabra: Psychology.
<https://doi.org/10.1525/collabra.33267>
- Loftus, G. R. (1993). A picture is worth a thousand p values: On the irrelevance of hypothesis testing in the microcomputer age. *Behavior Research Methods, Instruments, & Computers*, 25, 250-256. <https://doi.org/10.3758/BF03204506>
- Lykken, D. T. (1968). Statistical significance in psychological research. *Psychological Bulletin*, 70, 151-159. <https://doi.org/10.1037/h0026141>
- Lykken, D. T. (1991). What's wrong with psychology anyway? Dans: Cicchetti, D. et Grove, W. M. (eds), *Thinking clearly about psychology [Vol. 1]: Matters of public interest*. Minneapolis: Minnesota University Press.
- McShane, B. B., et Gal, D. (2016). Blinding us to the obvious? The effect of statistical training on the evaluation of evidence. *ManagementScience*, 62(6), 1707-1718.
<https://doi.org/10.1287/mnsc.2015.2212>
- McShane, B. B., Gal, D., Gelman, A., Robert, C. et Tackett, J. L. (2019). Abandon Statistical Significance, *The American Statistician*, 73(sup1), 235-245.
<https://doi.org/10.1080/00031305.2018.1527253>

- 1 Melton, A. W. (1962). Editorial. *Journal of Experimental Psychology*, 64(6), 553-557.
2 <https://doi.org/10.1037/h0045549>
- 3 Mitroff, S. R., & Biggs, A. T. (2014). The ultra-rare-item effect: Visual search for
4 exceedingly rare items is highly susceptible to error. *Psychological Science*, 25(1),
5 284-289. <https://doi.org/10.1177/0956797613504221>
- 6 Murdoch, D. J., Tsai, Y.-L., et Adcock, J. (2008). P-values are random variables. *The*
7 *American Statistician*, 62(3), 242-245. <https://doi.org/10.1198/000313008X332421>
- 8 Nahm, F. S. (2017). What the *P* values really tell us. *The Korean journal of pain*, 30(4),
9 241-242. <https://doi.org/10.3344/kjp.2017.30.4.241>
- 10 Nester, M. R. (1996). An applied statistician's creed. *Applied Statistics*, 45(4), 401-410.
11 <https://doi.org/10.2307/2986064>
- 12 Neyman, J., & Pearson, E. S. (1928). On the use and interpretation of certain test criteria
13 for purposes of statistical inference: Part I. *Biometrika*, 20A(1/2), 175–240.
14 <https://doi.org/10.2307/2331945>
- 15 Neyman, J. and Pearson, E. S. (1933). On the problem of the most efficient tests of
16 statistical hypotheses. *Philosophical Transactions of the Royal Society A*, 231(694-
17 706), 289-337. <https://doi.org/10.1098/rsta.1933.0009>
- 18 Nuzzo, R. (2014) Statistical errors: P values, the ‘gold standard’ of statistical validity, are
19 not as reliable as many scientists assume. *Nature*, 506, 150-152.
20 <https://doi.org/10.1038/506150a>
- 21 Pashler, H., & Wagenmakers, E.-J. (2012) Editor's introduction to the special section on
22 replicability in psychological science: A crisis of confidence? *Perspectives on*
23 *Psychological Science*, 7(6), 528-530. <https://doi.org/10.1177/1745691612465253>
- 24 Rafi, Z. et Greenland, S. (2020). Semantic and cognitive tools to aid statistical science:
25 replace confidence and significance by compatibility and surprise. *BMC Med Res*
26 *Methodol*, 20, 244. <https://doi.org/10.1186/s12874-020-01105-9>
- 27 Ranstam J. (2012). Why the P-value culture is bad and confidence intervals a better
28 alternative. *Osteoarthritis and cartilage*, 20(8), 805-808.
29 <https://doi.org/10.1016/j.joca.2012.04.001>
- 30 Rosnow, R. L. et Rosenthal, R. (1989). Statistical Procedures and the Justification of
31 Knowledge in Psychological Science. *American Psychologist*, 44(10), 1276-1284.
32 <https://doi.org/10.1037/0003-066X.44.10.1276>
- 33 Sijtsma, K. (2009). On the use, the misuse, and the very limited usefulness of Cronbach's
34 alpha. *Psychometrika*, 74(1), 107-120. <https://doi.org/c76r6c>
- 35 Schuirmann, D.J. A comparison of the Two One-Sided Tests Procedure and the Power
36 Approach for assessing the equivalence of average bioavailability. *Journal of*
37 *Pharmacokinetics and Biopharmaceutics* 15, 657–680 (1987).
38 <https://doi.org/10.1007/BF01068419>
- 39 Sijtsma K. (2016). Playing with Data--Or How to Discourage Questionable Research
40 Practices and Stimulate Researchers to Do Things Right. *Psychometrika*, 81, 1-15.
41 <https://doi.org/10.1007/s11336-015-9446-0>
- 42 Sullivan, G.M., & Feinn, R. (2012). Using effect size-or why the *p* value is not enough.
43 *Journal of Graduate Medical Education*, 4(3), 279-82.
44 <https://doi.org/10.4300/JGME-D-12-00156.1>
- 45 Trafimow, D. et Marks, M. (2015) Editorial. *Basic and Applied Social Psychology*, 37(1),
46 1-2. <https://doi.org/10.1080/01973533.2015.1012991>

- 1 Trafimow, D., Amrhein, V., Areshenkoff, C. N. et 51 autres auteurs (2018). Manipulating
2 the Alpha Level Cannot Cure Significance Testing. *Frontiers in Psychology*,
3 9(699), 1-7. <https://doi.org/10.3389/fpsyg.2018.00699>
4 Tufte, E. R. (2001). The visual display of quantitative information. Graphics Press.
5 Wasserstein, R.L. et Lazar, N.L. (2016) The ASA Statement on p-Values: Context,
6 Process, and Purpose. *The American Statistician*, 70(2), 129-133.
7 <https://doi.org/10.1080/00031305.2016.1154108>
8 Wasserstein, R.L., Schirm, A.L., & Lazar, N.A. (2019) Moving to a world beyond
9 “ $p < 0.05$ ”, *The American Statistician*, 73(S1), 1-19,
10 [10.1080/00031305.2019.1583913](https://doi.org/10.1080/00031305.2019.1583913)
11

Tableau 1 : Résultats à un test t pour échantillons indépendants où $\mu_1 - \mu_2 = 2,5$, $ET_1 = ET_2 = 10$ avec six tailles d'échantillon totales

$n = 20$	$t(18) = 0,559, p = 0,583$
$n = 200$	$t(198) = 1,768, p = 0,079$
$n = 248$	$t(246) = 1.969, p = 0.0501$
$n = 250$	$t(248) = 1,976, p = 0,0492$
$n = 600$	$t(598) = 3,062, p = 0,002$
$n = 2000$	$t(1998) = 5,590, p < 0,001$

Note. n représente la taille des deux groupes combinés, lesquels sont égaux.

Tableau 2: Types de tailles d'effet, des plus simples et intuitives aux plus abstraits et complexes

Familles	Exemples
Tailles d'effets brutes	la moyenne, l'écart entre deux moyennes, la variance ou l'écart type
Tailles d'effets standardisées	le d de Cohen, le η^2 -partiel, la corrélation
Tailles d'effets pondérées	la valeur t , la valeur F
Tailles d'effet uniformisées	la valeur p

Tableau 3 : conversion de quelques valeurs p en LFB

p	0,1	0,05	0,01	0,005	0,001
LFB	1,60	2,44	8,13	13,9	52,9



Figure 1 : Exemple de photomontage référant à une valeur p significative

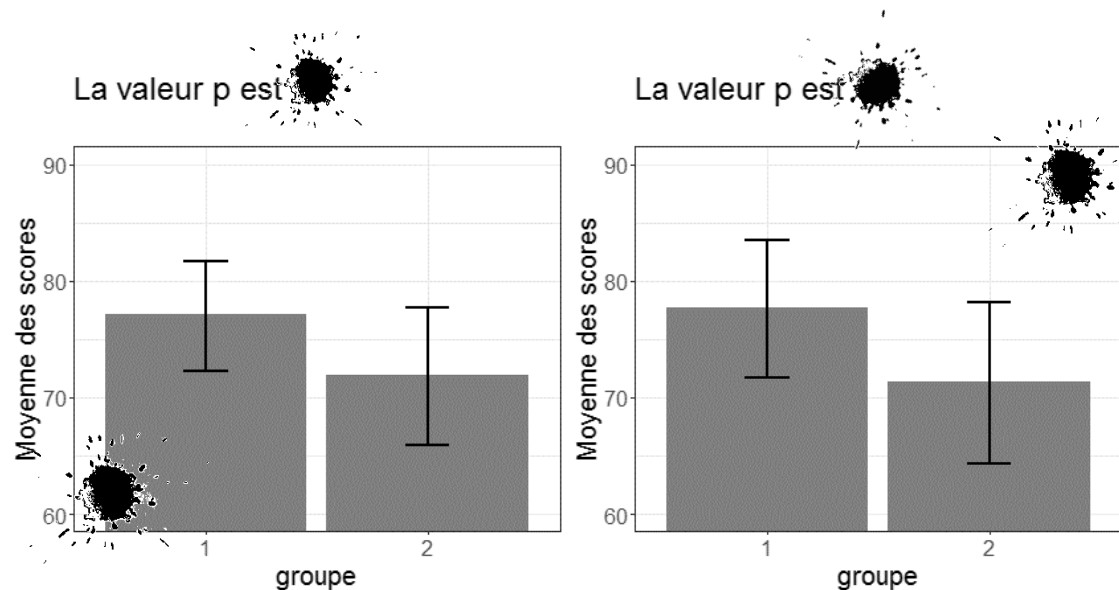


Figure 2: Le test de la tache de café. Benoit est très enthousiaste par ses résultats significatifs alors que Raphaëlle est déçue de ses résultats non significatifs. Leurs graphiques en main (moyennes et intervalles de confiance à 95%; il y a 50 participants par groupes dans les deux études), ils vont voir Guy pour recueillir son avis. Benoit et Raphaëlle sont tellement absorbés qu'ils ne voient pas ce dernier sortir de son bureau avec sa tasse de café noir et ils s'entrechoquent. Les deux graphiques s'envolent dans les airs et se tachent de café du côté des résultats. Pouvez-vous dire quel graphique (celui de gauche ou celui de droite) provient des analyses de Raphaëlle? (voir <https://osf.io/prab2/> pour les données et le code utilisés pour réaliser ce graphique, XXXXX).

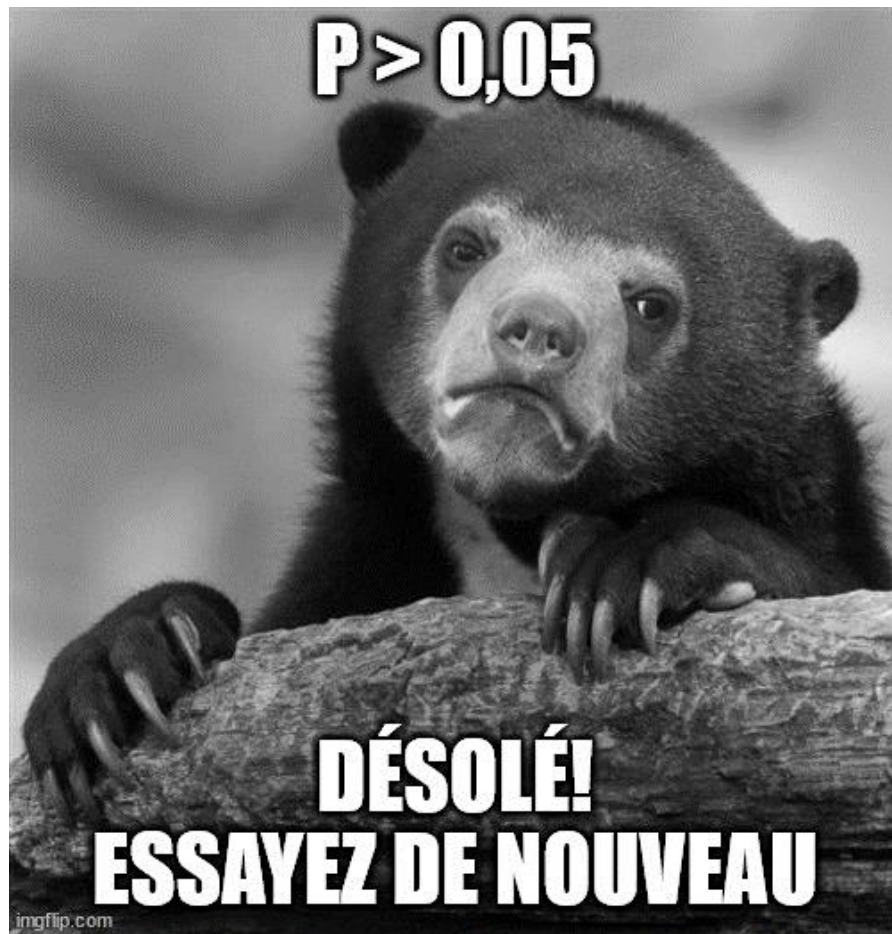


Figure 3 : Désolé! Essayez de nouveau

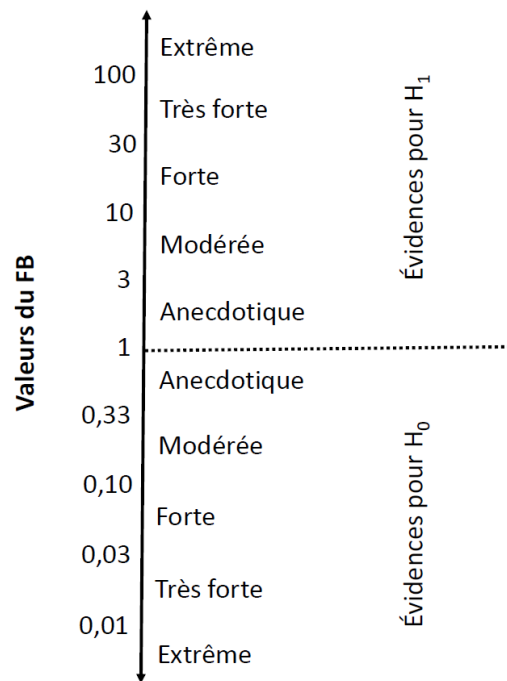


Figure 4 : Barème pour interpréter le facteur de Bayes pour H_0 et H_1 (traduit de la Figure 1 de Kelter, 2020)