

n° 2013-47

**Des spécificités de l'approche
bayésienne et de ses
justifications en statistique
inférentielle**

C. P. ROBERT¹

December 2013

Les documents de travail ne reflètent pas la position du CREST et n'engagent que leurs auteurs.
Working papers do not reflect the position of CREST but only the views of the authors.

¹ Université Paris-Dauphine, CEREMADE et CREST (ENSAE). Email : xian@ceremade.dauphine.fr

Des spécificités de l’approche bayésienne et de ses justifications en statistique inférentielle*

CHRISTIAN P. ROBERT

Résumé: Ce chapitre vise à établir les fondements de l’approche bayésienne en statistique inférentielle, ses racines historiques et ses justifications philosophiques, ainsi qu’à présenter des illustrations de sa mise en oeuvre pratique.

Keywords and phrases: inférence bayésienne, statistique inférentielle, loi a priori, probabilités inverses, subjectivité.

1. Introduction

La statistique bayésienne est une approche de la statistique inférentielle qui propose une réponse unitaire et globale, dans le cadre paramétrique comme dans le cadre non-paramétrique. On peut légitimement se demander pourquoi une telle spécificité est nécessaire, d’autant qu’elle s’accompagne de débats philosophiques particulièrement virulents et de positions militantes pouvant parfois évoquer des dérives sectaires, attitudes qui ne retrouvent pas dans les autres approches de la statistique inférentielle.

Le principe de l’inférence bayésienne se résume assez simplement : étant donné un modèle statistique, l’ensemble des constituants (paramètres et/ou fonctions) inconnus de ce modèle est traité comme une variable aléatoire—potentiellement de dimension infinie—, donc munie d’une loi de probabilité, et l’ensemble des réponses inférentielles se fonde sur la loi de cette variable aléatoire, conditionnellement aux données. Ce qui fait la beauté de cette approche et explique en partie son attractivité est que la démarche inférentielle est alors quasi-automatique, étant donné cette loi. La contrepartie, que certains pourraient qualifier de beauté vénéneuse!, est que la modélisation en variable aléatoire exige un choix de loi de probabilité, choix subjectif opéré par le statisticien bayésien car ne reposant pas directement sur l’observation. Cet aspect de l’approche bayésienne concentre la majorité des critiques de la communauté statistique, une critique secondaire (et plus répandue dans la communauté du *machine learning*) étant l’obligation faite au statisticien de construire intégralement le modèle statistique sur les données, ce qui empêcherait le traitement de structures complexes et/ou de grandes tailles.

*. Christian P. Robert, CEREMADE, Université Paris-Dauphine, 75775 Paris cedex 16, France xian@ceremade.dauphine.fr. Recherche financée en partie par l’Agence Nationale de la Recherche (ANR, 212, rue de Bercy 75012 Paris) sur le contrat 2012–2015 ANR-11-BS01-0010 “Calibration” et par une chaire senior de l’Institut Universitaire de France.

Ce chapitre présente brièvement le développement historique de l'approche bayésienne, depuis Bayes et Laplace jusqu'à nos jours, ainsi que les motivations philosophiques et méthodologiques qui la sous-tendent. La seconde partie présente quelques éléments de mise en œuvre au travers d'exemples simples. Nous dirigeons le lecteur vers Robert (2005) pour une couverture plus complète (en français) de cette méthodologie spécifique, de nombreux autres ouvrages étant disponibles en anglais.

Quelques avertissements au lecteur sont de mise quant au contenu de ce chapitre : ne pouvant y définir correctement les concepts de probabilité nécessaires, je suppose que les lecteurs sont suffisamment familiers avec ceux-ci, à un niveau de fin de Licence (3^{ème} année d'université), pour manipuler les concepts de probabilité conditionnelles. Par ailleurs, ces mêmes lecteurs devront se référer à leur source favorite en ce qui concerne les distributions usuelles (Wikipédia, bien sûr, ou l'annexe de Robert, 2005).

2. Un peu d'histoire sur le développement de la spécificité bayésienne, des origines naturelles, à la quasi-extinction post-kolmogorienne, et au renouvellement modélisateur

2.1. Des probabilités inverses comme définition de la statistique

Le concept de "statistique bayésienne" part du néologisme "bayésien", tiré du nom de Thomas Bayes,¹ qui introduisit le théorème qui porte à présent son nom dans un article posthume de 1763, il y a 250 ans. Ce théorème exprime une probabilité conditionnelle en termes des probabilités conditionnelles inverses,

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)},$$

ce qui a valu à la statistique bayésienne d'être appelée probabilités inverses pendant plus d'un siècle, de Laplace (1812) à Keynes (1920), avant que Fisher² n'introduise le terme "bayésien" (Fienberg, 2006).³ Bien que la formule (ou théorème) de Bayes soit une conséquence directe de la définition des probabilités

1. Thomas Bayes (1702?-1761) était un prêtre presbytérien non-conformiste, membre de la Royal Society, qui publia à titre posthume un *Essay towards solving a Problem in the Doctrine of Chances*. On connaît très peu de détails sur sa vie et le seul portrait de lui dont on dispose demeure incertain. Voir McGrayne (2011) pour une introduction et Dale (1999) pour une étude plus profonde sur les contributions de Bayes à la théorie qui porte à présent son nom.

2. Ronald Fisher, statisticien et généticien anglais, est à l'origine de la notion de vraisemblance. Féroce critique des approches alternatives à la sienne, et en particulier de la perspective bayésienne, il proposa à la fin de sa carrière la notion de statistique fiducière qui, s'appuyant sur des quantités pivotales, ressemblait formellement à une modélisation bayésienne non-informative.

3. Pierre Simon de Laplace a contribué à formaliser et à généraliser la mise en œuvre des probabilités inverses, bien plus que Thomas Bayes. Cette approche aurait donc mérité de s'appeler laplacienne plutôt que bayésienne. Notons également que, si Keynes a été formé à la statistique suivant des principes bayésiens, il rédige son traité de 1921 dans un esprit assez critique, sans pour autant proposer une alternative constructive (Robert, 2011).

(et densités) conditionnelles, son application à des problématiques statistiques, où une observation x dépend d'un paramètre inconnu θ est effectivement appropriée, au sens où le contexte *inverse* ce qui est connu et ce qui est inconnu. Effectuer l'inversion pour obtenir l'information contenue dans x à propos de θ conduit à définir la loi *a posteriori*, loi conditionnelle de θ sachant x .

2.2. Des notions d'*a priori* et d'*a posteriori*, et sur l'illusion des paramètres aléatoires

Pour que cette loi *a posteriori* soit définie, le modèle doit non seulement contenir une loi des observations, de densité $f(x; \theta)$ ou $f(x|\theta)$ mais aussi une loi de probabilité sur le paramètre θ , de densité $\pi(\theta)$ et appelée loi *a priori*. Dans ce cas, la loi *a posteriori* s'obtient par la formule de Bayes,

$$\pi(\theta|x) = \frac{f(x|\theta)\pi(\theta)}{\int_{\Theta} f(x|\theta)\pi(\theta) d\theta}.$$

Dans ce contexte, où θ est traité comme une variable aléatoire, la loi des observations apparaît comme une loi conditionnelle à la valeur du paramètre, plutôt que comme une loi indicée par θ comme dans l'approche classique.

Cette distinction entre paramètre inconnu (mais fixe) et paramètre aléatoire peut à la fois paraître fondamentale et sembler condamner l'approche bayésienne comme inappropriée vis-à-vis de la compréhension moderne—élaborée par Kolmogorov et ses successeurs—de la notion de probabilité. La distinction faite entre probabilité comme stabilisation des fréquences (loi des grands nombres) et probabilité comme quantification d'un degré d'incertitude met en lumière l'imprécision et la subjectivité liées à la seconde approche. C'est certainement la perspective qu'adopta Fisher très rapidement dans sa carrière et la raison de sa querelle avec Jeffreys au cours des années 1930 (Robert et al., 2009). Bien que fondée elle-aussi sur la notion de vraisemblance,

$$\ell(\theta|x) = f(x|\theta),$$

qui reflète également un principe d'inversion, l'approche de Fisher refuse l'interprétation de $f(x|\theta)$ comme densité conditionnelle et la probabilisation de θ . Il est cependant abusif de faire de cet aspect de l'analyse bayésienne autre chose qu'une source de discussion philosophique. En effet, le passage de la notion de paramètre *inconnu* à la notion de paramètre *aléatoire* est incompatible avec la plupart des expériences en particulier dans les sciences physiques. La modélisation statistique suppose au contraire l'existence d'un paramètre fixe θ , sur lequel elle vise à obtenir une information aussi précise que possible. La constante de Hubble, la vitesse de la lumière, le coefficient de diffraction d'un prisme sont autant d'exemples où le concept d'aléa ne fait pas sens. Cependant, ce que propose l'approche bayésienne se situe dans un autre plan qui ne remet pas en cause cette modélisation de la réalité. La loi *a posteriori* est utilisée comme un nouvel (et efficace) outil de résumé de l'information disponible sur θ ,

sans véritablement remettre en cause l'existence de ce paramètre inconnu et *non aléatoire*. En d'autres termes, la loi a posteriori est un outil de représentation de l'information disponible sur θ une fois les observations obtenues, elle reflète l'incertitude inhérente aux données et à l'expérimentateur plutôt qu'un aléa physique sur ce paramètre.

Cette distinction entre outil de représentation et véritable aléa a presque causé l'extinction de l'approche bayésienne lorsque, au tournant du siècle, et comme dans de nombreuses autres branches des mathématiques, Kolmogorov et ses collègues probabilistes ont formalisé la théorie des probabilités. Il devenait alors difficile de donner un sens mathématique à la loi a priori si elle ne correspond pas à un véritable phénomène aléatoire mais plutôt à une traduction de ce que l'expérimentateur est prêt à parier sur les valeurs possibles du paramètre.⁴

Cet aspect subjectif de la loi a priori, résultant du choix de l'expérimentateur, peut apparaître comme un élément réducteur de la perspective bayésienne, mais il permet de refléter les informations (et leur degré d'imprécision) dont dispose cet expérimentateur et donc conduit à une inférence plus précise et plus riche, de ce fait. Que deux expérimentateurs adoptent deux lois a priori différentes ne devrait pas porter plus à critique que le fait que deux expériences (donc deux séries d'observations différentes) conduisent à des fonctions de vraisemblance différentes. Il est fondamental d'observer que l'approche bayésienne ne dispose pas d'une *seule loi a priori*, une sorte de Graal qu'il conviendrait d'obtenir après de longues recherches! Cette vision, souvent observée dans la littérature appliquée, est du même ordre que celle qui attribue au paramètre θ un véritable caractère aléatoire. (Bien entendu, dans certains contextes comme ceux des modèles à effets aléatoires ou de la prévision, certaines composantes de θ sont effectivement considérées comme aléatoires, mais elles le sont aussi dans l'approche classique de ces modèles.) Dans l'approche bayésienne, comme dans dans la plupart des autres écoles de statistique, parler de la "vraie" valeur du paramètre fait sens.

Comme nous le verrons dans les sections suivantes, disposer d'une loi a posteriori permet de construire des procédures inférentielles (tests, intervalles de confiance, densités prédictives) de la manière la plus naturelle possible, ce qui explique entre autres la persistance de cette approche, contre vents et marées (McGrayne, 2011), depuis 250 ans.

2.3. Une justification par le principe de vraisemblance

Tandis que la vision de Fisher, fondée sur la vraisemblance, n'admet pas une extension vers une modélisation bayésienne, il existe un principe de vraisemblance, formalisé par Birnbaum (1962), dont découle comme implémentation première la méthodologie bayésienne. Bien que ce principe soit régulièrement remis en cause (Mayo, 2010), il sous-tend suffisamment cette approche pour que nous le rappelions brièvement ici.

4. C'est l'argument du *Dutch Book*, voir Berger (1985).

Le *principe de vraisemblance* impose à deux séries d'observations \mathbf{x}_1 et \mathbf{x}_2 sur le même paramètre θ avec des fonctions de vraisemblance proportionnelles, donc telles que, pour tout θ ,

$$\ell(\theta|\mathbf{x}_1) \propto \ell(\theta|\mathbf{x}_2),$$

de conduire à la même inférence sur ce paramètre θ . Cette situation se produit par exemple lorsqu'on compare un échantillonnage binomial, $m \sim \mathcal{Bin}(n, p)$, à un échantillonnage binomial négatif, $n \sim \mathcal{Neg}(m, p)$. Le nombre d'essais total n et le nombre de succès m sont les mêmes, mais l'aléa portant sur deux parties différentes, les fonctions de vraisemblance sont proportionnelles,

$$\begin{aligned} \ell_1(p|n) &= \binom{n}{m} p^m (1-p)^{n-m} \\ &\propto \binom{n-1}{m-1} p^m (1-p)^{n-m} = \ell_2(p|m). \end{aligned}$$

Bien que cet exemple puisse paraître très artificiel, il sous-tend la classe des problèmes de règles d'arrêt, où un échantillon (x_1, \dots, x_n) est de taille n aléatoire, la loi de n ne dépendant pas du paramètre θ . Le principe de vraisemblance implique alors qu'une inférence distinguant un échantillon (x_1, \dots, x_n) iid d'un échantillon (x_1, \dots, x_n) produit par une règle d'arrêt contredit le principe de vraisemblance. Plus généralement, des méthodologies fondées sur des propriétés fréquentistes comme les p -values ne s'accordent pas avec ce principe.

Birnbaum (1962) a démontré que le principe de vraisemblance découlait logiquement de la conjonction de deux principes, le principe de conditionalité et le principe d'exhaustivité, à savoir que, (i) si deux expériences sont possibles pour mesurer θ , seule compte l'expérience effectivement réalisée, et que (ii) l'inférence ne doit dépendre que des statistiques exhaustives.⁵ Comme le discutent Berger and Wolpert (1988), la mise en œuvre du principe de vraisemblance conduit naturellement à une construction bayésienne (qui ne dépend effectivement que de la fonction de vraisemblance), même s'il existe des méthodologies alternatives (comme le test du rapport de vraisemblance) respectant le principe de vraisemblance et couvrant certains aspects de l'inférence.

3. Sur la construction des lois a priori, sur leur calibration et des réponses aux critiques afférentes

Comme signalé ci-dessus, la loi a priori est choisie par l'expérimentateur et témoigne à un degré ou un autre d'un choix (comme peut l'être la sélection de la procédure inférentielle voire de la loi des observations). Ces lois sont souvent choisies dans des classes de lois standard pour faciliter leur utilisation, même après l'avènement de techniques de simulation plus puissantes. L'impact de ce choix sur la réponse inférentielle est non-nulle, même si elle disparaît à mesure

5. La réfutation de cette démonstration par Mayo (2010) semble découler d'une définition tautologique de la notion d'inférence.

que la taille de l'échantillon grandit, du fait de la consistance de l'approche bayésienne dans un grand nombre de situations. Il peut être évalué de manière analytique ou numérique, mais aussi comparé à des solutions dites de référence (Jeffreys, 1939; Berger and Bernardo, 1992), décrites ci-dessous.

3.1. Des lois conjuguées comme éléments de base de la modélisation a priori

Dans le cadre de lois d'observations standard comme la loi normale, il existe des familles de lois a priori facilitant le calcul de la loi a posteriori. Bien que leur motivation soit surtout fondée sur cette simplification (plutôt que sur une caractéristique particulière de l'information a priori), elles n'en constituent pas moins un élément de base dans la construction effective de lois a priori plus personnelles. Par ailleurs, elles sont l'élément constitutif du logiciel BUGS (Lunn et al., 2000, 2010), car elles lui permettent de construire automatiquement l'algorithme de Gibbs approprié.⁶

Les lois conjuguées sont naturellement associées aux lois de familles exponentielles :

$$f(x|\theta) = h(x) \exp\{\theta \cdot R(x) - \psi(\theta)\},$$

où θ est un paramètre de dimension d et R une fonction à valeurs dans \mathbb{R}^d , $\psi(\theta)$ permettant la normalisation de la densité de probabilité. Ainsi, ces densités sont associées aux lois a priori de la forme

$$\pi(\theta) \propto \exp\{\theta \cdot \rho - \lambda\psi(\theta)\},$$

où $\lambda > 0$ et $\rho \in \mathbb{R}^d$ sont contraints par l'intégrabilité de la fonction, puisque

$$\pi(\theta|x) \propto \exp\{\theta \cdot [\rho + R(x)] - [\lambda + 1]\psi(\theta)\}.$$

La loi a posteriori est donc définie par un changement de paramètres, de λ à $\lambda + 1$, et de ρ à $\rho + R(x)$. Un point important est que ces lois conjuguées ne peuvent exister *que* dans le cas des familles exponentielles (Robert, 2005). Elles sont donc disponibles pour un grand nombre de lois classiques, comme la loi normale à un ou deux paramètres, les lois binomiale, binomiale négative, de Poisson, gamma à un ou deux paramètres, Dirichlet, de Wishart... Par exemple, pour des observations x_1, \dots, x_n normales de loi $\mathcal{N}(\theta, \sigma^2)$, où σ est connu, la loi conjuguée est aussi normale $\theta \sim \mathcal{N}(\mu, \tau^2)$. En effet, la loi de θ a posteriori est la loi normale

$$\theta|x_1, \dots, x_n \sim \mathcal{N}(\{\tau^{-2} + n\sigma^{-2}\}^{-1}\{\tau^{-2}\mu + n\sigma^{-2}\bar{x}_n\}, \{\tau^{-2} + n\sigma^{-2}\}^{-1}).$$

Les lois conjuguées ayant une forme *imposée* par la loi des observations, elles ne sont pas à même de refléter toute sorte d'information a priori, tout en demandant la spécification des hyperparamètres λ et ρ . Une extension efficace et

6. Nous rappelons d'une part que *BUGS* signifie *Bayesian analysis Using Gibbs Sampling* et d'autre part que le physicien Josiah Willard Gibbs n'a rien à voir avec les algorithmes de simulation conditionnelles qui portent à présent son nom. Ces algorithmes ont simplement été utilisés pour la première fois sur des champs de Gibbs (Geman and Geman, 1984).

robuste de ces lois est fournie par les mélanges, discrets ou continus, obtenus par l'introduction de lois sur les paramètres λ et ρ , en particulier parce qu'elles autorisent à leur tour des implémentations par l'algorithme de Gibbs (Diaconis and Ylvisaker, 1979).

3.2. Sur les lois de Jeffreys et de l'impossibilité de trouver la loi la moins informative

Dans une situation où l'information a priori n'est pas disponible, et où l'expérimentateur renacle à faire un choix, il est tentant de chercher à proposer une loi reflétant ce manque d'information. Hélas (ou pas!), il n'existe pas de loi "la moins informative" et on peut seulement, au mieux, définir une procédure automatique de construction d'une loi a priori de référence, à l'aune de laquelle les lois a priori subjectives peuvent être évaluées. Il s'agit donc d'une convention et non d'une optimalité quelconque au titre de la faible information a priori. Même si ces lois sont souvent appelées non-informatives, elles contiennent néanmoins des informations sur le paramètre et ne peuvent prétendre représenter une complète ignorance (Kass and Wasserman, 1996).

Une illustration de cette impossibilité est fournie par le principe de la raison insuffisante de Laplace (1812) : par extension du cas d'un ensemble fini de paramètres, Laplace propose d'utiliser une loi uniforme dans toute situation non-informative. Si l'espace des paramètres n'est pas compact, cela exige l'emploi d'une mesure de Lebesgue en lieu et place d'une loi de probabilité, ce qui ne pose pas problème en soi du moment que la loi a posteriori est définie, et surtout la solution dépend de la paramétrisation adoptée, puisque la loi uniforme ne "résiste" pas à un changement de variable. Bien que cette solution reste celle adoptée tout au long du 19^{ème} siècle, elle sert de focus aux critiques naissantes sur le paradigme bayésien.

Il faut en fait attendre les années 1930 et Harold Jeffreys⁷ pour voir apparaître une perspective globale sur le choix des lois de référence. Ironiquement, cette solution emprunte à Fisher en ce qu'elle est fondée sur l'information du même nom. Si l'information de Fisher associée à un modèle est définie par

$$I(\theta) = \mathbb{E}_{\theta} \left[\frac{\partial \ell}{\partial \theta^t} \frac{\partial \ell}{\partial \theta} \right]$$

la loi de Jeffreys associée est donnée par

$$\pi^*(\theta) \propto |I(\theta)|^{1/2}$$

soit donc la racine du déterminant de l'information. Plusieurs remarques sont de mise :

7. Harold Jeffreys fut à la fois mathématicien, statisticien, géophysicien et astronome. Son livre, *Theory of Probability*, demeure une référence sur la formalisation de l'approche bayésienne, écrite à une époque où celle-ci n'était plus si populaire, même si certains aspects du livre sont criticables sur le plan mathématique (Robert et al., 2009).

1. la solution de Jeffreys est invariante par reparamétrisation puisque le jacobien du changement de variable annule celui du changement d'information ;
2. la loi est directement dépendante de l'information, ce qui signifie que les régions où l'information est plus importante sont privilégiées, car les données y sont plus discriminantes ;
3. cette construction reproduit en général les lois invariantes sous l'action d'un groupe, les mesures de Haar ;
4. cette construction ne garantit en rien l'intégrabilité de π^* , ce qui signifie que les lois de Jeffreys seront souvent des mesures ;
5. les lois de Jeffreys dépendent de l'intégralité de la loi des observations, donc ne respectent pas le principe de vraisemblance.

Une illustration de ce dernier point est fournie par l'opposition entre la loi binomiale, $\mathcal{B}(n, \theta)$ pour qui la loi de Jeffreys est la loi $\mathcal{Be}(1/2, 1/2)$, et la loi binomiale négative, $\mathcal{Neg}(x, \theta)$, pour qui la loi de Jeffreys est la mesure (impropre) de densité $1/\theta\sqrt{1-\theta}$.

Comme cette construction universelle de lois de référence peut donner naissance à des “monstres”, par exemple dans le cas de l'estimation de $\|\theta\|^2$ quand $x \sim \mathcal{N}(\theta, I_p)$ et p grand, il existe des déterminations de lois de référence qui distinguent entre paramètres d'intérêt et paramètres de nuisance avant de construire des lois de Jeffreys conditionnelles et marginales sur les deux groupes. Nous renvoyons le lecteur à Berger and Bernardo (1992) et Bernardo and Smith (1994) pour des entrées sur cette extension qui demeure peu utilisée à ce jour. Voir aussi Kass and Wasserman (1996) pour une revue de certains principes de détermination des lois “non-informatives”. En lien avec ces principes, insistons une nouvelle fois sur le fait que, de même qu'il n'existe pas “une” loi a priori unique qu'il faudrait découvrir, il n'existe pas non plus “une” seule loi a priori non-informative. Il s'agit bien de poser une référence, qui peut servir à la fois à analyser des problèmes sans information visible et à comparer différentes lois a priori, si besoin. (Cette référence se trouve faire défaut dans le cadre des tests, à moins d'accepter quelques compromis avec le paradigme bayésien, Robert, 2005)

3.3. Sur l'utilisation de lois en dimension infinie et de l'apparition de restaurants ethniques

Nous avons mentionné dans l'Introduction que l'approche bayésienne permettait également de traiter de problèmes dans un cadre non-paramétrique. Cela signifie par exemple que, à partir d'un échantillon iid x_1, \dots, x_n de loi inconnue⁸ F , une loi a posteriori peut être construite sur F . Les lois a priori qui sous-tendent ces constructions sont donc des lois sur des espaces fonctionnels.

8. Notons que cette hypothèse est peu restrictive en ce qu'elle correspond *in fine* à celle d'échangeabilité sur la loi du n -uplet, suivant la représentation de Bruno de Finetti. Voir par exemple Bernardo and Smith (1994).

Par exemple, la loi de Dirichlet *fonctionnelle*, $Dir(\alpha, G_0)$, est définie à partir d'un coefficient de précision $\alpha > 0$ et d'une moyenne a priori G_0 , loi de probabilité. Cette loi (Ferguson, 1974) généralise la loi de Dirichlet sur le simplexe $Dir_d(\beta_1, \dots, \beta_d)$ au sens où, pour toute partition $\{A_1, \dots, A_d\}$ de l'espace des observations, si $G \sim Dir(\alpha, G_0)$, alors

$$(F(A_1), \dots, F(A_d)) \sim Dir_d(\alpha G_0(A_1), \dots, \alpha G_0(A_d)).$$

La loi a posteriori correspondante est conjuguée,

$$F|x_1, \dots, x_n \sim Dir(\alpha + n, G_0 + \hat{F}_n),$$

où \hat{F}_n dénote la distribution empirique associée à l'échantillon x_1, \dots, x_n , ce qui signifie que la variable aléatoire n'est pas absolument continue par rapport à la mesure de Lebesgue, même si G_0 l'est. En particulier, si $F \sim \mathcal{D}(\alpha, G_0)$, la loi conditionnelle de x_1 sachant (x_2, \dots, x_n) est de la forme

$$\frac{\alpha}{\alpha + n - 1} F_0 + \frac{1}{\alpha + n - 1} \sum_{i=2}^n \delta_{x_i}.$$

L'évaluation complète de la loi a posteriori de F exige l'emploi d'outils de simulation qui ne seront pas abordés ici (voir par exemple Holmes et al., 2002 et Hjort et al., 2010).

Une représentation équivalente du processus de Dirichlet est appelée *processus du restaurant chinois* pour la raison suivante : générer un échantillon de taille n suivant le processus de Dirichlet est équivalent à simuler un flux d'arrivées de clients dans un restaurant chinois stylisé. Chaque client de ce restaurant s'assied à une table déjà occupée avec une probabilité en proportion du nombre de clients présents à cette table et à une nouvelle table avec probabilité proportionnelle à α . D'autres versions du processus de Dirichlet sont fondées sur des successions de partitions de l'intervalle $[0, 1]$ appelées "stick breaking" en raison de leur construction récurrente (Sethuraman, 1994). Ces différentes interprétations sont utilisées soit au niveau théorique, pour valider la convergence des estimateurs résultants (van der Vaart, 1998; Yang and Barron, 1999; Ghosal et al., 2008), soit au niveau computationnel, pour simuler ces processus (Hjort et al., 2010; Müller and Mitra, 2013).

De nombreuses extensions de cette loi a priori existent dans la littérature, fondées sur différents processus comme les processus gaussiens ou le processus de Lévy. Des représentations par mélanges permettent en particulier de dépasser le défaut des processus de Dirichlet de ne simuler que des lois à support discret. Une extension de la représentation du restaurant chinois s'appelle le buffet indien (!) et permet le mélange de plusieurs caractéristiques au lieu d'imposer une partition (Weiss et al., 2006).

Enfin, notons que les outils de modélisation d'espaces de fonctions comme les bases d'ondelettes autorisent une autre forme de représentation bayésienne de l'inférence fonctionnelle (Vidakovic, 1999; Clyde and George, 2000).

4. De la mise en œuvre des principes bayésiens

Nous décrivons dans cette partie quelques mises en œuvre de l'inférence bayésienne dans des modèles standard. Il s'agit en grande partie d'illustrations et nous renvoyons le lecteur par exemple à Robert (2005) pour une perspective plus complète.

4.1. De la loi a posteriori comme pivot de l'inférence

Si la loi a priori $\pi(\theta)$ est l'élément central qui concentre l'essentiel des critiques sur la perspective bayésienne, la loi a posteriori $\pi(\theta|x)$ est l'élément central permettant de conduire l'inférence bayésienne. Une fois cette loi construite, et le modèle accepté (voir ci-dessous), il n'est plus besoin de considérer un autre aspect des données! En effet, les diverses procédures inférentielles sont toutes construites (automatiquement) à partir de cette loi : estimateurs, variations, intervalles de confiance, tests, choix de variables, choix de modèle, prévisions, etc. Certains défendent en fait la définition de la loi a posteriori comme ultime but de l'inférence bayésienne, considérant que les procédures ci-dessus ne sont que des résumés qui dégradent le contenu informatif de la loi a posteriori. Bien qu'une loi de probabilité ne puisse en effet se résumer à certains de ses moments, il est néanmoins pertinent de fournir des réponses aux questions des clients et des décideurs, pour lesquels la loi a posteriori demeure un concept abstrait. Cette perspective pragmatique est particulièrement pertinente dans les problèmes à paramètres de grande dimension, dont la plupart correspondent à des paramètres de nuisance.⁹

La loi a posteriori $\pi(\theta|x)$ étant construite, on peut dériver (analytiquement ou non) les estimateurs ponctuels que sont la moyenne et de la médiane a posteriori de n'importe quelle transformation $\varphi(\theta)$, ainsi que les évaluations de leur variabilité fournies par les écart-types a posteriori. De plus, la loi a posteriori induit des lois a posteriori marginales pour toute transformation $\varphi(\theta)$, déduites par projection du modèle probabiliste joint, $\varphi \sim \pi(\varphi|x)$. De ces lois marginales peuvent se déduire des intervalles ou des régions de confiance naturelles, les régions de plus forte densité a posteriori

$$\mathcal{C}(x) = \{\varphi; \pi(\varphi|x) > \tau\},$$

dont la borne τ peut se calibrer en fonction du taux de couverture α souhaité pour cette région. (Elle dépendra alors de x .) Contrairement aux intervalles fondés sur l'approximation normale des "deux sigma", ces régions ont une couverture exacte et peuvent être asymétriques. Elles sont par contre dépendantes de la paramétrisation choisie (ou, ce qui est équivalent, de la mesure utilisée pour mesurer les volumes, voir Druilhet and Marin, 2007). Il est même possible

9. Répétons ici la remarque que la prise en compte de la distinction entre paramètres d'intérêt et paramètres de nuisance a conduit Berger and Bernardo (1992) à proposer des lois de référence reproduisant cette distinction et éliminant les paramètres de nuisance par une intégration suivant une loi de Jeffreys conditionnelle.

de discuter de la vraisemblance d'une hypothèse comme $H_0 : \varphi = 0$ en examinant si la valeur $\varphi = 0$ appartient à la région HPD, bien que cette perspective ne soit pas conseillée puisqu'elle ne prend pas en compte l'action résultant d'un rejet de H_0 .

Comme simple illustration, prenons l'exemple d'un phénomène binaire, à valeurs dans $\{0, 1\}$, observé durant n répétitions indépendantes et identiquement distribuées (i.i.d.). Le nombre de 1, X , est alors modélisé par une loi binomiale $\mathfrak{B}(n, p)$ où p est la probabilité d'obtenir une valeur 1. Si on associe à ce paramètre p une loi de Jeffreys,

$$\pi(p) = 1/\sqrt{p(1-p)},$$

qui correspond à une loi Béta $\mathfrak{B}e(1/2, 1/2)$ une fois normalisée, la construction de la loi a posteriori est immédiate :

$$\begin{aligned} \pi(p|x) &\propto \pi(p) \times \ell(p|x) \\ &\propto p^x (1-p)^{n-x} / \sqrt{p(1-p)} \\ &\propto p^{x-1/2} (1-p)^{n-x-1/2}, \end{aligned}$$

qui correspond à une loi Béta $\mathfrak{B}e(x + 1/2, n - x + 1/2)$ (toujours après normalisation). Si, par exemple, $n = 232$ et $x = 117$, la loi a posteriori est la loi Béta $\mathfrak{B}e(117.5, 115.5)$, de mode et moyenne $117.5/233 = 0.504$. L'intervalle de crédibilité à queues égales, valant $[0.440, 0.558]$ pour un niveau de crédibilité $\alpha = 0.95$, est plus facile à calculer que la région HPD.¹⁰ Si à présent nous cherchons à tester l'hypothèse nulle et ponctuelle $H_0 : p = 1/2$, le facteur de Bayes en faveur de H_0 (voir Section 4.3) est donné par

$$\mathfrak{B}_{01}(x) = 1/2^n \left/ \frac{B(x + 1/2, n - x + 1/2)}{B(1/2, 1/2)} \right.,$$

où $B(\alpha, \beta)$ dénote la constante de normalisation de la densité de la loi $\mathfrak{B}e(\alpha, \beta)$. Dans le même exemple numérique que ci-dessus, ce facteur de Bayes vaut $\pi(1/2|x)/\pi(1/2) = 18.95$, donc exprime un soutien assez important en faveur de H_0 .

Même si les estimateurs ci-dessus sont des produits de la loi a posteriori, il est légitime de se demander comment les comparer et quel estimateur choisir comme "meilleur résumé". À l'exception de cas spécifiques, il n'existe pas de réponse à cette question sans passer par une modélisation de la décision prise, qui donne un sens au terme "meilleur".

4.2. De l'évaluation des procédures par des fonctions de coût et de la dualité avec les lois a priori

Il n'est guère difficile de se convaincre qu'il n'existe pas d'estimateur optimal à l'exception du trivial $\hat{\theta}(x) = \theta$. Suivant les perspectives adoptées pour classer

10. Une résolution numérique plus poussée de cette dernière conduit à... $[0.440, 0.558]$, soit donc le même intervalle (pour cette précision) que la solution symétrique!

les estimateurs, certains seront préférables à d'autres mais très très rarement meilleurs que tous les autres. C'est du moins le cas dans le cadre fréquentiste où coexistent en général des classes d'estimateurs faiblement optimaux. À l'opposé, la perspective bayésienne permet de définir une forme plus forte d'optimalité et d'aboutir à une solution *unique*.

La dérivation de cette propriété repose sur la notion de fonction de perte ou de coût, empruntée à la théorie des jeux, $L(d, \theta)$, qui évalue l'erreur ou la pénalité résultant de la décision d lorsque la véritable valeur du paramètre est θ . Par exemple, si d est la décision d'estimer θ , la formalisation de l'erreur peut être

$$L(d, \theta) = \sum_{i=1}^n |\theta_i - d_i| \quad \text{ou} \quad L(d, \theta) = \max_{i=1, \dots, n} |\theta_i - d_i|.$$

Même si des problèmes véritables peuvent conduire à une détermination unique de la fonction de coût, dictée par des impératifs financiers par exemple, il est plus fréquent qu'on doive adopter des fonctions de coût abstraites comme celles ci-dessus ou la plus traditionnelle, déjà adoptée par Legendre et Laplace,

$$L_2(d, \theta) = \sum_{i=1}^n (\theta_i - d_i)^2,$$

dont le principal intérêt est de fournir des solutions explicites (voir ci-dessous).

Étant donné un problème modélisé par une fonction à valeurs réelles $L(d, \theta)$, un estimateur δ propose une décision $\delta(x)$ (ou estimation) pour chaque observation x . La comparaison des estimateurs ne peut se faire au travers de $L(\delta(x), \theta)$, puisque x est aléatoire et θ est inconnu. Tandis que l'approche classique évalue les estimateurs au travers de l'erreur moyenne $\mathbb{E}_\theta[L(\delta(X), \theta)]$ ou risque, l'approche bayésienne de la théorie de la décision consiste à prendre l'erreur moyenne par rapport aux paramètres,

$$\varrho(\delta, x) = \int L(\delta(x), \theta) \pi(\theta|x) d\theta.$$

L'intérêt fondamental de cette perspective est qu'elle permet à la fois d'éliminer le paramètre inconnu et de conditionner par rapport à l'observation x . Les estimations $\delta(x)$ sont alors toutes comparables à x donné, ce qui permet d'obtenir la meilleure décision. Un *estimateur de Bayes* est ainsi défini comme la fonction qui associe à tout x la décision

$$\delta^\pi(x) = \arg \min_d \int L(d, \theta) \pi(\theta|x) d\theta.$$

Par exemple, le coût L_2 ci-dessus conduit à un estimateur de Bayes égal à la moyenne a posteriori, $\mathbb{E}[\theta|x]$. Cet estimateur classique est cependant peu recommandé dans le cas de lois a posteriori multimodales, comme dans le cas des mélanges de distribution (Lee et al., 2009). Sous certaines hypothèses, dont la convexité stricte de la fonction de coût et l'absolue continuité de la loi a priori,

l'estimateur résultant est unique. Il a de plus la propriété de minimiser le risque intégré,

$$\int \mathbb{E}_\theta[L(\delta(X), \theta)] d\theta,$$

ce qui lui permet de plus de bénéficier de propriétés fréquentistes d'optimalité comme la minimaxité et l'admissibilité (Robert, 2005). Bien que cette approche de la construction des procédures statistiques optimales produise des arguments additionnels pour la justification de l'approche bayésienne (les théorèmes de classes complètes de Wald, 1950 démontrant ainsi que les seuls estimateurs admissibles sont les estimateurs de Bayes), le recours à une perspective décisionnelle n'est guère suivie dans les ouvrages récents de statistique bayésienne (voir par exemple Gelman et al., 2003). Mon point de vue est qu'il est quelque peu regrettable de ne pas prendre en compte les aspects décisionnels, car ils conduisent à une formulation rigoureuse du choix des estimateurs et illustrent la dualité entre loi a priori et fonction de coût. En effet, dans la minimisation de

$$\int L(\delta(x), \theta)\pi(\theta|x) d\theta,$$

seul compte le produit entre fonction de coût et loi a priori. Le transfert de masse entre L et π n'a donc aucun effet sur la décision finale. Plus fondamentalement, connaître les régions de l'espace des paramètres où une erreur est la plus dommageable est similaire à favoriser dans la loi a priori cette région. C'est d'ailleurs une manière d'aboutir à la loi de Jeffreys à partir d'une fonction de coût intrinsèque (Robert et al., 2009).

4.3. De la spécificité des tests en inférence bayésienne, de l'opposition aux p-values, et des deux versions du paradoxe de Lindley

Le cas particulier des tests statistiques, si souvent utilisés en statistique classique pour démontrer des "différences significatives", met en lumière une distinction majeure entre cette approche classique et l'approche bayésienne, distinction qui les rend incompatibles sur ce plan. Dans l'approche classique (des tests), il s'agit d'identifier des événements improbables, c'est à dire des observations incompatibles avec une certaine loi de probabilité (ou une famille de lois). Par exemple, dans l'illustration binomiale ci-dessus, on recherche les valeurs de X les plus improbables lorsque $p = 1/2$. Dans cette approche, l'hypothèse alternative y joue un rôle mineur voire négligeable, servant au mieux à définir la zone des statistiques extrêmes. L'approche bayésienne vise à une résolution complète en modélisant les deux hypothèses simultanément, permettant d'une part la résolution du problème décisionnel du choix de l'hypothèse la plus probable (choix que certains bayésiens récuse) et d'autre part la construction d'une loi a posteriori sur les paramètres du modèle retenu. Le choix d'une hypothèse comme $H_0 : p = 1/2$ se fait donc *contre* ou relativement à une hypothèse alternative comme $H_1 : p > 1/2$ et non pas dans un absolu mal défini et évitant la

modélisation des contraires (et des décisions à prendre en cas de rejet de l'hypothèse nulle). Cette opposition entre approche classique et approche bayésienne est souvent mal comprise, d'une part à cause du jugement relatif associé à la seconde et d'autre part du fait de la validation fréquentiste du taux d'erreur de première espèce des tests de Neyman-Pearson. Ce dernier est en fait très souvent interprété à tort comme la probabilité de l'hypothèse nulle, ce qui n'a pas de sens en dehors d'une perspective bayésienne.¹¹ Au delà de l'opposition philosophique, et des interprétations différentes des résultats, les deux principes conduisent aussi à des oppositions claires dans les décisions, en particulier parce que les tests classiques tendent à rejeter plus fréquemment les hypothèses nulles.

La procédure de décision bayésienne est fondée sur la probabilité a posteriori de l'hypothèse nulle

$$\pi(H_0|x) = \frac{\int_{H_0} \pi(\theta)\ell(\theta|x) d\theta}{\int \pi(\theta)\ell(\theta|x) d\theta}.$$

ou, de manière (presque) équivalente,¹² sur le facteur de Bayes (Jeffreys, 1939)

$$\mathfrak{B}_{01}(x) = \frac{\int_{H_0} \pi(\theta)\ell(\theta|x) d\theta}{\int_{H_1} \pi(\theta)\ell(\theta|x) d\theta} \bigg/ \frac{\int_{H_0} \pi(\theta) d\theta}{\int_{H_1} \pi(\theta) d\theta}.$$

Ce dernier se compare à 1 comme la probabilité a posteriori se compare à 1/2. Dans l'exemple binomial, nous avons ainsi vu une valeur de $\mathfrak{B}_{01}(x) = 18.95$ qui nous permet de conclure en faveur de H_0 . (La probabilité a posteriori serait alors 0.95.) Dans un cadre classique, la p -value est la probabilité de dépasser $x = 117$ sous l'hypothèse nulle, soit $\mathbb{P}(X \geq 117|H_0) = 0.42$, qui s'avère beaucoup moins favorable à H_0 .

L'opposition mentionnée ci-dessus est des plus claires dans une étude menée par Berger and Sellke (1987) : ils démontrent que la borne inférieure des probabilités a posteriori est supérieure à la p -value (les deux quantités sont comparables d'un point de vue décisionnel, en utilisant une fonction de coût 0–1, voir Robert, 2005, Chapitre 5). Cela signifie que pour toute loi a priori donnant le même poids aux deux hypothèses, l'hypothèse nulle H_0 sera toujours jugée plus improbable par l'approche classique. Berger and Sellke (1987) remarquent par exemple que, en simulant les données et paramètres suivant une loi a priori spécifique, dans 20% des cas où l'hypothèse nulle est rejetée, elle est en fait correcte.

Plus clairement encore, le paradoxe de Lindley (1957) exprime cette opposition : dans le modèle normal de moyenne θ avec n observations, lorsque $H_0 : \theta = 0$, pour une p -value donnée, $p(\bar{x}_n) = \alpha$, sous une loi a priori normale arbitraire, la probabilité a posteriori de l'hypothèse nulle tend vers 1 avec la taille de l'échantillon n . Les interprétations de ce paradoxe abondent, allant du rejet des p -values (car pourquoi faudrait-il garder cette p -value ou le seuil d'ac-

11. L'approche bayésienne des tests est souvent qualifiée de "holmesienne" en référence à la citation de Sherlock Holmes sur la sélection de l'hypothèse la moins improbable : "When you have eliminated the impossible, whatever remains, however improbable, must be the truth".

12. Le "presque" est dû au fait que le facteur de Bayes ne dépend plus des probabilités a priori des deux hypothèses en jeu.

ception constant ?) à celui des probabilités a posteriori (Spanos, 2013), voire les deux (Sprenger, 2013).¹³

Une autre spécificité des tests bayésiens est de nécessiter une attention particulière envers les lois impropres. D'un point de vue strictement mathématique, ces lois ne peuvent être employées car elles introduisent une constante de normalisation arbitraire dans le calcul de la loi a posteriori. Cela conduit à la seconde forme du paradoxe de Lindley (1957) : étant donné une observation x arbitraire d'une loi normale, la probabilité a posteriori de l'hypothèse nulle $H_0 : \theta = 0$ tend vers 1 avec la variance a priori sous l'hypothèse alternative.

Il existe de nombreuses réponses à ce dilemme, de l'interdiction globale des lois impropres (DeGroot, 1970) à diverses stratégies de validation croisée (ou non-croisée) utilisant une partie des données pour contruire une véritable loi (a priori ? a posteriori ? a mediori ?) et le reste des données pour effectuer le test (Berger and Pericchi, 1996), voire utilisant *toutes* les données pour contruire une loi et les réutilisant une seconde fois pour le test (Aitkin, 1991, 2010; Gelman et al., 2013). Je ne mentionnerai pas ici les alternatives sur des "critères d'information" comme le BIC (pour *Bayesian Information Criterion*) et le DIC (pour *Deviance Information Criterion*), qui fournissent des classements de modèles ou d'hypothèses en comparaison en dehors des cadres décisionnel *et* bayésien.¹⁴

5. Envoi

Ce bref chapitre n'a pas abordé de nombreux points et modèles où l'approche bayésienne fait sens, soit pour des raisons de modélisation (comme l'émergence des modèles graphiques à la fin des années 1980, définis uniquement par des lois conditionnelles, voir par exemple Lauritzen, 1996), soit pour des raisons de complexité générique des modèles (comme dans le cadre de modèles hiérarchiques rencontrés en marketing comme en sélection animale, en climatologie comme en épidémiologie), soit encore grâce à sa démarche intégrée permettant de mêler description stochastique et impératifs de décision (comme par exemple dans la construction de plans d'expérience dans Müller et al., 2004). Je n'ai pas mentionné non plus les nouvelles capacités de l'analyse bayésienne à la frontière entre statistique et machine learning, comme le classificateur BART de Chipman et al. (2010), techniques qui, associées avec des méthodologies non-paramétriques aussi rapides, autorisent le traitement de grands jeux de données, la norme de l'ère du "Big Data".

Bien que la statistique bayésienne ait principalement grandi dans un terreau philosophique et polémique, de Bayes et Price s'opposant à Hume, à Jeffreys (1939) contre Fisher, Keynes (1920) rejetant Laplace (Robert, 2011), &tc., il est assez paradoxal que l'émergence d'une composante bayésienne significative dans la statistique (ou science des données ?) moderne soit surtout dû à sa

13. Comme je le souligne dans Robert (2013), l'opposition est naturelle et doit être relativisée par l'existence de solutions convergentes dans les deux approches.

14. Ces critères sont néanmoins populaires car (i) ne nécessitant pas de loi a priori pour BIC et (ii) disponible directement sur le logiciel BUGS pour DIC (Lunn et al., 2010).

maniabilité face à des modèles complexes et à l'émergence de techniques de calcul par ordinateur finalement assez simples à comprendre et implémenter pour que des disciplines utilisatrices de la statistique s'en saisisse. On peut le regretter du point de vue des fondations de la discipline, mais la dissémination des idées bayésiennes dans les autres sciences et au delà est aussi à même de fournir de nouvelles perspectives et propositions à même d'assurer la continuation de la discipline dans les décennies à venir : *From practice stems theory!*

Références

- AITKIN, M. (1991). Posterior Bayes factors (with discussion). *J. Royal Statist. Society Series B*, **53** 111–142.
- AITKIN, M. (2010). *Statistical Inference : A Bayesian/Likelihood approach*. CRC Press, Chapman & Hall, New York.
- BERGER, J. (1985). *Statistical Decision Theory and Bayesian Analysis*. 2nd ed. Springer-Verlag, New York.
- BERGER, J. and BERNARDO, J. (1992). On the development of the reference prior method. In *Bayesian Statistics 4* (J. Berger, J. Bernardo, A. Dawid and A. Smith, eds.). Oxford University Press, London, 35–49.
- BERGER, J. and PERICCHI, L. (1996). The intrinsic Bayes factor for model selection and prediction. *J. American Statist. Assoc.*, **91** 109–122.
- BERGER, J. and SELLKE, T. (1987). Testing a point-null hypothesis : the irreconcilability of significance levels and evidence (with discussion). *J. American Statist. Assoc.*, **82** 112–122.
- BERGER, J. and WOLPERT, R. (1988). *The Likelihood Principle (2nd edition)*, vol. 9 of *IMS Lecture Notes — Monograph Series*. 2nd ed. IMS, Hayward.
- BERNARDO, J. and SMITH, A. (1994). *Bayesian Theory*. John Wiley, New York.
- BIRNBAUM, A. (1962). On the foundations of statistical inference. *J. American Statist. Assoc.*, **57** 269–306.
- CHIPMAN, H., GEORGE, E. and MCCULLOCH, R. (2010). Bart : Bayesian additive regression trees. *Annals of Applied Statistics*, **4** 266–298.
- CLYDE, M. and GEORGE, E. I. (2000). Flexible empirical Bayes estimation for wavelets. *J. Royal Statist. Society Series B*, **62** 681–698.
- DALE, A. I. (1999). *A History of Inverse Probability*. Springer-Verlag, New York. (Second edition.).
- DEGROOT, M. (1970). *Optimal Statistical Decisions*. McGraw-Hill, New York.
- DIACONIS, P. and YLVIKAKER, D. (1979). Conjugate priors for exponential families. *Ann. Statist.*, **7** 269–281.
- DRUILHET, P. and MARIN, J.-M. (2007). Invariant hpd credible sets and map estimators. *Bayesian Analysis*, **2(4)** 681–692.
- FERGUSON, T. (1974). Prior distributions in spaces of probability measures. *Ann. Statist.*, **2** 615–629.
- FIENBERG, S. (2006). When did Bayesian inference become “Bayesian“? *Bayesian Analysis*, **1(1)** 1–40.

- GELMAN, A., CARLIN, J., STERN, H., DUNSON, D., VEHTARI, A. and RUBIN, D. (2013). *Bayesian Data Analysis*. 3rd ed. Chapman and Hall, New York, New York.
- GELMAN, A., CARLIN, J., STERN, H. and RUBIN, D. (2003). *Bayesian Data Analysis*. 2nd ed. Chapman and Hall, New York, New York.
- GEMAN, S. and GEMAN, D. (1984). Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Trans. Pattern Anal. Mach. Intell.*, **6** 721–741.
- GHOSAL, S., LEMBER, J. and VAN DER VAART, A. (2008). Nonparametric Bayesian model selection and averaging. *Electronic J. of Statistics*, **2** 63–89.
- HJORT, N., HOLMES, C., MÜLLER, P. and WALKER, S. (2010). *Bayesian nonparametrics*. Cambridge University Press.
- HOLMES, C., DENISON, D., MALLICK, B. and SMITH, A. (2002). *Bayesian methods for nonlinear classification and regression*. John Wiley, New York.
- JEFFREYS, H. (1939). *Theory of Probability*. 1st ed. The Clarendon Press, Oxford.
- KASS, R. and WASSERMAN, L. (1996). Formal rules of selecting prior distributions : a review and annotated bibliography. *J. American Statist. Assoc.*, **91** 343–1370.
- KEYNES, J. (1920). *A Treatise on Probability*. Macmillan and Co., London.
- LAPLACE, P. (1812). *Théorie Analytique des Probabilités*. Courcier, Paris.
- LAURITZEN, S. (1996). *Graphical Models*. Oxford University Press, Oxford.
- LEE, K., MARIN, J.-M., MENGENSEN, K. and ROBERT, C. (2009). Bayesian inference on mixtures of distributions. In *Perspectives in Mathematical Sciences I : Probability and Statistics* (N. N. Sastry, M. Delampady and B. Rajeev, eds.). World Scientific, Singapore, 165–202.
- LINDLEY, D. (1957). A statistical paradox. *Biometrika*, **44** 187–192.
- LUNN, D., THOMAS, A., BEST, N. and SPIEGELHALTER, D. (2000). WinBUGS – a Bayesian modelling framework : concepts, structure, and extensibility. *Statist. Comput.*, **10** 325–337.
- LUNN, D., THOMAS, A., BEST, N. and SPIEGELHALTER, D. (2010). *The BUGS Book : A Practical Introduction to Bayesian Analysis*. Chapman & Hall/CRC Press.
- MAYO, D. (2010). An error in the argument from conditionality and sufficiency to the likelihood principle. In *Error and inference : recent exchanges on experimental reasoning, reliability, and the objectivity and rationality of science* (D. Mayo and A. Spanos, eds.), chap. 7 (III). Cambridge University Press, 305–314.
- MCGRAYNE, S. (2011). *The Theory that Would Not Die*. Yale Univ Press, New Haven, CT.
- MÜLLER, P., PARMIGIANI, G., ROBERT, C. and ROUSSEAU, J. (2004). Optimal sample size for multiple testing : the case of gene expression microarrays. *J. American Statist. Assoc.*, **99** 990–1001.
- MÜLLER, P. and MITRA, R. (2013). Bayesian nonparametric inference – why and how. *Bayesian Anal.*, **8** 269–302.
- ROBERT, C. (2005). *Le Choix Bayésien*. Springer-Verlag, Paris.

- ROBERT, C. (2011). Re-reading Keynes' treatise on probability. *International Statistical Review*, **79** 1–15.
- ROBERT, C. (2013). On the Jeffreys–Lindley's paradox. *Philosophy of Science*. (To appear).
- ROBERT, C., CHOPIN, N. and ROUSSEAU, J. (2009). Theory of Probability revisited (with discussion). *Statist. Science*, **24(2)** 141–172 and 191–194.
- SETHURAMAN, J. (1994). A constructive definition of Dirichlet priors. *Statistica Sinica*, **4** 639–650.
- SPANOS, A. (2013). Who should be afraid of the Jeffreys–Lindley paradox? *Philosophy of Science*, **80** 73–93.
- SPRENGER, J. (2013). Testing a precise null hypothesis : The case of Lindley's paradox. *Philosophy of Science*. (To appear).
- VAN DER VAART, A. (1998). *Asymptotic Statistics*. Cambridge University Press.
- VIDAKOVIC, B. (1999). *Statistical Modeling by Wavelets*. Wiley, New York.
- WALD, A. (1950). *Statistical Decision Functions*. John Wiley, New York.
- WEISS, Y., SCHÖLKOPF, B. and PLATT, J. (eds.) (2006). *Infinite Latent Feature Models and the Indian Buffet Process*, vol. 18. Cambridge, MA : MIT Press.
- YANG, Y. and BARRON, A. (1999). Information-theoretic determination of minimax rates of convergence. *Ann. Statist.*, **27** 1564–1599.