

Modélisation et statistique bayésienne computationnelle

Notes de cours MAJ

`nicolas.bousquet@sorbonne-universite.fr`

30 mars 2023

Master 2, Sorbonne Université, 2023



Résumé

Ce cours a pour objectif de présenter d'une part les principales méthodologies de modélisation bayésienne appliquées à des problèmes d'aide à la décision en univers risqué, et d'autre part les principales méthodes de calcul inférentiel permettant l'enrichissement de l'information utile, en fonction de l'emploi et de la nature des modèles. Il nécessite les pré-requis suivants : notions fondamentales de probabilités et statistique, introduction aux statistiques bayésiennes, méthodes de Monte-Carlo, calcul scientifique en `R` ou/ et en `Python`. Tout au long du cours, des liens avec l'apprentissage statistique (*machine learning*) sont présentés.

Ce document évolue au fil du temps, et comporte parfois des coquilles ou quelques contresens qui peuvent m'échapper. Certains de mes étudiants, par leurs remarques et parfois leur aide pour déceler ces coquilles, par leurs demandes d'éclaircissement, contribuent notablement à améliorer son propos et sa fluidité. Qu'ils en soient particulièrement remerciés, et plus spécialement Paul Liautaud (M2A 2022). Enfin, ce document fait parfois quelques emprunts graphiques à des ouvrages par ailleurs recommandés.

Table des matières

1	Notations	6
2	Introduction et rappels	9
2.1	Modélisation, inférence et décision statistique	9
2.2	Cadre statistique paramétrique	9
2.3	Estimation statistique classique ("fréquentiste")	10
2.3.1	Rappel des principes	11
2.3.2	Difficultés pratiques, théoriques et conceptuelles	11
2.4	Principes de la statistique bayésienne	13
2.4.1	Paradigme	13
2.4.2	Fondations théoriques	14
2.4.3	Plan du cours	15
2.5	Liens avec le <i>machine learning</i>	16
2.6	Quelques lectures conseillées	17
3	Éléments de théorie de la décision	18
3.1	Existence d'une fonction de coût	18
3.2	Supériorité des estimateurs de Bayes sur les estimateurs fréquentistes	21
3.3	Choix d'une fonction de coût	21
3.4	Coûts intrinsèques	23
3.5	Mode <i>a posteriori</i> (MAP)	24
3.6	Sélection de modèle et facteur de Bayes	24
3.7	TP : Création d'un système d'alerte pour la circulation routière	28
4	Propriétés fondamentales du cadre bayésien	29
4.1	Prédiction (prévision)	29
4.2	Propriétés asymptotiques	29
4.3	Régions de crédibilité et régions HPD	30
5	Compréhension et représentation de l'information incertaine	33
5.1	Une vision subjectiviste de la théorie bayésienne	33
5.2	Théories de la connaissance incertaine	33
5.3	Une vision plus claire de la statistique bayésienne	33
5.4	Validité des priors informatifs <i>via</i> la logique probabiliste de l'information incertaine	34
6	Modélisation <i>a priori</i>	38
6.1	Priors objectifs régularisant (priors peu ou "non informatifs")	38
6.1.1	Priors de Laplace (prior uniforme) et de Jeffreys	39
6.1.2	Prior de référence de Berger-Bernardo	41
6.1.3	Priors coïncidants ou concordants	42
6.2	Priors objectifs informatifs	43
6.2.1	Maximum d'entropie	43
6.2.2	Priors conjugués et famille exponentielle	45
6.3	Priors hiérarchiques	50
6.3.1	Un exemple utile : les <i>Latent Gaussian Models</i> (LGM)	51
6.4	Convergence des priors informatifs vers des priors régularisant	53
6.5	Démarches critiques d'éllicitation	55
6.5.1	Détecter et limiter les conflits entre prior et données	55
6.5.2	Fusionner plusieurs priors	61
6.6	TP : Un exemple complet dans un cadre de fiabilité industrielle	63
6.7	L'importance du prior en <i>deep learning</i>	64

7	Incorporation de connaissance <i>a priori</i>	65
7.1	Notion de crédibilité	65
8	Méthodes de calcul bayésien	67
8.1	Introduction	67
8.1.1	Principe de la simulation indirecte	68
8.2	Méthodes d'échantillonnage dans la loi <i>a posteriori</i>	69
8.2.1	Rappel : approches par inversion et transformations simples	70
8.2.2	Simulation multidimensionnelle	71
8.2.3	Algorithmes d'acceptation-rejet (AR)	72
8.2.4	Algorithmes d'échantillonnage préférentiel ou d'importance (IS)	73
8.2.5	Méthodes de Monte Carlo par Chaînes de Markov (MCMC)	75
8.2.6	Échantillonneur de Gibbs et approches hybrides	81
8.2.7	Un résumé de ces premières méthodes	84
8.3	Méthodes d'échantillonnage accélérées	85
8.3.1	Réduction de variance par utilisation des corrélations négatives (variables antithétiques)	85
8.3.2	Variables de contrôle	86
8.3.3	Rao-Blackwellisation	86
8.3.4	Amélioration de lois instrumentales par dynamique de Langevin / hamiltonienne	87
8.4	Méthodes particulières	90
8.5	Méthodes variationnelles	90
8.5.1	Principe fondamental	90
8.5.2	Principes d'usage	91
8.6	Méthodes d'échantillonnage sans vraisemblance (ABC)	92
8.7	Vérification <i>a posteriori</i>	92
	ANNEXES	95
A	Rappels : concepts et outils fondamentaux de l'aléatoire	96
A.1	Problèmes unidimensionnels	96
A.2	Familles de modèles paramétriques	98
A.3	Cas multidimensionnels	102
A.4	Processus aléatoires et stationnarité	103
A.5	Modélisations probabiliste et statistique	104
A.6	Contrôle de l'erreur de modélisation	105
B	Descriptif de quelques modèles statistiques utiles	111
B.1	Lois discrètes	111
B.2	Lois continues	112
C	Éléments sur la simulation pseudo-aléatoire	113
D	Rappels sur les chaînes de Markov	115
E	Calcul bayésien avec OpenBUGS et JAGS	117
E.1	Contexte de développement	117
E.2	Un exemple "fil rouge" : le modèle bêta-binomial	117
E.3	Fonctionnement résumé d'OpenBUGS	118
E.4	Quelques détails supplémentaires concernant OpenBUGS	120
E.5	Liste des distributions de probabilités disponibles	120
E.6	Noeuds logiques et indexation	121
E.7	Pièges à éviter	121
E.8	Utilisation d'OpenBUGS avec R	122
E.9	Détails sur JAGS : exemple de script	123
E.10	D'autres outils (R/Python)	123

1 Notations

La définition des notations suivantes sera rappelée à leur première occurrence dans le document, et elles seront réutilisées par la suite sans rappel obligatoire. D’une manière générale, les variables aléatoires (v.a.) seront notées en majuscules, les réalisations de ces variables en minuscules. Les vecteurs et matrices sont indiqués en gras, à la différence des scalaires.

NOTATIONS GÉNÉRALES

X	variable aléatoire d’étude, unidimensionnelle ou multidimensionnelle
$\mathbb{P}(\cdot)$	mesure de probabilité usuelle
$\mathcal{B}(A)$	tribu (σ –algèbre) des boréliens sur un espace A
$P(A)$	ensemble des parties de A
$\mathbb{1}_{\{x \in A\}}$	fonction indicatrice
\emptyset	ensemble vide
F_X	fonction de répartition de X
f_X	fonction de densité de probabilité de \mathbf{X}
$F_X(\cdot \theta)$	fonction de répartition de X , paramétrée par le vecteur θ
$f_X(\cdot \theta)$	fonction de densité de probabilité de X , paramétrée par θ
$\ell(x_1, \dots, x_n \theta)$	vraisemblance statistique des observations conditionnelle au vecteur θ
$\pi(\theta)$	densité <i>a priori</i> (bayésienne) sur le vecteur θ
$\pi(\theta x_1, \dots, x_n)$	densité <i>a posteriori</i> (bayésienne) sur le vecteur θ sachant un échantillon d’observations x_1, \dots, x_n
$\Pi(\theta)$	fonction de répartition <i>a priori</i>
$\Pi(\theta x_1, \dots, x_n)$	fonction de répartition <i>a posteriori</i>
$\text{sign}(x)$	signe de x
$\text{Supp}(f)$	support de la densité f
X^T	transposée de X
$[x]$	partie entière de X

NOTATIONS GÉNÉRALES (SUITE)

$\mathbb{E}_X[\cdot]$	espérance selon la loi de X (le X peut être ôté si pas d'ambiguïté)
$\mathbb{V}_X[\cdot]$	variance selon la loi de X
$\mathbb{C}ov_{\mathbf{X}}[\cdot]$	matrice de covariance selon la loi de \mathbf{X}
\mathbb{R}	ensemble des réels
\mathbb{N}	ensemble des entiers naturels
L^2	espace des fonctions de carré intégrable
\mathcal{C}	notation générique pour une classe de régularité fonctionnelle
$\langle \cdot, \cdot \rangle$	produit scalaire canonique
A^T	transposée de A
$\text{tr}(A)$	trace de A
$\text{diag}(A)$	vecteur diagonal de A
$ A $	déterminant de A
∇X	gradient de X
$\mathbf{0}_d$	vecteur nul de dimension d
$X_1 \vee X_2$	vecteur de composantes maximales deux à deux
$\xrightarrow{\mathcal{L}}$	convergence en loi
$\xrightarrow{\mathbb{P}}$	convergence en probabilité
$\xrightarrow{p.s.}$	convergence presque sûre
\log	logarithme népérien (\ln)
$\exp(\cdot)$	exponentielle
cste	valeur constante
<i>resp.</i>	respectivement

NOTATIONS ET FONCTIONS DE RÉPARTITION DE LOIS STATISTIQUES

Bernoulli $\mathcal{B}(p)$	$\mathbb{P}(X = 1) = 1 - \mathbb{P}(X = 0) = p$
Binomiale $\mathcal{B}(N, p)$	$\mathbb{P}(X \leq k) = \sum_{i=0}^k \frac{i!(n-i)!}{n!} p^i (1-p)^{n-i}$
Poisson $\mathcal{P}(\lambda)$	$\mathbb{P}(X \leq k) = \sum_{i=0}^k \frac{\lambda^i}{i!} \exp(-\lambda)$
Normale centrée réduite $\mathcal{N}(0, 1)$	$F_X(x) = \Phi(x)$
Gaussienne $\mathcal{N}(\mu, \sigma^2)$	$F_X(x) = \Phi\left(\frac{x-\mu}{\sigma}\right)$
Exponentielle $\mathcal{E}(\lambda)$	$F_X(x) = 1 - \exp(-\lambda x)$
Bêta $\mathcal{B}_e(a, b)$	$\mathbb{P}(X \leq x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1} \mathbb{1}_{\{0 \leq x \leq 1\}}$
Gamma $\mathcal{G}(a, b)$	$F_X(x) = \frac{\gamma(a, bx)}{\Gamma(a)}$ avec $\gamma(a, x) = \int_0^x t^{a-1} \exp(-t) dt$
Inverse gamma $\mathcal{IG}(a, b)$	$F_X(x) = \frac{\Gamma(a, b/x)}{\Gamma(a)}$ avec $\Gamma(a, x) = \int_x^\infty t^{a-1} \exp(-t) dt$
χ_k^2 (Chi-2)	$F_X(x) = \frac{\gamma(k/2, x/2)}{\Gamma(k/2)}$
Student $\mathcal{S}_t(k)$	$F_X(x) = \frac{1}{\sqrt{k\pi}} \frac{\Gamma(\frac{k+1}{2})}{\frac{k}{2}} \int_{-\infty}^x \left(1 + \frac{t^2}{k}\right)^{-\frac{k+1}{2}} dt$

Voir également l'Annexe B pour des précisions sur les modèles fréquemment rencontrés durant le cours.

2 Introduction et rappels

2.1 Modélisation, inférence et décision statistique

Afin d'aborder sereinement ce cours, rappelons que la Statistique (avec une majuscule) peut être vue comme une théorie de la description d'un phénomène incertain, perçu au travers de données $x_n = (x_1, \dots, x_n)$, décrites comme des observations d'une variable X vivant dans un espace Ω . Cette incertitude du phénomène est fondamentalement supposée aléatoire ; c'est-à-dire que l'incertitude sur les valeurs que prend X ne peut pas être réduite à 0 même si le nombre n d'observations tend vers $+\infty$.

La distribution probabiliste à l'origine de ce caractère aléatoire est notée \mathcal{P} , et l'objectif premier de la Statistique est donc d'inférer sur \mathcal{P} à partir de x_n .

Le second objectif est de pouvoir mener une prévision (ou "prédiction") d'une répétition future du phénomène. Le troisième objectif est de prendre une décision ayant des conséquences mesurables, sur la base de l'étude du phénomène.

Remarque 1 Une intelligence artificielle (IA) dite connexioniste (qui se fonde sur l'exploitation des structures de corrélation dans des données) agglomère ces trois objectifs en fournissant une réponse finale à la prise de décision (troisième objectif). Comprendre le comportement d'une telle IA (par exemple en vue de l'étude de sa robustesse puis sa certification) nécessite donc de comprendre les fondations en modélisation et en inférence de la Statistique, et ses liens avec la théorie de la décision.

La modélisation du phénomène consiste en une interprétation réductrice faite sur \mathcal{P} par le biais d'une approche statistique qui peut être :

- non-paramétrique, qui suppose que l'inférence doit prendre en compte le maximum de complexité et à minimiser les hypothèses de travail, en ayant recours le plus souvent à l'estimation fonctionnelle ;
- paramétrique, par laquelle la distribution des observations x_n est représentée par une fonction de densité $f(x|\theta)$ où seul le paramètre θ (de dimension finie) est inconnu.

Ce cours s'intéresse uniquement au cas de l'approche statistique paramétrique. On considèrera en effet en permanence un nombre n fini (et parfois restreint) d'observations, qui ne peut en théorie servir qu'à estimer un nombre fini de paramètres. L'évaluation des outils inférentiels paramétriques peut d'ailleurs être faite avec un nombre fini d'observations.

La section suivante résume brièvement le cadre de la statistique paramétrique. Une revue des concepts fondamentaux de l'aléatoire est donnée en Annexe A, ceux-ci n'étant pas rappelés durant le cours.

2.2 Cadre statistique paramétrique

Pour formaliser la description faite précédemment, et fixer les notations pour le reste du cours, on décrit X comme une variable évoluant dans un espace mesuré et probabilisé

$$(\Omega, \mathcal{A}, \mu, \mathcal{P})$$

où :

1. Ω est l'espace d'échantillonnage des $X = x$, soit l'ensemble de toutes les valeurs possibles prises par X ;
2. la tribu (ou σ -algèbre) \mathcal{A} est la collection des événements (sous-ensembles de Ω) mesurables par μ ;
3. μ est une mesure positive dominante sur (Ω, \mathcal{A}) .
4. \mathcal{P} est une famille de distributions de probabilité dominée par μ , que suit X .

Définition 1 (Domination) Le modèle $P \in \mathcal{P}$ est dit dominé s'il existe une mesure commune dominante μ tel que P admet une densité par rapport à μ ¹

$$f(X) = \frac{dP(X)}{d\mu}.$$

De manière générale, on travaillera avec $\Omega \subset \mathbb{R}^d$ avec $d < \infty$ et des échantillons de réalisations $x_n = (x_1, \dots, x_n)$ de X . La mesure dominante μ sera Lebesgue (cas continus) ou Dirac (cas discrets). Enfin, \mathcal{A} sera très généralement / classiquement choisie comme la tribu des boréliens

$$\mathcal{A} = \mathcal{B}(\mathbb{R}^d) = \sigma \left(\{ \otimes_{i=1}^d]a_i, b_i]; a_i < b_i \in \mathbb{R} \} \right).$$

Dans le cadre paramétrique, on supposera que \mathcal{P} peut se définir par

$$\mathcal{P} = \{ \mathbb{P}_\theta; \theta \in \Theta \subset \mathbb{R}^p \}$$

où $p < \infty$. De plus, on notera généralement $f(\cdot|\theta)$ la densité (ou fonction de masse) induite par la dérivée de Radon-Nikodym de P_{p_θ} :

$$\frac{d\mathbb{P}_\theta}{d\mu} = f(X|\theta)$$

et parfois, lorsque X sera unidimensionnelle ($d = 1$), nous utiliserons aussi la notation classique $F(x|\theta)$ pour désigner la fonction de répartition $P_{p_\theta}(X \leq x)$. Par la suite, on parlera indifféremment de la variable aléatoire

$$X \sim f(x|\theta)$$

ou de son observation $x \sim f(x|\theta)$, et on parlera plus généralement de loi en confondant P_{p_θ} et $f(\cdot|\theta)$. Enfin, la notation μ sera généralement induite dans les développements techniques :

$$\mathbb{P}_\theta(X < t) = \int_{\Omega} f(x) \mathbb{1}_{\{x < t\}} dx.$$

Remarque 2 Suivant l'usage classique, les variables et processus aléatoires sont décrits par des majuscules, tandis que leurs réalisations sont décrits par des minuscules. On notera souvent v.a. pour variable aléatoire.

Nous retrouverons et utiliserons abondamment la notion de *vraisemblance* statistique $f(\mathbf{x}_n|\theta)$, définie dans un cadre paramétrique comme la densité jointe des observations $\mathbf{x}_n = (x_1, \dots, x_n)$ sachant le paramètre θ . Lorsque les données sont *indépendantes et identiquement distribuées* (iid) selon $f(\cdot|\theta)$, alors

$$f(\mathbf{x}_n|\theta) = \prod_{i=1}^n f(x_i|\theta).$$

D'autres formes de vraisemblance existent, notamment lorsque les données sont bruitées, censurées, etc. Voir Annexe A pour des rappels sur ces principaux concepts.

Remarque 3 (Statistique bayésienne non paramétrique) Jusqu'à présent, θ est considéré comme appartenant à un espace Θ de dimension finie. On peut étendre la statistique bayésienne à Θ un ensemble comme $[0, 1]^{\mathbb{R}}$ (l'ensemble des distributions sur $[0, 1]$) ou encore l'ensemble des probabilités sur \mathbb{R} . Ces deux espaces ne sont pas dominés par μ . C'est le principe fondateur de la statistique non paramétrique (au sens où le paramètre n'a pas de dimension finie).

2.3 Estimation statistique classique ("fréquentiste")

(ou fréquentielle en meilleur français)

1. Pour des mesures σ -finies et de part le théorème de Radon-Nykodim, ceci est équivalent à être absolument continue par rapport à μ

2.3.1 Rappel des principes

L'inférence statistique consiste à estimer "les causes à partir des effets". Ces *causes* sont réduites, dans le cadre paramétrique, au paramètre θ du mécanisme générateur des données que représente la distribution \mathbb{P}_θ . Les *effets* sont naturellement les données observées $\mathbf{x}_n = (x_1, \dots, x_n)$. De ce fait, dans un cadre paramétrique, l'inférence consiste à produire des règles d'estimation de θ à partir de \mathbf{x}_n . Dans ce cadre classique, θ **est supposé inconnu, mais fixe** (et à Θ n'est pas conféré la structure d'un espace probabilisé).

Les règles d'estimation les plus courantes, fondées sur de l'optimisation de critère (M —estimation, telles la maximisation de la vraisemblance

$$\hat{\theta}_n(\mathbf{x}_n) = \arg \max_{\theta} \log f(\mathbf{x}_n | \theta)$$

ou les *estimateurs des moindres carrés*), par *moments*, par des combinaisons linéaires de statistiques d'ordre (L —estimation, en général moins robuste), etc. sont nombreuses et doivent faire l'objet d'une sélection. Voir Annexe A.6 pour quelques rappels.

Pour mener cette sélection, les estimateurs sont comparés en fonction de différents critères, comme le biais, la rapidité de convergence vers la valeur supposée "vraie" θ_0 du paramètre, et d'autres différentes propriétés asymptotique (telle la nature de la loi d'un estimateur $\hat{\theta}_n(\mathbf{X}_n)$, qui est une variable aléatoire dont la loi dépend de celle des X .

D'une manière générale, si l'on note $\hat{\theta}_n = \hat{\theta}_n(\mathbf{X}_n)$ tout estimateur classique de θ , à de rares exceptions près la validité de ce choix d'estimateur est dépendante du caractère *reproductible* et *échangeable* des données x_1, \dots, x_n conditionnellement à θ .

Définition 2 (Échangeabilité.) Les données x_1, \dots, x_n sont dites *échangeables* si, pour toute permutation $\sigma : \mathbb{N}^n \rightarrow \mathbb{N}^n$, la loi jointe $f(x_{\sigma(1)}, \dots, x_{\sigma(n)})$ est indépendante de σ .

Cette validité, donc en général fondée sur des critères asymptotiques ($n \rightarrow \infty$), s'exprime en termes de *région de confiance* (cf. Annexe 34)

$$\mathbb{P} \left(\hat{\theta}_n - \theta \in A_\alpha \right) = 1 - \alpha.$$

En général, la distribution \mathbb{P} de l'estimateur est inconnue pour $n < \infty$, elle est le plus souvent approximée asymptotiquement via un théorème de convergence en loi, tel que :

$$\text{si } x_1, \dots, x_n \text{ sont iid } \quad \Sigma_n^{-1/2} \left(\hat{\theta}_n - \theta_0 \right) \xrightarrow{\mathcal{L}} \mathcal{Q}$$

où $\mathbb{E}_{\mathcal{Q}}[X] = 0$ et $\mathbb{V}_{\mathcal{Q}}[X] = 1$. Ici Σ_n est lui-même un estimateur consistant de la matrice de covariance de $\hat{\theta}_n$, et le résultat précédent est issu de l'usage de la méthode Delta de dérivation des lois d'estimateur, ainsi que du théorème de Slutsky de composition des convergences.

Remarque 4 On utilise souvent le terme d'inférence en machine learning pour désigner la tâche de prévision (prediction) d'un modèle appris, et l'entraînement la phase d'estimation de ce modèle. En ce sens, le mot inférer est tout aussi valide, car il signifie "aller des principes vers la conclusion".

2.3.2 Difficultés pratiques, théoriques et conceptuelles

Ce *paradigme*² forme, depuis les travaux de Fisher, Neyman et Pearson dans la première moitié du XXème siècle, le socle théorique de la majeure partie des études statistiques. Il n'est pas cependant sans poser quelques problèmes :

- (a) Tout d'abord, les difficultés rencontrées sont **pratiques** : face à de petits échantillons, le cadre asymptotique ne tient plus : la comparaison des estimateurs doit alors reposer sur des critères non asymptotiques³, et on perd l'usage des résultats de la convergence en loi et ses dérivées (ex : production des régions de confiance). De même, la plupart des résultats utiles pour mener des tests statistiques (voir Annexe A.2) deviennent inutilisables.

2. Modèle censé être cohérent d'un univers scientifique, faisant l'objet d'un consensus.

3. Parmi ces critères, les inégalités de concentration (Markov, Bienaymé-Chebychev, Bernstein, etc.) se révèlent fondamentales.

(b) Des difficultés peuvent aussi être **théoriques**.

1. Ainsi, pour de nombreux modèles complexes, tels les modèles à espace d'états (qui font partie des modèles à données latentes), tels que les modèles de population, la dimension de Θ peut augmenter linéairement avec le nombre de données. Dans ce cas, la théorie asymptotique classique n'a plus de sens.

EXEMPLE 1. *On considère une population suivie annuellement, n étant le nombre d'années de mesure. A chaque année est associée un paramètre spécifique de renouvellement de la population. La dimension augmente donc linéairement avec le nombre de donnée, si aucune réduction de dimension (par exemple via des covariables connues) n'est effectuée.*

2. Plus fondamentalement, l'utilisation d'un estimateur fréquentiste peut contredire le principe fondamental de la statistique inférentielle :

Définition 3 (Principe de vraisemblance) *L'information (= l'ensemble des inférences possibles) apportée par une observation x sur θ est entièrement contenue dans la fonction de vraisemblance $\ell(\theta|x) = f(x|\theta)$. De plus, si x_1 et x_2 sont deux observations qui dépendent du même paramètre θ , telle qu'il existe une constante c satisfaisant*

$$\ell(\theta|x_1) = c\ell(\theta|x_2) \quad \forall \theta \in \Theta,$$

alors elles apportent la même information sur θ et doivent conduire à la même inférence.

Exercice 1 (Adapté de [38]) *Soient (x_1, x_2) deux réalisations aléatoires. Nous disposons de deux candidats pour la loi jointe de ces observations : $x_i \sim \mathcal{N}(\theta, 1)$ ou encore*

$$g(x_1, x_2|\theta) = \pi^{-3/2} \frac{\exp\left\{-\frac{(x_1 + x_2 - 2\theta)^2}{4}\right\}}{1 + (x_1 - x_2)^2}.$$

Quel est l'estimateur du maximum de vraisemblance de θ dans chacun des cas ? Que constate-on ?

3. Citons également le fait que l'estimateur du maximum de vraisemblance (EMV), considéré généralement comme le plus efficace (atteignant la borne de Cramer-Rao et asymptotiquement sans biais dans la plupart des cas), peut ne pas exister ou être unique.

EXEMPLE 2. *Modèles à paramètre de position, modèles de mélange...*

Par ailleurs, l'usage de l'EMV pose un autre problème, qui contredit le principe de vraisemblance : les régions de confiance de la forme (*test du rapport de vraisemblance*)

$$\mathcal{C} = \left\{ \theta; \frac{\ell(\theta|x)}{\ell(\hat{\theta}|x)} \geq c \right\}$$

qui sont les plus petites asymptotiquement, ne dépendront pas uniquement de la fonction de vraisemblance si la borne c doit être choisie de manière à obtenir un niveau de confiance α .

4. Une dernière difficulté théorique posée par les estimateurs fréquentiels apparaît lorsqu'on cherche à mener une *prévision*. Considérons en effet Soit $\mathbf{X}_n = (X_1, \dots, X_n) \stackrel{iid}{\sim} f(\cdot|\theta)$. On cherche à prévoir le plus précisément possible ce que pourrait être le prochain tirage X_{n+1} . Dans l'approche classique, on utilise

$$f(X_{n+1}|X_1, \dots, X_n, \hat{\theta}_n) = \frac{f(X_1, \dots, X_n, X_{n+1}|\hat{\theta}_n)}{f(X_1, \dots, X_n|\hat{\theta}_n)}$$

et ce faisant on utilise deux fois les données et on risque de sous-estimer les incertitudes (intervalles de confiance) en renforçant arbitrairement la connaissance.

- (c) Enfin, les difficultés peuvent être **d'ordre conceptuel**. En effet, le sens donné à une probabilité est, dans la statistique bayésienne, celui d'une *limite de fréquence*, et la notion de *confiance* est uniquement fondée sur la répétabilité des expériences peut ne pas être pertinente.

EXEMPLE 3. *Le premier pari d'une course de chevaux ?*

En prévision, nous souhaiterions connaître parfaitement l'incertitude sur le mécanisme générateur de X , mais c'est une tâche impossible en pratique. Dans de nombreux contextes, toute variable aléatoire est la représentation mathématique d'une grandeur soumise à deux types d'incertitude :

1. \mathbb{P}_θ représente la partie *aléatoire* du phénomène considéré ;
2. l'estimation de θ souffre d'une incertitude *épistémique*, réductible si de l'information supplémentaire (données) est fournie (typ. : données).

L'approche classique des statistiques souffre donc de difficultés qui limitent son usage à des situations généralement restreintes à l'asymptotisme. Elle constitue en fait une *approximation* d'un paradigme plus vaste, celui de la *statistique bayésienne*, qui permet notamment de *correctement appréhender la gestion des incertitudes en estimation, prévision, et en aide à la décision*.

Remarque 5 (Écriture fiduciaire) *L'écriture fiduciaire $\ell(\theta|x) = f(x|\theta)$ a été proposée au début du XXème siècle pour témoigner du fait qu'on cherche à mesurer l'éventail des valeurs possibles de θ sachant l'observation des x_i . Toutefois, il s'agissait d'une confusion entre la définition d'un estimateur statistique et celle d'une véritable variable aléatoire nécessitant l'ajout d'une mesure dominante sur θ . Il vaut mieux ne pas l'utiliser pour ne pas oublier le sens statistique d'une vraisemblance (loi jointe des données).*

2.4 Principes de la statistique bayésienne

2.4.1 Paradigme

Le paradigme de la statistique bayésienne paramétrique part du principe que le **vecteur θ est une variable aléatoire**, vivant dans un espace probabilisé (on utilisera généralement $(\Theta, \Pi, \mathcal{B}(\Theta))$).

En reprenant la formulation *L'inférence statistique consiste à estimer "les causes à partir des effets"* au § 2.3.1, cela revient à associer X aux effets, et θ aux causes, et d'"estimer ces causes" par la mise à jour de la distribution (mesure) $\Pi(\Theta)$ via la *règle de Bayes* :

Si C (cause) et E (effet) sont des événements tels que $P(E) \neq 0$, alors

$$\begin{aligned} P(C|E) &= \frac{P(E|C)P(C)}{P(E|C)P(C) + P(E|C^c)P(C^c)} \\ &= \frac{P(E|C)P(C)}{P(E)} \end{aligned}$$

Il s'agit d'un principe d'*actualisation*, décrivant la mise à jour de la vraisemblance de la cause C de $P(C)$ vers $P(C|E)$.

Ce paradigme a historiquement été proposé par Bayes (1763) puis Laplace (1795), qui ont supposé que l'*incertitude sur θ* pouvait être décrite par une distribution de probabilité Π de densité $\pi(\theta)$ sur Θ , appelée *loi a priori*. On notera en général

$$\theta \sim \pi$$

Formulation en densité. Sachant des données \mathbf{x}_n , la mise à jour de cette loi *a priori* s'opère par le conditionnement de θ à \mathbf{x}_n ; on obtient la *loi a posteriori*

$$\pi(\theta|\mathbf{x}_n) = \frac{f(\mathbf{x}_n|\theta)\pi(\theta)}{\int_{\Theta} f(\mathbf{x}_n|\theta)\pi(\theta) d\theta} \quad (1)$$

Définition 4 Un modèle statistique bayésien est constitué d'un modèle statistique paramétrique (ou vraisemblance) $f(x|\theta)$ et d'une mesure a priori $\pi(\theta)$ pour les paramètres.

En conséquence, là où la statistique classique s'attache à définir des procédures d'estimation ponctuelle de θ , la statistique bayésienne va s'attacher à définir des procédures d'estimation de la loi *a posteriori* $\pi(\theta|\mathbf{x}_n)$.

Exercice 2 (Bayes (1763)) Une boule de billard Y_1 roule sur une ligne de longueur 1, avec une probabilité uniforme de s'arrêter n'importe où. Supposons qu'elle s'arrête à la position θ . Une seconde boule Y_2 roule alors n fois dans les mêmes conditions, et on note X le nombre de fois où Y_2 s'arrête à gauche de Y_1 . Connaissant X , quelle inférence peut-on mener sur θ ?

Exercice 3 (Loi gaussienne / loi exponentielle) Soit une observation $x \sim \mathcal{N}(\theta, \sigma^2)$ où σ^2 est connu. On choisit a priori

$$\theta \sim \mathcal{N}(m, \rho\sigma^2)$$

Quelle est la loi *a posteriori* de θ sachant x ? Même question en supposant que $X \sim \mathcal{E}(\lambda)$ et

$$\lambda \sim \mathcal{G}(a, b).$$

Définition 5 (Loi impropre) Une "loi impropre" est une mesure a priori σ -finie qui vérifie $\int_{\Theta} \pi(\theta) d\theta = \infty$.

La mesure de Lebesgue sur un ouvert est un exemple de loi impropre. Le choix de manier ce type de mesure peut sembler étrange, mais ce choix peut s'avérer en fait particulièrement intéressant. Par exemple, travailler avec une loi normale centrée à grande variance pour approcher une "loi uniforme sur \mathbb{R} " peut être précieux. Une telle loi *a priori* n'a cependant d'intérêt que si la loi *a posteriori* correspondante existe. On se limitera donc aux lois impropres telles que la loi marginale soit bien définie :

$$m_{\pi}(x) = \int_{\Theta} f(x|\theta) d\pi(\theta) < \infty$$

Exercice 4 (Loi uniforme généralisée) Soit $X \sim \mathcal{N}(\mu, \sigma^2)$ et $d\pi(\mu) = d\mu$ (mesure de Lebesgue). Que vaut $m_{\pi}(x)$?

Exercice 5 (Loi d'échelle) Soit $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ et $\pi(\mu, \sigma) = 1/\sigma$ avec $\Theta = \mathbb{R} \times \mathbb{R}_*^+$. Que vaut $m_{\pi}(x_1, \dots, x_n)$? La mesure $\pi(\mu, \sigma)$ peut-elle être utilisable ?

Dans le cas où π est une mesure impropre σ -finie, on considère $\pi^*(\theta) = c\pi(\theta)$ où c est une constante arbitraire. Elle doit être sans influence pour l'usage du modèle bayésien. On peut facilement voir que c'est bien le cas dans le calcul *a posteriori* (exercice), puisqu'elle apparaît aussi bien au numérateur qu'au dénominateur de l'expression (1) : on a bien

$$d\pi^*(\theta|X) = d\pi(\theta|X).$$

Ainsi, l'usage de lois impropres *a priori* est justifié si la loi *a posteriori* est propre⁴ car cette dernière ne dépend pas de la constante multiplicative c inconnue. C'est à rapprocher du principe de vraisemblance énoncé précédemment.

2.4.2 Fondations théoriques

Les fondations théoriques de la statistique bayésienne seront progressivement investiguées durant le cours, notamment en lien avec la section consacrée à la théorie de la décision (§ 3), mais il est important de connaître un premier résultat, dû originellement à De Finetti. Il s'agit d'un *théorème de représentation*, c'est-à-dire un théorème qui permet de justifier un choix de représentation probabiliste des variations de θ dans Θ .

4. C'est-à-dire intégrable : une loi de probabilité qui mesure les informations une fois les données connues.

Théorème 1 (De Finetti (1931)) Soit X_1, \dots, X_n, \dots une séquence échangeable de variables aléatoires binaires (0-1) de probabilité jointe P . Alors il existe une mesure de probabilité unique $\pi(\theta)$ telle que

$$P(X_1 = x_1, \dots, X_n = x_n, \dots) = \int_{\Theta} f(x_1, \dots, x_n, \dots | \theta) \pi(\theta) d\theta$$

où $f(x_1, \dots, x_n | \theta)$ est la vraisemblance d'observations iid de Bernoulli (également notée $\ell(\theta | x_1, \dots, x_n, \dots)$).

Nous admettons ce théorème ainsi que ses nombreux dérivés. En effet, il a été généralisé successivement par Hewitt, Savage (1955), Diaconis, Freedman (1980) pour l'ensemble des distributions discrétisées puis continues.

Selon ce théorème, la modélisation bayésienne apparaît comme une modélisation statistique naturelle de *variables corrélées mais échangeables*. L'existence formelle d'une mesure *a priori* (ou *prior* dans la suite de ce cours) $\pi(\theta)$ est assurée en fonction du mécanisme d'échantillonnage, qui apparaît dès lors comme une simplification d'un mécanisme par essence mal connu ou inconnu.

Un autre théorème fondamental qui nous permet de justifier l'usage du cadre bayésien est le *théorème de Cox-Jaynes*, qui sera introduit plus tard dans le cours (Section 5). Il est fondé sur une *axiomatique de la représentation de l'information* et il constitue aujourd'hui à la fois une autre façon de défendre le choix de la théorie des probabilités et l'un des théorèmes fondamentaux de l'intelligence artificielle.

Le prior correspond donc à une mesure d'information incertaine à propos de θ , et (comme on le verra) un *prior probabiliste* pour certains théoriciens des probabilités. Cette probabilisation de θ va permettre de répondre de façon pratique :

- à la nécessité de *satisfaire le principe de vraisemblance* ;
- à la nécessité de *tenir compte de toutes les incertitudes épistémiques* s'exprimant sur θ , en particulier dans un objectif de *prévision* ;
- de distinguer ces incertitudes de l'incertitude *aléatoire*, intrinsèque au modèle $f(\cdot | \theta)$;
- à la possibilité d'intégrer de la connaissance *a priori* sur le phénomène considéré, autre que celle apportée par les données \mathbf{x}_n ;
- à la nécessité de faire des choix de modèles en évitant les difficultés des tests statistiques classiques ;
- l'invariance $\pi(\theta | \mathbf{x}_n) = \pi(\theta)$ permet en outre d'identifier des problèmes d'*identifiabilité* du modèle d'échantillonnage $X \sim f(x | \theta)$

2.4.3 Plan du cours

Ce cours va considérer successivement plusieurs aspects du choix et de la mise en œuvre du cadre statistique bayésien. Il cherche à fournir les éléments nécessaires pour répondre aux questions fondamentales suivantes :

- Quand le paradigme bayésien est-il préférable ?** Hors du contexte spécifique des petits échantillons, pour lesquels la statistique classique apporte des réponses limitées, cette question revient d'abord à comprendre que la statistique bayésienne est d'abord une *théorie de la décision, centrale en apprentissage statistique* et dans la formalisation du travail du statisticien. Le cadre décisionnel proposé par la statistique bayésienne améliore la vision fréquentielle du monde, et s'accorde avec elle lorsque l'information apportée par les données augmente. Ces deux aspects sont considérés dans les Sections 3 et 4.
- Comment construire une ou plusieurs mesures *a priori* $\pi(\theta)$?** Cette partie importante du cours est traitée plusieurs sections. La section 5 propose d'abord de formuler les principes généraux de compréhension et de représentation probabiliste de l'information incertaine. Sur la base de ces principes, issus d'une axiomatique, la section 6 proposera un panorama des méthodes et outils de la modélisation bayésienne.

(b) **Comment faire du calcul bayésien ?** La mise en oeuvre concrète des outils et méthodes de la statistique bayésienne suppose de pouvoir manipuler les lois *a posteriori* $\pi(\theta|\mathbf{x}_n)$. Les méthodes par simulation (échantillonnage) et les approches par approximation variationnelle font aujourd'hui partie des outils courants pour ce faire. Elles seront abordées dans la section 8.

2.5 Liens avec le *machine learning*

Dans une optique de *régression supervisée*, le paradigme du *machine learning* propose de produire un estimateur (ou *prédicteur*) de la fonction inconnue $g : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$ (plus généralement vers un espace euclidien de dimension d_2) telle que

$$Y = g(X)$$

à partir de couples connus $\mathbf{z}_n = (x_i, y_i)_{1 \leq i \leq n}$, où chaque x_i est un ensemble de d_1 *covariables* et chaque y_i est un *label* de dimension d_2 . La recette est la suivante :

1. Faire un choix g_θ pour "mimer" g en régression, tel que

$$\mathbb{E}[Y|X] = g_\theta(X)$$

- Dans un problème de régression linéaire, θ (noté généralement β) est le vecteur des coefficients de la régression).
 - Si g_θ est un réseau de neurones d'architecture choisie, alors θ constitue un vecteur de paramètres structurant pour ce réseau (poids, biais, nombre de neurones par couche, éventuellement les choix de fonctions d'activation, etc.).
2. Décider d'une *fonction de coût*⁵ souvent définie comme la somme d'un regret quadratique et d'une pénalité

$$L(\theta|\mathbf{z}_n) = \sum_{i=1}^n \|y_i - g_\theta(x_i)\|_2^2 + \text{pen}(\theta) \quad (2)$$

où $\text{pen}(\theta)$ dépend de la complexité du problème.

3. Définir l'estimateur $\hat{\theta}_n$ par

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} L(\theta|\mathbf{z}_n) \quad (3)$$

et choisir une méthode pour minimiser la fonction de coût (exemple : rétropropagation du gradient).

On peut alors réécrire l'équation (3) de la façon suivante :

$$\begin{aligned} \hat{\theta}_n &= \arg \max_{\theta \in \Theta} \{-L(\theta|\mathbf{z}_n)\}, \\ &= \arg \max_{\theta \in \Theta} \log \{f_g(\mathbf{z}_n|\theta)\pi(\theta)\}, \\ &= \arg \max_{\theta \in \Theta} \log \pi(\theta|\mathbf{z}_n), \\ &= \arg \max_{\theta \in \Theta} \pi(\theta|\mathbf{z}_n) \end{aligned}$$

où $f_g((\mathbf{z}_n|\theta))$ est une vraisemblance de forme gaussienne de \mathbf{z}_n et

$$\pi(\theta) \propto \exp(-2\text{pen}(\theta)).$$

Le cadre bayésien explique le sens d'une pénalisation comme celui d'une transformation d'une mesure *a priori*, et l'optimisation en *machine learning* consiste à estimer le mode d'une distribution *a posteriori* (calcul simplificateur de la véritable inférence, qui serait celle de la loi $\pi(\theta|\mathbf{z}_n)$ toute entière).

EXEMPLE 4. La régression lasso propose un choix de pénalisation $\text{pen}(\theta) = \lambda\|\theta\|_1$, qui correspond à l'action d'un prior $\pi(\theta) \propto \exp(-2\lambda\|\theta\|_1)$. De même, la régularisation ridge est similaire à l'action d'un prior $\pi(\theta) \propto \exp(-2\lambda\|\theta\|_2^2)$.

5. On retrouvera ce terme plus tard dans la partie du cours consacré à la théorie de la décision (Section 3).

2.6 Quelques lectures conseillées

Ce cours s'inspire de plusieurs ouvrages et résultats publiés ces dernières années. L'étudiant intéressé par une vision générale du cadre pourra approfondir les aspects théoriques à partir de l'ouvrage de référence [38]. Une démarche plus appliquée de la statistique bayésienne bénéficie d'une présentation pédagogique dans l'ouvrage [31]. Les aspects computationnels historiques sont au coeur des ouvrages de référence [39, 26]. Le cadre décisionnel de la théorie bayésienne, dans un contexte d'usage concret (et relié à l'industrie), fait l'objet de l'article (français) [19].

L'article de revue récent [47] offre enfin une vision générale du cadre statistique bayésien, et complète utilement les lectures précédentes.

3 Éléments de théorie de la décision

L'objectif général de la plupart des études inférentielles est de fournir une *décision* au statisticien (ou au client) à partir du phénomène modélisé par $X \sim f(x|\theta)$ (dans le cadre paramétrique). Il faut donc exiger un *critère d'évaluation* des procédures de décision qui :

- prenne en compte les conséquences de chaque décision
- dépende des paramètres θ du modèle, càd du *vrai état du monde (ou de la nature)*.

Un autre type de décision est d'*évaluer* si un nouveau modèle descriptif est compatible avec les données expérimentales disponibles (*choix de modèle*). Le critère en question est habituellement nommé **fonction de coût**, **fonction de perte** ou **utilité** (opposé du coût).

EXEMPLE 5. Acheter des capitaux selon leurs futurs rendement θ , déterminer si le nombre θ des SDF a augmenté depuis le dernier recensement...

Formellement, pour le modèle $X \in \{\Omega, \mathcal{B}, \{\mathbb{P}_\theta, \theta \in \Theta\}\}$ on définit donc trois espaces de travail :

- Ω = espace des observations x ;
- Θ = espace des paramètres θ ;
- \mathcal{D} = espace des décisions possibles d .

En général, la décision $d \in \mathcal{D}$ demande d'évaluer (*estimer*) une *fonction d'intérêt* $h(\theta)$, avec $\theta \in \Theta$, estimation fondée sur l'observation $x \in \Omega$. On décrit alors \mathcal{D} comme l'ensemble des fonctions de Θ dans $h(\Theta)$ où h dépend du contexte :

- si le but est d'estimer θ alors $\mathcal{D} = \Theta$;
- si le but est de mener un test, $\mathcal{D} = \{0, 1\}$.

3.1 Existence d'une fonction de coût

La *théorie de la décision* suppose alors que :

- chaque décision $d \in \mathcal{D}$ peut être évaluée et conduit à une *récompense* (ou *gain*) $r \in \mathcal{R}$
- l'espace \mathcal{R} des récompenses peut être *ordonné totalement* :
 - (1) $r_1 \preceq r_2$ ou $r_2 \preceq r_1$;
 - (2) si $r_1 \preceq r_2$ et $r_2 \preceq r_3$ alors $r_1 \preceq r_3$;
- l'espace \mathcal{R} peut être étendu à l'espace \mathcal{G} des distributions de probabilité dans \mathcal{R} ;
 - les décisions peuvent être alors partiellement aléatoires
- la relation d'ordre \preceq peut être étendue sur les **moyennes** des récompenses aléatoires (*et donc sur les distributions de probabilité correspondantes*) ;
 - il existe au moins un ordre partiel sur les gains (même aléatoires) et un gain optimal.

Ces axiomes expriment une certaine **hypothèse de rationalité du décideur**. Ils impliquent l'existence d'une **fonction d'utilité** $U(r)$ permettant de trier les gains aléatoires. Cette utilité ne dépend en fait que de θ et de d : on la note donc $U(\theta, d)$. Elle peut être vue comme une *mesure de proximité* entre la décision proposée d et la vraie valeur (inconnue) θ .

Définition 6 On appelle fonction de coût ou fonction de perte une fonction L mesurable de $\Theta \times \mathcal{D}$, telle que

$$L(\theta, d) = -U(\theta, d),$$

à valeurs réelles positives :

$$L : \Theta \times \mathcal{D} \longrightarrow \mathbb{R}^+.$$

La fonction de coût est définie selon le problème étudié et constitue l'armature d'un problème de décision statistique (qui comprend notamment les problèmes d'estimation).

EXEMPLE 6. On considère le problème de l'estimation de la moyenne θ d'un vecteur gaussien

$$x \sim \mathcal{N}_p(\theta, \Sigma)$$

où Σ est une matrice diagonale connue avec pour éléments diagonaux σ_i^2 ($i = 1, \dots, p$). Dans ce cas $\mathcal{D} = \Theta = \mathbb{R}^p$ et d représente une évaluation de θ . S'il n'y a pas d'information additionnelle disponible sur ce modèle, il paraît logique de choisir une fonction de coût qui attribue le même poids à chaque composante, soit un coût de la forme

$$\sum_{i=1}^p L\left(\frac{x_i - \theta_i}{\sigma_i}\right) \quad \text{avec } L(0) = 0.$$

Par normalisation, les composantes avec une grande variance n'ont pas un poids trop important. Le choix habituel de L est le coût **quadratique** $L(t) = t^2$.

Dans un contexte de gain aléatoire, l'approche fréquentiste propose de considérer le coût moyen ou *risque fréquentiste*. Pour une fonction de coût quadratique, le risque fréquentiste est souvent appelé *risque quadratique*. On appelle $\delta : \Omega \mapsto \mathcal{D}$ minimisant un risque un estimateur et $\delta(x)$ une estimation.

Définition 7 (Risque fréquentiste) Pour $(\theta, \delta) \in \Theta \times \mathcal{D}$, le risque fréquentiste est défini par

$$R(\theta, \delta) = \mathbb{E}_\theta [L(\theta, \delta(x))] = \int_{\Omega} L(\theta, \delta(x)) f(x|\theta) dx$$

où $\delta(x)$ est la règle de décision = attribution d'une décision connaissant l'observation x .

Cette définition du risque n'est pas sans poser problème. En effet :

- le critère évalue les procédures d'estimation selon leurs *performances à long terme* et non directement pour une observation donnée ;
- on suppose tacitement que le problème sera rencontré de nombreuses fois pour que l'évaluation en fréquence ait un sens

$$R(\theta, \delta) \simeq \text{coût moyen sur les répétitions};$$

- ce critère n'aboutit pas à un *ordre total* sur les procédures de construction d'estimateur.

Exercice 6 Soient x_1 et x_2 deux observations de la loi définie par

$$P_\theta(x = \theta - 1) = P_\theta(x = \theta + 1) = 1/2 \quad \text{avec } \theta \in \mathbb{R}$$

Le paramètre d'intérêt est θ (donc $\mathcal{D} = \Theta$) et il est estimé par δ sous le coût

$$L(\theta, \delta) = 1 - \mathbb{1}_\theta(\delta)$$

appelé coût 0-1, qui pénalise par 1 toutes les erreurs d'estimation quelle que soit leur magnitude (grandeur). Soit les estimateurs

$$\begin{aligned}\delta_1(x_1, x_2) &= \frac{x_1 + x_2}{2}, \\ \delta_2(x_1, x_2) &= x_1 + 1, \\ \delta_3(x_1, x_2) &= x_2 - 1.\end{aligned}$$

Calculez les risques $R(\theta, \delta_1)$, $R(\theta, \delta_2)$ et $R(\theta, \delta_3)$. Quelle conclusion en tirez-vous ?

L'approche bayésienne de la théorie de la décision considère que le coût $L(\theta, d)$ doit plutôt être moyenné sur tous les états de la nature possibles. Conditionnellement à l'information x disponible, ils sont décrits par la loi *a posteriori* $\pi(\theta|x)$. On définit donc le coût moyenné *a posteriori*, ou *risque a posteriori*, qui est l'erreur moyenne résultant de la décision d pour un x donné.

Définition 8 (Risque *a posteriori*)

$$R_P(d|\pi, x) = \int_{\Theta} L(\theta, d) \pi(\theta|x) d\theta.$$

On peut enfin définir le risque fréquentiste intégré sur les valeurs de θ selon leur distribution *a priori*. Associant un nombre réel à chaque estimateur δ , ce risque induit donc une *relation d'ordre total* sur les procédures de construction d'estimateur. Il permet donc de définir la notion d'estimateur bayésien (ou estimateur de Bayes).

Définition 9 (Risque intégré) À fonction de coût (perte) donnée, le risque intégré est défini par

$$R_B(\delta|\pi) = \int_{\Theta} \int_{\Omega} L(\theta, \delta(x)) f(x|\theta) dx \pi(\theta) d\theta.$$

Définition 10 (Estimateur bayésien et risque de Bayes) Un estimateur de Bayes associé à une distribution *a priori* π et une fonction de coût L est défini par

$$\delta^\pi = \arg \min_{\delta \in \mathcal{D}} R_B(\delta|\pi)$$

la valeur $r(\pi) = R_B(\delta^\pi|\pi)$ est alors appelée **risque de Bayes**.

Le résultat suivant peut être obtenu par interversion d'intégrales (théorème de Fubini). *Modulo* un peu de machinerie technique, on peut montrer que celui-ci reste vrai même si $\int_{\Theta} \pi(\theta) d\theta = \infty$ (mesure *a priori* non informative) à condition que $\int_{\Theta} \pi(\theta|x) d\theta = 1$.

Théorème 2 Pour chaque $x \in \Omega$,

$$\delta^\pi(x) = \arg \min_{d \in \mathcal{D}} R_P(d|\pi, x). \quad (4)$$

Un corollaire est le suivant : s'il existe $\delta \in \mathcal{D}$ tel que $R_B(\delta|\pi) < \infty$, et si $\forall x \in \Omega$ l'équation (4) est vérifiée, alors $\delta^\pi(x)$ est un estimateur de Bayes.

3.2 Supériorité des estimateurs de Bayes sur les estimateurs fréquentistes

Le risque minimax est le coût fréquentiste minimum dans le cas le moins favorable (l'écart entre θ et δ , c'est-à-dire l'erreur d'estimation, est maximal(e)).

Définition 11 (Risque minimax) On définit le risque minimax pour la fonction de coût L par

$$\bar{R} = \inf_{\delta \in \mathcal{D}} \sup_{\theta \in \Theta} R(\theta, \delta) = \inf_{\delta \in \mathcal{D}} \sup_{\theta \in \Theta} \mathbb{E}_{\theta} [L(\theta, \delta(x))].$$

Théorème 3 Le risque de Bayes est toujours plus petit que le risque minimax

$$R = \sup_{\pi} r(\pi) = \sup_{\pi} \inf_{\delta \in \mathcal{D}} R_B(\delta|\pi) \leq \bar{R}.$$

Si elle existe, une distribution *a priori* π^* telle que $r(\pi^*) = R$ est appelée *distribution a priori la moins favorable*. Ainsi, l'apport d'information *a priori* $\pi(\theta)$ ne peut qu'améliorer l'erreur d'estimation, même dans le pire des cas.

Définition 12 (Inadmissibilité d'un estimateur) Un estimateur δ_0 est dit inadmissible s'il existe un estimateur δ_1 qui domine δ_0 au sens du risque fréquentiste, c'est-à-dire si

$$R(\theta, \delta_0) \geq R(\theta, \delta_1) \quad \forall \theta \in \Theta$$

et $\exists \theta_0$ tel que $R(\theta_0, \delta_0) > R(\theta_0, \delta_1)$. Sinon, il est dit admissible.

Théorème 4 Si un estimateur de Bayes δ^π associé à une mesure *a priori* π (probabiliste ou non) est tel que le risque $R(\theta, \delta^\pi) < \infty$ et si la fonction $\theta \mapsto R(\theta, \delta)$ est continue sur Θ , alors δ^π est admissible.

Théorème 5 Si un estimateur de Bayes δ^π associé à une mesure *a priori* π (probabiliste ou non) et une fonction de coût L est unique, alors il est admissible.

Notons que les critères de minimaxité et d'admissibilité sont éminemment *fréquentistes* (car construits à partir du risque fréquentiste). Selon ces critères fréquentistes, les estimateurs de Bayes font mieux ou au moins aussi bien que les estimateurs fréquentistes :

- leur risque minimax est toujours égal ou plus petit ;
- ils sont tous admissibles (si le risque de Bayes est bien défini).

Les estimateurs de Bayes, plus généralement, sont souvent optimaux pour les concepts fréquentistes d'optimalité et devraient donc être utilisés même lorsque l'information *a priori* est absente. On peut ignorer la signification d'une distribution *a priori* tout en obtenant des estimateurs corrects d'un point de vue fréquentiste.

3.3 Choix d'une fonction de coût

La fonction de coût L est l'élément fondamental du choix d'un estimateur. Le choix dépend du contexte décisionnel et s'écrit souvent sous la forme

$$L = \text{Coût financier, etc.} - \text{Bénéfice}.$$

Une alternative, lorsqu'il est difficile de la construire, est de faire appel à des *fonctions de coût usuelles, mathématiquement simples et de propriétés connues*. L'idée est simplement de construire une "distance" usuelle entre $\theta \in \Theta$ et $d \in \mathcal{D}$ permettant une bonne optimisation (convexe par exemple).

EXEMPLE 7. Fonction de coût quadratique Soit $\mathcal{D} = \Theta$. On pose

$$L(\theta, \delta) = \|\theta - \delta\|^2. \quad (5)$$

Cette fonction de coût constitue le critère d'évaluation le plus commun. Elle est convexe (mais pénalise très (trop) fortement les grands écarts peu vraisemblables). Elle est justifiée par sa simplicité, le fait qu'elle permet de produire des estimateurs de Bayes intuitifs, et qu'elle peut être vue comme issue d'un développement limité d'un coût symétrique complexe.

Proposition 1 *L'estimateur de Bayes associé à toute loi a priori π et au coût (5) est l'espérance (moyenne) de la loi a posteriori $\pi(\theta|\mathbf{x}_n)$*

La fonction de coût absolu, également convexe, croît plus lentement que le coût quadratique et ne surpénalise pas les erreurs grandes et peu vraisemblables.

EXEMPLE 8. **Fonction de coût absolu (Laplace 1773)** Soit $\mathcal{D} = \Theta$ et $\dim \Theta = 1$. On pose

$$L(\theta, \delta) = |\theta - \delta| \quad (6)$$

ou plus généralement une fonction linéaire par morceaux

$$L_{c_1, c_2}(\theta, \delta) = \begin{cases} c_2(\theta - \delta) & \text{si } \theta > \delta \\ c_1(\delta - \theta) & \text{sinon} \end{cases} \quad (7)$$

Proposition 2 *L'estimateur de Bayes associé à toute loi a priori π et au coût (7) est le fractile $c_1/(c_1 + c_2)$ de la loi a posteriori $\pi(\theta|\mathbf{x}_n)$. En particulier, la médiane de la loi a posteriori est l'estimateur de Bayes lorsque $c_1 = c_2$ (qui sont donc des coûts associés à la sous-estimation et la surestimation de θ).*

La fonction de coût 0-1, non quantitative, est utilisé dans l'approche statistique classique pour construire des test d'hypothèse.

EXEMPLE 9. **Fonction de coût 0-1**

$$L(\theta, \delta) = \begin{cases} 1 - \delta & \text{si } \theta \in \Theta_0 \\ \delta & \text{sinon} \end{cases} \quad (8)$$

Le risque fréquentiste associé est

$$R(\theta, \delta) = \mathbb{E}_\theta[L(\theta, \delta(x))] = \begin{cases} P_\theta(\delta(x) = 0) & \text{si } \theta \in \Theta_0 \\ P_\theta(\delta(x) = 1) & \text{sinon} \end{cases}$$

Proposition 3 *L'estimateur de Bayes associé à toute loi a priori π et au coût 0-1 est*

$$\delta^\pi = \begin{cases} 1 & \text{si } \Pi(\theta \in \Theta_0|\mathbf{x}_n) > \Pi(\theta \notin \Theta_0|\mathbf{x}_n) \\ 0 & \text{sinon} \end{cases}$$

Ainsi, l'estimation bayésienne permet d'accepter une hypothèse (nulle) $H_0 : \theta \in \Theta_0$, si c'est l'hypothèse la plus probable *a posteriori*, ce qui est une réponse intuitive.

Une variante du test 0-1 est le test de Neyman-Pearson qui permet de distinguer risques de première et de deuxième espèce :

$$L(\theta, d) = \begin{cases} 0 & \text{si } d = \mathbb{1}_{\Theta_0} \\ a_0 & \text{si } \theta \in \Theta_0 \text{ et } d = 0 \\ a_1 & \text{si } \theta \notin \Theta_0 \text{ et } d = 1 \end{cases} \quad (9)$$

qui donne l'estimateur bayésien

$$\delta^\pi(x) = \begin{cases} 1 & \text{si } \Pi(\theta \in \Theta_0|x) > a_1/(a_0 + a_1), \\ 0 & \text{sinon.} \end{cases}$$

Ainsi, l'hypothèse nulle est rejetée quand la probabilité *a posteriori* de H_0 est trop petite. Il est cependant délicat de choisir les poids a_0 et a_1 sur des considérations d'utilité.

Plus généralement, ce résultat permet d'illustrer une différence majeure entre statistique classique et statistique bayésienne. L'approche classique (dite de Fisher-Neyman-Pearson) suppose qu'on puisse définir une statistique de test dont la loi, sous l'hypothèse nulle H_0 , est indépendante du paramètre estimé sous H_0 . Ce faisant, la seule décision que l'on peut prendre avec une bonne certitude est de refuser H_0 . Cette dissymétrie entre H_0 et toute autre hypothèse alternative H_1 n'existe pas dans le cadre bayésien : celui-ci émet un prior sur chaque modèle en compétition, puis compare les modèles selon leur probabilité d'explication des données disponibles *a posteriori*. Cette approche semble plus séduisante d'un point de vue intuitif et opérationnel. Voir également § 3.6 pour plus de détails.

3.4 Coûts intrinsèques

On peut enfin chercher à trouver des fonctions de coûts qui restent invariantes par *transformation monotone inversible* sur les données (action d'un C^1 -difféomorphisme sur Ω). On obtient ce faisant des fonctions de coûts définies à partir de *distances* ou de *divergences* D entre distributions

$$L(\theta, d) = D(f(\cdot|\theta) \parallel f(\cdot|d)).$$

Ci-dessous, quelques distances ou divergences usuelles entre des densités $(f_\theta, f_{\theta'})$ de fonctions de répartition $(F_\theta, F_{\theta'})$, qui induisent des fonctions de coût intrinsèques, sont présentées.

1. Distance de Kolmogoroff-Smirnoff :

$$d_{KS}(f_\theta, f_{\theta'}) = \sup_x |F_\theta(x) - F_{\theta'}(x)|$$

2. Distance L^1 :

$$d_1(f_\theta, f_{\theta'}) = \int |f_\theta(x) - f_{\theta'}(x)| dx \quad (2.1)$$

$$= 2 \sup_A |P_\theta(A) - P_{\theta'}(A)| \quad (2.2)$$

3. Distance de Hellinger :

$$d_H(f_\theta, f_{\theta'}) = \left(\int (\sqrt{f_\theta(x)} - \sqrt{f_{\theta'}(x)})^2 dx \right)^{\frac{1}{2}}$$

4. Pseudo-distance⁸ de Kullback-Liebler :

$$K(f_\theta, f_{\theta'}) = \int f_\theta(x) \log \frac{f_\theta(x)}{f_{\theta'}(x)} dx$$

Avec l'inégalité de Jensen, on prouve l'inégalité $K(f_\theta, f_\delta) \geq 0$. De plus, $K(f_\theta, f_\delta) = 0$ si et seulement si $f_\theta = f_{\theta'}$ μ -presque sûrement.

5. Distance L^2 :

$$d_2(f_\theta, f_{\theta'}) = \int (f_\theta(x) - f_{\theta'}(x))^2 dx$$

Ceci peut s'utiliser si les densités sont de carré intégrable.

Rappels sur la divergence de Kullback-Leibler (KL)

La divergence KL a un sens issu de la théorie de l'information (qui sera rappelé plus loin dans le cours). Rappeler ses principales propriétés peut être utile. Soit $\pi_1(\theta)$ et $\pi_2(\theta)$ deux densités de probabilité définies sur un même espace Θ , π_1 étant absolument continue par rapport à $\pi_2(\theta)$. Alors

$$\begin{aligned} \text{KL}(\pi_1 \parallel \pi_2) &= \mathbb{E}_{\pi_1} \left[\log \frac{\pi_1(\theta)}{\pi_2(\theta)} \right], \\ &= \int_{\Theta} \pi_1(\theta) \log \frac{\pi_1(\theta)}{\pi_2(\theta)} d\theta. \end{aligned}$$

Proposition 4 Si $\pi_1(\theta)$ et $\pi_2(\theta)$ sont propres, alors $KL(\pi_1||\pi_2) \geq 0$ et vaut 0 si et seulement si $\pi_1 = \pi_2$ (presque partout)

La preuve de ce résultat est fondée sur l'application de l'inégalité de Jensen à la fonction $\pi_1 \rightarrow \pi_1 \log \pi_1$, qui est différentiable deux fois sur l'espace des densités de probabilités, et dont la différentielle seconde vaut $1/\pi_1 > 0$. Elle est donc convexe. Alors $KL(\pi_1||\pi_2) \geq -\log 1 = 0$.

Exercice 7 Lorsqu'on fait un choix de fonction de coût $L(\theta, \delta)$ dans un ensemble $U : \Theta \times \mathcal{D} \rightarrow \Lambda \in \mathbb{R}^+$, on commet une erreur par rapport à la meilleure fonction de coût possible pour le problème. On peut donc proposer un estimateur bayésien de cette fonction de coût en introduisant une fonction de coût sur les fonctions de coût $L(\theta, \delta)$:

$$\begin{aligned} \tilde{L} : \Theta \times U \times \mathcal{D} &\rightarrow \mathbb{R}^+ \\ (\theta, \ell, \delta) &\mapsto \tilde{L}(\theta, \ell, \delta). \end{aligned}$$

Quel est l'estimateur bayésien de $\tilde{L}(\theta, \ell, \delta)$ sous un coût quadratique, lorsque $L(\theta, \delta)$ est elle-même quadratique ?

3.5 Mode a posteriori (MAP)

L'estimateur du mode *a posteriori*, ou MAP, est défini par

$$\delta^\pi(\mathbf{x}_n) = \arg \max_{\theta \in \Theta} \pi(\theta|\mathbf{x}_n).$$

Cet estimateur, contrairement aux précédents, n'est pas issu de la minimisation d'une fonction de coût (il n'est donc pas bayésien *stricto sensu*) mais peut être vu comme la limite d'estimateurs bayésiens.

Il correspond à un maximum de vraisemblance (MV) pénalisé (voir § 2.5) et souffre donc en général des mêmes inconvénients que le MV, en particulier une certaine instabilité d'estimation ponctuelle. Par ailleurs, à la différence du MV, il est en général non invariant par reparamétrisation. Cette gêne décisionnelle mène à le déconseiller formellement, ou du moins à s'en méfier, même si ce type d'estimateur est couramment privilégié par les praticiens du *machine learning*.

3.6 Sélection de modèle et facteur de Bayes

La sélection de modèle bayésien est un choix particulier de décision. D'une façon générale, supposons vouloir tester deux hypothèses de modèles, M_0 et M_1 , l'une contre l'autre, et pouvoir assigner à ces modèles des probabilités *a priori* non nulles

$$\Pi(M_0), \Pi(M_1)$$

d'expliquer des données X non encore observées. Typiquement, on peut vouloir :

- pour un même modèle de vraisemblance $f(X|\theta)$, tester *a posteriori* deux sous-ensembles différents de $\Theta : \Pi(\theta \in \Theta_0|X) > 0$ et $\Pi(\theta \in \Theta_1|X) > 0$;
- plus généralement tester un couple $(f_0(x|\theta_0), \pi_0(\theta_0))$ contre un couple $(f_1(x|\theta_1), \pi_1(\theta_1))$, avec $\theta_i \in \Theta_i$.

Le facteur de Bayes est le *rapport de la vraisemblance marginale* de ces deux hypothèses concurrentes :

$$B_{01}(X) = \frac{P(X|M_0)}{P(X|M_1)} \quad (10)$$

où

$$\begin{aligned} P(X|M_i) &= \int_{\Theta_i} f_i(X|\theta_i) \pi_i(\theta_i) d\theta_i, \\ &= \frac{\Pi(M_i|X)P(X)}{\Pi(M_i)} \end{aligned}$$

où $P(X)$ représente la loi inconnue des données. Le rapport (10) peut alors se simplifier en

$$B_{01}(X) = \frac{\Pi(M_0|X) \Pi(M_1)}{\Pi(M_1|X) \Pi(M_0)}. \quad (11)$$

Le facteur de Bayes est donc une transformation bijective de la probabilité *a posteriori*, qui a fini par être l'outil le plus utilisé pour choisir un modèle bayésien. Lorsque les deux modèles sont également probables *a priori*, alors $\frac{\Pi(M_1)}{\Pi(M_0)} = 1$ et le rapport de Bayes est simplement le rapport de leurs probabilités *a posteriori*.

Le cas le plus fréquemment rencontré en pratique est celui où l'hypothèse de modèle M_i se réduit à $\theta \in \Theta_i$. Il amène à la définition suivante.

Définition 13 *Le facteur de Bayes associé au problème du choix entre $\theta \in \Theta_0$ et $\theta \in \Theta_1$ est le rapport des probabilités a posteriori des hypothèses nulle et alternative sur le rapport a priori de ces mêmes hypothèses*

$$B_{01}(X) = \left(\frac{\Pi(\theta \in \Theta_0|X)}{\Pi(\theta \in \Theta_1|X)} \right) / \left(\frac{\Pi(\theta \in \Theta_0)}{\Pi(\theta \in \Theta_1)} \right) \quad (12)$$

qui se réécrit comme le pendant bayésien du rapport de vraisemblance en remplaçant les vraisemblances par les marginales (les vraisemblances intégrées sur les *a priori*) sous les deux hypothèses

$$B_{01}(X) = \frac{\int_{\Theta_0} f(X|\theta) \pi_0(\theta) d\theta}{\int_{\Theta_1} f(X|\theta) \pi_1(\theta) d\theta} = \frac{f_0(X)}{f_1(X)}$$

Sous le coût généralisé (9), en posant

$$\gamma_0 = \Pi(\theta \in \Theta_0) \quad \text{et} \quad \gamma_1 = \Pi(\theta \in \Theta_1).$$

Ainsi l'hypothèse H_0 est acceptée si

$$B_{01}(x) > (a_1 \gamma_1) / (a_0 \gamma_0).$$

Dans la pratique, on utilise souvent des échelles logarithmiques pour faire une sélection de modèle (échelle de Jeffreys améliorée par Kass et Raftery) :

- (i) si $\Lambda = \log_{10} B_{10}(\mathbf{x}_n)$ varie entre 0 et 0.5, la certitude que H_0 est fautive est faible ;
- (ii) si $\Lambda \in [0.5, 1]$, cette certitude est substantielle ;
- (iii) si $\Lambda \in [1, 2]$, elle est forte ;
- (iv) si $\Lambda > 2$, elle est décisive.

Malgré le côté heuristique de l'approche, ce genre d'échelle reste très utilisé.

Exercice 8 Soit $X \sim \mathcal{B}(\theta)$ (loi de Bernoulli) avec $\Theta = [0, 1]$. Soit M_0 un modèle défini par $\{\theta = 1/2\}$ et M_1 un modèle défini par un θ inconnu dans $[0, 1]$, avec $\pi_1(\theta) = \mathcal{U}[0, 1]$. Un échantillon de 200 tirages fournit 115 succès et 85 échecs. Au vu de ces données, quel modèle choisir ? Ce résultat diffère-t-il significativement d'un test fréquentiste ?

Remarque 6 Le calcul du facteur de Bayes n'est pas évident et demande le plus souvent de savoir simuler *a posteriori*.

Cas de l'estimation ponctuelle et des tests de significativité en régression

Dans la définition (12), on sous-entend que chaque alternative $\Pi(\theta \in \Theta_i) > 0$ sinon le facteur de Bayes n'est pas défini. Cela exclurait les situations fréquentes où Θ_i est de mesure nulle. Par exemple lorsqu'on veut tester $\Theta_0 = \{\theta_0\}$, ou pour mener un test de significativité pour les modèles de régression.

Considérons le cas suivant : on dispose d'un modèle de régression

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d$$

où $\beta_d \in \mathbb{R}$, et on veut mener un test de significativité sur $\theta = \beta_d$:

$$M_0 : \{\beta_d = 0\} \quad vs \quad M_1 : \{\beta_d \neq 0\}$$

On peut alors noter $\Theta_0 = \{\theta_0\} = \{0\}$ et $\Theta_1 = \mathbb{R}^*$. Il est clair que $\Pi(\theta \in \Theta_0) = 0$ car la mesure dominante est Lebesgue. Donc la formulation (12) n'est pas applicable. On peut alors reposer le test en redéfinissant plus généralement

$$M_0 : \{|\beta_d - \theta_0| = \epsilon\} \quad vs \quad M_0 : \{|\beta_d| \neq \epsilon\}$$

avec un ϵ tendant vers 0 par valeurs positives. Reste la difficulté majeure que ϵ est inconnu et dépend du contexte de l'étude.

Plus généralement, et pour mieux formaliser les choses, il convient alors d'introduire une masse de Dirac δ_{θ_0} en θ_0 et de considérer un poids $\rho_\epsilon \in [0, 1]$ à ϵ fixé, valant

$$\rho_\epsilon = \Pi(\theta || \theta - \theta_0 | < \epsilon) = \Pi(\theta \in \Theta_0)$$

Le mélange

$$\pi(\theta) = \rho_\epsilon \delta_{\theta_0} + (1 - \rho_\epsilon) \pi_1(\theta),$$

désigne la loi *a priori* qui généralise $\Pi(\theta \in \Theta_0)$ et $\Pi_1(\theta) = \Pi(\theta \in \Theta_1)$, la densité de cette dernière loi étant absolument continue par rapport à la mesure de Lebesgue. Dans ce cas, l'application de la formule (12) donne directement

$$\begin{aligned} B_{01}(X) &= \frac{\rho_\epsilon f(X|\theta_0)}{(1 - \rho_\epsilon) \int_{\theta \in \Theta_1} f(X|\theta) \pi_1(\theta) d\theta} \left(\frac{1 - \rho_\epsilon}{\rho_\epsilon} \right), \\ &= \frac{f(X|\theta_0)}{\int_{\theta \in \Theta_1} f(X|\theta) \pi_1(\theta) d\theta}. \end{aligned}$$

(on retrouve bien le rapport des lois marginales).

Cas d'un ou plusieurs priors impropres

Une remarque importante concerne le cas où π_0 ou π_1 n'est pas un *prior propre* (mesure non intégrable). Dans ce cas, B_{01} n'est alors pas défini de manière unique. En effet, si une loi est impropre (par exemple π_0) alors elle est définie à une constante multiplicative près. Pour $\pi_0^*(\theta) = c\pi_0(\theta)$, alors le facteur de Bayes est lui aussi multiplié par c : $B_{01}^* = cB_{01}$ et par conséquent les ordres de grandeur de B_{01} n'ont plus de sens : il n'est plus possible de comparer ses valeurs à une échelle dédiée.

EXEMPLE 10. Soit $X \sim \mathcal{N}(\theta, 1)$ et $\pi(\theta) = c \neq 0$. On considère le test suivant : $H_0 : \theta = 0$ vs $H_1 : \theta \neq 0$. Dans ce cas, le facteur de Bayes est

$$B_{01}(X = x) = \frac{\exp(-x^2/2)}{c\sqrt{2\pi}}$$

c étant inconnu, le facteur de Bayes ne peut avoir d'interprétation.

Pour pallier ce problème, des démarches ont été inventées, qui visent à séparer l'échantillon X en sous-échantillons et transformant les priors impropres en posteriors propres (mais faiblement informatifs) : les *facteurs de Bayes fractionnels* ou les *facteurs de Bayes intrinsèques*. Ils ont de bonnes propriétés asymptotiques. Toutefois, dans les dernières années, la recherche se tourne plutôt vers des approches par *mélanges de modèles bayésiens*, qui permettent de remplacer la sélection d'un sous-modèle via le facteur de Bayes par la sélection d'un sous-modèle via les poids du mélange. Ce type d'approche permet de se débarrasser, sous certaines conditions, des problèmes posés par les priors impropres.

Exercice 9 Sélection de modèle discret avec des priors impropres.

1. Pour des données discrètes x_1, \dots, x_n , on considère un modèle de Poisson $\mathcal{P}(\lambda)$ ou une loi binomiale négative $\mathcal{NB}(m, p)$ avec les a priori

$$\begin{aligned}\pi_1(\lambda) &\propto 1/\lambda \\ \pi_2(m, p) &= \frac{1}{M} \mathbb{1}_{\{1, \dots, M\}}(m) \mathbb{1}_{[0, 1]}(p)\end{aligned}$$

Peut-on sélectionner l'un des deux modèles ?

2. Si on remplace $\pi_1(\lambda)$ par un a priori **vague**

$$\pi_1(\lambda) \equiv \mathcal{G}(\alpha, \beta)$$

avec $\alpha(\beta)$ ou/et $\beta(\alpha) \rightarrow 0$, peut-on de nouveau résoudre le problème ?

3.7 TP : Création d'un système d'alerte pour la circulation routière

On s'intéresse à un évènement routier $X = x$ relevé par un système de détection vivant dans l'espace χ de dimension finie. Ce système de détection peut prédire des évènements répétés du type "un animal sur la voie", "accrochage", "accident", "bouchon"... La question est de déterminer si, à chaque fois qu'un évènement routier x est collecté, il est utile qu'une intervention de secours soit menée.

Nommons θ une variable indiquant la gravité de l'évènement. Cette variable a des valeurs dans les ensembles disjoints Θ_0 (incidents sans gravité) et Θ_1 (accidents nécessitant possiblement une intervention). On suppose disposer d'un échantillon labélisé $\mathbf{e}_n = (\mathbf{x}_n, \theta_n)$.

Questions.

1. Lorsqu'une observation x apparaît, comment prévoir θ ?
2. Comment peut-on en déduire une alarme efficace ?

4 Propriétés fondamentales du cadre bayésien

4.1 Prédiction (prévision)

Le contexte du problème de la prédiction est le suivant : les observations X sont identiquement distribuées selon P_θ , qui est absolument continue par rapport à une mesure dominante μ . Il existe donc une fonction de densité conditionnelle $f(\cdot|\theta)$. Par ailleurs on suppose que θ suit une loi a priori π . *Mener une prévision* consiste alors, à partir de n tirages observés x_1, \dots, x_n , de déterminer le plus précisément possible ce que pourrait être le tirage suivant X_{n+1} .

Dans l'approche fréquentiste, on calcule dans les faits $f(x_{n+1}|x_1, \dots, x_n, \hat{\theta}_n)$, puisqu'on ne connaît pas θ et qu'on doit l'estimer : on utilise donc deux fois les données (une fois pour l'estimation de θ , et une nouvelle fois pour la prévision). En règle générale, ceci amène à sous-estimer les intervalles de confiance.

La stratégie du paradigme bayésien consiste à intégrer la prévision suivant la loi courante *a posteriori* sur θ et ce, afin d'avoir la meilleure prévision compte-tenu à la fois de notre savoir et de notre ignorance sur le paramètre. La loi prédictive s'écrit ainsi :

$$f(X_{n+1}|x_1, \dots, x_n) = \int_{\Theta} f(X_{n+1}|x_1, \dots, x_n, \theta) \pi(\theta|x_1, \dots, x_n) d\theta$$

qui s'écrit plus simplement, lorsque *sachant* θ les tirages sont iid :

$$f(x_{n+1}|x_1, \dots, x_n) = \int_{\Theta} f(x_{n+1}|\theta) \pi(\theta|x_1, \dots, x_n) d\theta.$$

Ainsi, le prédicteur de X_{n+1} sous le coût quadratique est

$$\mathbb{E}[X_{n+1}|x_1, \dots, x_n] = \int_{\Omega} x f(x|x_1, \dots, x_n) dx.$$

4.2 Propriétés asymptotiques

Les approches classique et bayésienne de la modélisation et de la décision statistique aboutissent à des résultats similaires à l'asymptotisme, et les principaux théorèmes classiques connaissent leur pendant bayésien. Ainsi, le théorème central limite "classique" devient le théorème de Bernstein-von Mises dans le cadre bayésien (on l'appelle également *théorème central limite bayésien* par abus de langage). Afin de comparer les deux approches, on doit d'abord définir ce que signifie "vraie valeur θ_0 du paramètre θ ".

Notons $\tilde{f}(x)$ la "vraie loi" inconnue des données. Si on fait maintenant le choix d'une loi paramétrique $X \sim f(x|\theta_0)$ (ou mécanisme génératif), alors la loi $f(x|\theta_0)$ doit être la plus proche possible de $\tilde{f}(x)$. Cette notion de proximité est généralement définie de la façon suivante.

Définition 14 Soit $\tilde{f}(x)$ la loi inconnue des données. On définit θ_0 par

$$\theta_0 = \arg \min_{\theta \in \Theta} KL(\tilde{f}(x) || f(x|\theta))$$

où KL est la divergence de Kullback-Leibler. On notera par la suite plus simplement ce terme $KL(\theta)$.

Théorème 6 Consistance Si $f(\cdot|\theta)$ est suffisamment régulière et identifiable, soit si $\theta_1 \neq \theta_2 \Rightarrow f(x|\theta_1) \neq f(x|\theta_2) \forall x \in \Omega$, alors pour tout échantillon \mathbf{x}_n iid

$$\pi(\theta|\mathbf{x}_n) \xrightarrow{p.s.} \delta_{\theta_0}.$$

Par ailleurs, si $g : \Theta \rightarrow \mathbb{R}$ est mesurable et telle que $\mathbb{E}[g(\theta)] < \infty$, alors sous les mêmes hypothèses

$$\lim_{n \rightarrow \infty} \mathbb{E}[g(\theta)|X_1, \dots, X_n] = g(\theta) \text{ p.s.}$$

Un résultat utile, intermédiaire entre la consistance et la convergence en loi (Théorème 9), est la convergence en probabilité.

Théorème 7 Si Θ est fini et discret et $\Pi(\theta = \theta_0) > 0$, alors pour tout échantillon iid $X_1, \dots, X_n | \theta \sim f(X|\theta)$,

$$\Pi(\theta = \theta_0 | X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 1.$$

Si Θ est continu, alors $\pi(\theta_0|x)$ vaut toujours 0 pour tout échantillon fini x , et on ne peut appliquer les outils menant au résultat précédent. Pour adapter cette preuve, il faut définir un voisinage V_{θ_0} qui est un ensemble ouvert de points de Θ à une distance maximum fixée de θ_0 (Θ étant un espace métrique).

Théorème 8 Si Θ est un ensemble compact et si V_{θ_0} est tel que $\Pi(\theta \in V_{\theta_0}) > 0$ avec

$$\theta_0 = \arg \min_{\theta \in \Theta} KL(\theta)$$

alors

$$\Pi(\theta \in V_{\theta_0} | x_1, \dots, x_n) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 1.$$

Le théorème de Bernstein-von Mises suppose l'existence de l'information de Fisher I_θ . Il n'existe pas d'ensemble de conditions de régularité minimal nécessaire pour l'existence de I_θ ; cependant, la plupart des auteurs s'accordent sur les conditions suffisantes suivantes d'existence, de positivité et de continuité dans un sous-espace de Θ :

- $f(x|\theta)$ est absolument continue en θ ;
- sa dérivée doit exister pour tout $x \in \Omega$.

Alors

$$I_\theta = \mathbb{E} \left[\left(\frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2 \right] = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta) \right]$$

si $\log f(x|\theta)$ est deux fois différentiable en θ .

Théorème 9 Normalité asymptotique (**Bernstein-von Mises**) Soit I_θ la matrice d'information de Fisher du modèle $f(\cdot|\theta)$ et soit $g(\theta)$ la densité de la gaussienne $\mathcal{N}(0, I_\theta^{-1})$. Soit $\hat{\theta}_n$ le maximum de vraisemblance. Alors, dans les conditions précédentes,

$$\int_{\Theta} \left| \pi \left(\sqrt{n} \{ \theta - \hat{\theta}_n \} | \mathbf{x}_n \right) - g(\theta) \right| d\theta \rightarrow 0.$$

4.3 Régions de crédibilité et régions HPD

Soit $x \sim f(\cdot|\theta)$ une (ou plusieurs) observations.

Définition 15 Région α -crédible Une région A de Θ est dite α -crédible si $\Pi(\theta \in A|x) \geq 1 - \alpha$.

Notons que le paradigme bayésien permet une nouvelle fois de s'affranchir d'un inconvénient de l'approche fréquentiste. Rappelons qu'au sens fréquentiste, A est une région de confiance $1 - \alpha$ si, en refaisant l'expérience (l'observation d'un $X \sim f(\cdot|\theta)$) un nombre de fois tendant vers ∞ ,

$$P_\theta(\theta \in A) \geq 1 - \alpha.$$

Une région de confiance n'a donc de sens que pour un très grand nombre d'expériences tandis que la définition bayésienne exprime que la probabilité que θ soit dans A au vu des celles déjà réalisées est plus grande que $1 - \alpha$. Il n'y a donc pas besoin ici d'avoir recours à un nombre infini d'expériences pour définir une région α -crédible, seules comptent les expériences effectivement réalisées.

Remarque 7 On distingue bien ici la probabilité "fréquentiste" P_θ de la probabilité bayésienne Π . Dans le premier cas, l'aléatoire concerne la région A , qui est un estimateur statistique dépendant d'un estimateur classique $\hat{\theta}(X_1, \dots, X_n)$ et θ est considéré comme fixe. Dans le second cas, c'est bien θ qui est aléatoire.

Il y a une infinité de régions α -crédibles. Il est donc logique de s'intéresser à la région qui a le volume minimal. Le volume étant défini par $\text{vol}(A) = \int_A d\mu(\theta)$, si $\pi(\theta|x)$ est absolument continue par rapport à une mesure de référence μ .

Définition 16 Région HPD. $A_{\alpha,\pi}$ est une région HPD (highest posterior density) si et seulement si

$$A_{\alpha,\pi} = \{\theta \in \Theta, \pi(\theta|x) \geq h_\alpha\}$$

où h_α est défini par

$$h_\alpha = \sup_h \{\Pi(\theta|\pi(\theta|x) \geq h, X) \geq 1 - \alpha\}.$$

$A_{\alpha,\pi}$ est parmi les régions qui ont une probabilité supérieure à $1 - \alpha$ de contenir θ (et qui sont donc α -crédibles) et sur lesquelles la densité *a posteriori* ne descend pas sous un certain niveau (restant au dessus de la valeur la plus élevée possible).

Théorème 10 $A_{\alpha,\pi}$ est parmi les régions α -crédibles celle de volume minimal si et seulement si elle est HPD.

Exercice 10 Soit x_1, \dots, x_n des réalisations iid de loi $\mathcal{N}(\mu, \sigma^2)$. On choisit la mesure *a priori* (non probabiliste) jointe

$$\pi(\mu, \sigma^2) \propto 1/\sigma^2.$$

1. Déterminez la loi *a posteriori* jointe $\pi(\mu, \sigma^2|x_1, \dots, x_n)$
2. Déterminez la loi *a posteriori* marginale $\pi(\mu|x_1, \dots, x_n)$
3. Calculez la région HPD de seuil α pour μ et comparez-la à la région de confiance fréquentiste, de même seuil, qu'on pourrait calculer par l'emploi du maximum de vraisemblance.
4. Déterminez la loi *a posteriori* marginale $\pi(\sigma^2|x_1, \dots, x_n)$; le calcul de la région HPD est-il simple ?

Remarque. "Déterminez" signifie indiquer si la loi appartient à une famille connue, par exemple largement implémentée sur machines. La connaissance des lois gamma, inverse gamma et Student est peut-être nécessaire pour répondre aux questions.

Les régions HPD sont à manier avec précaution, car elles ne sont pas indépendantes de la paramétrisation.

Exercice 11 Soit $A_{\alpha,\pi} = \{\theta \in \Theta, \pi(\theta|x) \geq h_\alpha\}$ une région HPD et soit

$$\eta = g(\theta)$$

un C^1 -difféomorphisme (bijection). On définit alors la région HPD correspondante pour $\pi(\eta|x)$:

$$\tilde{A}_{\alpha,\pi} = \left\{ \eta \in g(\Theta), \pi(\eta|x) \geq \tilde{h}_\alpha \right\}$$

- Sous quelle condition peut-on écrire que $\tilde{A}_{\alpha,\pi} = g(A_{\alpha,\pi})$?
- Illustrons cela en supposant $X \sim \mathcal{N}(\theta, 1)$ et $\pi(\theta) \propto 1$, puis en posant $\eta = \exp(\theta)$.

Nous pouvons comprendre pourquoi une région de confiance n'est pas invariante par reparamétrisation. En effet, cette région se définit comme une solution du problème de minimisation suivant :

$$A_{\alpha,\pi} = \arg \min_{A, \Pi(A|X) \geq 1-\alpha} \text{Vol}(A)$$

où $\text{Vol}(A) = \int_A d\mu(\theta)$. Or la mesure de Lebesgue n'est pas invariante par reparamétrisation. Une idée pour lever cette difficulté est donc logiquement d'abandonner la mesure de Lebesgue et de considérer pour une mesure s :

$$A_{\alpha,\pi,s} = \arg \min_{A, \Pi(A|X) \geq 1-\alpha} \int_A ds(\theta).$$

Calcul de régions HPD

Pour calculer les régions HPD, il y a plusieurs méthodes :

1. *Méthode analytique et numérique* : c'est ce qui a été fait lors de l'exemple précédent. Précisons une nouvelle fois que cette méthode ne peut s'appliquer que dans des cas assez rares.
2. *Méthode par approximation* : cette méthode peut être appliquée si le modèle est régulier. L'usage du théorème de Bernstein-von Mises permet d'approximer la loi *a posteriori* par une gaussienne. On retombe peu ou prou sur des régions HPD proches de celles du maximum de vraisemblance.
3. *Méthode par simulation*. En effet, une région α -crédible peut génériquement être estimée par les quantiles empiriques de la simulation *a posteriori* (voir plus loin).

Théorème 11 Supposons avoir un échantillon iid $\theta_1, \dots, \theta_m \sim \pi(\theta|x_1, \dots, x_n)$ avec $\theta \in \mathbb{R}$. Alors les intervalles de quantiles empiriques de la forme $[\theta^{(\alpha/2)}, \theta^{(1-\alpha/2)}]$ sont tels que

$$\Pi \left(\theta \in \left[\theta^{(\alpha/2)}, \theta^{(1-\alpha/2)} \right] | x_1, \dots, x_n \right) \xrightarrow{m \rightarrow \infty} 1 - \alpha.$$

Il n'est cependant pas garanti qu'une telle région soit HPD. Pour m grand, $\theta^{(\alpha/2)}$ s'approche du quantile d'ordre $\alpha/2$ de la loi *a posteriori*. Cette région n'est pas nécessairement HPD mais reste α -crédible. Cette méthode est particulièrement adaptée lorsque la loi *a priori* est unimodale. Il est toujours utile de représenter graphiquement les sorties pour fixer les idées. Enfin, il est aussi envisageable d'avoir recours à une estimation non paramétrique par noyaux.

5 Compréhension et représentation de l'information incertaine

5.1 Une vision subjectiviste de la théorie bayésienne

La fonction de coût et le processus décisionnel permettent de proposer une interprétation importante de la distribution *a priori*. Elle peut être comprise comme pari (personnel) fait sur l'éventualité d'un événement, et notamment un gain conditionné par l'occurrence du phénomène modélisé par $f(x|\theta)$. Cette interprétation subjective, proposée par de Finetti (1948), est certainement le point le plus critiqué de la démarche bayésienne, mais c'est aussi celui qui permet à l'application de cette théorie d'être ancrée dans le réel.

On peut en fait mieux appréhender cette interprétation subjectiviste en la reliant à l'histoire récente de l'axiomatique de la connaissance incertaine.

5.2 Théories de la connaissance incertaine

Dans l'histoire des théories de représentation mathématique de la connaissance incertaine, il existe essentiellement deux grandes écoles de pensée :

1. des **théories de la représentation qui s'adaptent** aux moyens variés, pour un humain, d'exprimer son opinion personnelle sur le comportement d'une variable d'ancrage X ou d'un paramètre perceptible θ (plus rare) ;
 - *Exemples* : théories extra-probabilistes : Dempster-Schafer, possibilités, logique floue ...
2. des **théories qui visent à établir des axiomes de rationalité** à propos des décisions sous-tendant l'expression d'une opinion : un expert est perçu comme un preneur de décision selon ces axiomes.

Une vision axiomatique de la représentation mathématisée de la connaissance incertaine, qui permet d'interpréter la théorie des probabilités comme une "bonne" façon de représenter cette connaissance (ou plutôt cette information), a été construite par Cox et Jaynes. Elle s'incarne dans le théorème de Cox-Jaynes, dont les versions successives, au cours du temps, sont devenus les théorèmes fondamentaux de l'intelligence artificielle. Ce théorème permet de donner un cadre plus robuste à la vision subjectiviste du sens d'un modèle *a priori*.

5.3 Une vision plus claire de la statistique bayésienne

En définitive, il apparaît après ces premiers chapitres que la statistique bayésienne est à la fois :

- une théorie de la description d'un phénomène incertain, où "incertitude" signifie "mélange d'aléatoire (incertitude non-réductible) et d'épistémique (incertitude réductible) ;
- une théorie de la décision, sous certains axiomes de rationalité.

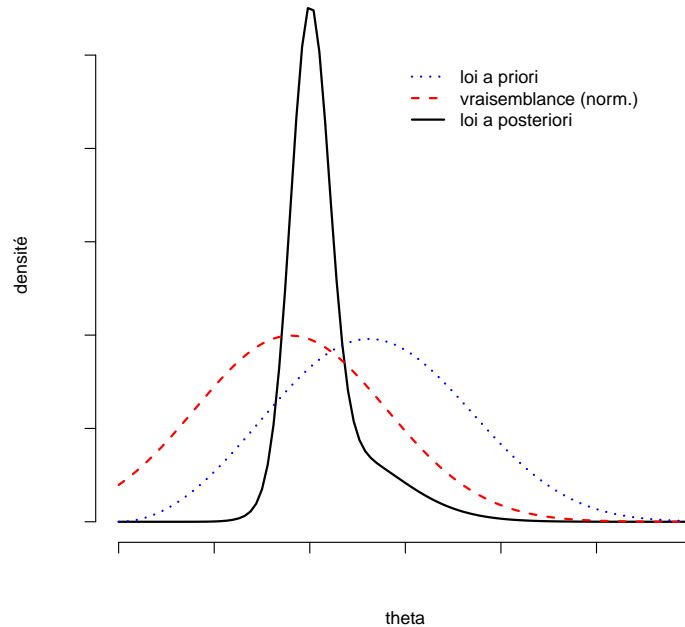
Sachant un modèle $f(x|\theta)$, le travail bayésien consiste donc à :

1. déterminer le coût associé aux décisions, $L(\theta, \delta)$;
2. éliciter ("construire") une loi *a priori* $\pi(\theta)$;
3. réaliser l'inférence *a posteriori* et produire un ou plusieurs estimateurs, voire faire un choix de modèle.

Notons qu'il y a redondance entre les deux premières étapes : présupposer l'existence d'une fonction de coût implique qu'une certaine information *a priori* sur le problème considéré est disponible.

La figure 14 illustre le positionnement classique de la vraisemblance statistique $f(\mathbf{x}_n|\theta)$, vue comme une fonction de θ , par rapport aux densités $\pi(\theta)$ et $\pi(\theta|\mathbf{x}_n)$; la mutualisation des sources d'information sur θ se traduit logiquement par une distribution *a posteriori* plus "piquée" – plus informative donc – que la loi *a priori*⁶.

6. Excepté dans les cas où les deux sources d'information sont en désaccord : voir [10].



380380

FIGURE 1 – Illustration par les densités du renforcement de l’information sur θ à partir de la loi *a priori* et la vraisemblance des données (vue comme fonction de θ et ici renormalisée).

Nature de la loi *a priori*. Qu’exprime la loi *a priori* $\pi(\theta)$? Un état d’information initial sur les valeurs possibles de θ , indépendant des observations \mathbf{x}_n . Encodé sous forme probabiliste, cet état d’information est donc susceptible de permettre l’ajout d’une connaissance réelle et sérieuse du phénomène exprimée autrement qu’au travers d’observations passées : expertise, prévisions de modèles physiques, etc. Il faut noter que l’information *a priori* sur θ n’est jamais inexistante : en effet, θ est un choix de paramétrisation du modèle de statistiques extrêmes choisi (type MAXB ou POT), et la structure de ce modèle est connue. En découlent des contraintes sur la structure de corrélation de $\pi(\theta)$ (cf. § ??). Nous conseillons au lecteur intéressé par une présentation didactique des possibles interprétations de $\pi(\theta)$ l’ouvrage [31] et l’article [32].

5.4 Validité des priors informatifs *via* la logique probabiliste de l’information incertaine

Un prior informatif est la représentation probabiliste d’une information *a priori*, généralement sur X (et non sur θ) *a piece of information whose the value of truth is justified by considerations independent on experiment on focus* (Pegny 2012). Cette information peut avoir pour origine :

- des résultats d’essais autres que les données (par exemple sur des maquettes) ;
- des spécifications techniques d’exploitation ;
- des bornes physiques ;
- des extractions de corpus référencés ;

- des connaissances exprimés par des experts humains.

Cette dernière information est souvent *incomplète*, toujours *incertaine*, à cause de la non-existence d'un système permettant *a priori* de vérifier si l'expertise est complète ou non ; mais aussi à cause de la non-existence d'un système suffisamment précis pour spécifier que $X = x_0$ exactement (sauf dans des cas rares et pathologiques).

Se pose alors la question de la pertinence de la théorie des probabilités pour représenter une information incertaine. Celle-ci offre de nombreux avantages pratiques, mais est-elle *auditable* ? La subjectivité qui semble inhérente à ce choix peut être un élément de critique fort et limiter la confiance dans les approches bayésiennes informatives. À cela, on peut tout d'abord tirer parti de nombreux travaux épistémologiques.

Hypothèse 1 (Lakatos, 1974). L'information sur l'état de la nature est cachée et partiellement révélée par une *théorie consensuelle* (au sens de Popper (1972) : *par décision mutuelle des protagonistes*) définissant l'*objectivité* [Gelman2015]. Dans ce cadre, la connaissance est un "filtrage" de l'information

Ce filtrage est produit par l'intervention de symboles, ou signes, afin de la *transmettre* ou de l'*implémenter*.

Hypothèse 2 (issu des neurosciences) [42, 41, 35, 15, 8, 5]. Face à des situations où de l'information incertaine est mobilisée, le raisonnement humain produit des inférences probabilistes. Les difficultés apparaissent au moment de l'explicitation de cette information par *langage interprétatif* \Rightarrow *expertise utilisable*.

Sous ces hypothèses, nous ne savons pas comment définir formellement la "déconvolution" retransformant la connaissance incertaine en information incertaine. Mais nous pouvons avoir des idées sur l'impact de l'ajout d'une connaissance incertaine mais utile sur la résolution du problème de détermination de X . Cet ajout se manifeste par un accroissement de l'information sur X – c'est-à-dire une *inférence* (mise à jour) :

- (a) cette inférence doit s'appuyer sur un principe de raisonnement ;
- (b) ce principe de raisonnement s'établit lui-même sur une *logique*, c'est-à-dire un ensemble de **règles formelles**.

Quelles propriétés souhaitons-nous pour cette logique ? Qu'elle permette de trier des assertions *atomiques* du type $X = x_0$ at chaque ajout d'information (*logique exclusive*). Ainsi, une situation initiale (*prémisse*) doit être moins informatif qu'une conclusion (*mise à jour*). De plus, cette logique doit permettre de représenter une information *incertaine* : on doit pouvoir trier plus d'assertions que celles qui sont simplement vraies ou fausses (*logique non-booléenne*). Construire une telle logique de l'information incertaine nécessite de définir les concepts suivants.

Définition 17 État d'information. Soit S_X un ensemble de propositions (assertions) atomiques du type $X = x_i$. L'ensemble B_X de toutes les propositions composées générées par

$$\begin{aligned} \neg X = x_i, \quad X = x_i \wedge X = x_j, \\ X = x_i \vee X = x_j, \quad X = x_i \Rightarrow X = x_j \\ \text{and} \quad X = x_i \Leftrightarrow X = x_j \end{aligned}$$

est appelé état d'information, avec $\text{Dom}(B_X) = \text{fermeture logique de } S_X$.

L'état d'information B_X résume l'information existante sur un ensemble de propositions portant sur X . Si la logique souhaitée précédemment existe, elle devrait guider la façon dont B_X évolue : il croît selon une certaine métrique quand l'information sur X croît. Pour mesurer cette croissance de l'information, il faut introduire le concept de plausibilité, qui lui-même permet le calcul propositionnel.

Définition 18 Plausibilité. *Considérons une proposition A on X . Sachant B_X , la plausibilité $[A|B_X]$ est un nombre réel, supérieurement borné par un nombre (fini ou infini) T .*

Supposer que la plausibilité existe est un axiome dit de *non-ambiguïté* est particulièrement important⁷. Il produit une hypothèse de *comparabilité universelle*⁸, qui a pour conséquence qu’une information additionnelle (pas forcément une connaissance) peut seulement faire croître ou décroître la plausibilité d’une proposition.

Définition 19 Consistance. B_X est consistant s’il n’existe aucune proposition A pour qui $[A|B_X] = T$ et $\neg[A|B_X] = T$.

Définition 20 Calcul propositionnel. *Le calcul propositionnel est un ensemble de règles permettant à tout domaine de problème de formuler des propositions utiles. Il s’exprime ainsi :*

- (i) Si $A = A'$ alors $[A|B_X] \Leftrightarrow [A'|B_X]$
- (ii) $[A|B_X, C_X, D_X] = [A|(B_X \wedge C_X), D_X]$
- (iii) Si B_X est consistant et $\neg[A|B_X] < T$, alors $A \cup B_X$ est consistant

- **Cohérence** : il existe une fonction non croissante S_0 telle que, pour tout x et tout B_X consistant

$$\neg[A|B_X] = S_0([A|B_X])$$

- **Densité** : l’ensemble $[S_0(T), T]$ admet un sous-ensemble non vide, dense et consistant.

Ce calcul propositionnel, qui repose sur l’axiome de non-ambiguïté, formalise donc des règles simples permettant de décrire les propriétés d’une logique de l’incertain. Celle-ci peut se vulgariser de la façon suivante :

1. *Règle de reproductibilité* : deux assertions équivalentes sur X ont la même plausibilité.
2. *Règle de non-contradiction* : s’il existe plusieurs approches aboutissant aux mêmes conclusions sur X , celles-ci ont la même plausibilité.
3. *Règle de consistance* : la logique ne peut formuler une conclusion contredite par les règles élémentaires de déduction (ex : transitivité).
4. *Règle d’intégrité* : la logique ne peut exclure une partie de l’information sur X pour parvenir à une conclusion sur X .
5. *Règle de monotonie* : la plausibilité d’une union non exclusive de deux assertions est au moins égale à la plus grande des plausibilités de chacune des assertions prises séparément.
6. *Règle de produit* : la plausibilité de l’intersection de deux assertions est au plus égale à la plus petite des plausibilités de chacune des assertions prises séparément.

On arrive dès lors à la question suivante : d’un point de vue pratique, cette logique est-elle équivalente à une théorie mathématique connue de représentation de l’information ? Le théorème suivant, absolument fondamental pour justifier l’usage de la théorie des probabilités pour représenter un réel incertain, est initialement dû à Cox (1946) et Jaynes (1954). Il a été étendu plus rigoureusement par Paris (1994), Van Horn (2003), Dupré and Tipler (2009) (entre autres) et finalisé par Terenin and Draper (2015). Il indique que sous les axiomes précédents, une plausibilité se comporte exactement comme une probabilité.

7. les différences entre logique probabiliste et logique non probabiliste (ou *extra-probabiliste*) sont issues de l’accord ou du désaccord avec cette hypothèse. Jaynes (1954) justifie toutefois la validité de cette hypothèse sur des bases pragmatiques.

8. Cette hypothèse est notamment motivée lorsque nous parlons de quantités X possédant une *signification physique et prenant une unique valeur à chaque instant* (étant donnée, possiblement, une précision de mesure finie).

Théorème 12 Cox-Jaynes *Sous les axiomes précédents, il existe une fonction \mathbb{P} , croissante et continue, telle que pour toute proposition A, C et tout ensemble B_X consistant,*

(i) $\mathbb{P}([A|B_X]) = 0$ si et seulement si A est fausse étant donnée l'information sur X

(ii) $\mathbb{P}([A|B_X]) = 1$ si et seulement si A est vraie étant donnée l'information sur X

(iii) $0 \leq \mathbb{P}([A|B_X]) \leq 1$

(iv) $\mathbb{P}([A \wedge C|B_X]) = \mathbb{P}([A|B_X])\mathbb{P}([C|A, B_X])$

(v) $\mathbb{P}(\neg[A|B_X]) = 1 - \mathbb{P}([A|B_X])$

On en conclut que tout système de raisonnement plausible, sous les hypothèses précédentes, est isomorphe à la théorie des probabilités.

Ce théorème est donc particulièrement fondamental en IA. Il est particulièrement résistant à la variation de certains axiomes. Ainsi, Goertzel (2013) a prouvé que si la règle de consistance était affaiblie, alors les plausibilités se comportent approximativement comme des probabilités. On en conclut La théorie des probabilités est pertinente pour représenter les incertitudes sur un sujet exploré par un systèmes cognitif implicite (humain ou artificiel) qui pourrait ne pas être complétement consistant.

6 Modélisation *a priori*

Un grand intérêt de la théorie bayésienne est sa cohérence et sa méthodologie unifiée. Ainsi, donner les lois *a priori* et *a posteriori* ainsi que la fonction de perte suffit pour déterminer, entre autres, un estimateur optimal et des régions α -crédibles. Le choix de la loi *a priori* π est donc crucial et aussi important que celui de la fonction de coût. Avec beaucoup d'observations, le comportement asymptotique peut guider ce choix (*approche bayésienne dite empirique*) mais sinon il est nécessaire de le justifier avec précision. Il existe principalement deux méthodes : l'approche **objectiviste**, qui s'appuie sur des règles formelles dites d'élicitation⁹, et l'approche **subjectiviste**. Elles sont toutes les deux présentées successivement dans ce cours. Par ailleurs, différentes méthodologies visant à mener des examens critiques des choix *a priori* sont décrites.

6.1 Priors objectifs régularisant (priors peu ou “non informatifs”)

Comme on l'a vu précédemment dans le cadre des méthodes usuelles d'apprentissage statistique, l'action d'une mesure *a priori* peut être assimilable à celle d'un terme de pénalisation de la vraisemblance : ce prior cherche à tenir compte des degrés de liberté (la complexité) du modèle de vraisemblance pour contre-balancer le phénomène de surapprentissage et fournir, associé à la vraisemblance, une sorte de diagnostic sur la capacité du modèle à représenter les informations.

Ainsi, la seule information *a priori* qui est utilisée à ce stade est le choix de la vraisemblance elle-même (sa structure algébrique). À cette structure n'est *a priori* attaché aucun choix de paramétrisation θ spécifique.

EXEMPLE 11. Une durée de vie X pourra être représentée par une loi de Weibull, dont la densité peut s'écrire indifféremment selon la paramétrisation suivante

$$f(x|\eta, \beta) = \frac{\beta}{\eta^\beta} x^{\beta-1} \exp(-(x/\eta)^\beta) \mathbb{1}_{x \geq 0}$$

ou

$$f(x|\mu, \beta) = \mu\beta x^{\beta-1} \exp(-\mu x^\beta) \mathbb{1}_{x \geq 0}.$$

L'information possible *a priori* s'exerçant sur X est indifférente à ce choix de paramétrisation.

En effet, θ n'a en soi généralement pas de sens physique, biologique, etc qui imposerait une paramétrisation précise. Il s'agit d'un choix subjectif condensant un *état caché de la nature* permettant de manipuler un modèle de génération de données, qui est lui-même artificiel (simplificateur) vis-à-vis de la réalité.

En ne s'attachant à définir un prior sur θ ou n'importe quelle reparamétrisation $g(\theta)$ (g étant bijectif) qu'à partir de l'information disponible minimale “la vraisemblance d'une donnée potentielle est ainsi”, différentes règles d'élicitation ont été proposées. Comme très généralement ces règles aboutissent à produire des priors $\pi(\theta)$ qui sont σ -finies sur Θ mais qui ne sont pas des mesures de probabilité, soit telles que

$$\int_{\Theta} \pi(\theta) d\theta = \infty,$$

on les nomme usuellement “non informatives”. C'est cependant faux, et il n'existe pas de mesure $\pi(\theta)$ qui n'apporte pas d'information supplémentaire. Plus prudemment, on devrait les qualifier de *faiblement informatives*, ou plus formellement de *priors non intégrables*. Rappelons que de telles mesures n'ont d'intérêt que si leur loi *a posteriori* est propre (cf. Définition 5 et page 26).

Il faut toutefois les différencier d'approches portant souvent le même nom, dans la littérature, qui consiste à choisir $\pi(\theta)$ comme une distribution arbitraire “morale de variabilité très grande” : par exemple, choisir

9. L'élicitation de l'*a priori* est le travail d'encodage probabiliste d'une connaissance incertaine (méconnaissance) voire d'une incertitude complète sur l'état de la nature, au travers d'une distribution *a priori* $\Pi(\theta)$ de densité $\pi(\theta)$. Ce terme français (et non anglais) vient du latin *elicere* qui signifie “tirer de, faire sortir, arracher, obtenir (*ex aliquo verbum elicere*)”. En anglais, le verbe *to elicit* signifie *to get, to draw out*.

une gaussienne dont l'échelle de variance est bien supérieure à celle de la moyenne. On qualifiera plutôt de "lois *a priori* vagues" de telles lois, qui cherchent à minimiser des apports subjectifs (typiquement le choix d'une espérance pour $\pi(\theta)$), et qui peuvent donner une impression de fausse sécurité (minimiser ces apports).

Si elles sont (faiblement) informatives, les lois impropres produites par des systèmes de règles formelles (et donc auditables) doivent être comprises comme des lois de référence ou choisies par défaut, auxquelles on peut avoir recours quand toute information *a priori* autre que la structure même de la vraisemblance est absente.

6.1.1 Priors de Laplace (prior uniforme) et de Jeffreys

Commençons cette section par un exercice d'analyse critique.

Exercice 12 Soit $\theta \in [1, 2]$ le paramètre d'un modèle $X \sim f(\cdot|\theta)$. On suppose ne connaître rien d'autre sur X ou θ . Considérons le choix de la loi *a priori* $\theta \sim \mathcal{U}[1, 2]$. Choisissons alors de reparamétriser le modèle en décrivant une loi *a priori* sur $1/\theta \in [1/2, 1]$. Peut-on faire un choix également uniforme, sachant qu'on n'a ajouté ni retranché aucune information ?

L'exemple précédent met en lumière la nécessité, si l'on ne souhaite pas ajouter d'information autre que la nature de la vraisemblance, de prêter attention à la *propriété d'invariance par reparamétrisation* (cf. Définition 22). Ce principe (expliqué plus bas) n'a pas été compris par Laplace, qui en 1773 considérait que les probabilités ("les chances") s'appliquait d'abord aux variables dénombrables.

Exercice 13 (Laplace, 1793). Une urne contient un nombre n de cartes noires et blanches. Si la première sortie est blanche, quelle est la probabilité que la proportion θ de cartes blanches soit θ_0 ?

Laplace suppose ici que tous les nombres de 2 à $n - 1$ sont équiprobables comme valeurs de θn , donc que θ soit *uniformément distribué* sur $2/n, \dots, (n - 1)/n$. Il applique alors ce qui est connu comme le *principe de raison insuffisante*.

Définition 21 Principe de raison insuffisante. En l'absence d'information, tous les événements élémentaires sont équiprobables, et le même poids doit être donnée à chaque valeur du paramètre, ce qui débouche automatiquement sur une distribution *a priori* uniforme $\pi(\theta) \propto 1$.

Très clairement, ce principe s'applique uniquement à des univers Θ finis. Par ailleurs, il débouche sur une incohérence en termes de partitionnement des événements équiprobables :

- si $\theta = \{\theta_1, \theta_2\}$ alors $\pi(\theta_1) = \pi(\theta_2) = 1/2$;
- si on détaille plus, avec $\theta = \{\theta_1, \omega_1, \omega_2\}$, alors $\pi(\theta_1) = 1/3$.

En d'autres termes, comme l'illustre également l'exemple 12, de façon contre-intuitive, la loi uniforme $\pi(\theta) \propto 1$ qui découle du principe de raison insuffisante mène à des paradoxes et des incohérences, et doit généralement être considérée comme une loi subjective, attachée à un choix particulier de paramétrisation.

On peut donc proposer une première propriété souhaitable pour des *a priori* que l'on veut qualifier de faiblement (ou minimalement) informatifs : ils doivent être *invariants par certaines reparamétrisations*.

Définition 22 Principe d'invariance par reparamétrisation. Si on passe de θ à $\eta = g(\theta)$ par une bijection g , l'information *a priori* reste inexistante et ne devrait pas être modifiée.

Comment faire ? On peut d'abord remarquer que

$$\pi^*(\eta) = |Jac(g^{-1}(\eta))| \pi(g^{-1}(\eta)) = \left| \det \frac{\partial \eta}{\partial \theta} \right| \pi(g^{-1}(\eta))$$

qui (en général) n'est pas constante si $\pi(\theta) = 1$, ce qui confirme le choix discutable de la loi uniforme.

EXEMPLE 12. Soit la reparamétrisation $\eta = -\log(1 - \theta)$. Supposons que $\pi(\theta) = \mathbb{1}_{[0,1]}$ (loi uniforme). Alors η suit une loi exponentielle $\mathcal{E}(1)$.

Au-delà de ces exemples particuliers, on peut considérer des **classes de paramètres** :

1. **Paramètres de position.** On suppose pouvoir écrire $f(x|\theta) = f(x - \theta)$. Alors la famille f est *invariante par translation* : si $x \sim f$, alors $y = x - x_0 \sim f \forall x_0$. Dans ce cas, une exigence d'invariance découlant du principe en Définition 22 est que $\pi(\theta)$ soit invariante par translation elle aussi :

$$\pi(\theta) = \pi(\theta - \theta_0) \quad \forall \theta_0.$$

Cette règle aboutit à une *loi uniforme* sur Θ .

2. **Paramètres d'échelle.** Si on peut écrire $f(x|\theta) = \frac{1}{\theta} f(x/\theta)$ avec $\theta > 0$, alors la famille f est *invariante par changement d'échelle* : $y = x/\theta_0 \sim f \forall \theta_0 > 0$. Dans ce cas, la loi *a priori* invariante par changement d'échelle satisfait $\pi(A) = \pi(A/c)$ pour tout ensemble mesurable $A \in]0, +\infty[$ et $c > 0$

$$\pi(\theta) = \frac{1}{c} \pi\left(\frac{\theta}{c}\right)$$

et implique

$$\pi(\theta) \propto 1/\theta$$

Dans ce deuxième cas, la mesure invariante n'est plus constante.

Une résolution générale du problème de l'invariance par n'importe quelle reparamétrisation bijective a été ainsi proposée par Jeffreys (1946). Il implique de pouvoir, pour le modèle $X|\theta \sim f(x|\theta)$, définir et manipuler la matrice d'information de Fisher $I(\theta)$ (cf. § A.6). Rappelons qu'en notant $\theta \in \Theta \subset \mathbb{R}^d$, l'élément $(i, j) \in \{1, \dots, k\}^2$ de I_θ est (sous des conditions de régularité suffisantes)

$$I_{ij}(\theta) = -\mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(x|\theta) \right].$$

Définition 23 Prior de Jeffreys. Si $I(\theta)$ existe, alors on définit le prior de Jeffreys par

$$\pi(\theta) \propto \sqrt{\det I(\theta)}.$$

Théorème 13 Propriété d'invariance du prior de Jeffreys. Soit $\pi_\theta(\theta)$ le prior de Jeffreys pour la paramétrisation θ , et soit $\eta = g(\theta)$ n'importe quelle reparamétrisation bijective de θ . Alors

$$\pi_\eta(\eta) \propto \sqrt{\det I(\eta)}.$$

Le prior de Jeffreys vérifie donc le principe d'invariance (intrinsèque) proposé par la Définition 22 pour n'importe quelle reparamétrisation.

Exercice 14 Soit $X|\theta$ de loi exponentielle $\mathcal{E}(\theta)$. Quel est le prior de Jeffreys ?

La règle de Jeffreys de construction d'un prior objectif n'aboutit pas forcément à une mesure non intégrable, comme l'illustre l'exemple suivant.

Exercice 15 Prior de Jeffreys pour une loi binomiale. Soit x un nombre de boules tirées dans une urne en contenant n avec probabilité θ . Alors $x \sim \mathcal{B}(n, \theta)$. Calculer la mesure a priori de Jeffreys sur θ . Mener un calcul bayésien sur un échantillon simulé. Peut-on dire que cette mesure est défavorable a priori ?

Exercice 16 Prior de Jeffreys pour une loi de Weibull. Considérons deux paramétrisations classiques de la loi de Weibull \mathcal{W} :

$$\begin{aligned} f(x|\eta, \beta) &= \frac{\beta}{\eta} \left(\frac{x}{\eta}\right)^{\beta-1} \exp\left(-\left\{\frac{x}{\eta}\right\}^\beta\right) \mathbb{1}_{\{x \geq 0\}}, \\ f(x|\mu, \beta) &= \beta \mu x^{\beta-1} \exp(-\mu x^\beta) \mathbb{1}_{\{x \geq 0\}}. \end{aligned}$$

Calculer chaque prior de Jeffreys correspondant, et vérifier la cohérence des résultats grâce à la règle de changement de variable.

On peut apporter une justification complémentaire au fait de définir une mesure *a priori* comme une fonction de l'information de Fisher : $I(\theta)$ est largement accepté comme un indicateur de la quantité d'information apportée par le modèle (ou l'observation) sur θ (Fisher, 1956); de plus, $I(\theta)$ mesure la capacité du modèle à discriminer entre θ et $\theta + / - d\theta$ via la pente moyenne de $\log f(x|\theta)$. Enfin, favoriser les valeurs de θ pour lesquelles $I(\theta)$ est grande équivaut à minimiser l'influence de la loi *a priori*.

Ainsi, l'*a priori* de Jeffreys est l'une des meilleures techniques automatiques pour obtenir des lois non-informatives. Il est le plus souvent impropre, sauf pour des modèles pour lesquels Θ est borné ou/et discret (cf. exercice 15). Précisons qu'il est généralement utilisé en dimension 1, où il permet d'obtenir des estimateurs bayésiens similaires au maximum de vraisemblance. En effet, en multidimensionnel, le prior de Jeffreys peut mener à des incohérences ou des paradoxes.

Exercice 17 Problème de Neyman-Scott. On considère le problème suivant

$$x_i \sim \mathcal{N}(\mu_i, \sigma^2) \text{ pour } i = 1, \dots, n$$

où les x_i sont indépendants. Soit $\theta = (\mu_1, \dots, \mu_n, \sigma)$. Calculez le prior de Jeffreys $\pi^J(\theta)$ puis l'espérance a posteriori de σ^2 . Est-ce un estimateur consistant ?

6.1.2 Prior de référence de Berger-Bernardo

Une variante de l'approche précédente, mieux adaptée au cadre multidimensionnel, a été proposée sous le nom de *reference prior* par Bernardo (1979) et Berger (1992). L'idée est de reposer sur un critère d'objectivité. On regarde ici la loi *a priori* qui apporte le moins d'information par rapport à l'information que peuvent apporter les données. Cette notion de variation de l'information a été synthétisée par l'entropie de Shannon (cf. § 6.2.1) ou, sous une forme équivalente, par la divergence de Kullback-Leibler dont on définit ci-dessous l'usage dans ce cadre.

Définition 24 Divergence prior-posterior de Kullback-Leibler. La divergence de Kullback-Leibler entre a posteriori et a priori

$$KL(\pi, \mathbf{x}_n) = \int_{\Theta} \pi(\theta|\mathbf{x}_n) \log \frac{\pi(\theta|\mathbf{x}_n)}{\pi(\theta)} d\theta$$

mesure l'information apportée par les données observées \mathbf{x}_n sur le modèle, indépendamment de la paramétrisation θ .

L'approche repose donc sur la maximisation de $KL(\pi, \mathbf{x}_n)$ en π pour des données \mathbf{x}_n pouvant être typiquement observées : elles sont générées par la vraisemblance (marginale)

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

Soit donc le critère à maximiser en π :

$$J_n(\pi) = \int KL(\pi, \mathbf{x}_n) m_\pi(\mathbf{x}_n) d\mathbf{x}_n$$

qui possède la propriété suivante lorsque l'information de Fisher $I(\theta)$ est bien définie, avec $d = \dim(\Theta)$:

$$J_n(\pi) \xrightarrow{n \rightarrow \infty} \frac{d}{2} \log \frac{n}{2\pi} - \int_{\Theta} \pi(\theta) \log \frac{\pi(\theta)}{\sqrt{\det I(\theta)}} d\theta + \mathcal{O}_p(1).$$

Dans le cas d'un modèle régulier, on constate que $J_n(\pi)$ est maximal quand $n \rightarrow \infty$ lorsque $\pi(\theta) \propto \sqrt{\det I(\theta)}$. Pour éviter de dépendre d'une taille n arbitraire, Berger et Bernardo ont donc proposé, dans la forme définitive de cette approche, de faire tendre n vers l'infini.

Définition 25 Prior de référence (Berger-Bernardo). *Le prior est défini formellement par*

$$\pi^* = \arg \min_{\pi \in \mathcal{D}} \lim_{n \rightarrow \infty} J_n(\pi) \quad (13)$$

où \mathcal{D} est l'ensemble des mesures positives sur Θ .

En dimension 1, on retombe par construction sur le prior de Jeffreys. En dimension supérieure à 1, cette approche permet de résoudre des problèmes d'inconsistance *a posteriori* (tels que celui constaté dans l'exercice 17) en suivant la stratégie de construction suivante, proposée par Bernardo :

- On sépare θ en (θ_1, θ_2) où θ_1 est un *paramètre d'intérêt* et θ_2 un *paramètre de nuisance*.
- On raisonne alors séquentiellement en écrivant $\pi(\theta) = \pi(\theta_2|\theta_1)\pi(\theta_1)$ et en choisissant pour $\pi(\theta_1)$ le prior de Jeffreys (le caractère séquentiel pouvant être généralisé si on peut hiérarchiser l'ensemble des paramètres unidimensionnels composant θ , par intérêt croissant).
- On calcule $\pi(\theta_2|\theta_1)$ la plus objective possible par la règle (13) (puis, si on a θ_3 , on reproduit cela sur $\pi(\theta_3|\theta_1, \theta_2)$, etc.).

L'invariance par reparamétrisation est maintenue à l'intérieur des groupes (intérêt et nuisance), sous certaines conditions de bijectivité. Cependant, il est clair ce raisonnement n'est pas purement objectif parce que donner plus d'importance à un paramètre qu'à un autre relève d'un choix qui peut être subjectif. Cette méthodologie d'éllicitation est en général délicate à mettre en oeuvre.

EXEMPLE 13. *Dans un problème de Markov à état cachés où l'on dispose d'observations bruitées*

$$X_t^* = X_t + \epsilon(\theta_2)$$

d'un phénomène réel non connu (ex : une dynamique de population)

$$X_{t+1} = g(X_t, \theta_1),$$

le paramètre θ_2 correspond souvent à un biais et une variance d'observation. On cherche surtout à estimer le paramètre θ_1 qui gouverne la dynamique et dont la valeur permet de simuler le phénomène pour une analyse prévisionnelle. Dans ce cas, la hiérarchie entre θ_1 et θ_2 est dictée par l'usage que l'on souhaite faire du modèle.

6.1.3 Priors coïncidants ou concordants

Le but ici est de trouver une loi *a priori* concernant le paramètre θ qui se rapproche le plus possible de la méthode de choix fréquentiste. Rappelons quelques exemples d'une région de confiance ou α -crédible. Elle peut être par exemple un intervalle unilatéral $\{\theta \leq \theta_d^{(x)}\}$ ou bien bilatéral $\{\theta_{\alpha,1} \leq \theta \leq \theta_{\alpha,1}\}$. Il peut s'agir aussi de région HPD.

On cherche alors π tel que $\forall \theta \in \Theta, \Pi(\theta \in C|X) = 1 - \alpha$ (égalité entre la région de confiance fréquentiste et la région de crédibilité. C'est en général impossible. On va alors chercher r_m le plus petit possible tel que

$$\forall \theta \in \Theta, \forall \alpha \in]0, 1[, \Pi(\theta \in C|X) = 1 - \alpha + \mathcal{O}(r_m).$$

Une telle loi *a priori* est alors dite *coincidente* ou *concordante* à l'ordre r_m . En anglais, on parle de *coverage matching prior*.

De tels priors, possédant de fortes propriétés de recouvrement fréquentiste *a posteriori*, peuvent être exprimés comme des solutions d'équation différentielle. On peut montrer qu'à l'ordre 1, on retombe sur le prior de Jeffreys. Notons que ces propriétés de recouvrement servent généralement à discriminer entre plusieurs *a priori* faiblement informatifs pour un problème donné.

6.2 Priors objectifs informatifs

6.2.1 Maximum d'entropie

Si l'on possède des informations partielles du type *contraintes linéaires en π* :

$$\int_{\Theta} g_i(\theta) \pi(\theta) d\theta = c_i, \quad i = 1, \dots, M,$$

on cherche la loi la moins informative sous ces contraintes, seules informations dont on dispose (outre la vraisemblance). Pour comparer le caractère informatif de mesures informatives, il est nécessaire d'avoir recours à un critère d'information. L'entropie de Shannon, issue de la théorie du signal et à l'origine de la première théorie de l'information¹⁰, permet de définir ce niveau d'informativité. Nous l'introduisons ci-dessous d'abord dans un cas discret, puis nous passons au continu.

Le concept d'entropie. À l'origine du concept d'entropie, on tente de résoudre un problème de recherche (tri) d'une information discrétisée. Soit N sites numérotés de 1 à N où une information recherchée peut être présente. On suppose ne pouvoir poser que des questions à réponse binaire (oui ou non). Une première stratégie consiste à visiter tous les sites et poser la question de présence ou d'absence de l'information : on doit donc poser N questions. Une meilleure stratégie est *dichotomique* : si $\exists Q_2 \in \mathbb{N}$ tel que $N = 2^{Q_2}$, on range les sites en 2 parties et on pose la question d'appartenance au premier ou au second groupe. En itérant ce procédé, on peut trier en $Q_2 = \log_2 N$ questions.

Ici, Q_2 s'interprète comme le nombre de bits (*binary digits*) nécessaire à l'écriture de N en base 2, c'est-à-dire la longueur du mot à utiliser pour coder N dans un alphabet de 2 caractères. Si on imagine maintenant un autre alphabet de c caractères, le nombre minimal de questions sera $Q_c = \log_c N = \frac{\log_2 N}{\log_2 c} = \frac{\log N}{\log c}$. Si enfin on suppose inconnue la taille de l'alphabet (donc la nature des questions posées), le nombre de questions à poser pour identifier un site parmi N est, à une constante multiplicative près,

$$Q = \log N \quad (\text{logarithme naturel}).$$

Généralisons à présent : on suppose qu'il existe une partition de k aires géographiques et que chaque aire contienne $N_i = N \times p_i$ sites, avec $i = 1, \dots, k$

- il suffit de poser $Q'_i = \log N_i = \log p_i + \log N$ questions pour trier l'aire i ;
- en moyenne sur l'ensemble des aires, on trie avec $Q' = \sum_{i=1}^k p_i Q'_i$ questions.

Le fait de savoir en probabilité dans quelle aire est l'information réduit le nombre de questions à poser en moyenne de la quantité

$$\Delta Q = Q - Q' = - \sum_{i=1}^k p_i \log p_i, \quad (14)$$

10. Une seconde théorie de l'information, qui s'intéresse à la sémantique du message plutôt que sa longueur, est due à Kolmogorov.

quantité positive et maximale quand $p_i = 1/k$. Moins la distribution de probabilité $\Pi = (p_1, \dots, p_k)$ est *informative*, plus cette quantité est grande. Ainsi, la formulation (14) offre une mesure de la quantité d'incertitude dans une distribution de probabilité d'événements discrets, notion qui sera ensuite étendue à des cas continus. Cette généralisation est produite par les considérations suivantes :

- le cas discret correspond à une partition fine de Θ en k intervalles dont l'étendue individuelle tend vers 0 ;
- le résultat dépend de la mesure de partitionnement sur Θ ;
- l'entropie doit être invariante par tout changement de variable $\theta \mapsto \nu(\theta)$.

Définition 26 Entropie de Shannon-Kullback L'entropie d'une variable aléatoire (continue) décrite par sa distribution de probabilité $\pi(\theta)$ est

$$\mathcal{H}(\pi) = - \int_{\Theta} \pi(\theta) \log \frac{\pi(\theta)}{\pi_0(\theta)} d\theta \quad (\text{entropie de Kullback})$$

où $\pi_0(\theta)$ est une mesure (positive) de référence sur Θ . S'il s'agit d'une mesure de probabilité, elle représente l'ignorance de la valeur θ sur Θ .

Appliquée au cas discret, l'entropie est toujours positive. Ce n'est plus forcément le cas dans le cas continu. Lorsque $\pi_0(\theta)$ est une mesure intégrable (mesure de probabilité), on voit que $\mathcal{H}(\pi) = -KL(\pi, \pi_0)$ et elle est donc finie et négative. Dans tous les cas, elle est toujours maximale en $\pi(\theta) = \pi_0(\theta)$. La preuve repose sur la démonstration des propriétés de la divergence de Kullback-Leibler (KL ; cf. § 3.4).

Le choix de $\pi_0(\theta)$ est assez déterminant. S'il est usuel de choisir une mesure uniforme, préférer des priors régularisants (Jeffreys, etc.) permet de limiter la dépendance aux choix de mesure dominante sur Θ .

Théorème 14 Prior par maximum d'entropie. Le problème

$$\pi^*(\theta) = \arg \max_{\pi \in \mathcal{P}} - \int_{\Theta} \pi(\theta) \log \frac{\pi(\theta)}{\pi_0(\theta)} d\theta$$

dans l'ensemble \mathcal{P} des mesures positives, sous M contraintes linéaires

$$\int_{\Theta} g_i(\theta) \pi(\theta) d\theta = c_i, \quad i = 1, \dots, M,$$

a pour solution unique presque partout

$$\pi(\theta) \propto \pi_0(\theta) \exp \left(\sum_{i=1}^M \lambda_i g_i(\theta) \right)$$

où les λ_i sont des réels.

Exercice 18 Exemple industriel. On cherche la distribution de la profondeur X d'un défaut de fabrication dans une enceinte d'acier difficile d'accès. À partir d'anciennes données mesurées sur des aciers différents, on suppose connaître un modèle exponentiel pour $X|\theta \sim \mathcal{E}(\theta)$ de densité $f(x|\theta) = \theta^{-1} \exp(-x/\theta) \mathbb{1}_{x \geq 0}$. Les mesures ultrasonores pour estimer la distribution de X sont cependant coûteuses et ont besoin d'être calibrées a priori. D'où les questions posées successivement à un expert en métallurgie :

1. Pouvez-vous préciser la profondeur moyenne θ_e d'un défaut de fabrication dans cette coulée ?
2. Pouvez-vous préciser un écart-type σ_e pour cette profondeur en général ?

(Remarque : la question à l'expert est ici quelque peu idéalisée. Il faut en pratique plutôt passer par des questions sur X et les connecter à $\pi(\theta)$ via la loi prédictive a priori.

Exercice 19 On considère les contraintes suivantes sur une mesure a priori sur $\theta \geq 0$: pour $\beta > 0$,

$$\mathbb{E}[\theta^\beta] = 1 \quad (15)$$

$$\mathbb{E}[\log \theta] = -\frac{\gamma}{\beta}. \quad (16)$$

où γ est la constante d'Euler ($\simeq 0,577$). Calculer la mesure de maximum d'entropie $\pi(\theta)$ relativement à une mesure $\pi^J(\theta) \propto \theta^{-1}$. Quelles sont les conditions pour que cette mesure soit une vraie mesure de probabilité ?

Indication : si $X \sim \mathcal{G}(a, b)$, alors $\mathbb{E}[\log X] = \psi(a) - \log b$ où ψ est la fonction digamma, dérivée logarithmique de la fonction gamma :

$$\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)},$$

et $\psi(1) = -\gamma$.

6.2.2 Priors conjugués et famille exponentielle

Le principe de maximisation d'entropie peut aussi être appliqué à X conditionnellement à θ , et il permet de déboucher sur la famille paramétrée suivante.

Définition 27 Famille exponentielle. Soient $(C, h) : \Theta \times \Omega \mapsto \mathbb{R}_+^2$, et $(R, T) : \Theta \times \Omega \mapsto \mathbb{R}^k \times \mathbb{R}^k$. La famille des distributions de densité

$$f(x|\theta) = C(\theta)h(x) \exp \{R(\theta) \cdot T(x)\}$$

est dite famille exponentielle de dimension k . Si R est linéaire¹¹, et lorsque $\Theta \subset \mathbb{R}^k$ et $\Omega \subset \mathbb{R}^k$, on peut écrire plus simplement

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

où ∇ désigne l'opérateur gradient, avec

$$\begin{aligned} \mathbb{E}_\theta[X] &= \nabla \psi(\theta), \\ \text{cov}(X_i, X_j) &= \frac{\partial^2 \psi}{\partial \theta_i \partial \theta_j}(\theta) \end{aligned}$$

On parle alors de forme naturelle (ou canonique) de la famille exponentielle (NEF).

Une propriété importante de la famille exponentielle est la suivante : $T(x)$ est une *statistique exhaustive* (vectorielle) de x . Rappelons quelques résultats fondamentaux sur la notion de statistique exhaustive (ou *suffisante*).

Définition 28 Statistique exhaustive. Si $x \sim f(x|\theta)$, une statistique T de x est exhaustive si la distribution de x conditionnellement à $T(X)$ ne dépend pas de θ .

On peut également caractériser une statistique exhaustive dans le cadre bayésien en utilisant la définition de Berger-Wolpert :

11. Si ce n'est pas le cas, on procède à une reparamétrisation.

Définition 29 Statistique exhaustive (Berger-Wolpert). Si $x \sim f(x|\theta)$ et si $z = T(x)$, alors z est une statistique exhaustive si et seulement si pour tout a priori π sur θ , $\pi(\theta|x) = \pi(\theta|z)$.

Définition 30 Critère de factorisation de Neyman-Fisher. La statistique $T = T(x_1, \dots, x_n)$ est exhaustive pour θ si et seulement si il existe deux fonctions (h, g) telle que la vraisemblance de l'échantillon $(x_1, \dots, x_n) \sim f(x|\theta)$ puisse s'écrire

$$f(x_1, \dots, x_n|\theta) = h(x_1, \dots, x_n)g(T, \theta).$$

L'exhaustivité de T permet en fait de caractériser complètement la famille exponentielle, comme le formalise le lemme suivant.

Lemme 1 Pitman-Koopman. Si une famille $f(\cdot|\theta)$ à support constant admet une statistique exhaustive de taille fixe à partir d'une certaine taille d'échantillon, alors $f(\cdot|\theta)$ appartient à la famille exponentielle.

EXEMPLE 14. La loi de Dirichlet appartient à la NEF. Sa densité est

$$f(x|\theta) = \frac{\Gamma(\sum_{i=1}^k \theta_i)}{\prod_{i=1}^k \Gamma(\theta_i)} \prod_{i=1}^k x_i^{\theta_i-1} \mathbb{1}_{\{S_k(x)\}},$$

définie sur le simplexe $S_k(x) = \left\{ x = (x_1, \dots, x_k); \sum_{i=1}^k x_i = 1, x_i > 0 \right\}$.

EXEMPLE 15. Soit $\mathbf{x}_n = (x_1, \dots, x_n) \sim \mathcal{N}_p(\mu, \sigma^2 I_p)$. Alors la distribution jointe satisfait

$$f(\mathbf{x}_n|\theta) = C(\theta)h(\mathbf{x}_n) \exp \left(n\bar{x} \cdot (\mu/\sigma^2) + \sum_{i=1}^n \|x_i - \bar{x}\|^2 (-1/2\sigma^2) \right)$$

avec $\theta = (\mu, \sigma)$, et la statistique $(\bar{x}, \sum_{i=1}^n \|x_i - \bar{x}\|^2)$ est exhaustive pour tout $n \geq 2$.

EXEMPLE 16. La loi de Weibull, de densité

$$f(x|\mu, \beta) = \mu\beta x^{\beta-1} \exp(-\mu x^\beta)$$

avec $(\mu, \beta) > 0$, n'appartient pas à la NEF.

Supposons donc que $X|\theta$ suive une loi construite par maximum d'entropie, de densité de forme :

$$f(x|\theta) = \exp \left(\sum_{j=1}^L T_j(x) d_j(\theta) \right).$$

Si, de plus, la loi a priori $\pi(\theta)$ est également construite par maximum d'entropie :

$$\pi(\theta) \propto \nu(\theta) \exp \left(\sum_{i=1}^M \lambda_i g_i(\theta) \right),$$

alors, sachant l'échantillon $\mathbf{x}_n = (x_1, \dots, x_n)$, la loi a posteriori est de la même forme structurelle que $\pi(\theta)$:

$$\pi(\theta|\mathbf{x}_n) \propto \nu(\theta) \exp \left(\sum_{i=1}^M \lambda_i g_i(\theta) + \sum_{j=1}^L \left[\sum_{k=1}^n T_j(x_k) \right] d_j(\theta) \right).$$

L' a priori est alors dit **conjugué**.

Définition 31 Prior conjugué. Une loi a priori $\pi(\theta)$ appartenant à une famille (hyper)paramétrée $\Pi_\omega = \{\pi(\theta|\omega), \omega \in \mathbb{R}^q\}$ est dite conjuguée si, pour tout échantillon x_1, \dots, x_n , la loi (éventuellement impropre) $\pi(\theta|x_1, \dots, x_n)$ appartient également à Π_ω .

Conjuguée naturelle. Si l'on utilise l'écriture canonique de la famille exponentielle, on pose

$$f(x|\theta) = h(x) \exp(\theta \cdot x - \psi(\theta))$$

alors la mesure *a priori* générée automatiquement par

$$\pi(\theta|a, b) = K(a, b) \exp(\theta \cdot a - b\psi(\theta))$$

lui est *conjuguée* (naturelle), et la mesure *a posteriori* sachant une donnée x est

$$\pi(\theta|a + x, b + 1).$$

$K(a, b)$ est la constante de normalisation

$$K(a, b) = \left[\int_{\Theta} \exp(\theta \cdot a - b\psi(\theta)) \right]^{-1}$$

qui est finie si $b > 0$ et $a/b \in \mathring{N}$. Quelques cas importants de priors conjugués sont présentés sur la figure 2.

$f(x \theta)$	$\pi(\theta)$	$\pi(\theta x)$
Normale $\mathcal{N}(\theta, \sigma^2)$	Normale $\mathcal{N}(\mu, \tau^2)$	$\mathcal{N}(\varrho(\sigma^2\mu + \tau^2x), \varrho\sigma^2\tau^2)$ $\varrho^{-1} = \sigma^2 + \tau^2$
Poisson $\mathcal{P}(\theta)$	Gamma $\mathcal{G}(\alpha, \beta)$	$\mathcal{G}(\alpha + x, \beta + 1)$
Gamma $\mathcal{G}(\nu, \theta)$	Gamma $\mathcal{G}(\alpha, \beta)$	$\mathcal{G}(\alpha + \nu, \beta + x)$
Binomiale $\mathcal{B}(n, \theta)$	Bêta $\mathcal{B}e(\alpha, \beta)$	$\mathcal{B}e(\alpha + x, \beta + n - x)$
Binomiale Négative $\mathcal{N}eg(m, \theta)$	Bêta $\mathcal{B}e(\alpha, \beta)$	$\mathcal{B}e(\alpha + m, \beta + x)$
Multinomiale $\mathcal{M}_k(\theta_1, \dots, \theta_k)$	Dirichlet $\mathcal{D}(\alpha_1, \dots, \alpha_k)$	$\mathcal{D}(\alpha_1 + x_1, \dots, \alpha_k + x_k)$
Normale $\mathcal{N}(\mu, 1/\theta)$	Gamma $\mathcal{G}a(\alpha, \beta)$	$\mathcal{G}(\alpha + 0.5, \beta + (\mu - x)^2/2)$

FIGURE 2 – Priors conjugués pour des choix de vraisemblance NEF (tiré de C.P. Robert, 2006).

Intérêt de la conjugaison.

1. **C'est pratique.** Le premier intérêt de la conjugaison est généralement le côté pratique : on sait directement résoudre le problème du calcul *a posteriori*, et généralement on connaît des algorithmes pour manipuler la famille Π_ω .
2. **C'est logique.** On peut considérer que l'information $x \sim f(x|\theta)$ transformant $\pi(\theta)$ en $\pi(\theta|x)$ est limitée, donc elle ne devrait pas entraîner une modification de *toute* la *structure* de $\pi(\theta)$, mais simplement de ses *hyperparamètres* :

$$\pi(\theta) = \pi(\theta|\delta) \Rightarrow \pi(\theta|x) = \pi(\theta|\delta + s(x)).$$

Cette modification devrait être de dimension finie, et un changement plus radical de $\pi(\theta)$ est peu acceptable.

3. **Elle offre une interprétation claire.** En effet, si un prior conjugué existe et qu'il est exponentiel, tel

$$\pi(\theta|x_0, m) \propto \exp\{\theta \cdot x_0 - m\psi(\theta)\}$$

alors l'espérance *a priori prédictive* est

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|\theta]] = \mathbb{E}[\nabla\psi(\theta)] = \frac{x_0}{m}$$

et l'espérance *a posteriori prédictive*, sachant un échantillon $\mathbf{x}_n = (x_1, \dots, x_n)$, est

$$\mathbb{E}[X|\mathbf{x}_n] = \frac{x_0 + n\bar{x}}{m + n}. \quad (17)$$

Autrement dit, m a le sens d'une *taille d'échantillon virtuelle*, offrant une indication de la "force" informative de l'*a priori* (d'un expert, etc.). Ce résultat bénéficie d'une réciprocité forte offerte par le théorème de Diaconis-Ylvisaker.

Théorème 15 (Diaconis-Ylvisaker). Si la mesure de référence est continue par rapport à la mesure de Lebesgue, alors

$$\mathbb{E}[X|\mathbf{x}_n] = \frac{x_0 + n\bar{x}}{m + n} \Rightarrow \pi(\theta|x_0, m) \propto \exp\{\theta \cdot x_0 - m\psi(\theta)\}.$$

Exercice 20 Cas multinomial. On considère $X|\theta$ suivant une loi multinomiale avec $X = (X_1, \dots, X_d)$ et $\theta = (\theta_1, \dots, \theta_d)$ tel que $0 \leq \theta_i \leq 1$ et $\sum_{i=1}^d \theta_i = 1$:

$$P(X_1 = k_1, \dots, X_d = k_d|\theta) = \frac{n!}{k_1! \dots k_d!} \theta_1^{k_1} \dots \theta_d^{k_d}.$$

Montrer que la loi de Dirichlet est conjuguée pour cette vraisemblance.

Exercice 21 Loi inverse Wishart. Soient des observations $x_1, \dots, x_n \sim \mathcal{N}_p(\mu, \Sigma)$, de loi jointe

$$f(x_1, \dots, x_n|\theta = (\mu, \Sigma)) \propto (\det \Sigma)^{-n/2} \exp\left(-\frac{1}{2} [n(\bar{x}_n - \mu)^T \Sigma^{-1}(\bar{x}_n - \mu) + \text{tr}(\Sigma^{-1} S_n)]\right)$$

avec $S_n = \sum_{i=1}^n (x_i - \bar{x}_n)(x_i - \bar{x}_n)^T$. On suppose prendre *a priori*

$$\begin{aligned} \mu|\Sigma &\sim \mathcal{N}_p\left(\mu_0, \frac{1}{n_0}\Sigma\right) \\ \Sigma &\sim \mathcal{IW}(\alpha, V) \end{aligned}$$

la loi de Wishart inverse \mathcal{IW} étant définie (sur l'espace des matrices symétriques non indépendantes de rang d) par la densité

$$f(x) = \frac{|V|^{\alpha/2}}{2^{\alpha d/2} \Gamma_d(\alpha/2)} |x|^{-\frac{\alpha+d+1}{2}} \exp\left\{-\frac{1}{2} \text{tr}(Vx^{-1})\right\}$$

où Γ_d est la fonction gamma multivariée. Le prior est-il conjugué ? En dimension 1, à quelle loi se réduit-il ?

Exercice 22 Soit $X \sim \mathcal{N}(\theta, \theta)$ avec $\theta > 0$.

1. Déterminer la loi *a priori* de Jeffreys $\pi^J(\theta)$
2. Établir si la loi de X appartient à la famille exponentielle et construire les lois *a priori* conjuguées sur θ .
3. Utiliser la propriété de linéarité des espérances des familles exponentielles pour relier les hyperparamètres des lois conjuguées à l'espérance de θ .

Modèles de mélange. La conjugaison peut s'étendre à des modèles de mélange. En effet, soit $\mathcal{F} = \{\pi(\theta|a, b) = K(a, b) \exp(\theta \cdot a - b\psi(\theta))\}$ la famille conjuguée naturelle de la famille exponentielle

$$f(x|\theta) = C(\theta)h(x) \exp(\theta \cdot x).$$

Alors l'ensemble des mélanges de N lois conjuguées

$$\mathcal{F}_N = \left\{ \sum_{i=1}^N \omega_i \pi(\theta|a_i, b_i); \sum_{i=1}^N \omega_i = 1, \omega_i > 0 \right\}$$

est aussi une famille conjuguée. *A posteriori*, on a

$$\pi(\theta|x) = \sum_{i=1}^N \omega'_i(x) \pi(\theta|a_i + 1, b_i + x)$$

avec

$$\omega'_i(x) = \frac{\omega_i K(a_i, b_i) / K(a_i + 1, b_i + x)}{\sum_{j=1}^N \omega_j K(a_j, b_j) / K(a_j + 1, b_j + x)}.$$

Les mélanges d'*a priori* conjugués peuvent alors être utilisés comme *base pour approcher une loi a priori quelconque*, au sens où la distance de Prohorov entre une loi et sa représentation par un mélange dans \mathcal{F}_N peut être rendue arbitrairement petite. Cela fournit un argument fort en faveur de l'utilisation des lois conjuguées.

Définition 32 Distance de Prohorov. La distance de Prohorov entre deux mesures π et $\tilde{\pi}$ est définie par

$$D^p(\pi, \tilde{\pi}) = \inf_A \{ \epsilon; \pi(A) \leq \tilde{\pi}(A^\epsilon) + \epsilon \}$$

où l'infimum est pris sur les ensembles boréliens de Θ et où A^ϵ indique l'ensemble des points distants de A d'au plus ϵ .

Théorème 16 Pour toute loi *a priori* π sur Θ , $\forall \epsilon > 0$, on peut trouver N et $\tilde{\pi} \in \mathcal{F}_N$ tel que

$$D^p(\pi, \tilde{\pi}) < \epsilon.$$

Au-delà de la famille exponentielle. Le mécanisme de conjugaison n'est pas réservé aux lois de la famille exponentielle. Par exemple, la loi de Pareto avec $\alpha > 0$ connu, et $\theta > 0$, de densité

$$f(x|\theta) = \alpha \frac{\theta^\alpha}{x^{\alpha+1}} \mathbb{1}_{] \theta, \infty[}(x),$$

admet un *a priori* conjugué, qui est Pareto sur $1/\theta$.

Exercice 23 Considérons les deux lois uniformes de densité

$$\begin{aligned} f(x|\theta) &= \frac{\mathbb{1}_{[-\theta, \theta]}(x)}{2\theta} \\ f(x|\theta) &= \frac{\mathbb{1}_{[C\theta, \theta]}(x)}{\theta}. \end{aligned}$$

Admettent-elles des priors conjugués ?

6.3 Priors hiérarchiques

Pour des raisons liées à la modélisation des observations ou à la décomposition de l'information *a priori*, le modèle bayésien $(f(x|\theta), \pi(\theta))$ peut être défini comme *hiérarchique* : $\pi(\theta)$ est décomposé en plusieurs lois conditionnelles

$$\begin{aligned}\pi(\theta|\theta_1, \dots, \theta_k) &= \pi_1(\theta|\theta_1) \cdots \pi_2(\theta_1|\theta_2) \cdots \pi_k(\theta_{k-1}|\theta_k) \cdot \pi_{k+1}(\theta_k) \\ \text{et } \pi(\theta) &= \int_{\Theta_1 \times \dots \times \Theta_k} \pi(\theta|\theta_1, \dots, \theta_k) d\theta_1 \dots d\theta_k.\end{aligned}$$

EXEMPLE 17. *Modèle linéaire à effets aléatoires :*

$$\begin{aligned}y|\theta &\sim \mathcal{N}_p(\theta, \Sigma_1), \\ \theta|\beta &\sim \mathcal{N}_p(X\beta, \Sigma_2)\end{aligned}$$

souvent utilisé en génétique animale pour différencier l'influence d'éléments fixes (ex : lignée, race, année) de celle de facteurs aléatoires (ex : nb de femelles dans une lignée).

Conditionnement. Le conditionnement peut apparaître par

- des dépendances statistiques naturelles ;
- l'appel à des variables latentes décrivant un mécanisme caché ;
- des grandeurs stochastiques jouant un rôle de forçage.

En général, on ne va guère plus loin que deux ou trois niveaux de hiérarchie. En incluant l'information *a priori* aux niveaux les plus élevés, l'approche bayésienne hiérarchique permet en général de gagner en robustesse.

Exemples de conditionnements naturels. Ces conditionnements peuvent apparaître à partir de résultats statistiques théoriques.

EXEMPLE 18. *La courbe de von Bertalanffy*

$$L(t|\theta) = L_\infty(1 - \exp(-g(t, \delta)))$$

est fréquemment utilisée pour produire une **clé âge-longueur**, en modélisant l'accroissement en longueur d'un organisme vivant (ex : arbre, poisson...). On note $\theta = (L_\infty, \delta)$ le vecteur des paramètres inconnus. Supposons disposer de données de capture-recapture, c'est-à-dire des couples d'observation $\{l^*(t_i), l^*(t_i + \Delta_i)\}$ tel que

$$\begin{aligned}l^*(t_i) &= L(t_i|\theta) \exp(\epsilon_1), \\ l^*(t_i + \Delta_i) &= L(t_i + \Delta_i|\theta) \exp(\epsilon_2)\end{aligned}$$

où (ϵ_1, ϵ_2) sont des bruits de mesure (générant donc une vraisemblance). Les estimations par maximum de vraisemblance de L_∞ étant très sensibles à la taille des données, on cherche donc à produire un prior sur L_∞ . On peut tirer parti du théorème de Pickands en comprenant que L_∞ peut être assimilé à une longueur maximale qu'un être vivant peut atteindre, en moyenne sur toutes les observations possibles. Posons alors $L_\infty^* = L_\infty \exp(\epsilon)$ la longueur maximale observée. Soit \bar{L} la longueur moyenne. On obtient alors une justification pour :

1. établir une forme *a priori* pour $\pi(L_\infty)$ (la forme de ϵ étant fixée) ;
2. conditionner ce prior par rapport à $\bar{L} \Leftrightarrow$ utiliser une approche bayésienne hiérarchique.

Théorème 17 Pickands (1975). Quand \bar{L} grandit, la distribution de $L_\infty^* | \bar{L} = l$ est une loi de Pareto généralisée :

$$P(L_\infty^* < x | L_\infty^* > \bar{L}, \sigma, \mu) = 1 - \left(1 + \mu \left(\frac{x - \bar{L}}{\sigma}\right)\right)^{-1/\mu}$$

EXEMPLE 19. Soit X_t un nombre d'individus dans une population. Soit $\theta = \theta_{t,t+1}$ la probabilité de survie entre t et $t + 1$. La vraisemblance de ces données markoviennes peut être définie par le noyau de transition

$$X_{t+1} | X_t, \theta_{t,t+1} \sim \mathcal{B}(X_t, \theta_{t,t+1}) \text{ (loi binomiale)}$$

On peut alors écrire

$$\theta_{t,t+1} = \prod_{i=0}^{M+1} \theta_{t+i/M, t+(i+1)/M}$$

donc, de par le Théorème de la Limite Centrale, quand $1 \ll M$,

$$\log(\theta_{t,t+1}) \sim \mathcal{N}(\mu_t, \sigma_t^2)$$

avec $\mu_t < -\sigma_t^2/2$ tel que $\mathbb{E}[\theta_{t,t+1}] \in [0, 1]$ qui est une contrainte de forme sur le niveau hiérarchique $\pi(\mu_t, \sigma_t)$.

Graphes acycliques orientés. On peut représenter les causalités et dépendances probabilistes sous forme de graphes acycliques orientés (DAG = *direct acyclic graphs*, cf. figures 3-4). Ce type de représentation peut s'avérer utile pour l'encodage sous certains logiciels (famille BUGS) et pour s'assurer de la cohérence interne d'un modèle bayésien hiérarchique. Dans ce type de représentation, les informations connues (typiquement : les données) sont représentées sous la forme de rectangles) tandis que les variables aléatoires sont placées dans des ellipses (ex : figure 5).

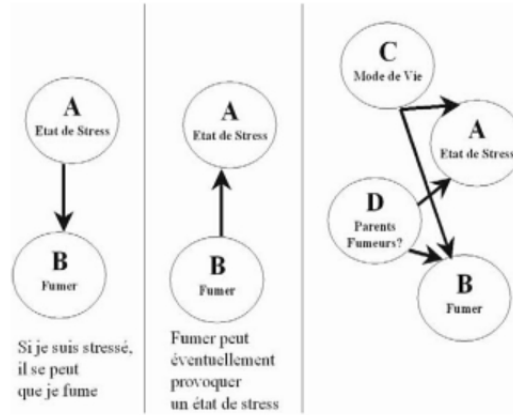


FIGURE 3 – DAG tiré de Parent et Bernier (2007).

6.3.1 Un exemple utile : les *Latent Gaussian Models* (LGM)

On note Y_i la variable réponse, reliée à μ_i par

$$Y_i = \mathcal{N}(\mu_i, \sigma^2)$$

où μ_i est tel que

$$g(\mu_i) = \eta_i = \alpha + \sum_j f^j(u_{i,j}) + \sum_k \beta_k Z_{k,i} + \varepsilon_i.$$

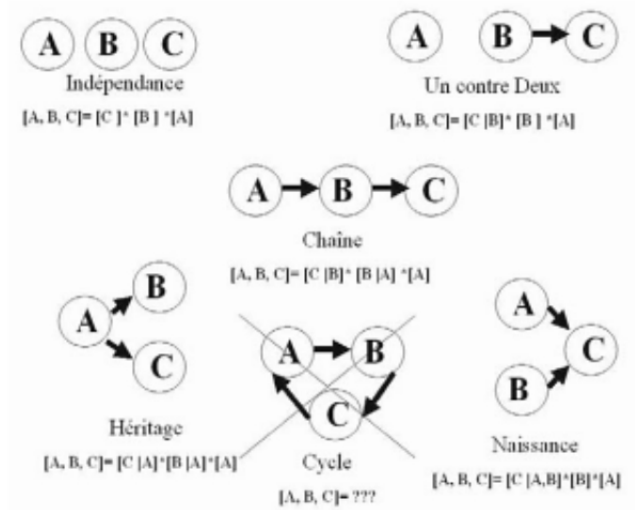


FIGURE 4 – Relations de dépendance conditionnelles possibles entre trois variables aléatoires. DAG tiré de Parent et Bernier (2007).

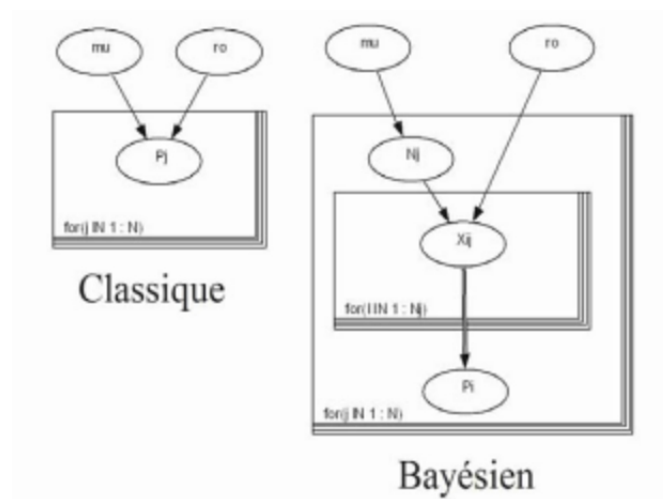


FIGURE 5 – DAG incluant des variables latentes. DAG tiré de Parent et Bernier (2007).

6.4 Convergence des priors informatifs vers des priors régularisant

On peut interpréter, et édifier comme une règle de bon sens permettant de juger le bien-fondé d'un prior informatif, le fait qu'un *a priori* régularisant, typiquement impropre, soit décrit comme la limite d'une suite de priors informatifs *en un certain sens* (de moins en moins informatifs). Ce "sens" a fait l'objet de nombreuses spéculations au cours du temps, et a été définitivement formalisé dans les travaux de Bioche (2015).

EXEMPLE 20. Soit $X \sim \mathcal{E}(\lambda)$ et $\pi(\lambda) \equiv \mathcal{G}(a, b)$. On sait alors que la loi *a posteriori* de λ sachant un échantillon iid x_1, \dots, x_n est $\lambda \sim \mathcal{G}(a + n, b + \sum_{i=1}^n x_i)$. Les hyperparamètres (a, b) ont le sens respectif d'une taille et d'une somme d'échantillon virtuel. "Moralement", faire tendre (a, b) vers 0 revient à annuler l'information apportée par cet échantillon virtuel. On voit alors que $\forall \lambda > 0$, $\pi(\lambda)$ "tend" vers $1/\lambda$, qui est bien le prior de Jeffreys.

Ainsi, Wallace (1959) a proposé une première démarche s'appuyant sur la convergence en tout point :

Proposition 5 Si π est une densité *a priori* impropre, alors il existe une suite de densités *a priori* propres $\{\pi_n\}_n$ engendrant une suite d'*a posteriori* $\{\pi_n(\cdot|x)\}_n$ telle que pour tout $\theta \in \Theta$ et pour tout x , tout x ,

$$\lim_{n \rightarrow \infty} \pi_n(\theta|x) = \pi(\theta|x).$$

Ce résultat reste vrai si $\{\pi_n\}_n$ est une suite de densités telle qu'il existe une constante K et une suite $\{a_n\}_n$ telle que, pour tout θ ,

$$\lim_{n \rightarrow \infty} a_n \pi_n(\theta) = \pi(\theta) \quad \text{et} \quad a_n \pi_n(\theta) \leq K \pi(\theta).$$

Comme l'a fait remarquer Stone en 1965, ce type d'approche est *rétrospective* : le jeu de données est fixé avant tout. Or on souhaite pouvoir caractériser le comportement de $\pi(\theta)$ *avant* l'occurrence de données. D'autres notions de convergence de mesures ont été pour cela étudiées.

1. La *convergence en probabilité* proposée par Stone (1965) vers des mesures impropres *relativement invariantes* continues et telles que $\pi(\theta_1, \theta_2) = \pi(\theta_1)\pi(\theta_2)$. Cette approche nécessite d'introduire des suites de mesures *a priori* obtenues par troncature (suite croissante de compacts sur Θ).
2. La *convergence en variation totale* considérée par Head et Sudderth (1989), définie par

$$\|\pi_n(\theta) - \pi(\theta)\| = \sup_{\mathcal{F}} |\pi_n(\theta) - \pi(\theta)|$$

où \mathcal{F} est une σ -algèbre de sous-ensembles de Θ .

3. La *convergence en entropie relative* étudiée par Berger et ses co-auteurs en 2009.

Mais ces approches supposent de nombreuses hypothèses et n'offrent pas une vision complète de la convergence. En particulier, elles nécessitent souvent de travailler sur des suites de sous-ensembles de Θ compacts, des priors tronqués, et les lois *a posteriori* limites ne sont pas forcément cohérents avec les mesures régularisantes obtenues par des règles formelles.

Ainsi, Bioche (2015) est la première à proposer une vision *prospective* ou intrinsèque (préalable à l'occurrence de données) complète, par le biais de la *convergence vague* des mesures de Radon strictement positives (soit des mesures strictement positives finies sur les compacts). Ce mode de convergence est équivalent à la notion de *convergence étroite* pour les mesures bornées (et donc les mesures de probabilité).

Définition 33 Convergence vague de mesures de Radon. Soit $\{\mu_n\}_n$ et μ des mesures (de Radon). La suite $\{\mu_n\}_n$ converge vaguement vers μ si, pour toute fonction h continue à support compact,

$$\lim_{n \rightarrow \infty} \int h d\mu_n = \int h d\mu.$$

Définition 34 Convergence étroite de mesures bornées. Soit $\{\mu_n\}_n$ et μ des mesures bornées. La suite $\{\mu_n\}_n$ converge étroitement vers μ si, pour toute fonction h continue bornée,

$$\lim_{n \rightarrow \infty} \int h d\mu_n = \int h d\mu.$$

À partir de ces ingrédients, on peut définir la q -convergence vague de mesures positives.

Définition 35 q -convergence vague de mesures positives. Une suite de mesures positives $\{\mu_n\}_n$ converge q -vaguement vers une mesure positive μ s'il existe une suite de réels positifs $\{a_n\}_n$ telle que $\{a_n \mu_n\}_n$ converge vaguement vers μ .

On obtient alors le résultat suivant.

Proposition 6 Si $(\pi_n)_n$ est une suite de priors qui converge q -vaguement vers π , et si $\theta \rightarrow f(x|\theta)$ est une fonction continue sur Θ et non nulle, alors $\pi_n(\cdot|x)$ converge q -vaguement vers $\pi(\cdot|x)$.

Dans les cas discrets, si μ et μ_n sont définies sur $\Theta = \{\theta_i\}_{i \in I}$, la convergence q -vague est équivalente à : $\forall i \in I$,

$$\lim_{n \rightarrow \infty} a_n \mu_n(\theta_i) = \mu(\theta_i).$$

Dans les cas continus, on peut obtenir le résultat suivant.

Proposition 7 Soient μ et μ_n des mesures a priori sur Θ . Supposons que :

1. il existe une suite de réels positifs $\{a_n\}_n$ tel que la suite $\{a_n \mu_n\}_n$ converge ponctuellement vers μ ;
2. pour tout ensemble compact K , il existe un scalaire M et $N \in \mathbb{N}$ tels que, pour tout $n > N$,

$$\sup_{\theta \in K} a_n \mu_n(\theta) < M.$$

Alors $\{\mu_n\}_n$ converge q -vaguement vers μ .

Exercice 24 Soit $\Theta = \mathbb{N}$ et $\Pi_n = \mathcal{U}(\{0, 1, \dots, n\})$ la distribution uniforme discrète sur le compact discret $\{0, \dots, n\}$. Prouver que $\{\Pi_n\}_n$ converge q -vaguement vers la mesure de comptage.

Exercice 25 Soit $\Theta = \mathbb{R}$ et $\Pi_n = \mathcal{N}(0, n)$. Prouver que $\{\Pi_n\}_n$ converge q -vaguement vers la mesure de Lebesgue.

6.5 Démarches critiques d'élicitation

6.5.1 Détecter et limiter les conflits entre prior et données

Si rien formellement n'empêche de chercher à utiliser un prior $\theta \rightarrow \pi(\theta)$ et une vraisemblance qui, vue comme une fonction $\Theta \rightarrow \mathbb{R}^+$ ($\theta \rightarrow f(x_1, \dots, x_n|\theta)$), peuvent être "éloignés" l'un de l'autre dans Θ (au sens où les régions de haute densité *a priori* favorisent d'autres zones que $\theta \rightarrow f(x_1, \dots, x_n|\theta)$), ce type de situation doit être en pratique détecté. On parle de **conflit entre prior et données** (même si ces deux objets ne vivent pas sur le même espace), qui peut aboutir à un calcul de posterior dont les régions HPD ne correspondent à aucune des zones privilégiées par les deux sources d'information. C'est donc un paradoxe qui a pour origine l'absence de règles d'élicitation du prior en cohérence avec la vraisemblance.

Afin de mieux appréhender ce problème, considérons l'exemple suivant. Soit un échantillon $\mathbf{x}_n \sim \mathcal{N}(\mu, \sigma^2)$. On suppose connaître σ , mais μ est inconnu. On place l'*a priori* conjugué $\mu \sim \mathcal{N}(m, \rho\sigma^2)$. Peut-on émettre une règle simple de cohérence entre $\pi(\mu)$ et la vraisemblance des données ?

On peut tenter de répondre à ce problème en utilisant des principes de statistique classique (fréquentiste), via une approche fondée sur la notion de *prior predictive check*.

Prior predictive check. En reprenant l'exemple précédent, on voit que $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$ est une **statistique exhaustive** de l'échantillon. Sous une hypothèse iid. des X_i , la loi de la variable aléatoire associée \bar{X}_n est, conditionnellement à (μ, σ^2)

$$\bar{X}_n | \mu, \sigma^2 \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right).$$

Intégrée sur $\pi(\mu)$, la loi *a priori prédictive* de \bar{X}_n est

$$\bar{X}_n | \sigma^2 \sim \mathcal{N}\left(m, \sigma^2\left(\frac{1}{n} + \rho\right)\right).$$

Alors, *a priori* et *prédictivement a priori*,

$$Z_n = \frac{(\bar{X}_n - m)^2}{\sigma^2(\frac{1}{n} + \rho)} \sim \chi_1^2.$$

Il y a donc incohérence entre *a priori* et vraisemblance des données \mathbf{x}_n si la valeur *observée* de Z_n , à partir de \bar{x}_n , est une *valeur extrême* de la distribution du χ_1^2 . On cherche donc à "tester" le caractère extrême de la valeur observée de Z_n , ce qui revient à devoir définir de façon subjective (comme un niveau de *p-value*) une frontière. Ce choix peut avoir une forte influence, comme en témoignent les figures 6 à 10.

Une autre problème subsiste : comment justifier le choix d'une statistique "résumée", pour comparer sa valeur observée avec sa distribution *a priori* prédictive, plutôt qu'une autre ? Enfin, remarquons qu'il n'est pas évident que la loi sous θ de la statistique choisie (ici, Z_n) dépende complètement de toutes les dimensions de θ . Dans le pire de cas, cette statistique peut être **ancillaire/pivotale**, ce qui signifie qu'on ne peut l'utiliser pour détecter un conflit. Et hors de cas simples il n'est pas évident de vérifier qu'une statistique est bien ancillaire (ou pivotale).

Définition 36 Statistique pivotale ou ancillaire. Une statistique est appelée pivotale ou ancillaire par rapport au paramètre θ lorsque sa loi est indépendante de ce paramètre.

Data-Agreement Criterion. Pour remédier à la vision "fréquentiste" du problème de la détection de conflit, mais aussi pour évacuer le problème du choix de la statistique à tester, et enfin pour éviter le problème supplémentaire posé par les statistiques ancillaires, d'autres critères de détection ont été proposés. C'est le cas du critère DAC décrit ci-dessous.

Ce critère vise non pas à tester un désaccord entre sources d'information sur X via un choix de statistique vivant dans l'espace des X , mais à tester un désaccord quand les sources d'information sont "projetées" dans l'espace Θ . Il repose sur les hypothèses suivantes, dépendante d'un choix de prior "non informatif" pour l'expérience, noté π^J :

1. π^J est toujours en accord avec les données \mathbf{x}_n ;
2. $\pi^J(\theta|\mathbf{x}_n)$ est considéré comme une *loi a priori "parfaite"*, correspondant à un expert fictif parfaitement en accord avec les données.

L'idée principale est que la divergence KL

$$D\{\pi^J(\theta|\mathbf{x}_n) \parallel \pi(\theta)\}$$

formule un regret informationnel dû au choix de states $\pi(\theta)$ alors que le choix de référence serait $\pi^J(\theta|\mathbf{x}_n)$. Quand $D\{\pi^J(\theta|\mathbf{x}_n) \parallel \pi(\theta)\}$ est large, cela signifie un désaccord entre l'information *a priori* et l'information apportée par les données sur θ . Supposons alors que $D\{\pi^J(\theta|\mathbf{x}_n) \parallel \pi(\theta)\} > D\{\pi^J(\theta|\mathbf{x}_n) \parallel \pi^J(\theta)\}$, et que $\pi(\theta)$ est plus informatif que $\pi^J(\theta)$. Nécessairement, le prior et les données sont en conflit. Cela permet de proposer le critère suivant.

Définition 37 Data Agreement Criterion (DAC ; 2007).

$$DAC^J(\pi|\mathbf{x}_n) = \frac{D\{\pi^J(\theta|\mathbf{x}_n) \parallel \pi(\theta)\}}{D\{\pi^J(\theta|\mathbf{x}_n) \parallel \pi^J(\theta)\}}.$$

et π et \mathbf{x}_n sont dits en conflit si $DAC^J(\pi|\mathbf{x}_n) > 1$.

Une conséquence immédiate de cette définition est que lorsqu'on a affaire à un prior hiérarchique, on peut calculer DAC de façon séparée à chaque niveau de la hiérarchie et obtenir le DAC pour le prior complet directement.

Proposition 8 Priors hiérarchiques. Soit $\pi(\theta) = \pi(\theta_1|\theta_2)\pi(\theta_2)$ et $\pi^J(\theta) = \pi^J(\theta_1|\theta_2)\pi^J(\theta_2)$. Soit $\tilde{\pi}_1(\theta) = \pi(\theta_1|\theta_2)\pi^J(\theta_2)$ et $\tilde{\pi}_2(\theta) = \pi^J(\theta_1|\theta_2)\pi(\theta_2)$.

Alors

$$DAC(\pi|\mathbf{x}_n) = DAC(\tilde{\pi}_1|\mathbf{x}_n) + DAC(\tilde{\pi}_2|\mathbf{x}_n) - 1.$$

Proposition 9 Fusion de priors (voir § 6.5.2). Soit

$$\pi(\theta) \propto \prod_{i=1}^M \pi_i(\theta)^{\omega_i}.$$

Alors

$$DAC^J(\pi|\mathbf{x}_n) \leq \sum_{i=1}^M \omega_i DAC^J(\pi_i|\mathbf{x}_n). \quad (\text{arguments de convexité simples})$$

Proposition 10 Les priors réguliers ne sont pas en conflit avec \mathbf{x}_n quand $n \rightarrow \infty$:

$$\mathbb{E}[DAC^J(\pi|\mathbf{x}_n)] \rightarrow 1.$$

Deux situations peuvent se présenter au statisticien :

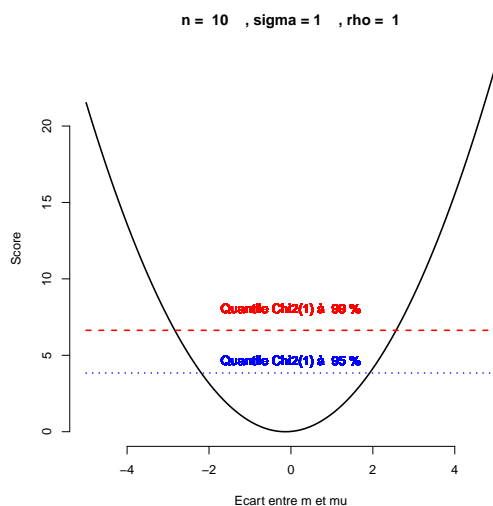


FIGURE 6 – Influence du choix d'un seuil définissant le caractère extrême.

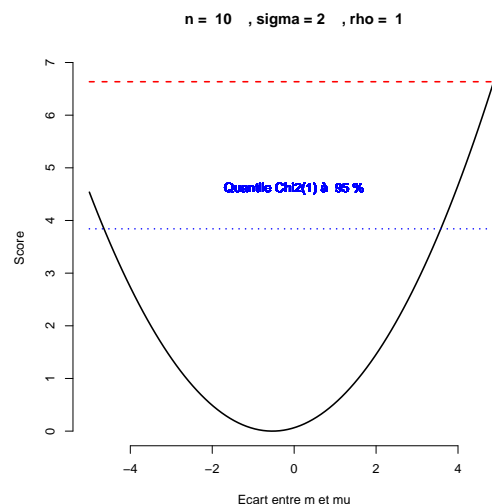


FIGURE 7 – Influence d'un σ plus large.

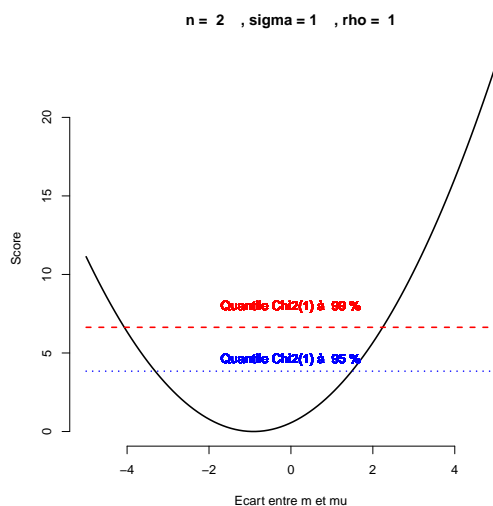


FIGURE 8 – Influence d'une très faible taille d'échantillon.

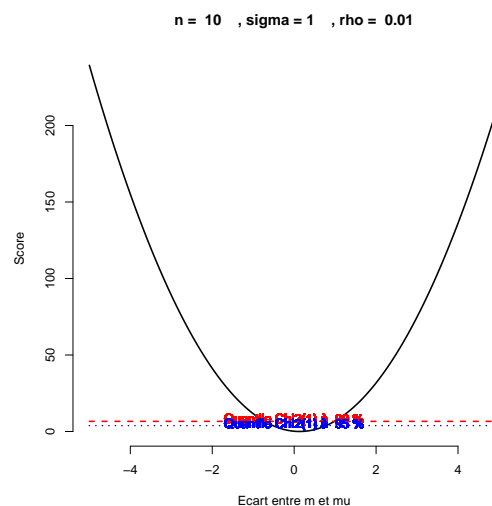


FIGURE 9 – Influence d'une expertise extrêmement précise.

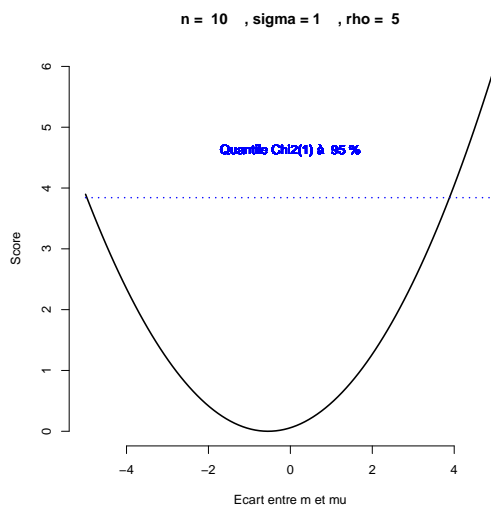


FIGURE 10 – Influence d'une expertise extrêmement vague.

1. *Cas idéal.* Θ est borné et / ou $\mathcal{M}(\theta)$ est discret. Alors π^J est souvent **propre** :

$$\int_{\Theta} \pi^J(\theta) d\theta < \infty.$$

2. *Cas problématique.* Si π^J est impropre (ce qui arrive le plus souvent), le dénominateur de DAC est défini à une *constante additive près inconnue*. Dans ce cas, DAC ne peut pas être utilisé en l'état, et doit être adapté.

EXEMPLE 21. *Modèle gaussien borné. Supposons avoir*

$$\begin{aligned} \mathbf{x}_n &\stackrel{i.i.d.}{\sim} \mathcal{N}(\theta, 1) \\ \theta &\sim \mathcal{N}(\mu_0, \sigma_0) \text{ restreint à } D = [-15, 15] \end{aligned}$$

$\pi^J(\theta)$ étant la loi uniforme sur D . Dans ce cas, le calcul est aisé, et on peut par exemple tracer l'évolution de DAC en fonction de μ_0 sur la figure 11.

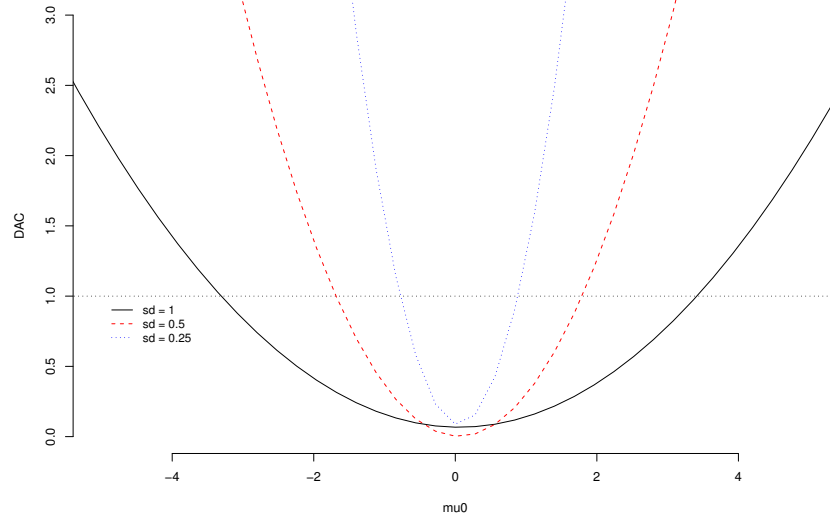


FIGURE 11 – Exemple 6.5.1. Vraie valeur $\theta_0 = 0$, taille $n = 5$

Quand π^J est impropre, on a les mêmes difficultés conceptuelles que lorsqu'on souhaite calculer un *facteur de Bayes*

$$B_{J,\pi}(\mathbf{t}_n) = \frac{\int_{\Theta} \mathcal{L}(\mathbf{x}_n|\theta) \pi(\theta) d\theta}{\int_{\Theta} \mathcal{L}(\mathbf{x}_n|\theta) \pi^J(\theta) d\theta}$$

qui est défini à une constante multiplicative inconnue près. En s'inspirant d'approches dites *intrinsèques*, proposées originellement par Berger et Perrichi (1996, 1998), on peut proposer une adaptation en usant de *petites quantités de données d'entraînement* $x(l) \subset \mathbf{t}_n$ pour remplacer le $\pi^J(\theta)$ impropre par un $\underbrace{\pi^J(\theta|x(l))}_{\text{"posterior prior"}}$ propre.

Le sous-échantillon $x(l)$ est appelé *échantillon minimal d'entraînement* (MTS).

Définition 38 Critère DAC intrinsèque. Soit $\mathbf{x}_n(-1) = \mathbf{x}_n / \{x(l)\}$ and ℓ la taille d'un MTS $x(l)$. Par un argument de validation croisée

$$DAC^{AIJ}(\pi|\mathbf{x}_n) = \frac{1}{L} \sum_{l=1}^L \frac{D\{\pi^J(\cdot|\mathbf{x}_n(-1)) \parallel \pi(\cdot|x(l))\}}{D\{\pi^J(\cdot|\mathbf{x}_n(-1)) \parallel \pi^J(\cdot|x(l))\}}.$$

La qualité de l'adaptation est généralement bonne quand $L \geq 10$ et $n/\ell \geq 8$. Il reste cependant coûteux à calculer (sauf, évidemment, dans les cas conjugués).

EXEMPLE 22. *Loi de Bernoulli. $X \sim \mathcal{B}_r(\theta)$ avec un prior de loi bêta. On choisit pour l'exemple $n = 20$, une valeur de simulation $\theta_0 = 0.7$, et un prior $\theta \sim \mathcal{B}_e$ sur $[0, 1]$, de moyenne μ_0 et pour un écart-type $\sigma = 0.2$ fixé. Dans ce cas conjugué pour lequel le prior π^J de Jeffreys est propre, on peut comparer les deux approches (figure 12).*

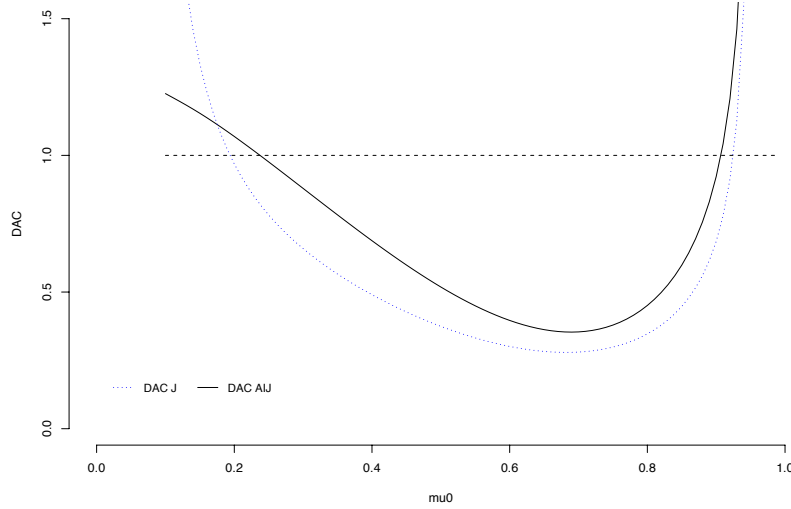


FIGURE 12 – Exemple 6.5.1 : comparaison du DAC et du DAC intrinsèque.

EXEMPLE 23. *Loi exponentielle. Soit $X \sim \mathcal{E}(\lambda)$ avec $n = 10$ et un EMV $\hat{\lambda} = 207$. On choisit naturellement un prior conjugué $\lambda \sim \mathcal{G}(a, a.te)$ et on décide de faire varier la "taille virtuelle a " et l'hyperparamètre t_e , qui correspond à la moyenne d'un échantillon virtuel (figure 13).*

Le critère DAC présente en outre un avantage pour la calibration. En effet, DAC^J et DAC^{AIJ} détectent des priors très biaisés et des priors non biaisés mais trop informatifs vis-à-vis des données \mathbf{x}_n . Si on se retrouve dans une situation où la seule information *a priori* dont nous disposons est une valeur centrale pour θ , on peut souhaiter calibrer raisonnablement la taille virtuelle m telle que

$$DAC^{AIJ}(m^*|\mathbf{x}_n) = 1.$$

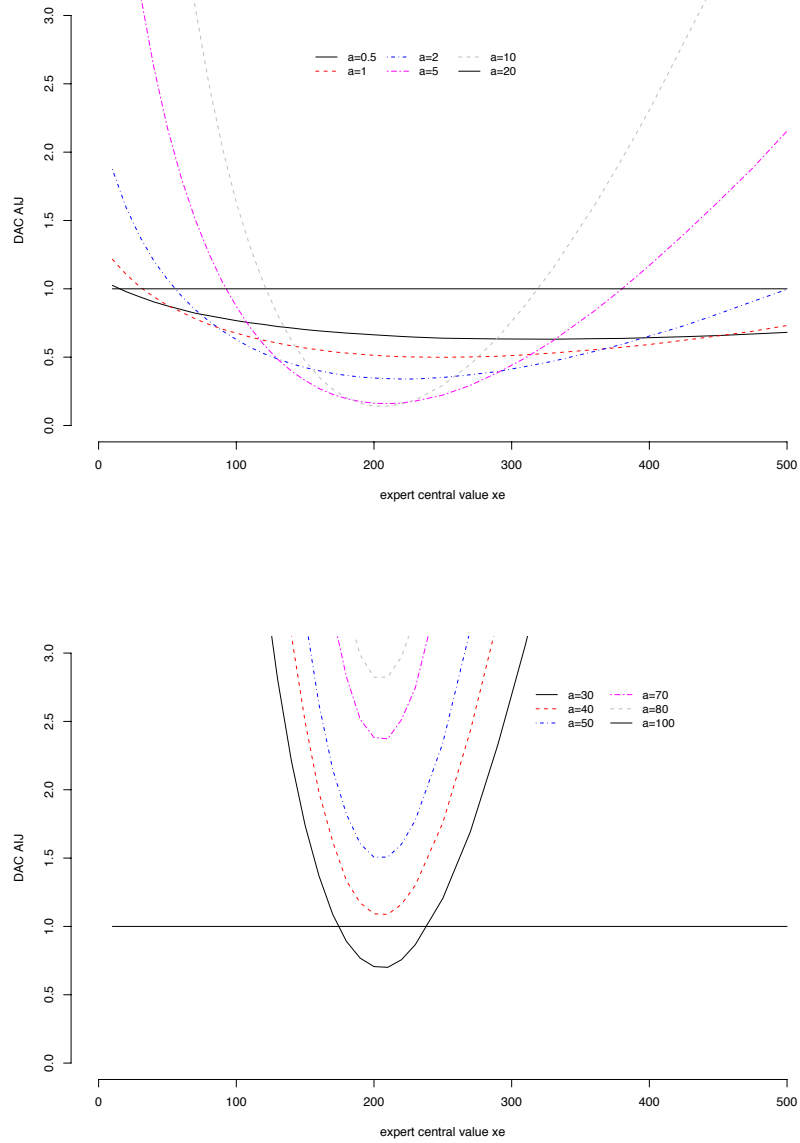


FIGURE 13 – Exemple 6.5.1 : DAC intrinsèques en fonction des variations des hyperparamètres.

6.5.2 Fusionner plusieurs priors

Dans de nombreux cas pratiques, on peut disposer de plusieurs *a priori* possibles $\pi_1(\theta), \dots, \pi_M(\theta)$.

EXEMPLE 24. *Réunions d'experts à la fin d'études pharmacologiques.*

Considérons d'abord la situation où les priors sont produits indépendamment les uns des autres. Une première idée est de proposer une **fusion linéaire pondérée** (ou moyenne arithmétique) :

$$\pi(\theta) = \sum_{i=1}^M \omega_i \pi_i(\theta)$$

avec $\sum_{i=1}^M \omega_i = 1$. Cette proposition n'est pas sans problème. En effet, le résultat produit peut être multi-modal, et contre-intuitif. De plus, cette approche n'est pas **externalement bayésienne** :

$$\pi(\theta | \mathbf{x}_n) \neq \sum_{i=1}^M \omega_i \pi_i(\theta | \mathbf{x}_n)$$

pour une ou plusieurs données \mathbf{x}_n . Cela pose un problème de cohérence.

Une autre approche est la **fusion logarithmique pondérée** (ou moyenne géométrique) :

$$\pi(\theta) = \frac{\prod_{i=1}^M \pi_i^\omega(\theta)}{\int_{\Theta} \prod_{i=1}^M \pi_i^\omega(\theta) d\theta}$$

avec $\sum_{i=1}^M \omega_i = 1$. Celle-ci est bien extérieurement bayésienne.

Toutefois, elle pose un autre problème : elle n'est pas *cohérente par marginalisation*. En fait, seule la fusion linéaire permet de respecter ce principe de cohérence.

Définition 39 Cohérence par marginalisation. Soit A et B tels que $A \cap B = \emptyset$ et $C = A \cup B$. Supposons avoir des systèmes experts indépendants produisant des estimateurs des probabilités des événements A et B . Pour chacun, on peut directement calculer $P(C)$ ou calculer séparément $P(A)$ puis $P(B)$. La modélisation de ces experts est cohérente par marginalisation si $P(C) = P(A) + P(B)$.

Mais la fusion logarithmique est globalement séduisante car elle peut s'expliquer en faisant appel à la théorie de l'information. Rappelons que la divergence de Kullback-Leibler

$$KL(\pi, \pi_i) = \int_{\Theta} \pi(\theta) \log \frac{\pi(\theta)}{\pi_i(\theta)}$$

exprime une perte en termes d'information lorsque le meilleur choix *a priori* π est remplacé par π_i .

Proposition 11 *Le minimiseur de la perte pondérée*

$$\pi^*(\theta) = \arg \min_{\pi} \sum_{i=1}^M \omega_i KL(\pi, \pi_i)$$

est l'a priori opérant la fusion logarithmique.

Notons cependant que la calibration des poids ω_i est un problème qui reste ouvert, malgré quelques réponses déjà proposées. Un dernier argument plaide pour ce choix de fusion. La famille exponentielle naturelle

est stable par ce type de fusion, et cette fusion correspond naturellement à celle réalisée par une inférence bayésienne croissante, indépendante de l'ordre d'arrivée des informations (ex : échantillons virtuels).

Il n'existe pas de démarche canonique pour traiter la situation où les priors sont produits dépendamment les uns des autres. Les approches par *copules* peuvent être utilisées si l'on connaît explicitement la structure de dépendance. On peut aussi considérer que les (systèmes) experts produisant les priors π_i sont des générateurs aléatoires dépendants de valeurs de θ , et une démarche de construction hiérarchique *a priori* peut alors être proposée.

Exercice 26 *Considérons M priors exponentiels*

$$\theta \sim \pi_i(\theta) = \lambda_i \exp(-\lambda_i \theta) \mathbb{1}_{\{\theta \geq 0\}}.$$

Quelle est la loi-fusion logarithmique ?

6.6 TP : Un exemple complet dans un cadre de fiabilité industrielle

Soit X la durée de vie d'un composant Σ , supposé tomber uniquement en panne par hasard. Le taux de défaillance λ de Σ est donc constant, ce qui implique $X \sim \mathcal{E}(\lambda)$. Il est courant de disposer d'un expert industriel familier de λ , avec qui le dialogue suivant peut être engagé. "Considérons une décision de management (remplacement) établie sur une valeur donnée $\bar{\lambda}$ (*différente de la vraie valeur inconnue* λ)

Pour un coût similaire $|\bar{\lambda} - \lambda|$, il y a 2 conséquences possibles au remplacement :

- soit C_1 le coût positif moyen d'être trop optimiste (d'avoir $\bar{\lambda} \leq \lambda$);
- soit C_2 le coût positif moyen d'être trop pessimiste (d'avoir $\bar{\lambda} > \lambda$).
- Pouvez-vous donner un estimé $\hat{\delta}$ du rapport des coûts moyens $\delta = C_2/C_1$?

L'axiome de rationalité dit que si l'expert n'est pas *averse au risque*, alors

$$\bar{\lambda} = \arg \min_x \underbrace{\int_0^\infty |x - \lambda| (C_1 \mathbb{1}_{\{x \leq \lambda\}} + C_2 \mathbb{1}_{\{x > \lambda\}}) \pi(\lambda) d\lambda}_{\text{fonction de coût intégrée sur toutes les valeurs possibles a priori du vrai } \lambda}$$

$$\text{Il s'ensuit que } \int_0^{\bar{\lambda}} d\Pi(\lambda) = \Pi(\lambda < \bar{\lambda}) = \frac{C_1}{C_1 + C_2}.$$

L'interprétation de la réponse de l'expert est que $1/(1 + \hat{\delta})$ est un estimé du quantile *a priori* d'ordre $\alpha = C_1/(C_1 + C_2)$. Avec $P(\lambda < \bar{\lambda}) = \frac{C_1}{C_1 + C_2} = \alpha$, on a :

- tant que les coûts sont équilibrés, un expert de plus en plus optimiste fournira un $\bar{\lambda}$ de plus en plus petit, car la durée moyenne avant la prochaine défaillance est

$$\mathbb{E}[X|\lambda] = \frac{1}{\lambda}.$$

- cependant l'expert s'exprime plutôt sur les coûts lorsqu'on lui fournit une valeur représentative de $\bar{\lambda}$:
 - plus l'expert est optimiste, plus le coût C_2 d'être optimiste (selon lui) est petit, donc α grandit vers 1 et

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|\lambda]] = \mathbb{E}[1/\lambda] \text{ augmente.}$$

- plus l'expert est pessimiste, plus le coût C_2 d'être optimiste augmente, donc α tombe vers 0 et

$$\mathbb{E}[X] = \mathbb{E}[1/\lambda] \text{ diminue.}$$

Quelle choix de loi *a priori* pouvons-nous proposer au décideur ?

6.7 L'importance du prior en *deep learning*

Les approches de *deep learning bayésien*, où les paramètres (typiquement les poids d'un réseau de neurones) sont supposés aléatoires, prennent de plus en plus d'importance dans la littérature scientifique, et commencent à être régulièrement déployés. Deux problèmes importants limitent encore leur usage :

1. Le calcul *a posteriori*, qui nécessite de nombreux calculs en très grande dimension ; ces limitations ont notamment mené à l'usage de techniques variationnelles pour approximer ce calcul. Le lecteur intéressé pourra regarder l'article de revue [14] pour s'informer sur les techniques de calcul des incertitudes bayésiennes pour l'usage des réseaux de neurones.
2. Le choix *a priori*, qui est souvent le parent pauvre de la démarche. En général, des lois gaussiennes standard sont utilisées. La revue [12] souligne cependant l'importance du choix des priors pour l'apprentissage profond bayésien et présente un aperçu des différents priors qui ont été proposés pour les processus gaussiens (profonds), les auto-encodeurs variationnels et les réseaux neuronaux bayésiens.

7 Incorporation de connaissance *a priori*

7.1 Notion de crédibilité

Le paradigme bayésien critique la notion de *confiance statistique classique* : celle-ci n'est pas réellement adaptée à une prise de décision fondée majoritairement sur les observations disponibles, car elle repose sur un comportement d'observations équivalentes aux observations réelles, mais répétées à l'infini en fonction du modèle choisi. Les arguments permettant de construire ces zones de confiance reposent en effet sur le comportement en loi des estimateurs $\hat{\theta}_n$ (cf. § A.6), qui est en général connu asymptotiquement (soit quand $n \rightarrow \infty$). C'est le cas des estimateurs utilisés précédemment pour les modèles de statistique extrême. Or, la seule information dont on dispose sans hypothèse de modèle est l'échantillon \mathbf{x}_n , et la taille n de ces données peut ne pas être assez grande pour s'assurer d'être "proche" de l'asymptotisme.

À cette notion la théorie bayésienne substitue celle de *crédibilité* : **conditionnellement aux données observées**, une zone de crédibilité estimée contient θ avec une certaine probabilité.

Le conditionnement de θ aux observations implique alors que θ est considéré non plus comme un vecteur de paramètres fixes et inconnus, mais comme une variable latente aléatoire dont la loi est notée génériquement Π (et sa densité π). L'existence d'une mesure de probabilité sur θ est attestée par le théorème de représentation de De Finetti et ses versions étendues [30] lorsque les observations sont simplement considérées comme *échangeables* (et donc potentiellement dépendantes) et non plus *iid*. Le cadre bayésien englobe donc le cadre classique. L'ouvrage [36] (chap. I.5) offre de nombreux détails sur l'existence de Π dans un cadre général.

Le procédé d'inférence consiste alors à opérer une *mise à jour de la loi* Π conditionnellement aux observations \mathbf{x}_n , via la règle de Bayes [38] : la densité dite *a priori* $\pi(\theta)$ est modifiée, par l'ajout d'information issue de la vraisemblance statistique $\ell(\mathbf{x}_n|\theta)$, en une densité *a posteriori* $\pi(\theta|\mathbf{x}_n)$ s'écrivant

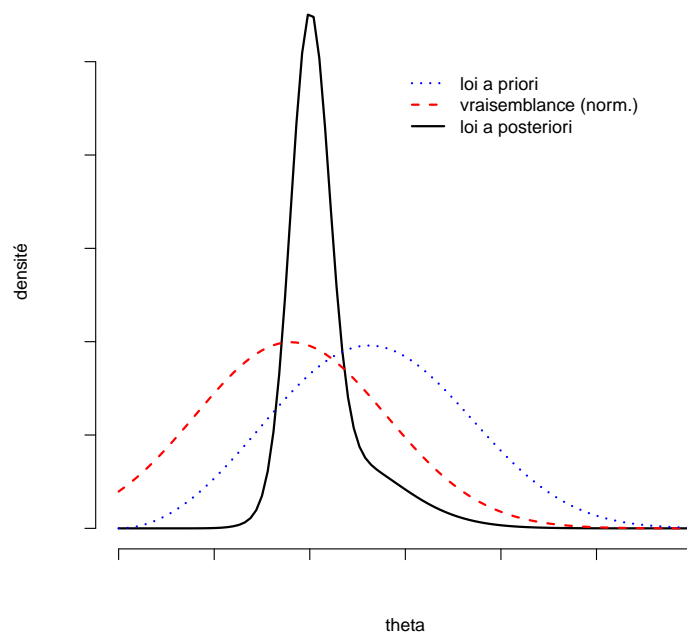
$$\pi(\theta|\mathbf{x}_n) = \frac{\pi(\theta)f(\mathbf{x}_n|\theta)}{\int_{\Theta} \pi(\theta)f(\mathbf{x}_n|\theta) d\theta}. \quad (18)$$

Ainsi, là où la statistique classique calcule un estimateur $\hat{\theta}_n$ (maximisant souvent la vraisemblance $f(\mathbf{x}_n|\theta)$) qui possède une loi connue asymptotiquement, la statistique bayésienne considère que la vraie nature de θ est mieux décrite par une autre loi statistique qui est définie à taille n fixée, selon l'expression (18). Des liens profonds existent entre ces deux approches ; en particulier, lorsque $n \rightarrow \infty$, les lois respectives de $\hat{\theta}_n$ et $\pi(\theta|\mathbf{x}_n)$ deviennent similaires. Ces liens sont présentés en détail dans l'ouvrage de référence [38].

La figure 14 illustre le positionnement classique de la vraisemblance statistique $f(\mathbf{x}_n|\theta)$, vue comme une fonction de θ , par rapport aux densités $\pi(\theta)$ et $\pi(\theta|\mathbf{x}_n)$; la mutualisation des sources d'information sur θ se traduit logiquement par une distribution *a posteriori* plus "piquée" – plus informative donc – que la loi *a priori*¹².

Nature de la loi *a priori*. Qu'exprime la loi *a priori* $\pi(\theta)$? Un état d'information initial sur les valeurs possibles de θ , indépendant des observations \mathbf{x}_n . Encodé sous forme probabiliste, cet état d'information est donc susceptible de permettre l'ajout d'une connaissance réelle et sérieuse du phénomène exprimée autrement qu'à travers d'observations passées : expertise, prévisions de modèles physiques, etc. Il faut noter que l'information *a priori* sur θ n'est jamais inexistante : en effet, θ est un choix de paramétrisation du modèle de vraisemblance, et la structure de ce modèle est connue. En découlent des contraintes sur la structure de corrélation de $\pi(\theta)$. Nous conseillons au lecteur intéressé par une présentation didactique des possibles interprétations de $\pi(\theta)$ l'ouvrage [31] et l'article [32].

12. Excepté dans les cas où les deux sources d'information sont en désaccord : voir [10].



380380

FIGURE 14 – Illustration par les densités du renforcement de l'information sur θ à partir de la loi *a priori* et la vraisemblance des données (vue comme fonction de θ et ici renormalisée).

8 Méthodes de calcul bayésien

8.1 Introduction

Nous nous plaçons dans le cadre à présent bien connu :

- Soit $\{X \sim f(\cdot|\theta), \pi(\theta)\}$ un modèle bayésien servant à prendre une décision $\delta \in \mathcal{D}$.
- Dans un cadre d'**analyse** (*a posteriori*), on a observé des données $\mathbf{x}_n = (x_1, \dots, x_n) \sim f(\cdot|\theta)$.
- La loi *a posteriori* $\pi(\theta|\mathbf{x}_n)$ décrit l'ensemble des incertitudes sur $\theta \in \Theta$, vecteur inconnu qui "paramétrise" l'état de nature.
- Toute décision δ peut être jugée par un *coût* $L(\theta, \delta)$, c'est-à-dire son écart par rapport à une décision idéale inatteignable, affecté par la distribution de probabilité *a posteriori* $\pi(\theta|\mathbf{x}_n)$.
- La *décision optimale* s'obtient en cherchant l'optimum de la fonction du coût moyen *a posteriori* (expected opportunity loss) : $\delta^\pi = \arg \min_{\delta \in \mathcal{D}} R_B(\delta|\pi)$ avec

$$R_B(\delta|\pi) = \int_{\Theta} L(\theta, \delta) d\Pi(\theta|\mathbf{x}_n) = \int_{\Theta} L(\theta, \delta) \pi(\theta|\mathbf{x}_n) d\theta$$

La question fondamentale du calcul bayésien est double : peut-on obtenir une expression explicite pour δ^π ? sinon, comment peut-on l'évaluer numériquement ?

En général, l'estimateur de Bayes n'a pas de caractère explicite, car

$$\pi(\theta|\mathbf{x}_n) = \frac{f(\mathbf{x}_n|\theta)\pi(\theta)}{\int f(\mathbf{x}_n|\theta)\pi(\theta)d\theta}$$

et le dénominateur n'est pas explicitement connu ($\pi(\theta|\mathbf{x}_n)$ est définie à une constante d'intégration près). On peut essayer de l'estimer par *intégration numérique* : les techniques classiques de Newton-Cotes ou Runge-Kutta mènent fréquemment à des instabilités numériques lorsque $\dim \Theta$ augmente. Toutefois, les méthodes de *quadrature bayésienne* permettent de pallier ce problème. Elles relèvent encore du domaine de la recherche, et seront abordées plus tard dans ce cours.

Fort généralement, le calcul d'estimateur de Bayes, mais aussi la sélection de modèle, la production de régions α -crédibles (par exemple), nécessitent de pouvoir **simuler** *a posteriori*. Par des techniques de Monte Carlo, fondée sur une capacité algorithmique de simulation pseudo-aléatoire (voir Annexe C), on peut ainsi mener les calculs suivants :

1. **Calcul d'une moyenne *a posteriori***. On peut estimer δ^π par un estimateur de Monte Carlo

$$\hat{\delta}_M^\pi = \frac{1}{M} \sum_{i=1}^M \theta_i$$

- Ingrédients : Loi Forte des Grands Nombres ($\hat{\delta}_M^\pi \xrightarrow{p.s.} \delta^\pi$) + Théorème Central Limite

2. **Estimation d'un quantile *a posteriori* d'ordre α** . On peut estimer δ^π par l'inversion de la fonction de répartition empirique *a posteriori*

$$\hat{\Pi}_M(\theta|\mathbf{x}_n) = \frac{1}{M} \sum_{i=1}^M \mathbb{1}_{\{\theta \leq \theta_i^*\}} \quad (\text{Théorème de Glivenko-Cantelli})$$

soit en prenant

$$\hat{\delta}_M^\pi = \begin{cases} \frac{1}{2} (\theta_{\alpha \cdot M}^* + \theta_{\alpha \cdot (M+1)}^*) & \text{si } \alpha \cdot M \text{ est entier} \\ \theta_{[\alpha \cdot M] + 1}^* & \text{sinon} \end{cases} \quad (\text{Théorème de Mosteller})$$

Un autre intérêt (et même une nécessité) apporté par la simulation *a posteriori* est le suivant : l'analyse prédictive. Plus précisément, plaçons-nous dans ce cadre : **simuler des réalisations de X** est une nécessité lorsqu'on s'intéresse au comportement Y d'un phénomène (par exemple physique) modélisé ainsi :

$$Y = g(X, \nu) + \epsilon$$

où :

- g est une fonction (ou un code de calcul) déterministe ;
- ν est un indice ou une variable indexant typiquement des conditions environnementales ;
- ϵ est un "bruit" stochastique qui modélise l'erreur entre la réalité du phénomène Y et la sortie de g .

Dans ce problème de *propagation d'incertitudes*, on cherche à reproduire un grand nombre de configurations de Y pour calculer (par exemple) la probabilité que Y dépasse un certain seuil.

EXEMPLE 25. Y représente une hauteur d'eau aval, g est un code hydraulique, X est un débit d'eau amont, ν caractérise le frottement de la rivière et ϵ tient compte de la méconnaissance du terrain, de la précision du code, etc.

Comment doit être simulé X en entrée de g ? La loi *prédictive* de densité

$$f(x|\mathbf{x}_n) = \int f(x|\theta)\pi(\theta|\mathbf{x}_n) d\theta$$

permet de simuler une prochaine observation x_{n+1} *crédible* sachant qu'on a déjà observé les \mathbf{x}_n . Si l'on cherche à simuler de façon crédible la succession de **deux** futures observations (x_{n+1}, x_{n+2}) , on doit procéder ainsi :

$$\begin{aligned} X_{n+1} &\sim f(x|\mathbf{x}_n), \\ X_{n+2} &\sim f(x|\mathbf{x}_n, x_{n+1}), \text{ etc.} \end{aligned}$$

Un algorithme simple de simulation repose donc sur la simulation de la loi *a posteriori*.

Remarque 8 Liens avec le *machine /deep learning*. Les techniques de simulation et d'échantillonnage, dans le cadre spécifique du *machine /deep learning*, servent notamment à :

- initialiser des algorithmes d'optimisation de métriques, fonctions de coût complexes... ;
- élaborer des modèles génératifs, tels les Generative Adversarial Networks (GAN) ;
- résoudre des problèmes de complétion de données manquantes ou en trop faible nombre (*data augmentation*) ;
- proposer des façons de sélectionner des mini-batches utiles parmi un ensemble de données trop grand pour être traité intégralement (problématique de plan d'expérience) .

8.1.1 Principe de la simulation indirecte

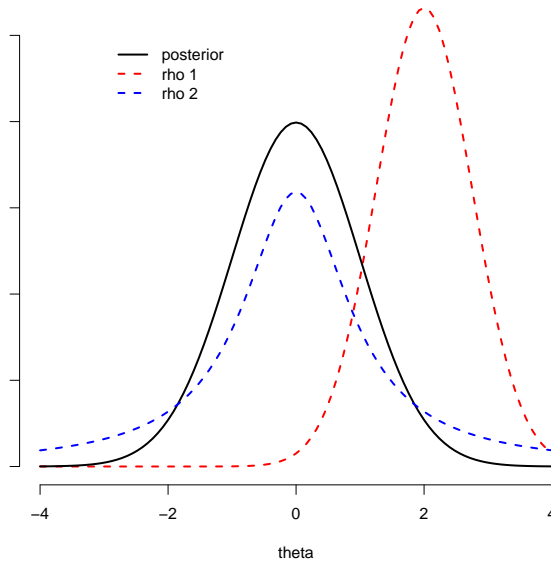
1. on simule un tirage θ_i suivant une *loi instrumentale* $\rho(\theta)$ (facile à simuler) ;
2. on utilise un *test* pour déterminer si θ_i aurait également pu être un *tirage plausible* de $\pi(\theta|\mathbf{x}_n)$.

Plus $\rho(\theta)$ est "proche" de $\pi(\theta|\mathbf{x}_n)$, plus ce test doit accepter les θ_i . Plus précisément :

1. La densité $\rho(\theta)$ doit être facilement simulable (ex : mélanges gaussiens si $\pi(\theta|\mathbf{x}_n)$ est multimodale...).
2. Le support¹³ de $\rho(\theta)$ contient nécessairement celui de $\pi(\theta|\mathbf{x}_n)$.

13. Le domaine de Θ où la densité est non nulle.

3. Les queues de $\rho(\theta)$ devraient être plus lourdes que celles de $\pi(\theta|\mathbf{x}_n)$.



4. Lorsque $\dim \Theta$ est petite (1 ou 2), on peut tracer $\pi(\theta|\mathbf{x}_n)$ à un coefficient près pour sélectionner une forme intéressante pour $\rho(\theta)$.

Un candidat logique peut parfois être la loi *a priori* $\pi(\theta)$, car elle respecte automatiquement la règle d'inclusion du support. Si l'*a priori* est très informatif par rapport aux données, l'*a posteriori* en sera proche. Une quantification de cette "force" relative d'information est donc pratique pour choisir $\rho(\theta)$. Toutefois, ce choix peut être délicat : si l'*a priori* est très large (peu informatif), alors

- il peut privilégier indûment des régions où la vraisemblance (comme fonction de θ) est nulle ou quasi-nulle ;
- il faudra beaucoup de tirages pour atteindre les régions HPD (de plus haute densité) *a posteriori*, ce qui entraînera un coût algorithmique très fort

Ce choix est aussi à proscrire si l'*a priori* privilégie des régions de Θ qui sont éloignées de celles privilégiées par les données. Une indication en faible dimension est de mesurer l'éloignement du mode *a priori* de θ et du maximum de vraisemblance $\hat{\theta}_n$.

8.2 Méthodes d'échantillonnage dans la loi *a posteriori*

Dans cette partie, on cherche donc à obtenir **indirectement** des tirages qui suivent (en général approximativement) la loi *a posteriori*. Citons quelques algorithmes classiques que nous étudierons (quasiment tous) dans ce cours :

1. algorithmes d'acceptation-rejet ;
2. échantillonnage d'importance (préférentiel) ;
3. méthodes de Monte Carlo par chaînes de Markov (MCMC) ;
4. filtrage particulière (pour les modèles à espace d'état).

Ces méthodes – et leurs hybrides – sont les outils actuels les plus puissants pour simuler des lois connues semi-explicitement (à une constante/une intégrale près). Ils connaissent plusieurs approches d'**accélération** dont nous discuterons également, et qui amènent par ailleurs à faire un lien avec les techniques classiques du *machine learning* : une technique de gradient stochastique, qui vise à estimer un mode *a posteriori* (et sous-entend en général que la loi *a posteriori* est approximativement gaussienne), peut être vue comme une forme dégénérée de MCMC adaptivement accélérée.

8.2.1 Rappel : approches par inversion et transformations simples

Rappelons avant de commencer que la méthode générique de simulation de $\theta \sim \pi(\theta|X)$ repose sur l'*inversion de la fonction de distribution* Π qui, en unidimensionnel, est la fonction de répartition.

Théorème 18 Si $U \sim \mathcal{U}[0, 1]$ et $\Pi(\theta|X)$ la fonction de répartition de $\theta|X$, alors $\Pi^{-1}(U|X)$ a la même loi que θ .

La preuve vient du fait que par définition, $\Pi(\Pi^{-1}(U|X) \leq \theta|X) = \Pi(I \leq \Pi(\theta|X)|X) = \Pi(\theta|X)$. Si Π n'est pas parfaitement croissante, on prend $\Pi^{-1}(u|X) = \inf\{\theta ; \Pi(\theta|X) \geq u\}$.

EXEMPLE 26.

- Loi binomiale $\mathcal{B}(n, p) : F_X(x) = \sum_{i \leq x} \binom{n}{i} p^i (1-p)^{n-i}$ et $F_X^{-1}(u)$ s'obtient numériquement.
- Loi exponentielle $\mathcal{E}(\lambda) : F_X(x) = 1 - \exp(-\lambda x)$ et $F_X^{-1}(u) = -\log(u)/\lambda$.
- Loi de Cauchy $\mathcal{C}(0, 1) : F_X(x) = \frac{1}{\pi} \arctan(x) + \frac{1}{2}$ et $F_X^{-1}(u) = \tan(\pi(u - 1/2))$.

La critique de cette approche est aisée : $\Pi^{-1}(\cdot|X)$ est rarement disponible, et ce "théorème" (plutôt un lemme) d'inversion ne s'applique qu'en dimension 1. Pour "contrer" ces problèmes, on peut proposer quelques transformations (voir ci-dessous), mais cela ne permet de régler que des cas particuliers.

Définition 40 Transformation de Box-Müller Pour la loi normale $\mathcal{N}(0, 1)$, si $X_1, X_2 \stackrel{iid}{\sim} \mathcal{N}(0, 1)$, alors

$$X_1^2 + X_2^2 \sim \chi_2^2, \quad \arctan(X_1/X_2) \sim \mathcal{U}([0, 2\pi]).$$

Comme χ_2^2 est identique à $\mathcal{E}(1/2)$, il vient par inversion :

$$X_1 = \sqrt{-2 \log(U_1)} \sin(2\pi U_2) \quad X_2 = \sqrt{-2 \log(U_1)} \cos(2\pi U_2).$$

Les lois de Student et de Fisher se déduisent naturellement de la loi normale et de la loi du chi-deux. La loi de Cauchy se déduit de la loi normale par la règle suivante : si $X_1, X_2 \sim \mathcal{N}(0, 1)$, alors $X_1/X_2 \sim \mathcal{C}(0, 1)$. La loi Beta $\mathcal{B}_e(\alpha, \beta)$, de densité

$$f_X(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1},$$

s'obtient à partir de la loi gamma par la règle suivante : si $X_1 \sim \mathcal{G}a(\alpha, 1)$, $X_2 \sim \mathcal{G}a(\beta, 1)$, alors

$$\frac{X_1}{X_1 + X_2} \sim \mathcal{B}_e(\alpha, \beta).$$

8.2.2 Simulation multidimensionnelle

Le cas de la simulation multidimensionnelle est réglé également en principe par la règle en cascade suivante :

Définition 41 Cascade rule. Supposons vouloir générer dans \mathbb{R}^p l'échantillon $(X_1, \dots, X_p) \sim f(x_1, \dots, x_p)$ dont les composantes ne sont pas nécessairement indépendantes. La densité jointe s'écrit alors

$$f(x_1, \dots, x_p) = f_1(x_1) \times f_{2|1}(x_2|x_1) \dots \times f_{p|-p}(x_p|x_1, \dots, x_{p-1}).$$

On peut donc en déduire la règle d'implémentation suivante :

Simuler pour $t = 1, \dots, T$

1 $X_1 \sim f_1(x_1)$

2 $X_2 \sim f_{2|1}(x_2|x_1)$

\vdots

$X_p \sim f_{p|-p}(x_p|x_1, \dots, x_{p-1})$

8.2.3 Algorithmes d'acceptation-rejet (AR)

Ce type d'algorithme permet de simuler de façon **exacte** et **indépendante** selon la loi *a posteriori*. Il repose sur l'hypothèse suivante sur $\rho(\theta)$:

$$0 < K = \sup_{\theta \in \Theta} \frac{f(\mathbf{x}_n|\theta)\pi(\theta)}{\rho(\theta)} < \infty.$$

Algorithme AR :

1. simulation indirecte : soit $\theta_i \sim \rho(\cdot)$
 2. test :
 - soit $U_i \sim \mathcal{U}[0, 1]$
 - si $U_i \leq \frac{f(\mathbf{x}_n|\theta_i)\pi(\theta_i)}{K\rho(\theta_i)}$ alors θ_i suit la loi $\pi(\theta|\mathbf{x}_n)$
-

(Preuve en cours)

Observons que la loi du nombre de tirages nécessaires selon $\rho(\theta)$ jusqu'à en accepter un suit la loi géométrique de probabilité $1/(K \cdot C)$ où C est la constante d'intégration inconnue

$$C = \int_{\Theta} f(\mathbf{x}_n|\theta)\pi(\theta) d\theta$$

donc $K \cdot C$ est l'espérance du nombre de tirages nécessaires avant l'acceptation. *Optimiser l'algorithme* revient donc à *diminuer* K .

Exercice 27 On suppose $X \sim \mathcal{N}(\theta, 1)$ et on suppose connaître un échantillon \mathbf{x}_n composé de :

- quelques observations x_1, \dots, x_{n-1} supposées iid.
- une pseudo-observation y qui est un cas-limite masquant (censurant) une observation x_n qui aurait dû être faite : $y < x_n$

A priori, on suppose $\theta \sim \mathcal{N}(\mu, 1)$. Pouvez-vous produire un algorithme d'AR qui génère des réalisations de la loi *a posteriori* de θ ?

Exercice 28 Soit un échantillon de loi gamma $x_1, \dots, x_n \stackrel{iid}{\sim} \mathcal{G}(a, \theta)$ où a est connu. On suppose $\pi(\theta) \equiv \mathcal{G}(c, d)$. Produisez une méthode par AR pour simuler la loi a posteriori $\pi(\theta|x_1, \dots, x_n)$ et vérifiez que les tirages obtenus sont bien issus de cette loi, par ailleurs explicite.

Remarque 9 On peut améliorer (faire baisser) le taux de rejet en encadrant la loi a posteriori entre 2 densités instrumentales (acceptation-rejet par enveloppe).

Le principe de l'AR est parfait en théorie, mais en pratique il est réservé aux cas simples ($\dim \Theta$ petite). De plus, cet algorithme est en général très coûteux en temps d'attente.

8.2.4 Algorithmes d'échantillonnage préférentiel ou d'importance (IS)

Ce type d'algorithme vise surtout à produire un estimateur consistant d'une quantité d'intérêt *a posteriori*, mais il peut être utilisé dans un but de produire un échantillonnage exact mais non indépendant de la loi-cible (approche SIR).

Algorithme IS :

1. Soit $(\theta_1, \dots, \theta_M)$ un tirage i.i.d. selon une densité instrumentale $\rho(\theta)$.
2. Soit $(\omega_1, \dots, \omega_M)$ les **poids d'importance** définis par

$$\omega_i \propto \frac{f(\mathbf{x}_n|\theta_i)\pi(\theta_i)}{\rho(\theta_i)}$$

et normalisés de façon à ce que leur somme fasse 1.

Théorème 19 Geweke 1989. Toute fonction prédictive

$$h(x|\mathbf{x}_n) = \int_{\Theta} h(x|\theta)\pi(\theta|\mathbf{x}_n) d\theta$$

(ex : $h = f$) peut être estimée de façon consistante, lorsque $M \rightarrow \infty$, par

$$\hat{h}(x|\mathbf{x}_n) = \sum_{i=1}^M \omega_i h(x|\theta_i).$$

L'approche *Sampling-Importance Resampling* (SIR), proposée par Rubin (1988), se fonde sur le résultat suivant :

Théorème 20 Les tirages

$$\tilde{\theta}_1, \dots, \tilde{\theta}_P \sim \mathcal{M}_{\text{multinomial}}(\theta_1, \dots, \theta_M | \omega_1, \dots, \omega_M)$$

suivent la loi a posteriori $\pi(\theta|\mathbf{x}_n)$.

L'heuristique de Rubin consiste à prendre $P < M/20$ pour diminuer la dépendance dans l'échantillon résé- mulé. On peut aussi ainsi estimer les caractéristiques de $\pi(\theta|\mathbf{x}_n)$ (moyenne, variance, etc.).

EXEMPLE 27. En reprenant une solution proposée pour l'exemple 27, par exemple $\rho(\theta) \equiv \mathcal{N}(\mu + \sum_{k=1}^{n-1} x_k, 1/n)$, les poids sont simplement proportionnels à

$$\omega_i \propto 1 - \Phi(y - \theta_i)$$

qu'on normalise en divisant le membre de droite par la somme des $1 - \Phi(y - \theta_i)$, $i = 1, \dots, M$. Les poids les plus hauts sont donc ceux pour lesquels $y \ll \theta_i$.

Remarquons que fort logiquement, plus la densité instrumentale $\rho(\theta)$ est "proche" de $\pi(\theta|\mathbf{x}_n)$, plus les poids sont équilibrés (donc meilleur est le rééchantillonnage). Voir figure 15 pour un résumé. Plutôt qu'une loi unique ρ , on peut mettre en place des algorithmes *adaptatifs* qui construisent itérativement une suite de densités $\{\rho_k(\theta)\}_k$ convergeant vers $\pi(\theta|\mathbf{x}_n)$, pour améliorer encore le rééchantillonnage. Une revue de tels algorithmes est proposée dans [26].

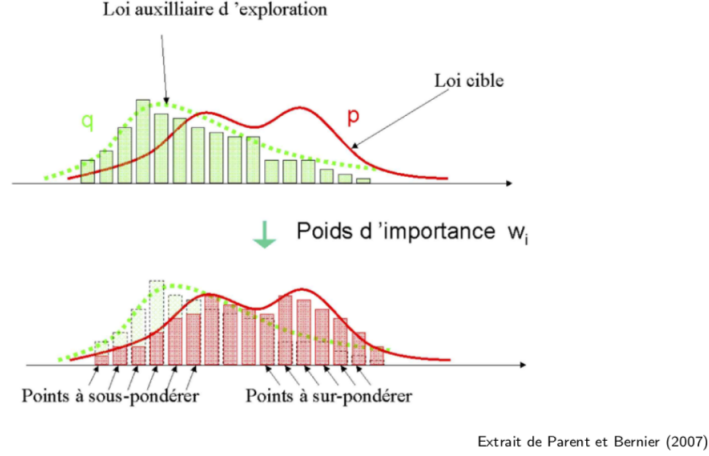


FIGURE 15 – Schématisation du principe de l'échantillonnage d'importance.

Notons le résultat suivant, important et dû à Rubinstein (1981), qui guide ces recherches d'algorithmes IS optimaux. Ce résultat n'est pas exploitable directement, car il revient à avoir déjà résolu que l'on cherche à résoudre, mais il sert à construire des lois instrumentales qui progressivement vont se rapprocher de cette optimalité.

Théorème 21 *Importance sampling optimal.* Soit l'estimateur de la fonction d'intérêt $h(\theta) \in \mathbb{R}$ par IS :

$$\hat{h}_M = \frac{1}{M} \sum_{i=1}^M \frac{\pi(\theta_i|x)}{\rho(\theta_i)} h(\theta_i) \rightarrow \mathbb{E}_\pi[h(\theta|X)] \text{ p.s.}$$

où les $\theta_i \stackrel{iid}{\sim} \rho(\theta)$. Alors le choix de ρ qui minimise la variance de l'estimateur \hat{h}_M est

$$\rho^*(\theta) = \frac{|h(\theta)|\pi(\theta|X)}{\int_{\Theta} |h(\theta)|\pi(\theta|X) d\theta}.$$

(Preuve en cours)

Un outil important, dès qu'on aborde le problème de la génération de données de même loi, mais corrélées, est la **taille d'échantillon effective**, notée en général ESS (*Effective Sample Size*). L'ESS relie la variance d'un estimateur de Monte Carlo idéal (échantillonnant directement dans la loi-cible) à la variance d'un estimateur fondé sur un échantillonnage corrélé, comme celui produit par l'IS, dans le cas où les deux estimateurs utilisent le même nombre de tirages. L'ESS mesure donc l'efficacité d'un algorithme d'échantillonnage corrélé.

Définition 42 Effective Sample Size. Soit un échantillonnage instrumental $(\theta_1, \dots, \theta_M)$ associé à des poids

d'importance normalisés $(\omega_1, \dots, \omega_M)$. Alors on définit

$$ESS = \left(\sum_{i=1}^M \omega_i^2 \right)^{-1}.$$

Lorsque l'échantillon produit est bien indépendant, les poids normalisés sont tous égaux à $1/M$, et donc $ESS = M^2/M = M$.

Remarque 10 Dans le cas des MCMC, la définition précédente nécessite d'être remaniée pour tenir compte de la corrélation dans l'échantillonnage obtenu.

Une autre propriété intéressante des techniques d'IS est de permettre de mener facilement certains types d'**analyse de sensibilité au prior**. Nous l'analysons dans l'exercice suivant.

Exercice 29 Considérons une fonction d'intérêt $h(\theta)$ que l'on cherche à résumer par un estimateur calculé sous un coût quadratique ; il s'agit donc de l'espérance a posteriori

$$h = \mathbb{E}_\pi[h(\theta)|x_1, \dots, x_n] = \int_{\Theta} h(\theta) \pi(\theta|x_1, \dots, x_n) d\theta \quad (19)$$

que l'on suppose pouvoir estimer simplement, de façon consistante, par Monte Carlo. Supposons vouloir modifier le prior : $\pi(\theta) \rightarrow \pi'(\theta)$, sans modifier le support, mais de façon à ce que la nouvelle loi a posteriori ne soit plus directement simulable. Peut-on (et sous quelles conditions) ne pas faire de calcul supplémentaire pour simuler le nouveau posterior $\pi'(\theta_1, \dots, x_n)$?

8.2.5 Méthodes de Monte Carlo par Chaînes de Markov (MCMC)

Le principe des MCMC est de partir d'un tirage d'une densité $\tilde{\pi}_0(\theta)$ arbitraire, puis de produire une *chaîne de Markov* de réalisations $\theta^{(1)}, \dots, \theta^{(M)}$ qui a pour loi **stationnaire** $\pi(\theta|\mathbf{x}_n)$.

Définition 43 Noyau de transition. Une chaîne de Markov homogène est déterminée par un noyau de transition, défini sur $\Theta \times \mathcal{B}(\Theta)$ à l'itération i par

$$\mathcal{K}(\theta|A) = P(\theta^{(i)} \in A | \theta^{(i-1)} = \theta) = \int_A \underbrace{\kappa(\theta, \tilde{\theta})}_{\text{densité de transition sur } \tilde{\theta}} d\tilde{\theta},$$

telle que $\mathcal{K}(\cdot|A)$ est mesurable $\forall A \in \mathcal{B}(\Theta)$. Cette notion de noyau généralise au cadre continu celle de matrice de transition d'un état à un autre dans un cadre discret.

Toute la structure d'une chaîne de Markov, que l'on considèrera toujours d'ordre 1 dans ce cours, dépend seulement du choix d'un noyau de transition et de l'état initial (ou la distribution initiale) de la chaîne, comme l'exprime la définition suivante.

Définition 44 Chaîne de Markov. Sachant un noyau de transition \mathcal{K} , une suite $\theta_0, \dots, \theta_n, \dots$ de variables aléatoires est une chaîne de Markov d'ordre 1 si, $\forall n \geq 0$, la distribution de θ_n conditionnelle à la σ -algèbre (filtration) générée par $\theta_{n-1}, \theta_{n-2}, \dots, \theta_0$ est la même que celle de $\theta_n | \theta_{n-1}$:

$$\pi(\theta_n \in \mathcal{A} | \theta_{n-1}, \theta_{n-2}, \dots, \theta_0) = \pi(\theta_n \in \mathcal{A} | \theta_{n-1}), = \mathcal{K}(\theta_{n-1} | \mathcal{A}).$$

Des éléments fondamentaux de théorie des chaînes de Markov sont rappelés en Annexe D. Nous en retenons surtout un résultat fondamental, le *théorème ergodique*. Celui-ci nous donne "le droit", sous certaines conditions, de mener des calculs de Monte Carlo à partir de chaînes de valeurs corrélées produites par une chaîne de Markov. On parlera alors de méthode de Monte Carlo par chaîne de Markov (MCMC).

Théorème 22 (Théorème ergodique). Si la chaîne de Markov $(\theta_n)_{n \geq 0}$ est récurrente positive^a, alors pour toute fonction $h : \Theta \rightarrow \mathbb{R}$ telle que $\mathbb{E}_\pi[|h| | x_1, \dots, x_n] < \infty$, alors

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n h(\theta_i) = \int_{\Theta} h(\theta) d\Pi(\theta | x_1, \dots, x_n).$$

Si de plus la chaîne $(\theta_n)_{n \geq 0}$ est réversible, alors

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (h(\theta_i) - \mathbb{E}_\pi[h | x_1, \dots, x_n]) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, \gamma^2)$$

avec

$$0 < \gamma^2 = \mathbb{E}_\pi[h^2(\theta_0) | x_1, \dots, x_n] + 2 \sum_{k=1}^{\infty} \mathbb{E}_\pi[h(\theta_0)h(\theta_k) | x_1, \dots, x_n] < \infty.$$

a. Plus précisément Harris-récurrente positive, cf. Annexe D.

Caractéristiques générales d'une MCMC. À l'itération i d'une MCMC, la densité de probabilité d'un θ simulé est

$$\tilde{\pi}_i(\theta) = \int_{\hat{\theta} \in \Theta} \tilde{\pi}_{i-1}(\hat{\theta}) \kappa(\hat{\theta}, \theta) d\hat{\theta}$$

et converge en loi vers une *unique densité stationnaire* $\tilde{\pi}_\infty(\theta)$, indépendamment de $\tilde{\pi}_0$, sous des conditions très générales de convergence et d'unicité :

- tout état (ou sous-ensemble) de Θ est accessible à partir de n'importe quel autre état (*irréductibilité*);
- le nombre minimal d'états intermédiaires est nul (*apériodicité*); cela se traduit, lorsque Θ est discret, par le fait que la chaîne ne peut pas boucler sur un ensemble d'états;
- l'espérance du temps de retour en n'importe quel état est fini (*récence positive*).

Les caractéristiques majeures d'une MCMC sont les suivantes (cf. figure 16) :

- le début de la chaîne (dit *temps de chauffe*) sert à explorer l'espace Θ et trouver les zones de **haute densité** *a posteriori*;
- on ne conserve que la seconde partie de l'ensemble des $\theta^{(i)}$ produits, qui suivent la distribution stationnaire (la chaîne "oublie" son état initial);
- la fréquence de visite de chaque état (ou sous-ensemble) de Θ est la même pour toute trajectoire MCMC;
- on ajoute souvent une étape de *rééchantillonnage* (SIR) ou de *décorrélation* des $\theta^{(i)}$ conservés pour obtenir un échantillon approximativement indépendant de $\tilde{\pi}_\infty(\theta)$.

Application au bayésien. Si on veut appliquer le principe des MCMC au bayésien, il faut que la loi stationnaire $\tilde{\pi}_\infty(\theta)$ soit la loi *a posteriori* $\pi(\theta | \mathbf{x}_n)$. Pour cela, le noyau \mathcal{K} doit être construit en fonction de la vraisemblance des données \mathbf{x}_n et de l'*a priori* $\pi(\theta)$. On peut prosaïquement réutiliser la structure de l'algorithme d'Acceptation-Rejet, en créant un noyau résultant du mélange de deux actions à l'itération i :

- on accepte un nouveau candidat-tirage de $\tilde{\pi}_i(\theta)$ avec une probabilité α_i ;
- on refuse et on conserve le tirage précédent dans la chaîne avec probabilité $1 - \alpha_i$.

Ce type d'algorithme a été formalisé et est connu sous le nom d'**algorithme de Hastings-Metropolis (HM)**. Sous certaines conditions de conditionnement explicite *a posteriori*, on peut accepter des candidats avec probabilité 1. Ceci donne une forme particulière de MCMC, connue sous le nom d'**algorithme de Gibbs**.

Algorithme HM. La structure de l'algorithme est la suivante.

Étape i :

1. Simuler selon une loi instrumentale $\tilde{\theta} \sim \rho(\theta|\theta^{(i-1)})$.

2. Calculer la probabilité

$$\alpha_i = \min \left\{ 1, \left(\frac{f(\mathbf{x}_n|\tilde{\theta})\pi(\tilde{\theta})}{f(\mathbf{x}_n|\theta^{(i-1)})\pi(\theta^{(i-1)})} \right) \cdot \left(\frac{\rho(\theta^{(i-1)}|\tilde{\theta})}{\rho(\tilde{\theta}|\theta^{(i-1)})} \right) \right\}.$$

3.
$$\left. \begin{array}{l} \text{Simuler } U \sim \mathcal{U}_{\text{unif}}[0, 1]. \\ \text{Si } U \leq \alpha_i \text{ choisir } \theta^{(i)} = \tilde{\theta}. \\ \text{Sinon choisir } \theta^{(i)} = \theta^{(i-1)}. \end{array} \right\} \text{ Accepter } \tilde{\theta} \text{ avec probabilité } \alpha_i.$$

Le noyau markovien est alors constitué d'un mélange d'un Dirac en $\theta^{(i-1)}$ et de la loi instrumentale, mélange pondéré par la probabilité de transition α_i . La partie continue du noyau de transition (de θ vers θ') s'écrit

$$p(\theta, \theta') = \alpha(\theta, \theta')\rho(\theta'|\theta)$$

avec $\alpha(\theta, \theta')$ la probabilité de transition

$$\alpha(\theta, \theta') = \min \left\{ 1, \left(\frac{f(\mathbf{x}_n|\theta')\pi(\theta')}{f(\mathbf{x}_n|\theta)\pi(\theta)} \right) \cdot \left(\frac{\rho(\theta|\theta')}{\rho(\theta'|\theta)} \right) \right\}.$$

On a

$$\pi(\theta) \times p(\theta, \theta') = \pi(\theta') \times p(\theta', \theta).$$

La chaîne MCMC produite est alors dite *réversible* et ceci suffit à montrer, si la chaîne est irréductible et apériodique, que :

- elle est ergodique ;
- la distribution des itérés $\theta^{(i)}, \dots, \theta^{(j)}$ de la chaîne converge en loi vers une loi-limite unique ;
- celle-ci est proportionnelle à $f(\mathbf{x}_n|\theta)\pi(\theta)$: il s'agit donc de $\pi(\theta|\mathbf{x}_n)$.

L'irréductibilité peut être facilement assurée par la contrainte suivante : le support de la loi instrumentale doit inclure le support de la loi-cible.

Le rapport de Metropolis fait intervenir le rapport des lois *a posteriori* (qui permet d'ôter la constante d'intégration inconnue) : à l'étape k , si $\tilde{\theta}_k$ se situe plus haut dans la zone de haute densité de $\pi(\theta|\mathbf{x}_n)$ que θ_{k-1} , ce rapport est plus grand que 1. Le rapport inverse des lois instrumentales et l'usage d'un tirage uniforme interdisent d'automatiser l'acceptation de ce nouveau point, en permettant à la chaîne de Markov d'explorer exhaustivement l'espace Θ .

Exercice 30 Soit X la variable "débit maximal de rivière". Elle est supposée suivre une loi des extrêmes (Gumbel) de densité

$$f(x|\theta) = \lambda\mu \exp(-\lambda x) \exp(-\mu \exp(-\lambda x)).$$

avec $\theta = (\mu, \lambda)$.



Considérons n observations $\mathbf{x}_n = (x_1, \dots, x_n)$ supposées iid selon cette distribution.

1. Comment s'écrit la vraisemblance ?
2. On considère l'a priori $\pi(\mu, \lambda) = \pi(\mu|\lambda)\pi(\lambda)$ avec

$$\begin{aligned}\mu|\lambda &\sim \mathcal{G}(m, b_m(\lambda)), \\ \lambda &\sim \mathcal{G}(m, m/\lambda_e)\end{aligned}$$

et $b_m(\lambda) = [\alpha^{-1/m} - 1]^{-1} \exp(-\lambda x_{e,\alpha})$. Ces hyperparamètres ont le sens suivant :

- $x_{e,\alpha}$ = quantile prédictif a priori d'ordre α :

$$P(X < x_{e,\alpha}) = \int P(X < x_{e,\alpha} | \mu, \lambda) \pi(\mu, \lambda) d\mu d\lambda = \alpha;$$

- m = taille d'échantillon fictif, associée à la "force" de la connaissance a priori $x_{e,\alpha}$;
- $1/\lambda_e$ = moyenne de cet échantillon fictif.

Pouvez-vous produire un algorithme de type MCMC qui permette de générer une loi jointe a posteriori pour (μ, λ) ?

Heuristique de progression du taux d'acceptation moyen α . Dans une chaîne MCMC produite par HM, la *stationarité* est l'atteinte par une chaîne d'un tirage stationnaire dans la loi *a posteriori*. La rapidité de convergence vers la stationarité est induite par le taux d'acceptation α . Au début de la MCMC, on cherche à *explorer l'espace* : α grand ($\simeq 0.5$). Si α est petit, la simulation est fortement dépendante du passé de la chaîne : l'exploration de l'espace est très lente. Si α reste grand, chaque chaîne évolue solitairement et elles risquent de se mélanger lentement. De différents travaux appliqués et théoriques, on a tiré une règle du pouce : un $\alpha = 0.25$ est souvent considéré, en pratique (en particulier lorsque $\dim \Theta$ est grande) comme un bon objectif de renouvellement à la stationnarité. Par ailleurs, la calibration de $\rho(\theta|\theta^{(i-1)})$ (en général, le choix de sa variance) peut être en général faite de façon **empirique** en "testant" le taux d'acceptation effectif.

Heuristique de choix d'une loi instrumentale $\rho(\theta|...)$. Dans le cas le plus simple, on choisit volontiers $\rho(\theta|\theta^{(i-1)}) = \rho(\theta)$ (*loi statique*). Mais une modélisation standard est de choisir $\rho(\theta|\theta^{(i-1)})$ centrée sur $\theta^{(i-1)}$, et donc seule la variance doit être calibrée (ou le coefficient de variation).

EXEMPLE 28. Marche aléatoire $\theta \sim \theta^{(i-1)} + \sigma \epsilon_i$ où $\epsilon_i \sim \mathcal{N}(0, 1)$.

À la différence du noyau, le caractère markovien de ρ peut être relâché : on peut construire des ρ **adaptatives** en utilisant tout le passé de la chaîne et non pas le dernier état connu $\theta^{(i-1)}$. Il existe une très vaste littérature à ce sujet, plutôt du domaine de la recherche que de la règle du pouce ou la “boîte à outils”. Là encore, on conseille de se référer à l’ouvrage [26].

Arrêt de chaîne MCMC. Une fois que le *temps de chauffe* est passée, la *phase ergodique* est atteinte. De nombreux diagnostics de convergence vers la stationarité ont été proposés [7] et *nécessitent d’avoir lancé plusieurs chaînes parallèles*. À la stationarité, ces chaînes parallèles se sont bien mélangées et ont “oublié” le passé de chacune. Les diagnostics sont surtout *visuels* : on regarde l’évolution du comportement d’une statistique informant sur la stabilité de la distribution des θ .

Parmi ces diagnostics de convergence, les statistiques de Gelman-Rubin (1992, cas 1D) et de Brooks-Gelman (1998, cas multidimensionnel) sont très standards : elles sont fondées sur la comparaison de variances inter et intra chaînes.

Définition 45 Diagnostic de Gelman-Rubin.

- Soit P trajectoires (chaînes) parallèles de longueur n (en pratique, $P = 3$.)
- Soit $\theta_k^{(i)}$ la $i^{\text{ème}}$ réalisation sur la trajectoire k .
- Soit B l’estimateur de la variance de θ inter-chaînes.

$$B = \frac{n}{P-1} \sum_{k=1}^P (\bar{\theta}_k - \bar{\theta})^2$$

avec

$$\bar{\theta}_k = \frac{1}{n} \sum_{i=1}^n \theta_k^{(i)} \quad \text{et} \quad \bar{\theta} = \frac{1}{P} \sum_{k=1}^P \bar{\theta}_k$$

- soit W l’estimateur de la variance de θ intra-chaînes (**ergodique**)

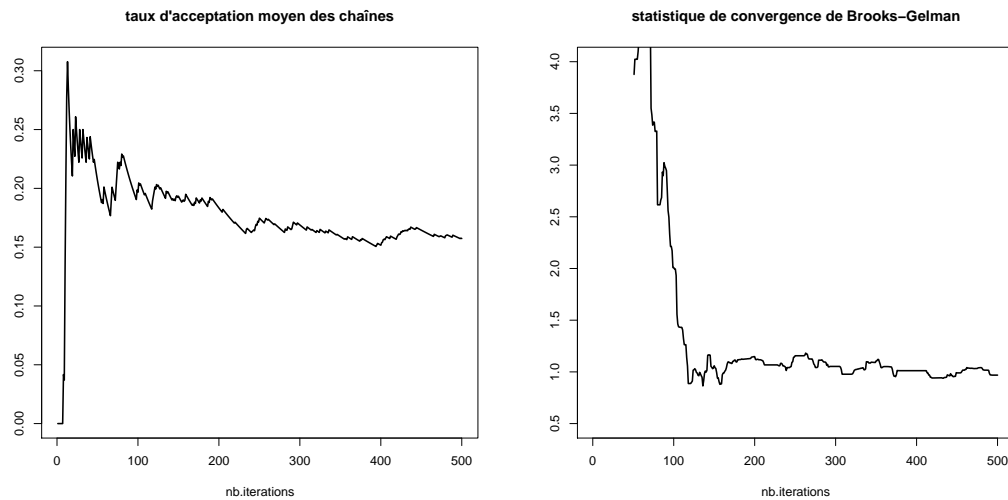
$$W = \frac{1}{P} \sum_{k=1}^P \left[\frac{1}{n-1} \sum_{i=1}^n (\theta_k^{(i)} - \bar{\theta}_k)^2 \right]$$

Alors, le rapport (statistique de Gelman-Rubin)

$$R = \frac{\frac{(n-1)}{n} W + \frac{1}{n} B}{W}$$

tends vers 1 par valeurs supérieures.

En reprenant l’exemple 30, on représente ci-dessous des exemples de trajectoires du taux d’acceptation et du diagnostic de Brooks-Gelman, qui généralise Gelman-Rubin en tenant compte de la covariance entre dimensions d’une même chaîne.



Décorrélation d'un échantillon MCMC. Soit M_c le nombre d'itérations d'une MCMC avant qu'on atteigne la stationnarité (*temps de chauffe*, cf. figure 16). En sortie de la MCMC, on obtient un échantillon de $M - M_c$ vecteurs $\theta^{(M-M_c+1)}, \dots, \theta^M$ qui suivent la loi stationnaire $\pi(\theta|\mathbf{x}_n)$.

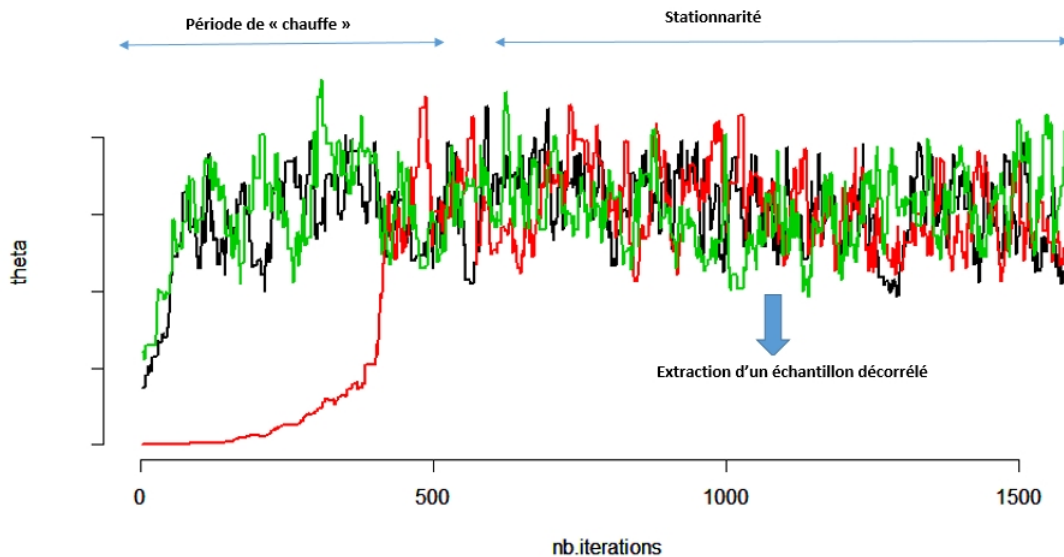


FIGURE 16 – Trois chaînes MCMC convergeant en parallèle vers la même loi-cible *a posteriori*. La période de *chauffe* correspond au nombre d'itérations nécessaire au mélange des chaînes, qui est un indicateur de l'atteinte de la stationnarité.

De par le caractère markovien de la MCMC, ces valeurs peuvent être très dépendantes le long d'une chaîne. Si on dispose de beaucoup de chaînes parallèles indépendantes, il suffit de prélever une valeur dans chacune... (mais c'est peu faisable en pratique).

Pour obtenir un échantillon décorrélé (qui offre une meilleure information sur les caractéristiques de la loi $\pi(\theta|\mathbf{x}_n)$), une bonne façon de faire repose sur l'usage de l'*autocorrélation*, de façon similaire au traitement des

séries temporelles.

On peut procéder comme suit :

Procédure de décorrélation.

1. On estime l'autocorrélation des éléments d'une chaîne :

$$\text{Aut}_{i,i+j} = \frac{\mathbb{E}[(\theta^{(i)} - \mathbb{E}[\theta|\mathbf{x}_n])(\theta^{(i+j)} - \mathbb{E}[\theta|\mathbf{x}_n])]}{\text{Var}[\theta|\mathbf{x}_n]}$$

à valeur dans $[-1, 1]$.

- À i fixé, $\text{Aut}_{i,i+j}$ tends vers 0 lorsque j augmente $\Leftrightarrow \theta^{(i+j)}$ devient de plus en plus décorrélé de $\theta^{(i)}$.
- On considère que cette décorrélation est effective lorsque l'estimateur de $\text{Aut}_{i,i+j}$ est un **bruit blanc gaussien**.
- On peut donc, en moyenne sur les i , estimer le nombre d'itérations nécessaire t pour obtenir 2 valeurs décorrélées de θ .

2. Sur chaque chaîne, on sélectionne le sous-échantillon (*thinning*)

$$\theta^{(M-P+1)}, \theta^{M-P+1+t}, \theta^{M-P+1+2t}, \dots$$

3. On baisse encore la dépendance des éléments de l'échantillon final en prélevant dans les chaînes indépendantes.

8.2.6 Échantillonneur de Gibbs et approches hybrides

Dans un cas où le paramètre $\theta = (\theta_1, \dots, \theta_d)$ est multidimensionnel (ce qui concerne notamment pour les modèles hiérarchiques), il est recommandé d'utiliser une algorithmique d'échantillonnage de Gibbs [39] qui tire parti du principe de *cascade rule* (§ 8.2.2).

Étant donné un échantillon \mathbf{x}_n , supposons disposer des lois *a posteriori* conditionnelles

$$\begin{aligned} \theta_1 | \theta_2, \dots, \theta_d, \mathbf{x}_n &\sim \pi(\theta_1 | \theta_2, \dots, \theta_d, \mathbf{x}_n), \\ \theta_2 | \theta_1, \theta_3, \dots, \theta_d, \mathbf{x}_n &\sim \pi(\theta_2 | \theta_1, \theta_3, \dots, \theta_d, \mathbf{x}_n), \\ &\dots \quad \dots \end{aligned}$$

Alors la chaîne de Markov de vecteurs

$$\theta^{(1)} = \begin{pmatrix} \theta_1^{(1)} \\ \dots \\ \theta_d^{(1)} \end{pmatrix}, \quad \theta^{(i)} = \begin{pmatrix} \theta_1^{(M)} \\ \dots \\ \theta_d^{(M)} \end{pmatrix}, \dots$$

produite par la simulation conditionnelle itérée converge également vers la loi *a posteriori*-cible $\pi(\theta|\mathbf{x}_n)$, sous des conditions très générales.

Lorsque les lois *a posteriori* conditionnelles ne sont elles-mêmes pas complètement explicites, une démarche hybride, dite de *Metropolis-Hastings-within-Gibbs*, consiste à faire appel à un test de Metropolis pour chaque dimension.

Remarque 11 La solution proposée pour l'exercice 30 est d'ailleurs un exemple d'algorithme de Gibbs en dimension 2, incluant une étape de Metropolis-Hastings pour l'estimation du paramètre λ .

L'algorithmique de Gibbs (hybride), qui exploite au maximum la structure conditionnelle des modèles hiérarchiques, est donc très générale, et elle est notamment particulièrement intéressante lorsqu'on traite des problèmes à données manquantes (ex : présentant des données censurées, ou des variables latentes comme les modèles de mélange...). Dans un cadre bayésien, il permet de considérer ces données comme des paramètres inconnus à simuler (*augmentation de données*). L'exemple suivant illustre ce procédé.

Exercice 31 (Retour à l'exercice 27). On suppose de nouveau connaître un échantillon $\mathbf{x}_n \sim \mathcal{N}(\theta, 1)$ composé de quelques observations x_1, \dots, x_{n-1} supposées iid de loi $\mathcal{N}(\theta, 1)$, et d'une pseudo-observation y qui est un cas-limite masquant (censurant) une observation x_n qui aurait dû être faite : $y < x_n$. On considère toujours $\theta \sim \mathcal{N}(\mu, 1)$ a priori. Pouvez-vous produire un algorithme d'échantillonnage par Gibbs qui génère des réalisations de la loi a posteriori de θ ?

La modélisation bayésienne par conditionnement peut fréquemment entraîner le mécanisme suivant :

1. On construit un *a priori* hiérarchique

$$\pi(\theta) = \pi(\theta_1 | \theta_2, \theta_3) \pi(\theta_2 | \theta_3) \pi(\theta_3)$$

avec des *a priori* non-informatifs

2. Ce conditionnement est souvent choisit pour tirer parti de *conjugaisons a posteriori* : les lois conditionnelles

$$\pi(\theta_1 | \mathbf{x}_n, \theta_2, \theta_3),$$

$$\pi(\theta_2 | \mathbf{x}_n, \theta_1, \theta_3),$$

$$\pi(\theta_3 | \mathbf{x}_n, \theta_1, \theta_2)$$

sont explicites, ce qui permet d'utiliser un algorithme de Gibbs.

Toutefois, même si ces lois *a posteriori conditionnelles* sont propres, la loi *jointe* peut ne pas l'être :

$$\int_{\Theta} \pi(\theta | \mathbf{x}_n) d\theta = \infty.$$

Il est donc indispensable de vérifier la *propriété a posteriori* avant de mettre en oeuvre un algorithme de Gibbs.

Exercice 32 Modèle à effets aléatoires autour d'une constante (Hobert-Casella). Pour $i = 1, \dots, I$ et $j = 1, \dots, J$, on considère

$$x_{ij} = \beta + u_i + \epsilon_{ij}$$

où $u_i \sim \mathcal{N}(0, \sigma^2)$ et $\epsilon_{ij} \sim \mathcal{N}(0, \tau^2)$. Ce type de modèle permet de représenter la distribution d'une caractéristique au sein d'une population, où β est une tendance moyenne, u_i correspond à une variation d'un groupe et ϵ_{ij} à une variation au sein d'un sous-groupe. On suppose choisir

$$\pi(\beta, \sigma^2, \tau^2) \propto \frac{1}{\sigma^2 \tau^2} \quad (\text{prior de Jeffreys}).$$

On note \mathbf{x}_{IJ} l'échantillon des données observées, \bar{x}_i la moyenne sur les j . On note \mathbf{u}_I l'échantillon manquant des u_1, \dots, u_I (reconstitué dans l'inférence).

1. Calculer les lois conditionnelles *a posteriori* de

$$U_i | \mathbf{x}_{IJ}, \beta, \sigma^2, \tau^2$$

$$\beta | \mathbf{x}_{IJ}, \sigma^2, \tau^2, \mathbf{u}_I$$

$$\sigma^2 | \mathbf{x}_{IJ}, \beta, \tau^2, \mathbf{u}_I$$

$$\tau^2 | \mathbf{x}_{IJ}, \beta, \sigma^2, \mathbf{u}_I$$

Ces lois sont-elles bien définies ?

2. Donner une formule (à un coefficient proportionnel près) pour la loi a posteriori jointe $\pi(\sigma^2, \tau^2 | \mathbf{x}_{\mathbf{IJ}})$. Comment se comporte-t-elle au voisinage de $\sigma = 0$, pour $\tau \neq 0$? Que pouvez-vous en déduire?
3. Mettre en place un algorithme de Gibbs permettant d'inférer sur $(\beta, \sigma^2, \tau^2)$. Que pouvez-vous dire sur la convergence des chaînes MCMC?

8.2.7 Un résumé de ces premières méthodes

On trouvera sur le tableau ci-dessous quelques informations qui résument l'usage et les particularités de ces algorithmes de calcul bayésien. Pour aider à faire un choix, voici quelques conseils reposant sur quelques questions essentielles.

- La loi-cible *a posteriori* est-elle proche d'un cas explicite (ou conjuguée) (ex : loi normale censurée) ?
- Si oui, que faudrait-il faire (typiquement, *simuler des données manquantes* \Rightarrow Gibbs)
- En multidimensionnel, a-t-on des propriétés de conjugaison conditionnelles (\Rightarrow Gibbs) ?
- Si nous n'avons aucune idée, peut-on trouver une loi $\rho(\theta)$ partageant certaines caractéristiques avec $\pi(\theta|\mathbf{x}_n)$?
 - Par exemple, on peut tenter de tracer $\pi(\theta|\mathbf{x}_n)$ (à un coefficient près) dans les cas unidimensionnels.
 - On peut également mener un calcul du mode *a posteriori* $\hat{\theta} = \text{maximum de la vraisemblance pondérée par l'a priori}$.
 - une idée de loi instrumentale typique est ainsi une loi $\rho(\theta)$ gaussienne $\mathcal{N}(\hat{\theta}, \sigma^2 I_d)$ avec σ calibré empiriquement.

	Acceptation - Rejet	Échant. d'importance	Métropolis - Hastings (MH)	Gibbs
<i>Contexte</i>				
Dimension de θ	1	grande	grande	grande
Nature simulation	iid	non-indep.*	non-indep. approx.*	non-indep. approx.*
Nature algo	itératif	statique	itératif	itératif
<i>Nb. itérations typ.</i>	quelques centaines	1	quelques dizaines de milliers	quelques milliers
<i>Implémentation</i>	aisée	aisée	calibration fine de $\rho(\theta)$ nécessaire	aisée
<i>Critère d'arrêt</i>	aucun	aucun	nécessaire	nécessaire
Risques	fort taux de rejet	poids mal équilibrés	mauvais mélange chaînes //	nécessite souvent couplage avec M-H
Temps de calcul	long	rapide	long	plutôt rapide

8.3 Méthodes d'échantillonnage accélérées

Les méthodes MCMC ont été et restent très utilisées pour simuler des valeurs *a posteriori* $\theta_1, \dots, \theta_n \stackrel{iid}{\sim} \pi(\theta|x) \propto f(x|\theta)\pi(\theta)$. Elles sont toutefois lentes pour des problèmes complexes, et la convergence vers la loi-limite (*a posteriori* dans un cadre bayésien) reste difficile à vérifier. En particulier, les *méthodes MCMC restent très lourdes* (voire impossibles) à utiliser dans les problèmes de type "Big Data".

EXEMPLE 29. *Modèles à états latents, tels les modèles de Markov cachés permettant de segmenter des images structurées de façon non supervisée...*

Ainsi, des avancées importantes pour **accélérer les algorithmes d'échantillonnage** ont été menées :

- en se fondant sur des techniques de *réduction de variance* ;
- en démarrant les chaînes de Markov *au plus proche du mode a posteriori* ;
- en *améliorant le choix des lois instrumentales*, inspirées par la physique énergétique (dynamique, particules, quantique...).

Nous en listons ci-dessous quelques-unes.

Soulignons cependant que dans un cadre d'usage de plus en plus soumis aux contraintes de grande dimensionnalité et d'apprentissage massif, dans la réalité des faits, on calcule rarement la distribution *posteriori*. On préfère se simplifier la vie et les calculs en approximant la loi distribution *posteriori* par une valeur unique (blackmode *a posteriori*). Il s'agit ni plus ni moins d'une *maximisation de vraisemblance pénalisée* (voir § 2.5 pour plus de détails). Dans ce contexte, les *techniques d'optimisation*, comme l'algorithme EM ou le Gradient Boosting, ont encore de beaux jours devant elle, et il est possible d'interpréter ces algorithmes comme des approches "dégénératives" du cadre MCMC.

8.3.1 Réduction de variance par utilisation des corrélations négatives (variables antithétiques)

Considérons deux échantillons iid $(\theta_1, \dots, \theta_m)$ et $(\tilde{\theta}_1, \dots, \tilde{\theta}_m)$ suivant $\pi(\theta|D)$ pour estimer

$$\mathcal{H} = \int_{\mathbb{R}} h(\theta) \pi(\theta|D) d\theta.$$

Soient

$$\hat{\mathcal{H}}_1 = \frac{1}{m} \sum_{i=1}^m h(\theta_i) \quad \text{et} \quad \hat{\mathcal{H}}_2 = \frac{1}{m} \sum_{i=1}^m h(\tilde{\theta}_i)$$

de moyenne m et de variance σ^2 . Alors, la variance de la moyenne vaut

$$\mathbb{V} \left[\frac{1}{2} (\hat{\mathcal{H}}_1 + \hat{\mathcal{H}}_2) \right] = \frac{\sigma^2}{2} + \frac{1}{2} \text{Cov}(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2)$$

Par conséquent, si les deux échantillons sont *négativement corrélés*,

$$\text{Cov}(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2) \leq 0$$

ils font mieux que deux échantillons de même taille. Ce type d'approche a été formalisé sous le nom de *réduction de variance par variables antithétiques* :

1. Si $\pi(\theta|D)$ symétrique autour de μ , prendre $\theta_i = 2\mu - \theta_i$;
2. Si $\theta_i = F^{-1}(U_i)$, prendre $Y_i = F^{-1}(1 - U_i)$;
3. Si $(A_i)_i$ est une partition de Θ , produire un échantillonnage *partitionné* (ou *stratifié*) en choisissant des θ_j dans chaque A_i (nécessite de connaître $\Pi(A_i)$).

8.3.2 Variables de contrôle

Soit

$$\mathcal{H} = \int_{\mathbb{R}} h(\theta) \pi(\theta|D) d\theta$$

à évaluer et

$$\mathcal{H}_0 = \int_{\mathbb{R}} h_0(\theta) \pi(\theta|D) d\theta$$

connue. On estime quand même \mathcal{H}_0 par $\hat{\mathcal{H}}_0$ (et \mathcal{H} par $\hat{\mathcal{H}}$). Dans ce cas, on peut définir un **estimateur combiné** :

$$\hat{\mathcal{H}}^* = \hat{\mathcal{H}} + \beta (\hat{\mathcal{H}}_0 - \mathcal{H}_0)$$

Proposition 12 $\hat{\mathcal{H}}^*$ est sans biais pour \mathcal{H} et

$$\mathbb{V} [\hat{\mathcal{H}}^*] = \mathbb{V} [\hat{\mathcal{H}}] + \beta^2 \mathbb{V} [\hat{\mathcal{H}}_0] + 2\beta \text{Cov} (\hat{\mathcal{H}}, \hat{\mathcal{H}}_0).$$

On peut alors faire un choix optimal de β , qui permet de diminuer la variance de l'estimateur :

$$\beta^* = - \frac{\text{Cov} (\hat{\mathcal{H}}, \hat{\mathcal{H}}_0)}{\mathbb{V} [\hat{\mathcal{H}}_0]}$$

avec

$$\mathbb{V} [\hat{\mathcal{H}}^*] = \mathbb{V} [\hat{\mathcal{H}}] (1 - \rho^2)$$

où ρ^2 est le coefficient de corrélation entre $\hat{\mathcal{H}}$ et $\hat{\mathcal{H}}_0$.

EXEMPLE 30. **Approximation de quantile.** Supposons vouloir évaluer

$$q = \Pi(\theta > a|D) = \int_a^\infty \pi(\theta|D) d\theta$$

par

$$\hat{q} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\theta_i > a\}}, \quad \theta_i \stackrel{iid}{\sim} \pi(\theta|D)$$

avec $\Pi(\theta > \mu|D) = 1/2$. La variable de contrôle

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\theta_i > a\}} + \beta \left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\theta_i > \mu\}} - \Pi(\theta > \mu|D) \right)$$

améliore \hat{q} si

$$\beta < 0 \quad \text{et} \quad |\beta| < 2 \frac{\Pi(\theta > a|D)}{\Pi(\theta > \mu|D)}.$$

8.3.3 Rao-Blackwellisation

On peut également vouloir tirer parti de l'inégalité

$$\mathbb{V} [\mathbb{E}[h(\theta)|Y]] \leq \mathbb{V} [h(\theta)],$$

aussi appelée *Théorème de Rao-Blackwell*, pour diminuer la variance d'un estimateur de Monte Carlo.

Proposition 13 Si $\hat{\mathcal{H}}$ est un estimateur sans biais de $\mathcal{H} = \mathbb{E}_\pi [h(\theta)]$, avec θ simulé à partir de la densité jointe $\tilde{\pi}(\theta, y)$ où

$$\int \tilde{\pi}(\theta, y) dy = \pi(\theta)$$

alors l'estimateur

$$\hat{\mathcal{H}}^* = \mathbb{E}_{\tilde{\pi}} [\hat{\mathcal{H}} | Y_1, \dots, Y_n]$$

domine $\hat{\mathcal{H}}(\theta_1, \dots, \theta_n)$ en variance (et est aussi sans biais).

L'algorithme EM, notamment, tire parti de ce procédé pour maximiser une vraisemblance de données incluant des données manquantes.

8.3.4 Amélioration de lois instrumentales par dynamique de Langevin / hamiltonienne

Les améliorations de lois instrumentales au sein des algorithmes de type MCMC sont généralement fondées sur des parallèles avec des problèmes de *diffusion / mouvement de système physique* (particulaire), régis par des équations aux dérivées partielles (EDP). L'idée générale est d'utiliser de l'information sur la densité cible au travers du passé de l'algorithme, et plus précisément l'information contenue dans le *gradient* $\nabla \log \pi(\theta|D)$. Le cadre général de la dynamique particulière de Langevin offre une formalisation de ce problème. Par la suite, pour simplifier les notations, on notera D les données disponibles.

Proposition 14 Metropolis-Adjusted Langevin Algorithm (MALA) *Discretisation d'une diffusion de Langevin de probabilité stationnaire $\pi(\theta|D)$:*

$$\rho(\theta, \cdot) = \mathcal{N}_d \left(\theta + \frac{h^2}{2} \nabla \log \pi(\theta|D), h^2 I_d \right)$$

pour un pas $h > 0$. Le choix de h relève de considérations sur la structure de $\text{Cov}(\theta|D)$.

De façon plus explicite, on veut simuler selon une loi non normalisée définie en termes de potentiel

$$\pi(\theta|D) \propto \exp(-U(\theta)).$$

Sous des conditions peu contraignantes, cette densité est l'unique mesure de probabilité invariante d'une *équation différentielle stochastique de Langevin*

$$d\theta_t = -(\theta_t)dt + \sqrt{2}dB_t$$

avec B_t un processus brownien (qui décrit le mouvement d'une particule soumise à une infinité de chocs en des temps très courts)

$$B_t \sim \mathcal{N}(0, t).$$

En général on ne peut résoudre exactement l'équation précédente : on s'appuie alors sur l'approximation produite par une discrétisation d'Euler-Maruyama

$$\theta_{k+1} = \theta_k - \gamma \nabla U(\theta_k) + \sqrt{2\gamma} Z_{k+1}.$$

Cette approche est l'*Unadjusted Langevin Algorithm* (ULA) et s'appelle aussi *Langevin Monte Carlo* (LMC). Il s'agit simplement d'un calcul de MAP par un algorithme de descente de gradient avec du bruit ajouté à chaque itération. Des résultats théoriques ont été obtenus récemment pour contrôler l'erreur d'approximation / taille d'échantillon et dimension (si $\pi(\theta|D)$ régulière), par Dumus et Moulines (2017). Cette approche marche bien pour le traitement de la grande dimension de θ en inférence bayésienne.

Toutefois, cette discrétisation induit du biais, qui peut être ôté par une étape d'acceptation-rejet de Metropolis-Hastings, ce qui donne la formulation MALA de l'algorithme. L'algorithmique MALA hérite donc des bonnes propriétés de convergence de ULA et affronte bien la grande dimension. Dans un tour d'horizon récent, Nemeth et Fearnhead (2019) ont récemment montré que le pas optimal pour MALA est grand, mais MALA coûte plus cher par itération que les approches ULA. Par ailleurs, les lois instrumentales fondées sur ULA ont un meilleur taux d'acceptation que les marches aléatoires MALA.

Approche MYULA. Quand $\pi(\theta|D)$ n'est pas régulière, on suppose que le potentiel peut s'écrire

$$U(\theta) = f(\theta) + g(\theta)$$

où f est convexe, continûment différentiable, et de gradient lipschitzien ; g est propre, convexe, et semi-continue à gauche. On peut alors remplacer g par son enveloppe dite de Moreau-Yosida (voir figure 17) :

$$\begin{aligned} g^\lambda(x) &= \min_{y \in \mathbb{R}^d} \{g(y) + (2\lambda)^{-1}\|x - y\|^2\} \\ \nabla g^\lambda(x) &= \lambda^{-1} (x - \text{prox}_g^\lambda(x)) \end{aligned}$$

avec

$$\text{prox}_g^\lambda(x) = \arg \min_{y \in \mathbb{R}^d} \{g(y) + (2\lambda)^{-1}\|x - y\|^2\}.$$

L'algorithme de *Moreau-Yosida ULA* (MYULA) s'écrit alors

$$\theta_{k+1} = \left(1 - \frac{\gamma}{\lambda}\right) \theta_k - \gamma \nabla f(\theta_k) + \frac{\gamma}{\lambda} \text{prox}_g^\lambda(\theta_k) + \sqrt{2\gamma} Z_{k+1}.$$

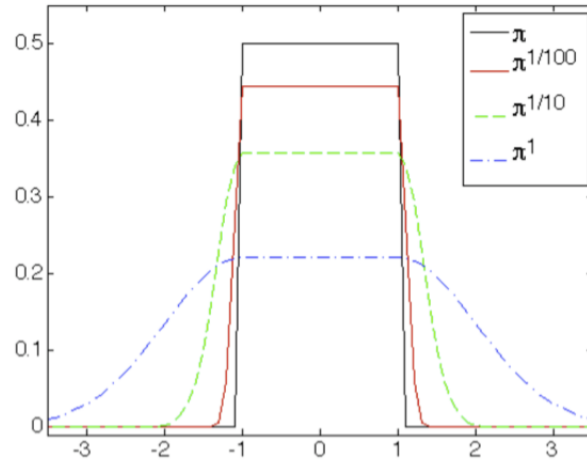


FIGURE 17 – Illustration issue de Nemeth et Fearnhead (2019).

En général, avec $\pi(\theta|D) \propto \pi(\theta) \prod_{i=1}^n f(x_i|\theta)$, on a

$$U = \sum_{i=0}^n U_i$$

avec $U_0(\theta) = -\log \pi(\theta)$ et $U_i(\theta) = -\log f(x_i|\theta)$. Le problème posé par cet algorithme est qu'une seule itération reste coûteuse. L'idée est d'utiliser des approches de *descente de gradient stochastique* (SGD) pour accélérer le calcul.

Descente de gradient stochastique (SGD). La version "Langevin" du SGD (*Stochastic Gradient Langevin Dynamics*, ou SGLD) s'écrit simplement

$$\theta_{k+1} = \theta_k - \gamma \left(\nabla U_0(\theta_k) + \frac{n}{p} \sum_{i \in S_{k+1}} \nabla U_i(\theta_k) \right) + \sqrt{2\gamma} Z_{k+1}.$$

Cette approche obtient un très faible coût de calcul par itération si $p \ll n$. Il s'agit donc, en résumé, d'un calcul de MAP par un algorithme de descente de gradient avec du bruit ajouté à chaque itération, à l'image de ce qui est massivement utilisé en apprentissage statistique. Des résultats théoriques prouvant la bonne convergence de ce type d'algorithme ont été obtenus récemment par Brosse *et al.* (2018).

Autres variantes. On peut hybrider le SGLD avec une approche par *variables de contrôle* si on peut estimer le mode de $\pi(\theta|D)$. De telles approches se nomment *SGLD Control Variate* (SGLDCV) ou *SGLD Fixed Point* (SGLDFP). Elles nécessitent très souvent de lancer un premier SGLD pour estimer le mode, puis opèrent un raffinement. D'autres approches s'hybridant avec des techniques d'échantillonnage d'importance (ou stratifié, par exemple) existent également.

Méthodes de Monte Carlo Hamiltoniennes (HMC). Les méthodes HMC, disponibles sous Python 3 (PYMC¹⁴) ou STAN¹⁵, sont des méthodes MCMC issues d'un parallèle entre deux problèmes :

- produire une dynamique pour la chaîne de Markov $\theta_i, \dots, \theta_j$ s'approchant de la loi visée $\pi(\theta|\mathbf{x}_n)$;
- prévoir le mouvement d'un système physique soumis à une *dynamique hamiltonienne*.

La dynamique hamiltonienne est une reformulation de la mécanique newtonienne selon laquelle la description d'un système est faite au travers de *coordonnées* (généralisées) et de *momentum* (quantité de mouvement), reliées par un **lagrangien** ; celui-ci exprime une différence entre énergie cinétique et énergie potentielle. Le lagrangien est une fonction des variables dynamiques qui permet d'écrire les équations de mouvement. La transformée de Legendre de ce lagrangien est nommé *hamiltonien*. Le tableau ci-dessous résume la correspondance algébrique entre les deux problèmes mentionnés :

Système physique	Applications MCMC
Position	θ
Énergie potentielle	$-\log \pi(\theta D)$
Momentum	Variables introduites artificiellement (variables gaussiennes en général)

L'idée générale de l'application aux MCMC est la suivante : à chaque pas de la MCMC, on met à jour les momentum et on produit une nouvelle trajectoire-candidate pour θ (et non un simple tirage) suivant une dynamique hamiltonienne (approche *leapfrog*).

Un récapitulatif rapide des méthodes fondées sur un parallèle avec des modèles de dynamique est présenté dans le résumé ci-dessous :

$$\zeta_{t+h} \approx \zeta_t - \frac{h}{2} [(\mathbf{D}(\zeta_t) + \mathbf{Q}(\zeta_t))\nabla H(\zeta_t) + \Gamma(\zeta_t)] + \sqrt{h}\mathbf{Z}$$

Algorithm	ζ	$H(\zeta)$	$\mathbf{D}(\zeta)$	$\mathbf{Q}(\zeta)$	
SGLD	θ	$U(\theta)$	\mathbf{I}	$\mathbf{0}$	
SG-RLD	θ	$U(\theta)$	$G(\theta)^{-1}$	$\mathbf{0}$	Riemannian SGLD (Fisher)
SG-HMC	(θ, ρ)	$U(\theta) + \frac{1}{2}\rho^\top \rho$	$\begin{pmatrix} 0 & 0 \\ 0 & \mathbf{C} \end{pmatrix}$	$\begin{pmatrix} 0 & -\mathbf{I} \\ \mathbf{I} & 0 \end{pmatrix}$	Hamiltonian MC
SG-RHMC	(θ, ρ)	$U(\theta) + \frac{1}{2}\rho^\top \rho$	$\begin{pmatrix} 0 & 0 \\ 0 & G(\theta)^{-1} \end{pmatrix}$	$\begin{pmatrix} 0 & -G(\theta)^{-1/2} \\ G(\theta)^{-1/2} & 0 \end{pmatrix}$	Riemannian Hamiltonian MC
SG-NHT	(θ, ρ, η)	$U(\theta) + \frac{1}{2}\rho^\top \rho + \frac{1}{2d}(\eta - A)^2$	$\begin{pmatrix} 0 & 0 & 0 \\ 0 & A \cdot \mathbf{I} & 0 \\ 0 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 0 & -\mathbf{I} & 0 \\ \mathbf{I} & 0 & \rho^\top/d \\ 0 & -\rho^\top/d & 0 \end{pmatrix}$	Nose-Hoover thermostat

14. https://colcarroll.github.io/hamiltonian_monte_carlo_talk/bayes_talk.html

15. <https://mc-stan.org>

Achevons cette section par une vision des possibilités d'application de ces méthodes en *machine learning* (ML) et en *deep learning* (DL) :

ML :

- Régression logistique (dimension $d = 123$) [Welling & Teh 2011]
- Débruitage d'image et déconvolution ($d = 256 \times 256$) [Durmus et al. 2018]
- Régression ($d \in [2, 90]$) $\times 256$) [Dubey et al. 2016]
- Factorisation de matrice ($d = 256 \times 140$) [Simsekli et al. 2016]

DL :

- Approches ensemblistes (dimension $d \in [100, 600]$) [Lakshminarayanan et al. 2017]
- Incertitude des poids ($d = 1200$) [Li et al. 2016]

8.4 Méthodes particulières

Le principe des méthodes particulières est de simuler N chaînes de Markov en parallèle, en éliminant celles qui restent loin du mode *a posteriori* et en multipliant (reproduisant) celles qui en sont proches (cf. figure 18).

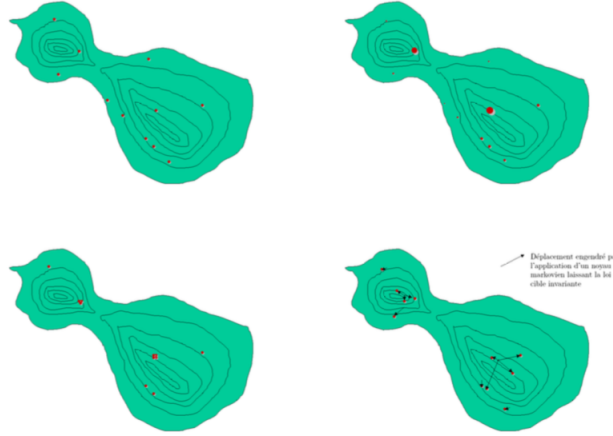


FIGURE 18 – Plusieurs étapes typiques d'une méthode particulière (extrait de Parent et Bernier, 2007).

8.5 Méthodes variationnelles

8.5.1 Principe fondamental

Les *méthodes d'inférence variationnelles* placent le problème d'inférence bayésienne dans un cadre d'optimisation déterministe qui approxime la distribution-cible $\pi(\theta|x) \propto f(x|\theta)\pi(\theta)$ avec une distribution plus simple $g(\theta|\hat{\lambda})$ à manier, en minimisant une divergence de Kullback-Leibler

$$g(\theta|\hat{\lambda}) = \arg \min_{\lambda} KL(\pi(\cdot|x) \| g(\cdot|\lambda)) \quad (20)$$

ou

$$g(\theta|\hat{\lambda}) = \arg \min_{\lambda} KL(g(\cdot|\lambda) \| \pi(\cdot|x)).$$

Notons qu'on n'écrit pas θ dans l'expression ci-dessus car KL est indépendant de tout choix de paramétrisation $h(\theta)$ avec h bijectif.

Remarque 12 Remarquons que si nous supposons posséder quand même un échantillon $\theta_1, \dots, \theta_n \stackrel{iid}{\sim} \pi(\theta|x)$, alors

$$\begin{aligned} \min_{\lambda} KL(\pi(\cdot|x) \| g(\cdot|\lambda)) &= \min_{\lambda} \mathbb{E}_{\pi(\cdot|x)} \left[\log \frac{\pi(\theta|x)}{g(\theta|\lambda)} \right] = \max_{\lambda} \mathbb{E}_{\pi(\cdot|x)} [\log g(\theta|\lambda)] \\ &\simeq \max_{\lambda} \frac{1}{M} \underbrace{\sum_{i=1}^M \log g(\theta_i|\lambda)}_{\text{log-vraisemblance}} \quad (\text{estimateur de Monte Carlo}). \end{aligned}$$

En d'autres termes, la meilleure approximation possible au sens de KL, sachant le choix de modèle $g(\theta|\lambda)$, est le modèle $g(\theta|\hat{\lambda})$ où $\hat{\lambda}$ est l'estimateur du maximum de vraisemblance (EMV) des données $\theta_1, \dots, \theta_n$. Cela permet de montrer d'ailleurs que l'EMV est invariant par reparamétrisation.

On s'attendrait, dans l'expression (20), à utiliser la divergence KL dans le sens précédent :

$$KL_1(g) = KL(\pi(\cdot|x) \| g(\cdot|\lambda))$$

car la loi supposée être la "bonne" (*target*) est $\pi(\theta|x)$: $KL_1(g)$ offre une *quantification de l'erreur informationnelle* résultant de l'approximation de $\pi(\theta|x)$ par $g(\theta|\lambda)$. Cependant, la théorie variationnelle bayésienne considère en général la divergence KL dans l'autre sens :

$$KL_2(g) = KL(g(\cdot|\lambda) \| \pi(\cdot|x))$$

En effet, cette approche est plus simple d'un point de vue calculatoire, et permet de préférer une approximation de $\pi(\cdot|x)$ pour lesquelles les régions de haute densité (de $g(\cdot|\lambda)$) sont les plus correctes. Elle permet en outre de construire des *bornes variationnelles* permettant de transformer le calcul de la *vraisemblance marginale* ou *évidence*

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

en un problème d'optimisation (voir Fox et Roberts (2012) pour une revue détaillée). En effet, d'après l'inégalité de Jensen, pour toute densité $g(\theta|\lambda)$ de support Θ ,

$$\begin{aligned} -KL(g(\cdot|\lambda) \| \pi(x, \cdot)) &= \int_{\Theta} g(\theta|\lambda) \log \frac{\pi(x, \theta)}{g(\theta|\lambda)} d\theta \leq \log \int_{\Theta} g(\theta|\lambda) \frac{\pi(x, \theta)}{g(\theta|\lambda)} \\ &\leq \log \int_{\Theta} \pi(x, \theta) = \log f(x) \end{aligned}$$

avec égalité si $g(\theta|\lambda) = \pi(\theta|x)$. De plus

$$\log f(x) = \underbrace{-KL(g(\cdot|\lambda) \| \pi(x, \cdot))}_{\text{énergie libre } \mathcal{F}(g)} + KL(g(\cdot|\lambda) \| \pi(\cdot|x)).$$

Si on maximise l'énergie libre $\mathcal{F}(g)$ en $g(\cdot|\lambda)$, on maximise une borne inférieure de $f(x)$, et plus celle-ci s'accroît, plus on se rapproche d'une situation où $KL(g(\cdot|\lambda) \| \pi(\cdot|x))$ est faible $\Rightarrow g(\theta|\lambda)$ devient géométriquement proche de la loi-cible $\pi(\theta|x)$, indépendamment de la paramétrisation θ .

8.5.2 Principes d'usage

Approximation champ moyen. L'approximation en champ moyen permet de faciliter la maximisation de $\mathcal{F}(g)$ lorsque le modèle est de Markov caché (ex : traitement d'images). Dans ce cas, $\theta = (\theta_0, z)$ où z est un ensemble de variables latentes z_1, \dots, z_M . On peut supposer par exemple adopter une approche *par séparabilité* ou *composite* simplifiant le problème :

$$g(\theta|\lambda) = g_0(\theta_0|\lambda_0) \prod_{i=1}^M g_z(z_i|\lambda_i).$$

Algorithme bayésien variationnel (VBEM). L'algorithme bayésien variationnel (Beal, 2003) maximise itérativement l'énergie libre $\mathcal{F}(g)$ par rapport aux distributions $g_z(z_i|\lambda_i)$ (étape VBE) puis $g_0(\theta_0|\lambda_0)$ (étape VBM). Dans ce cadre, les observations ne sont pas obligatoirement iid (comme en théorie champ moyen généralement). Cet algorithme est en fait réduit à l'algorithme EM si $g_0(\theta_0|\lambda_0)$ est la loi de l'estimateur du maximum de vraisemblance de la vraisemblance complète (les données latentes z étant alors connues). Toutefois, ce type d'approche est parfois hautement simplificateur : les lois obtenues ressemblent peu aux vrais posteriors.

L'approximation bayésienne variationnelle a donné cours à de nombreux travaux et a été initialement appliquée à différents modèles :

- modèles de mélange (Wang et Titterton);
- modèles à espace d'états (Wang et Titterton);
- modèles graphiques (Attias, Jordan, Beal et Ghahramani);
- réseaux de neurones (Titterton, etc.).

Ces dernières années, le développement d'algorithmes du type SVGD (*Stein Variational Gradient Descent* (Liu et al. 2016)), ce genre de technique s'est fortement répandu en optimisation de réseaux de neurones.

8.6 Méthodes d'échantillonnage sans vraisemblance (ABC)

Nous citons pour information les méthodes ABC (*Approximate Bayesian Computation*) forment une famille de méthodes s'attaquant à des problèmes bayésiens pour lesquels la vraisemblance est très difficilement (voire pas du tout) manipulable (calculable). Elles se fondent sur un principe très simple d'acceptation-rejet de θ dans un ensemble évolutif de lois instrumentales, le test statistique étant assuré par un choix de statistiques supposées représentatives.

De tels problèmes concernent typiquement des données spatio-temporelles complexes, des données excessivement bruitées et dépendantes ou encore des modèles stochastiques implicites (ex : modèles évolutifs de population). Dans tous les cas, nous partirons du postulat que si la vraisemblance n'est pas calculable, il est par contre aisé de simuler des variables synthétiques X à partir de tirages de θ .

8.7 Vérification *a posteriori*

La vérification *a posteriori* regroupe l'ensemble des techniques visant à faciliter le choix d'un modèle bayésien, en particulier lorsqu'on souhaite faire varier les choix *a priori*. La *validation* repose grossièrement sur des techniques de (contrôle *a posteriori* prédictif).

Définition 46 Contrôle *a posteriori* prédictif. Un contrôle *a posteriori* prédictif est un ensemble de comparaisons entre les données initiales x_1, \dots, x_n (au travers d'une statistique résumée éventuellement) et la loi prédictive *a posteriori*

$$f(x|x_1, \dots, x_n) = \int_{\Theta} f(x|\theta) \pi(\theta|x_1, \dots, x_n) d\theta.$$

Les contrôles *a posteriori* prédictifs sont utiles pour évaluer si votre modèle vous donne des prédictions "valables" sur la réalité - correspondent-elles ou non aux données observées? Il s'agit d'une phase utile de construction et de vérification de modèle; elle ne vous donne pas de réponse définitive sur si votre modèle est "ok" ou s'il est "meilleur" qu'un autre modèle, cependant, elle peut vous aider à vérifier si votre modèle fait sens.

Cependant, les contrôles *a posteriori* prédictifs impliquent une double utilisation des données, ce qui viole le principe de vraisemblance. Ils peuvent cependant être utilisés si l'utilisation se limite à des mesures de

divergence pour étudier l'adéquation du modèle, et non pour comparer et inférer des modèles (Meng 1994). Il en va ainsi du calcul des *posterior predictive p-values* en dimension 1 :

$$P(X < \tilde{x}_i | x_1, \dots, x_n) = \int_{\Theta} P(X < \tilde{x}_i | \theta) \pi(\theta | x_1, \dots, x_n) d\theta$$

qui fournissent des diagnostics visuels.

Afin d'obtenir des diagnostics plus fins, il faut comprendre l'idée que l'information d'une distribution est fondamentalement portée par l'espérance du log de sa densité (entropie relative). Cette densité correspond à la vraisemblance intégrée dans le cadre bayésien. Plus ce log sera élevé, plus le choix de modèle bayésien sera susceptible de bien décrire l'information apportée par les données.

Définition 47 elpd. On nomme *elpd* (expected (or mean) log predictive density for a new data point) la quantité

$$\begin{aligned} elpd &= \mathbb{E}_g [\log f(\tilde{x} | x_1, \dots, x_n)] = \int \left[\log \int_{\Theta} f(\tilde{x} | \theta) \pi(\theta | x_1, \dots, x_n) d\theta \right] g(\tilde{x}) d\tilde{x} \\ &\simeq \frac{1}{M} \sum_{j=1}^M \log \int_{\Theta} f(\tilde{x}_j | \theta) \pi(\theta | x_1, \dots, x_n) d\theta \quad \text{avec } \tilde{x}_j \stackrel{iid}{\sim} g \end{aligned}$$

où g le "vrai" modèle de production d'une donnée \tilde{x} .

Le *elpd* constitue en quelque sorte le critère idéal, qui ne peut jamais être atteint, car g reste inconnue. L'approche la plus naturelle pour estimer le *elpd* est une approche *out-of-bag* (OOB) (similaire à celle de l'apprentissage) : on sépare l'échantillon des X en un échantillon d'entraînement \mathbf{x}_k et un échantillon de test \mathbf{x}_{-k} . On peut alors définir les quantités suivantes.

Définition 48 lppd. On nomme *lppd* (log pointwise predictive density) la quantité suivante : pour $x_i \in \mathbf{x}_{-k}$

$$lppd = \log \prod_{x_i \in \mathbf{x}_{-k}} f(x_i | \mathbf{x}_k) = \sum_{x_i \in \mathbf{x}_{-k}} \log \int_{\Theta} f(x_i | \theta) \pi(\theta | \mathbf{x}_k) d\theta.$$

Cette quantité doit être estimée par le biais de l'estimateur ci-dessous.

Définition 49 clppd. On nomme *clppd* (computed log pointwise predictive density) la quantité suivante : pour $x_i \in \mathbf{x}_{-k}$

$$clppd = \sum_{x_i \in \mathbf{x}_{-k}} \log \frac{1}{S} \sum_{s=1}^S f(x_i | \theta_s)$$

où $\theta_s \stackrel{iid}{\sim} \pi(\theta | \mathbf{x}_k)$ (tirage a posteriori).

In fine, le *elcplpdd* ci-dessous fournit un estimateur calculable du *elpd*.

Définition 50 ecclppd. On nomme *ecclppd* (expected clppd) la quantité suivante, estimée par Monte Carlo OOB-CV :

$$ecclppd = \mathbb{E}_{\mathbf{x}_k, \mathbf{x}_{-k}} \left[\sum_{i=1}^n \sum_{x_i \in \mathbf{x}_{-k}} \log \frac{1}{S} \sum_{s=1}^S f(x_i | \theta_s) \right].$$

Cet estimateur peut être très coûteux en temps de calcul, mais ce dernier peut être diminué en menant des étapes d'importance sampling à partir de la loi *a posteriori*.

Il est intéressant de constater que le *elcplpdd* peut être relié à des problématiques de choix de modèle. Une statistique classique de choix de modèle est le critère AIC. Celui-ci n'est pas du tout bayésien et repose sur des considérations asymptotiques. Le critère DIC rend le format AIC "plus bayésien", mais le WAIC le supplante encore. À l'heure actuelle, le DIC et surtout le WAIC sont les critères de préférence pour sélectionner des modèles bayésiens.

Définition 51 *Aikake Information Criterion (AIC; 1973).* Soit k le nombre de paramètres estimés dans le modèle $f(x|\theta)$. La log-densité prédictive estimée au maximum de vraisemblance $\hat{\theta}_n$, à laquelle on soustrait k , fournit une information sur la façon dont l'estimation de k paramètres va accroître la finesse de prévision

$$AIC = -2elpd_{AIC}$$

avec

$$elpd_{AIC} = \log f(x_1, \dots, x_n | \hat{\theta}_n) - k.$$

Lorsqu'on utilise des des *a priori* sur θ et des structures hiérarchiques, cette approximation n'est pas suffisante.

Définition 52 Deviance Information Criterion (DIC; 2002). Il s'agit d'une généralisation du critère AIC à une modélisation hiérarchique, en utilisant l'estimateur bayésien $\tilde{\theta}_n = \mathbb{E}[\theta | x_1, \dots, x_n]$.

$$DIC = -2elpd_{DIC}$$

avec

$$elpd_{DIC} = \log f(x_1, \dots, x_n | \tilde{\theta}_n) - p_{DIC}$$

et

$$\begin{aligned} p_{DIC} &= 2 \left[\log f(x_1, \dots, x_n | \tilde{\theta}_n) - \mathbb{E}_{\theta \sim \pi(\theta | x_1, \dots, x_n)} [\log f(x_1, \dots, x_n | \theta)] \right] \\ &\simeq 2 \left[\log f(x_1, \dots, x_n | \tilde{\theta}_n) - \frac{1}{S} \sum_{s=1}^S \log f(x_1, \dots, x_n | \theta_s) \right] \end{aligned}$$

qui définit le nombre effectif de paramètres estimés dans le modèle

L'utilisation du DIC est donc très courante, car le calcul est stable et rapide. En utilisant un estimateur ponctuel (plug-in) de θ , il reste cependant d'essence très fréquentiste. Le WAIC est quant à lui un critère de sélection de modèle complètement bayésien.

Définition 53 Watanabe-Aikake information criterion (WAIC; 2010)

$$eclppd - p_{WAIC}$$

avec (dans sa version OOB-CV)

$$eclppd = \mathbb{E}_{\mathbf{x}_k, \mathbf{x}_{-k}} \left[\sum_{i=1}^n \sum_{x_i \in \mathbf{x}_{-k}} \log \frac{1}{S} \sum_{s=1}^S f(x_i | \theta_s) \right]$$

et

$$\begin{aligned} p_{WAIC} &= 2 \left[\log \mathbb{E}_{\theta \sim \pi(\theta | x_1, \dots, x_n)} [f(x_1, \dots, x_n | \theta)] - \mathbb{E}_{\theta \sim \pi(\theta | x_1, \dots, x_n)} [\log f(x_1, \dots, x_n | \theta)] \right] \\ &\simeq 2 \left[\log \frac{1}{S} \sum_{s=1}^S f(x_1, \dots, x_n | \theta_s) - \frac{1}{S} \sum_{s=1}^S \log f(x_1, \dots, x_n | \theta_s) \right]. \end{aligned}$$

ANNEXES

A Rappels : concepts et outils fondamentaux de l'aléatoire

Remarque 13 Pour faciliter la lecture et l'appropriation, cette annexe de rappels est illustrée par de nombreux exemples de phénomènes naturels dits extrêmes, telles des pluies diluviennes, des vents forts, etc. dont on cherche à modéliser le comportement.

La modélisation probabiliste d'un aléa X repose sur le caractère de *variable aléatoire* conféré à X , évoluant dans un ensemble d'échantillonnage Ω de dimension d . Puisque $\Omega \neq \emptyset$, les sous-ensembles de valeurs $\mathcal{A} \subset \Omega$ que peut parcourir X sont non vides, et ils présentent une certaine stabilité : l'union dénombrable de plusieurs \mathcal{A}_i est encore dans Ω , de même que le complémentaire de tout sous-ensemble \mathcal{A} .

Ces propriétés fondamentales permettent de "paver" (*mesurer*) l'ensemble Ω de façon à associer à toute observation (survenue) d'un événement $A \in \mathcal{A}$ une valeur numérique $\mathbb{P}(A)$. L'ensemble de ces valeurs numériques vit dans l'intervalle $[0, 1]$, et est tel que

$$\mathbb{P}(\Omega) = 1.$$

On parle alors, pour désigner \mathbb{P} , de *mesure de probabilité*.

La théorie des probabilités nomme le triplet $(\Omega, \mathcal{A}, \mathbb{P})$ *espace probabilisé*, l'ensemble Ω *univers* et \mathcal{A} *tribu* (ou σ -algèbre). En général, le choix de \mathcal{A} est l'ensemble des parties de Ω dont la mesure de Lebesgue peut être définie (cf. § A.1). Il n'est donc usuellement pas donné de précision, dans les problèmes appliqués, sur \mathcal{A} .

A.1 Problèmes unidimensionnels

Considérons tout d'abord le cas où $d = 1$. Si Ω est *discret* (par exemple si $\Omega = \{1, 2, 3, \dots\}$) ou *catégoriel*, et plus généralement si Ω est *dénombrable*, la distribution de probabilité est dite discrète et est déterminée par la *fonction de masse* probabiliste

$$f(x) = \mathbb{P}(X = x)$$

pour toute valeur $x \in \Omega$. Cependant, la très grande majorité des variables aléatoires considérées dans ce cours présente un caractère *continu*. En particulier, les valeurs prises par X (vitesse du vent, température, débit d'une rivière...) évoluent continûment – ce qui est indispensable pour appliquer la théorie des valeurs extrêmes – et Ω constitue généralement un sous-ensemble continu de \mathbb{R}^d , même si le dispositif de mesure est nécessairement limité, en pratique, par une précision donnée. Cette précision ne joue pas de rôle dans la construction du modèle probabiliste mais dans celui du modèle *statistique*, qui englobe le modèle probabiliste en établissant un lien direct avec des observations bruitées (voir § A.5). Dans la pratique, les deux modèles sont confondus quand le bruit d'observation est considéré comme négligeable.

Dans le cas continu, c'est-à-dire lorsque Ω n'est plus dénombrable, la distribution de probabilité peut être spécifiée par la *fonction de répartition*

$$F_X(x) = \mathbb{P}(X \leq x)$$

pour toute valeur $x \in \Omega$. Afin de satisfaire les axiomes des probabilités [21], cette fonction doit être croissante, et telle que, lorsque la dimension $d = 1$,

$$\begin{aligned} \lim_{x \rightarrow x_{\inf}} F_X(x) &= 0, \\ \lim_{x \rightarrow x_{\sup}} F_X(x) &= 1 \end{aligned}$$

où (x_{\inf}, x_{\sup}) sont les bornes inférieure et supérieure (éventuellement infinies) de Ω . Le cas multidimensionnel où $d > 1$ est précisé au § A.3. Toujours pour $d = 1$, l'équivalent de la probabilité discrète $f(x)$ dans le cas continu est fourni par la probabilité que X se situe entre les valeurs $x - a$ et $x + b$ (avec $a, b \geq 0$) :

$$\mathbb{P}(x - a \leq X \leq x + b) = F_X(x + b) - F_X(x - a).$$

Cette propriété pousse à définir, dans les cas où F_X est dérivable, la dérivée de F_X (dite de *Radon-Nikodym-Lebesgue*) définie comme le cas-limite $a = b = \epsilon \rightarrow 0$

$$f_X(x) = \frac{dF_X}{dx}(x),$$

appelée *densité de probabilité* de X , qui est donc telle que

$$F(x) = \int_{-\infty}^x f_X(u) du$$

et

$$\mathbb{P}(x-a \leq X \leq x+b) = \int_{x-a}^{x+b} f_X(u) du.$$

Nécessairement, $\int_{\Omega} f_X(u) du = 1$. Ainsi, toute distribution de probabilité continue, en dimension $d = 1$ (c'est aussi le cas en dimension $d > 1$) peut être représentée de façon équivalente (sous réserve de dérivabilité¹⁶) par sa fonction de répartition ou sa densité (figure 20).

Informellement, f_X peut être vue comme la limite de l'histogramme en fréquence des valeurs possibles de X , pour des classes de valeurs étroites (figure 20). Plus formellement, fonction de répartition et densité de probabilité doivent être interprétées comme des outils permettant d'opérer une *mesure* de la distribution des X relativement à une mesure de l'espace Ω .

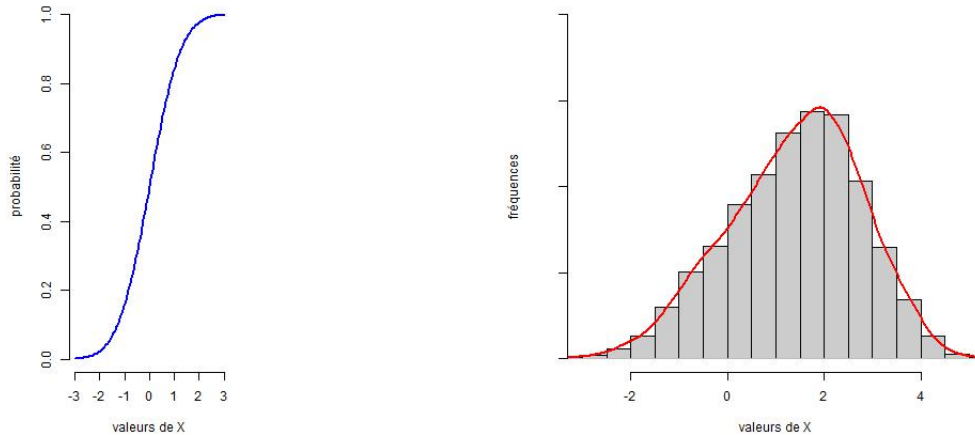


FIGURE 19 – Gauche : exemple de fonction de répartition. Droite : histogramme en fréquence de valeurs de X et densité de probabilité correspondante (courbe).

Considérons par exemple que $\Omega = I_1 \times I_2 \times \dots \times I_d$, où chaque I_k est un intervalle de \mathbb{R} (fermé, ouvert ou semi-ouvert), l'ensemble constituant un parallélépipède contenant toutes les valeurs de X pouvant être observées. Ce solide (ou cet espace) peut être décrit par un ensemble de mesures, par exemple son volume. La mesure de Lebesgue [23], notée μ_L , a été construite comme une mesure de référence permettant de décrire ce type d'espace de façon universelle et uniforme. Comme le volume, elle prend une valeur finie si Ω est compact. La densité f_X définit une autre mesure sur Ω , qui spécifie la forme de la distribution des X et permet de la différencier de l'uniformité. Il faut donc l'interpréter comme une mesure *relative* à celle de Lebesgue (ou *dominée* par la mesure de Lebesgue). Au lecteur intéressé par une introduction détaillée à la théorie de la mesure, nous suggérons les ouvrages [4] (pour une approche "ingénierie") et [22] (pour une vision plus mathématique).

16. Plus généralement de *différentiabilité* en dimension quelconque.

L'information incertaine transportée par les distributions de probabilité est très souvent résumée par des indicateurs statistiques particuliers : les *moments* d'ordre $k \in \mathbb{N}$, définis comme l'ensemble des valeurs moyennes de la variable X^k :

$$M_k = \mathbb{E}[X^k] = \int_{\Omega} x^k f_X(x) dx.$$

Si ceux-ci existent pour $k = 1$ et $k = 2$, ils permettent de définir l'*espérance* $\mathbb{E}[X]$ et la *variance*

$$\mathbb{V}[X] = \int_{\Omega} (x - \mathbb{E}[X])^2 f_X(x) dx.$$

L'espérance fournit une mesure de localisation moyenne de X dans la distribution f_X , tandis que $\mathbb{V}[X]$ est une mesure de la variabilité (ou dispersion) de f_X . L'*écart-type* de f_X , homogène à X , est défini par

$$\sigma_X = \sqrt{\mathbb{V}[X]}.$$

Alternativement, le *coefficient de variation* de X

$$\text{CV}[X] = \frac{\sigma_X}{\mathbb{E}[X]},$$

fournit une autre mesure *relative* de la variabilité ou dispersion de f_X (plus usuelle pour les ingénieurs). Enfin, on parlera de variable centrée-réduite si X est transformée en

$$X' = \frac{X - \mathbb{E}[X]}{\sigma_X},$$

d'espérance nulle et de variance unitaire.

A.2 Familles de modèles paramétriques

Rappelons quelques modèles probabilistes ou statistiques fondamentaux, qui interviennent très souvent dans les constructions plus élaborées qui seront décrites dans ce cours. Ces modèles seront, dans le cadre de ce cours, considérés *paramétriques*, c'est-à-dire descriptibles de façon exhaustive par un ensemble fini de paramètres.

La première raison de ce choix est liée au cadre d'étude : le comportement des extrêmes d'un échantillon aléatoire suit, sous certaines conditions théoriques, des lois paramétriques. C'est aussi le cas du comportement des estimateurs statistiques (§ A.6) obéissant à une loi des grands nombres.

Cet argument fondamental se renforce de la constatation suivante : lorsqu'on s'intéresse à ces comportements extrêmes, le nombre d'observations disponibles devient faible. Expliquer la production de ces observations par un mécanisme aléatoire déterminé par un nombre infini ou même simplement grand de paramètres (c'est-à-dire plus grand que le nombre de données) semble déraisonnable car la majeure partie de ces paramètres resteront inconnus, ou posséderont plusieurs valeurs possibles, et le modèle ainsi créé ne serait pas identifiable et utilisable.

Dans ce document, on notera très généralement θ ce vecteur de paramètres, qui évoluera donc dans un espace θ de dimension finie. Le conditionnement à θ du mécanisme de production aléatoire sera rappelé dans les notations des densités et fonctions de répartition : $f_X(x) = f(x|\theta)$ et $F_X(x) = F(x|\theta)$.

Lois

Dans un cadre discret, on peut s'intéresser à la survenue d'un événement ponctuel $Z > z_0$, où Z est, par exemple, un niveau d'eau maximal mensuel, et z_0 une hauteur de digue de protection. Supposons disposer d'un échantillon d'indicateurs $(\delta_1, \dots, \delta_n) \in \{0, 1\}^n$ valant chacun 1 si la crue ainsi définie survient, et 0 sinon. Faisons l'hypothèse que les δ_i sont indépendants et correspondent chacun au résultat d'un "essai de submersion" réussissant avec une même probabilité p . Si l'on note $X_n = \sum_{i=1}^n \delta_i$ le nombre total de "succès" parmi ces n essais, alors la fonction de masse probabiliste de X_n s'écrit

$$f(x) = \binom{n}{x} p^x (1-p)^{n-x}$$

pour $x \in \Omega = \{0, 1, 2, \dots, n\}$, et où

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}.$$

La variable aléatoire X_n est alors dite suivre la *loi binomiale* $\mathcal{B}(n, p)$ (figure 20). Dans le cadre d’une étude de risque, on s’attachera à estimer la probabilité de surverse p à partir de la statistique observée x_n .

La variable X_n dite de *comptage* définie ci-dessus peut être généralisée dans une perspective d’estimer l’occurrence d’événements survenant de façon aléatoire durant un laps de temps fixé (par exemple une année). Si on suppose que ces événements surviennent avec une fréquence moyenne unique $\lambda > 0$ dans cet intervalle de temps, alors la probabilité qu’il survienne exactement $X_n = x \in \Omega = \{0, 1, \dots, \infty\}$ occurrences est

$$f(x) = \frac{\lambda^x}{x!} \exp(-\lambda),$$

qui définit la fonction de masse probabiliste de la *loi de Poisson* d’espérance λ (figure 20). Celle-ci joue notamment un grand rôle dans l’établissement des lois statistiques associées aux observations historiques car elle permet de modéliser la survenue du nombre d’événements situés entre deux dates (par exemple séparés par plusieurs dizaines d’années) et non observés directement. Le lien technique entre la loi binomiale et la loi de Poisson s’exprime dans le lemme suivant :

LEMME 1. Si X_n suit une loi binomiale $\mathcal{B}(n, p)$ avec $p \ll 1$, alors la loi de X_n peut être approximée par la loi de Poisson d’espérance np lorsque $n \rightarrow \infty$.

Rappelons enfin, dans le cas continu, l’importance fondamentale de la *loi normale* $X \sim \mathcal{N}(\mu, \sigma^2)$, d’espérance μ et de variance σ^2 , et de densité de probabilité (pour $d = 1$ et $\Omega = \mathbb{R}$)

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right).$$

Celle-ci modélise un grand nombre de phénomènes, en particulier celui de la répartition de la moyenne d’un échantillon aléatoire (loi des grands nombres). La convergence en loi normale d’un estimateur statistique (cf. § A.6) constitue un type de résultat très classique (théorème de la limite centrale). La variable $(X - \mu)/\sigma$ suit la loi normale dite *centrée réduite* $\mathcal{N}(0, 1)$ (figure 20). On note usuellement par $\phi(\cdot)$ et $\Phi(\cdot)$ les densité et fonction de répartition de cette loi centrée réduite.

Tests statistiques

La démarche générale des tests consiste à rejeter ou ne pas rejeter (sans forcément accepter) une hypothèse statistique H_0 , dite *nulle*, en fonction d’un jeu de données \mathbf{x}_n . Par exemple, dans un cadre paramétrique cette hypothèse peut correspondre au choix spécifique d’une valeur $\theta = \theta_0$ dans une même famille $f(x|\theta)$ ou d’un domaine $\theta \in \theta_0$. Définir un test revient à définir une statistique

$$R_n = R(X_1, \dots, X_n)$$

qui est une variable aléatoire dont la loi \mathcal{F}_{R_n} est connue (au moins asymptotiquement, c’est-à-dire quand $n \rightarrow \infty$) lorsque l’hypothèse H_0 est vraie, et cette loi est indépendante de la *valeur de l’hypothèse*. (ex : indépendante de θ). Plus précisément, dans un cadre paramétrique où θ est testé, la loi \mathcal{F}_{R_n} ne doit pas dépendre de θ , et la variable R_n est dite *pivotale*. Lorsque R_n est défini indépendamment de θ , cette statistique est dite également *ancillaire*.

Le positionnement de la statistique *observée* $r_n = r(x_1, \dots, x_n)$ dans la loi \mathcal{F}_{R_n} a été définie par Fisher (1926; [11]) comme la probabilité p_{r_n} (dite *p-valeur* ou *p-value*) d’observer un événement plus “extrême” (plus petit ou plus grand) que r_n . Plus cette probabilité est faible, plus l’événement r_n est “loin” des valeurs de R_n de plus haute densité, et moins H_0 est probable (rappelons que la *p-valeur* n’est pas la probabilité que H_0

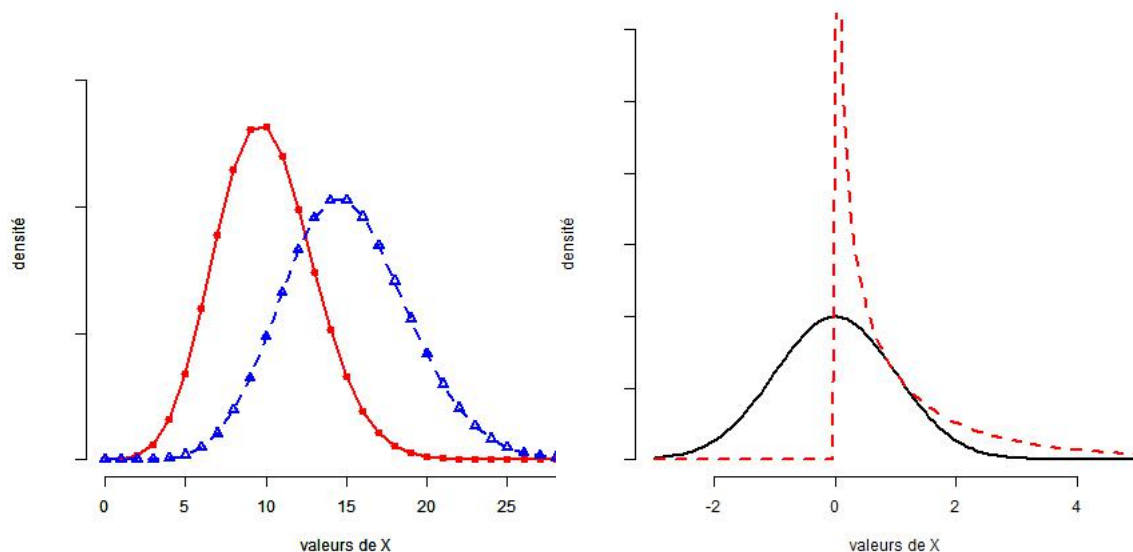


FIGURE 20 – Gauche : fonction de masse des lois discrètes binomiale $\mathcal{B}_n(100, 0.1)$ (carrés) et Poisson $\mathcal{P}(15)$ (triangles). Droite : densités de probabilité continues de la loi normale centrée réduite $\mathcal{N}(0, 1)$ (courbe pleine) et χ_1^2 .

soit vraie). En d'autres termes, si H_0 est fausse, r_n devrait être une valeur extrême de \mathcal{F}_{R_n} .

L'approche courante des tests, dite de *Neyman-Pearson* (1928 ; [24]), impose de fixer un *seuil de signification* $\alpha \ll 1$ définissant l'extrémalité et de comparer le quantile $q_{1-\alpha}$ de la loi \mathcal{F}_{R_n} avec p_{r_n} ; si $p_{r_n} < q_{1-\alpha}$, l'événement r_n est encore moins probable que α , et l'hypothèse H_0 doit être rejetée. Dans le cas contraire, cette hypothèse est plausible (mais pas forcément validée). La pratique courante dans l'ensemble des sciences expérimentales, là encore, est de fixer $\alpha = 5\%$ ou $\alpha = 1\%$, mais ces seuils arbitraires sont de plus en plus critiqués [28, 9], et il est actuellement recommandé [18, 3] de mener plusieurs tests et de tester des seuils α très faibles (ex : $\alpha \in [1\%, 5\%]$)

Dans de nombreux cas, la statistique R_n est choisie positive, afin de pouvoir définir simplement la p -valeur $p_{r_n} = \mathbb{P}(R_n > r_n)$.

EXEMPLE 31. **Test de Kolmogorov-Smirnov [46].** Disposant de l'estimateur empirique classique (cf. § A.6) $x \mapsto \hat{F}_n(x)$ de la fonction de répartition F d'un échantillon iid unidimensionnel x_1, \dots, x_n , défini par

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \leq x\}},$$

et d'un candidat F_0 pour F , on souhaite tester l'hypothèse $H_0 : F = F_0$. La statistique de test est définie par

$$R_n = \sqrt{n} \sup_{x \in \mathbb{R}} \left\| \hat{F}_n(x) - F_0(x) \right\|.$$

Sous H_0 et pour n grand, R_n suit approximativement la loi de Kolmogorov, définie par sa fonction de répartition

$$F_{KS}(x) = 1 - 2 \sum_{k=1}^{\infty} (-1)^{k+1} \exp(-2k^2 x^2) \quad \text{pour } x \in \mathbb{R}^+,$$

qui est généralement tabulée au sein des outils logiciels classiques.

Pour une classe importante de tests, dits du χ^2 (Chi-2), la statistique R_n est construite de façon à suivre loi du χ^2 avec $q \geq 1$ degrés de liberté

$$R_n \sim \chi_q^2$$

dont la densité est tracée sur la figure 20 pour $q = 1$. Les lois du χ^2 sont intrinsèquement liées aux lois normales par une relation quadratique. Par exemple, la somme des carrés de n variables $\mathcal{N}(0, 1)$ indépendantes suit une loi du χ_n^2 à n degrés de liberté. Les quantiles de cette loi sont fournis en pratique par des tables ou algorithmes spécifiques.

Puissance d'un test. Rappelons que deux procédures testant une même hypothèse H_0 ne sont pas forcément aussi pertinentes l'une que l'autre ; elles peuvent être comparées par leur *puissance*, c'est-à-dire leur probabilité respective de rejeter l'hypothèse nulle H_0 sachant qu'elle est incorrecte. Lorsqu'on utilise un test, il convient toujours de s'assurer que sa puissance est élevée, voire la meilleure possible [45]. Elle est définie par

$$1 - \beta$$

où β est nommée *erreur* ou *risque de deuxième espèce* - c'est-à-dire le risque d'accepter à tort l'hypothèse H_0 . L'erreur de deuxième espèce est équivalente à un *taux de faux positifs* dans une procédure de détection. Un exemple classique de test le plus puissant entre deux hypothèses simples $H_0 : \mathbb{P} = P_0$ et $H_1 : \mathbb{P} = P_1$ est le *test de rapport de vraisemblance* (Théorème de Neyman-Pearson), dit aussi test LRT (*likelihood ratio test*).

EXEMPLE 32. Test d'adéquation du χ^2 (cas discret) [49]. Soit $\mathbf{x}_n = (x_1, \dots, x_n)$ un échantillon de réalisations de X supposées iid dans un ensemble fini de valeurs $\{1, \dots, M\}$. On souhaite tester l'hypothèse nulle H_0 selon laquelle les probabilités que X prenne les valeurs 1 à M sont respectivement p_1, \dots, p_M avec $\sum_{k=1}^M p_k = 1$. On note alors

$$\hat{p}_k = \frac{1}{n} \sum_{j=1}^n \delta_{\{x_j=k\}}$$

où $\delta_{\{x_j=k\}} = 1$ si $x_j = k$ et 0 sinon. On définit alors

$$R_n = \sqrt{n \sum_{k=1}^M \frac{(\hat{p}_k - p_k)^2}{p_k}} \quad (21)$$

qui suit, sous l'hypothèse H_0 , une loi χ_{M-1}^2 .

Théorème 23 Test LRT (rapport de vraisemblance). Soit $\mathbf{X}_n = (X_1, \dots, X_n)$ un échantillon de variables aléatoires indépendantes et de même loi \mathbb{P} de densité f . On souhaite tester $H_0 : \mathbb{P} = P_0$ contre $H_1 : \mathbb{P} = P_1$. On nomme $L_i(\mathbf{X}_n) = \prod_{k=1}^n f_i(X_k)$ la vraisemblance statistique maximisée sous l'hypothèse $i \in \{0, 1\}$ (voir § A.6 pour une définition détaillée de la vraisemblance et sa maximisation). Soit

$$R_n = 2 \log \frac{L_1(\mathbf{X}_n)}{L_0(\mathbf{X}_n)}.$$

Alors, si P_0 désigne un modèle paramétré par θ tel que $\theta \in \theta_0$ et P_1 est spécifié par $\theta \notin \theta_0$, alors R_n suit asymptotiquement un mélange de mesures de Dirac et de lois du χ^2 dont le degré de liberté est égal ou inférieur au nombre de contraintes q imposées par l'hypothèse nulle.

De nombreuses précisions sur les mécanismes, les spécifications et les mises en garde sur l'interprétation des tests statistiques (tests paramétriques, non paramétriques, tests de conformité, d'adéquation, d'homogénéité, d'indépendance, d'association...) sont fournis dans [43] et [17]. Le cas spécifique des tests LRT est particulièrement détaillé dans [16]. Appliqués au cas spécifique des modèles d'extrêmes, le lecteur intéressé par une revue

générale pourra consulter avec profit l'article [27].

EXEMPLE 33. **Test LRT.** Dans le cas spécifique où θ_0 est dans l'intérieur strict de θ , alors

$$R_n \stackrel{n \rightarrow \infty}{\sim} \chi_q^2. \quad (22)$$

Considérons ainsi une loi normale $\mathcal{N}(\mu, \sigma)$ avec $\theta = (\mu, \sigma) \in \mathbb{R} \times \mathbb{R}_*^+$. On souhaite tester $H_0 : \mu = 0$ contre $H_1 : \mu \neq 0$. Une seule contrainte différencie les deux hypothèses, et $0 \in \mathbb{R}$. Donc $q = 1$ et le résultat (22) s'applique. Si on souhaite tester $H_0 : \mu = 0$ contre $H_1 : \mu > 0$, le domaine θ est alors restreint à $\mathbb{R}^+ \times \mathbb{R}_*^+$, et (22) doit être remplacé par

$$R_n \stackrel{n \rightarrow \infty}{\sim} \frac{1}{2}\delta_0 + \frac{1}{2}\chi_1^2.$$

A.3 Cas multidimensionnels

L'étude d'aléas conjoints nécessite de pouvoir généraliser les principaux concepts et notions décrits au § A.1. Soit $\mathbf{X} = (x_1, \dots, x_d)^T$ le vecteur des aléas considérés. La *fonction de répartition jointe* est définie par

$$F_X(\mathbf{x}) = \mathbb{P}(X_1 \leq x_1, \dots, X_d \leq x_d)$$

où $\mathbf{x} = (x_1, \dots, x_d)$. Lorsque les X_i sont des variables aléatoires continues, et en supposant F_X différentiable, la densité de probabilité jointe s'écrit

$$f_X(\mathbf{x}) = \frac{\partial^d F_X}{\partial x_1 \dots \partial x_d}(\mathbf{x}).$$

Alors, pour tout ensemble $\mathcal{A} \subset \Omega \subset \mathbb{R}^d$

$$\mathbb{P}(\mathbf{X} \in \mathcal{A}) = \int_{\mathcal{A}} f_X(\mathbf{u}) \, d\mathbf{u}.$$

En particulier, si $\Omega = \mathbb{R}^d$:

$$F_X(\mathbf{x}) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_d} f_X(\mathbf{u}) \, du_1 \dots du_d.$$

Chaque *densité marginale*, caractérisant X_i indépendamment des autres variables, s'obtient par intégration sur les autres composantes : si $\Omega = \bigotimes_{i=1}^d \Omega_i$, alors

$$f_{X_i}(x_i) = \iint_{\bigotimes_{j \neq i} \Omega_j} f_X(u_1, \dots, u_{i-1}, x_i, u_{i+1}, \dots, u_d) \, du_1 \dots du_d.$$

La notion de *covariance* permet de résumer la dépendance entre les X_i deux à deux :

$$\mathbb{C}ov(X_i, X_j) = \int_{\Omega_i} \int_{\Omega_j} (x_i - \mathbb{E}[X_i]) (x_j - \mathbb{E}[X_j]) f_{X_i, X_j}(x_i, x_j) \, dx_i dx_j$$

où $\mathbb{E}[X_i]$ est l'espérance marginale de X_i et f_{X_i, X_j} est la densité jointe bivariée de X_i et X_j , définie comme la marginale

$$f_{X_i, X_j}(x_i, x_j) = \int_{\bigotimes_{k \neq i, j} \Omega_k} f_X(\dots, u_{i-1}, x_i, \dots, u_{j-1}, x_j, \dots, u_d) \, du_1 \dots du_d.$$

La covariance généralise la notion de variance : $\mathbb{C}ov(X_i, X_i) = \mathbb{V}[X_i]$ (variance de la loi marginale de X_i). Dans la pratique, la loi multivariée est souvent résumée par son vecteur d'espérances $\mathbb{E}[\mathbf{X}] = (\mathbb{E}[X_1], \dots, \mathbb{E}[X_d])^T$ et sa matrice de variance-covariance

$$\Sigma = (\mathbb{C}ov(X_i, X_j))_{i,j}$$

ou sa *matrice de corrélation* $\Sigma' = (\rho_{i,j})_{i,j}$ définie par

$$\rho_{i,j} = \frac{\mathbb{C}ov(X_i, X_j)}{\sqrt{\mathbb{V}[X_i] \mathbb{V}[X_j]}}. \quad (23)$$

Chaque $\rho_{i,j}$ évolue entre -1 et 1 et fournit une information sur la dépendance *linéaire* entre les variables X_i et X_j . Toutefois, ce résumé est en général très incomplet. Par exemple, s'il y a indépendance entre X_i et X_j , alors $\text{Cov}(X_i, X_j) = 0$, mais la réciproque n'est pas toujours vraie. La matrice des coefficients de corrélation Σ' n'apporte une information exhaustive sur la structure de dépendance que dans des cas très précis, notamment lorsque \mathbf{X} est un vecteur gaussien, mais ne fournit pas en général une mesure réellement pertinente de cette dépendance. Il faut donc combattre la pratique bien établie d'accorder une confiance importante à cet indicateur [50].

Un cours spécifique doit préciser ce qui est entendu par *information exhaustive sur la structure de dépendance*, et fournir des outils plus adaptés au maniement des lois multivariées. Les premiers de ces outils sont les **copules**.

A.4 Processus aléatoires et stationnarité

Les lois apparaissant dans ce cours constituent un cas particulier des *processus aléatoires* (ou stochastiques) en temps discret¹⁷, qui définissent le comportement général d'une suite de variables aléatoires X_1, \dots, X_n . Ces variables ne sont plus obligatoirement considérées comme indépendantes et identiquement distribuées (*iid*). La loi f_{X_i} de chaque X_i peut varier selon i . Il peut aussi y avoir dépendance entre les X_i tout en conservant l'hypothèse d'une loi similaire pour chaque X_i . Dans ce dernier cas, le processus est alors dit *stationnaire*.

Définition 54 Stationnarité d'un processus. *Un processus aléatoire X_1, \dots, X_n est dit stationnaire si, pour tout ensemble d'entiers $\{k_1, \dots, k_s\}$ et pour tout entier m , les distributions de probabilité jointes de $(X_{k_1}, \dots, X_{k_s})$ et $(X_{k_1+m}, \dots, X_{k_s+m})$ sont identiques.*

Cette définition permet par exemple de caractériser les séries temporelles de façon plus appropriée que la mention *iid*. Le mécanisme stochastique définissant le processus aléatoire nécessite parfois d'être précisé. C'est en particulier vrai lorsqu'on étudie si ce processus converge vers un processus stationnaire lorsque n grandit.

On peut ainsi imaginer que X_1, \dots, X_n, \dots représentent des observations d'une température à des pas de temps très courts, et qu'il est souhaitable de pouvoir sélectionner des valeurs de températures stabilisées afin de calculer des grandeurs représentatives. Pour cela, il est nécessaire de pouvoir spécifier la distribution de probabilité de X_k conditionnelle à $X_{k-1}, X_{k-2}, \dots, X_1$ et d'utiliser une représentation par *chaîne de Markov*.

Définition 55 Chaîne de Markov. *Un processus aléatoire X_1, \dots, X_n, \dots est une chaîne de Markov d'ordre $r \in \mathbb{N}^*$ si, pour tout $i \geq r$,*

$$\mathbb{P}(X_i | X_{i-1}, \dots, X_1) = \mathbb{P}(X_i | X_{i-1}, \dots, X_{i-r}).$$

Si, de plus, $r = 1$ et que cette probabilité de transition ne dépend pas de i , le processus est dit homogène.

Les chaînes de Markov d'ordre 1 sont donc les plus aisées à spécifier, et constituent un outil de généralisation important des cas *iid* (un exemple est tracé sur la figure 21). Elles jouent également un grand rôle dans des cadres d'inférence et d'échantillonnage. Ainsi, un processus $\theta_1, \dots, \theta_n$ peut être construit comme un mécanisme d'exploration de l'espace θ , par exemple dans un cadre bayésien, et ce mécanisme d'exploration est très souvent construit en produisant une chaîne de Markov d'ordre 1, qui possède des propriétés de convergence vers un processus-limite stationnaire (on parle aussi de *distribution stationnaire*), dont les propriétés (espérance, variance, etc.) peuvent être estimées. Nous suggérons l'ouvrage de référence [39] au lecteur désireux d'explorer ce champ de la théorie des probabilités.

17. Les processus aléatoires en temps continu ne sont pas traités dans ce cours.

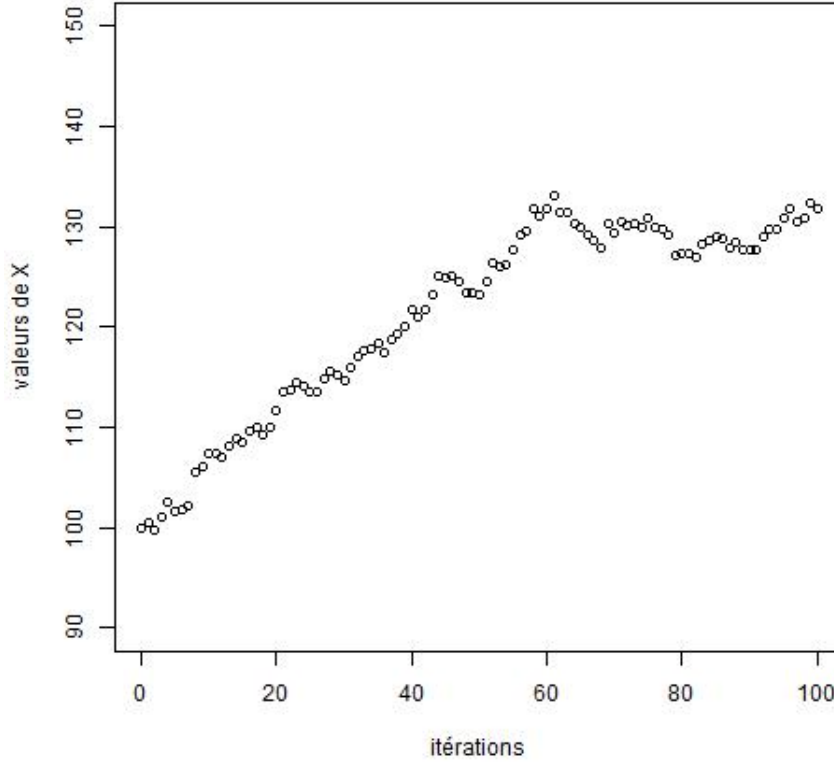


FIGURE 21 – Exemple d’une chaîne de Markov d’ordre 1 non stationnaire.

A.5 Modélisations probabiliste et statistique

Les termes de modélisations probabiliste et statistique sont souvent confondus, en particulier dans la littérature d’ingénierie. Cependant, ils possèdent des sens différents. Un modèle probabiliste décrit par sa densité de probabilité f_x est voué à représenter un phénomène (ex : physique) réel :

$$X \sim f_x$$

tandis que le modèle statistique traduit le fait qu’une ou plusieurs observations de X , notée(s) x^* , sont reliés à une réalisation réelle x de X par un dispositif de mesure : par exemple

$$x^* = x + \epsilon \quad (24)$$

où ϵ est un *bruit d’observation* dont la nature est aléatoire et qui est souvent supposé gaussien. On notera f_ϵ sa densité, qui est en général connue¹⁸. La connaissance de la relation (24) permet de définir la distribution de probabilité de la variable aléatoire X^* de réalisation x^* comme une *loi de convolution* de densité f_{x^*}

$$f_{x^*}(u) = \int f_x(u + y)f_\epsilon(y) dy$$

et cette loi détermine la *vraisemblance statistique* de l’observation x^* (cf. § A.6). Toutefois, il est essentiel pour la détermination de f_x , enjeu majeur de l’étude, que cette loi explique la majeure partie de la variabilité et

18. Notamment *via* les spécifications des constructeurs des dispositifs de mesure, ou par des tests répétés dans des conditions contrôlées.

des valeurs observées (celles de X^*). Souvent, on fera l'hypothèse que l'influence de ϵ est négligeable, ce qui revient à écrire

$$f_{x^*}(u) \simeq f_x(u) \quad \forall u \in \Omega,$$

et à confondre modèles probabiliste et statistique. Cette hypothèse n'est cependant pas toujours vérifiée en pratique, en particulier pour les observations historiques [37]. Le bruit affectant une mesure peut être important car cette dernière peut :

- ne pas être directe (par exemple, les mesures de pluie torrentielles ancestrales peuvent être reconstituées à partir d'études stratigraphiques [13]);
- être très imprécise (exemple : une crue datant du Moyen Âge a fait l'objet d'une chronique en termes qualitatifs (elle a emporté un pont, recouvert des champs...) ou quantitatif avec beaucoup d'incertitude (marque sur un mur de maison démolie depuis) [6, 33];
- souffrir d'un biais inconnu lié à un dispositif de mesure mal calibré (ou abîmé par l'aléa lui-même, surtout s'il est extrême) [2].

Même certaines mesures récentes peuvent souffrir d'un bruit potentiellement fort, car elles sont issues d'un calcul - et non d'une mesure directe - soumis à certaines incertitudes (voir également § ??).

A.6 Contrôle de l'erreur de modélisation

Convergence des modèles

La fiabilité des modèles probabilistes et statistiques repose sur une approximation du réel dont l'erreur peut être encadrée sous certaines hypothèses techniques. On distingue dans ce cours deux types d'approximation :

1. une approximation du comportement inconnu d'une grandeur X considérée comme aléatoire (par exemple le maximum d'un échantillon sur un intervalle de temps donné) par un comportement théorique (par exemple issu de la théorie statistique des valeurs extrêmes) permettant de quantifier et d'extrapoler ;
2. une approximation d'un modèle probabiliste théorique par un modèle statistique *estimé*, au sens où ce modèle théorique implique des paramètres *a priori* inconnus θ , qui seront quantifiés grâce aux observations réelles ; puisque ces observations x_1, \dots, x_n sont considérées comme des réalisations d'une variable aléatoire X , le paramètre *estimé* est également considéré comme une réalisation d'une autre variable aléatoire $\hat{\theta}_n$, définie comme un *estimateur statistique*.

La suite $(\hat{\theta}_n)_n$ constitue donc un premier processus stochastique, dont on souhaite qu'il approxime θ (paramètre fixe mais inconnu). L'ensemble des variables aléatoires $(X_n)_n$ produites alors par le modèle estimé forme un deuxième processus stochastique, dont on souhaite qu'il approxime le comportement réel X (variable aléatoire de loi inconnue).

Il est donc indispensable de vérifier que ces deux types d'approximation n'empêchent pas les *modèles statistiques estimés* - les outils concrets de l'étude - de fournir un diagnostic pertinent en termes de reproductibilité des observations, et n'entravent pas significativement leur emploi dans des études prévisionnelles. Une condition indispensable est d'avoir *convergence* entre modélisation théorique et réalité, puis entre modèle estimé et modélisation théorique. Cette convergence s'exprime sous la forme d'un écart entre les protagonistes, qui doit nécessairement diminuer lorsque la quantité d'information (c'est-à-dire le nombre d'observations n) s'accroît jusqu'à devenir nul lorsque $n \rightarrow \infty$.

Dans le monde probabiliste, cet écart est aléatoire, et il est donc possible qu'un écart soit nul sauf en un nombre k de situations données, formant un sous-ensemble de l'espace des événements Ω de mesure nulle. Typiquement, cet ensemble peut être formé d'un nombre fini de valeurs ponctuelles, ou d'éléments appartenant à la frontière de Ω ; en effet, dans le monde continu on sait que (sous des conditions d'indépendance)

$$\mathbb{P}(X \in \{x_1, \dots, x_m\}) = \sum_{i=1}^m \mathbb{P}(X = x_i)$$

et que $\mathbb{P}(X = x_i) = 0$ pour tout x_i (puisque X est continu). On parlera dans ce cas de nullité *presque sûre*.

La notion de *convergence presque sûre* s'en déduit assez naturellement : il s'agit de vérifier que la probabilité que la limite d'un processus stochastique $\hat{\theta}_n$ (ou X_n) corresponde à la cible θ (ou X) vaut 1 ; ou de façon équivalente, que l'écart entre la limite de ce processus et θ (ou X) est nul presque sûrement :

$$X_n \xrightarrow{p.s.} X \Leftrightarrow \mathbb{P}\left(\lim_{n \rightarrow \infty} X_n = X\right) = 1.$$

Cette notion de convergence est la plus forte et la plus courante en pratique pour démontrer le comportement attendu d'un processus aléatoire vers une variable aléatoire, éventuellement réduite à un vecteur (ou un scalaire). On parle également de *consistance forte*¹.

D'autres notions de convergence moins fortes, au sens où elles sont entraînées par la convergence presque sûre, sans réciprocity assurée, sont également très utilisées :

1. la convergence *en probabilité*

$$X_n \xrightarrow{\mathbb{P}} X \Leftrightarrow \forall \epsilon > 0, \lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \epsilon) = 0$$

joue un rôle important dans un grand nombre de démonstrations de convergences en loi, et implique également la convergence presque sûre d'une sous-suite de $(X_n)_n$; elle permet à X_n de s'écarter de X , mais de moins en moins significativement à mesure que n croît ;

2. la convergence *en loi*, qui est entraînée par la convergence en probabilité et qui constitue l'équivalent de la convergence simple¹⁹ dans le monde probabiliste

$$X_n \xrightarrow{\mathcal{L}} X \Leftrightarrow \lim_{n \rightarrow \infty} F_n(x) = F(x)$$

où (F_n, F) sont les fonctions de répartition de X_n et X , respectivement, pour tout x où F est continue. Cette notion de convergence ne caractérise pas les valeurs des processus stochastiques, mais uniquement les comportements aléatoires : celui de X_n ressemble de plus en plus à celui de X . On parle alors de *consistance faible*¹. Cette convergence caractérise notamment les statistiques de test (§ A.2).

D'autres notions de convergence (en norme L^p en particulier) sont également utilisées. Leur emploi est en général de s'assurer des convergences *déterministes* utiles, par exemple celles des espérances (moments), comme l'expriment les deux théorèmes suivants.

Théorème 24 Supposons que X_n converge en norme L^1 vers X dans $\Omega \in \mathbb{R}$:

$$\lim_{n \rightarrow \infty} \mathbb{E}[|X_n - X|] = 0.$$

Alors $\lim_{n \rightarrow \infty} \mathbb{E}[X_n] = \mathbb{E}[X]$.

Théorème 25 Supposons que $X_n \xrightarrow{\mathcal{L}} X$ avec $(X_n, X) \in \Omega^2$ avec $\Omega \subset \mathbb{R}$. Alors, pour toute fonction réelle, continue et bornée g (en particulier l'identité),

$$\lim_{n \rightarrow \infty} \mathbb{E}[g(X_n)] = \mathbb{E}[g(X)].$$

Un ensemble de résultats techniques permet de combiner ces différentes convergences et leurs transformations par des fonctions continues (*mapping theorem*) pour étudier des modèles complexes. Pour une exploration approfondie des notions évoquées dans ce paragraphe et leur généralisation dans un monde multidimensionnel, nous suggérons au lecteur l'ouvrage [48].

1. Bien que *stricto sensu*, la *consistance* est une propriété locale d'un estimateur, qui est induite par la convergence (possédant un sens global).

19. Au sens de la *fonction caractéristique* pour les spécialistes (théorème de continuité de Lévy [43]).

Estimation statistique classique

L'*inférence* est l'ensemble des méthodologies permettant de construire un ou plusieurs estimateurs de θ

$$\hat{\theta}_n = T(X_1, \dots, X_n)$$

où T est une fonction des variables aléatoires associées aux réalisations x_i du phénomène étudié. Comme indiqué précédemment, $\hat{\theta}_n$ est donc lui-même une variable aléatoire, et sa valeur *observée* $T(x_1, \dots, x_n)$ est appelée un *estimé*.

Principales propriétés des estimateurs statistiques Il existe une infinité d'estimateurs possibles pour un vecteur de paramètre θ , et il est donc indispensable de pouvoir opérer une sélection parmi eux. En statistique classique, les principales règles utilisées pour classer les estimateurs sont les suivantes :

1. *asymptotiquement* il doit y avoir *consistance* :

$$\hat{\theta}_n \xrightarrow{?} \theta$$

où ? représente, au mieux, la convergence presque sûre ;

2. l'*erreur quadratique*

$$\text{EQ}(\hat{\theta}_n) = \mathbb{E} \left[(\hat{\theta}_n - \theta)^T (\hat{\theta}_n - \theta) \right], \quad (25)$$

doit être la plus faible possible ; celle-ci peut s'écrire comme la somme du déterminant de la matrice de variance-covariance de $\hat{\theta}_n$, qui est une mesure de l'imprécision non-asymptotique de cet estimateur, et du carré du *biais*²⁰ de l'estimateur

$$\text{B}(\hat{\theta}_n) = \mathbb{E} \left[\hat{\theta}_n \right] - \theta,$$

que l'on peut définir comme l'*erreur non-asymptotique* en espérance. Ces deux termes ne peuvent être minimisés simultanément, et la minimisation de de (25) procède donc nécessairement d'un *équilibre biais-variance* . .

Remarquons que produire un estimateur $\hat{\theta}_n$ faiblement consistant pour un paramètre inconnu θ permet de produire un autre estimateur faiblement consistant sur n'importe quelle fonction $h(\theta)$ de ce paramètre, pourvu que h soit différentiable. Lorsque la loi de convergence est gaussienne, le procédé de dérivation permettant de le construire est connu sous le nom de *méthode Delta*.

Théorème 26 Méthode Delta multivariée [29]. Soit $\theta_1, \dots, \theta_n$ un processus stochastique dans \mathbb{R}^d et soit $g : \mathbb{R}^d \rightarrow \mathbb{R}^q$ une fonction différentiable et non nulle en θ . Notons $J_g(\theta)$ la jacobienne de g en θ . Supposons que $\sqrt{n}(\theta_n - \theta)$ converge en loi vers la loi normale multivariée $\mathcal{N}_d(\mathbf{0}_d, \Sigma)$, de moyenne le vecteur nul $\mathbf{0}_d$ en dimension d et de variance-covariance $\Sigma \in \mathbb{R}^{2d}$. Alors

$$\sqrt{n}(g(\theta_n) - g(\theta)) \xrightarrow{\mathcal{L}} \mathcal{N}_q(0, J_g^T(\theta) \Sigma J_g(\theta)).$$

Estimation des moindres carrés La classe des *estimateurs des moindres carrés* (EMC), qui cherchent à réaliser un compromis entre biais et variance, est donc naturellement définie par une règle du type

$$\hat{\theta}_n = \arg \min_{\hat{\theta}} \widetilde{\text{EQ}}(\hat{\theta}) \quad (26)$$

où $\widetilde{\text{EQ}}$ est une approximation empirique de EQ, construite comme une fonction de X_1, \dots, X_n . Les estimateurs ainsi produits possèdent souvent de bonnes propriétés de consistance, mais peuvent s'avérer sensibles aux choix du modèle et de la paramétrisation θ . Ainsi, il n'est pas évident que l'espérance et/ou la variance impliquées dans le critère (26) existent.

20. Un estimateur $\hat{\theta}_n$ dont l'espérance est égale à θ est dit *sans biais*.

Estimation par maximisation de vraisemblance

Principe de vraisemblance. Une règle plus générale est donc nécessairement fondée sur une représentation plus exhaustive, générique et toujours définie de l'information apportée par X_1, \dots, X_n sur le modèle paramétré par θ . Une telle représentation est la *vraisemblance statistique* ℓ , qui est définie (pour des X_i continus) comme la densité jointe des observations $X_i = x_i$ conditionnelle à θ . Ainsi, dans un cas où les observations x_i sont des réalisations iid :

$$\ell(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f_X(x_i | \theta). \quad (27)$$

La vraisemblance peut prendre des formes plus compliquées lorsque les réalisations ne sont pas indépendantes, non identiquement distribuées ou sont *manquantes* et ont été remplacées par des valeurs-seuils, par exemple parce que les limites du procédé de mesure ont été atteintes.

EXEMPLE 34. Mesure de la vitesse du vent. Certains vieux anémomètres ne peuvent mesurer la vitesse du vent au-delà d'une certaine valeur, et remplacent l'observation x_i qui aurait dû être faite par une vitesse maximale de vent mesurable, notée c . On parle alors d'observation statistique censurée à droite. Ce type d'observation partielle est fréquente en analyse de survie [25]. Le terme de densité $f(x_i)$ correspondant à une observation correcte est alors remplacé par la probabilité que la donnée manquante $P(X \geq c) = F(c)$ dans l'écriture de la vraisemblance (27), où F est la fonction de répartition de X .

L'exhaustivité de l'information portée par la vraisemblance constitue un principe fondamental de la théorie statistique classique. Alors, l'estimateur du maximum de vraisemblance (EMV)²¹

$$\hat{\theta}_n = \arg \max \ell(X_1, \dots, X_n | \theta) \quad (28)$$

définit la variable aléatoire dont la réalisation est la *valeur la plus probable* de θ ayant généré l'ensemble de réalisations $\{X_i = x_i\}_i$. Par le caractère générique de sa dérivation, sa signification et ses bonnes propriétés de consistance, il est l'estimateur statistique le plus courant et l'un des plus naturels.

EXEMPLE 35. Données censurées par intervalles. Très fréquemment, une donnée historique unidimensionnelle, ou mal mesurée, peut être simplement décrite comme une valeur manquante x_i entre deux bornes connues $x_{i,\min} < x_{i,\max}$. Le terme de densité $f(x_i)$ doit alors être remplacé dans la vraisemblance (27) par

$$P(x_{i,\min} \leq X \leq x_{i,\max} | x_{i,\min}, x_{i,\max}) \quad (29)$$

$$\begin{aligned} &= P(X \leq x_{i,\max} | x_{i,\max}) - P(X \leq x_{i,\min} | x_{i,\min}), \\ &= F(x_{i,\max}) - F(x_{i,\min}). \end{aligned} \quad (30)$$

Si l'on fait de plus l'hypothèse que les valeurs $(x_{i,\min}, x_{i,\max})$ sont des données elles-mêmes aléatoires (par exemple bruitées), décrites comme des réalisations de variables $(X_{i,\min}, X_{i,\max})$ de lois respectives $f_{i,\min}, f_{i,\max}$, le terme de vraisemblance (30) devient

$$\iint P(X_{i,\min} \leq X \leq X_{i,\max} | X_{i,\min} = y_1, X_{i,\max} = y_2) f_{i,\min}(y_1) f_{i,\max}(y_2) dy_1 dy_2.$$

Lorsque la donnée est multivariée, plusieurs situations peuvent se présenter : une ou plusieurs dimensions de X peuvent être censurées par intervalles, et des traitements approfondis doivent être menés pour obtenir des spécifications statistiques utiles (voir [20] pour les analyses de survie, et [40] pour le cas spécifique des extrêmes multivariés).

21. Pour des raisons de commodité, on remplace souvent ℓ par la log-vraisemblance $\log \ell$ dans la définition (28).

Théorème 27 Limite centrale pour l'EMV. Supposons que X_1, \dots, X_n soient indépendants et identiquement distribués. Soit q la dimension de θ . Alors, sous des conditions de régularité très générales (dites de Wald),

$$\hat{\theta}_n \xrightarrow{L} \mathcal{N}_q(\theta, I_\theta^{-1}) \quad (31)$$

où \mathcal{N}_q représente la loi normale multivariée en dimension q , de variance-covariance I_θ^{-1} , et I_θ est la matrice d'information de Fisher dont le terme $(i, j) \in \{1, \dots, q\}^2$ est défini par (sous ces mêmes conditions de régularité)

$$I_\theta^{(i,j)} = -\mathbb{E}_X \left[\frac{\partial \log \ell(X_1, \dots, X_n | \theta)}{\partial \theta_i \partial \theta_j} \right]. \quad (32)$$

Quelques informations supplémentaires sur la notion d'information et la matrice de Fisher sont indiquées au § A.6. Deux propriétés importantes de l'EMV sont d'être *asymptotiquement sans biais* et *fortement consistant et efficace asymptotiquement* : sa covariance asymptotique, fournie par l'inverse de la matrice de Fisher, est *minimale* pour tous les estimateurs sans biais de θ .

Information de Fisher La notion d'information a été proposée dans les années 1920 par le chercheur anglais Ronald A. Fisher (considéré comme le père de la statistique mathématique). La démarche de Fisher est la suivante : si l'on s'intéresse aux caractéristiques d'une population nombreuse (voire infinie, qui est le cas limite auquel on est ramené en permanence), on ne peut ni connaître ni traiter les informations trop abondantes relatives à chacun des individus qui la composent. Le problème devient donc d'être capable de décrire correctement la population au moyen d'indicateurs de synthèse pouvant être fournis par des échantillons issus de la population à étudier. Plus les données chiffrées que l'on peut extraire d'un échantillon représentent correctement la population de référence et plus l'information contenue dans cet échantillon doit être considérée comme élevée.

Partant de cette hypothèse, Fisher a défini techniquement l'information comme la valeur moyenne du carré de la dérivée du logarithme de la loi de probabilité étudiée. L'inégalité de Cramer permet alors de montrer que la valeur d'une telle information est proportionnelle à la faible variabilité – c'est-à-dire au fort degré de certitude – des conclusions qu'elle permet de tirer. Cette idée, qui est à la racine de toute la théorie de l'estimation et de l'inférence statistique, est exactement celle que l'on retrouvera vingt ans plus tard chez Shannon, exprimée cette fois en des termes non plus statistiques mais probabilistes.

Si X est un échantillon de densité de probabilité $f(x|\theta)$, on définit l'information de Fisher par

$$I_\theta = \mathbb{E} \left[\left(\frac{\partial \log f(X|\theta)}{\partial \theta} \right)^2 \right].$$

Dans le cas où la distribution de probabilité dépend de plusieurs paramètres, θ n'est plus un scalaire mais un vecteur. L'information de Fisher n'est plus définie comme un scalaire mais comme une matrice de covariance appelée matrice d'information de Fisher :

$$I_{\theta_i, \theta_j} = \mathbb{E} \left[\left(\frac{\partial \log f(X|\theta)}{\partial \theta_i} \right) \left(\frac{\partial \log f(X|\theta)}{\partial \theta_j} \right) \right]$$

Intervalle de confiance Dans la pratique, la loi asymptotique de $\hat{\theta}_n$ est à son tour estimée en remplaçant le terme inconnu I_θ par un estimateur consistant \hat{I}_n (en général $\hat{I}_n = I_{\hat{\theta}_n}$), ce qui permet de définir des *zones de confiance* $C_{\hat{\theta}_n, \alpha}$ associées à l'estimateur $\hat{\theta}_n$ telles que, lorsque n croît vers l'infini,

$$\mathbb{P}(\hat{\theta}_n \in C_{\hat{\theta}_n, \alpha}) = \alpha. \quad (33)$$

En particulier, lorsqu'on s'intéresse à une dimension spécifique θ_i , le théorème 27 permet de définir l'*intervalle de confiance (asymptotique)* $1 - \alpha$ associé à $\hat{\theta}_n$:

$$\mathbb{P} \left(\theta_i \in \left[\hat{\theta}_{n,i} - z_{\alpha/2} \sqrt{\sigma_{i,i}^2}, \hat{\theta}_{n,i} + z_{\alpha/2} \sqrt{\sigma_{i,i}^2} \right] \right) = 1 - \alpha, \quad (34)$$

où z_α est le quantile d'ordre α de la loi normale centrée réduite et $\sigma_{i,i}^2$ le terme diagonal (i, i) de l'*estimé* de l'inverse \hat{I}_n^{-1} .

Les équations (33) et (34) permettent d'évaluer la précision de l'estimation de θ à partir de l'échantillon x_1, \dots, x_n . Cependant, observons que la mesure de probabilité \mathbb{P} dans l'équation (33) concerne $\hat{\theta}_n$ et non θ (qui est inconnu mais fixe); une zone de confiance n'est donc pas définie par la probabilité $1 - \alpha$ que θ s'y situe, mais comme une zone où il y a *a priori* une très forte probabilité $1 - \alpha$ d'obtenir un *estimé* de θ . En simulant un grand nombre de fois des échantillons similaires à x_1, \dots, x_n , la distribution de ces estimés a $100(1 - \alpha)\%$ chances en moyenne de contenir la vraie valeur θ . L'intervalle de confiance sur une dimension i de θ vise donc à encadrer la vraie valeur θ_i avec une certaine probabilité reliée à la loi asymptotique de l'estimateur $\hat{\theta}_n$, qui présuppose que le modèle statistique est correct (et non selon une sorte de probabilité absolue, indépendante de tout modèle).

Tout comme l'EMC, l'EMV n'est pas toujours explicite et doit être en général calculé par des méthodes numériques. EMC et EMV peuvent ne pas être uniques pour des modèles complexes, et l'EMV peut aussi ne pas être défini (menant à une vraisemblance infinie). Toutefois ces cas restent rares dans le cadre de la théorie des valeurs extrêmes. À la différence de l'EMC, l'EMV est toujours invariant par reparamétrisation : l'EMV de $h(\theta)$ est $h(\hat{\theta}_n)$ pourvu que h soit une fonction bijective. Cette propriété est cruciale pour éviter des paradoxes et des inconsistances : si on remplace l'observation x par une transformation bijective $y = d(x)$, le modèle paramétré par θ est remplacé par un modèle paramétré par une transformation bijective $\theta' = h(\theta)$. Or l'information apportée par x et y est la même, et donc toute règle d'estimation $x \rightarrow \hat{\theta}(x)$ devrait être telle que $y = d(x) \rightarrow \hat{\theta}'(y) = h(\hat{\theta}(x))$.

Un dernier argument plaide en faveur de l'EMV : la vraisemblance maximisée constitue l'ingrédient fondamental de la plupart des techniques de *sélection de modèle* : assortie d'un facteur de pénalisation lié au nombre de degrés de liberté du modèle [44, 1], l'*estimé* de $\ell(x_1, \dots, x_n | \hat{\theta}_n)$ fournit un diagnostic utile, supplémentaire aux résultats de tests statistiques, pour évaluer la pertinence d'un modèle par rapport à un autre sur un même jeu de données. Nous renvoyons le lecteur intéressé par ce sujet à l'ouvrage spécialisé [34].

B Descriptif de quelques modèles statistiques utiles

Les tableaux suivants sont extraits d'un formulaire proposé par Aimé Lachal (Univ. Lyon).

B.1 Lois discrètes

<i>distribution</i>	<i>loi de probabilité</i>	$\mathbb{E}(X)$	$\text{var}(X)$	<i>fonction génératrice</i> $\mathbb{E}(z^X)$
Bernoulli	$\mathbb{P}(X = 0) = q, \mathbb{P}(X = 1) = p$ $q = 1 - p$	p	pq	$pz + q$
Binomiale $\mathcal{B}(n, p)$	$\mathbb{P}(X = k) = C_n^k p^k q^{n-k}$ $q = 1 - p, \quad k = 0, 1, \dots, n$	np	npq	$(pz + q)^n$
Poisson $\mathcal{P}(\lambda)$	$\mathbb{P}(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}$ $k = 0, 1, \dots$	λ	λ	$e^{\lambda(z-1)}$
Géométrique $\mathcal{G}(p)$	$\mathbb{P}(X = k) = pq^{k-1}$ $q = 1 - p, \quad k = 1, 2, \dots$	$\frac{1}{p}$	$\frac{q}{p^2}$	$\frac{pz}{1 - qz}$
Hypergéométrique $\mathcal{H}(N, n, p)$	$\mathbb{P}(X = k) = \frac{C_{Np}^k C_{Nq}^{n-k}}{C_N^n}$ $q = 1 - p$ $\max(0, n - Nq) \leq k \leq \min(Np, n)$	np	$npq \frac{N-n}{N-1}$	$\frac{C_{Nq}^n}{C_N^n} F(-n, -Np; Nq - n + 1; z)$
Binomiale négative	$\mathbb{P}(X = k) = C_{k+r-1}^{r-1} p^r q^k$ $q = 1 - p, \quad k = 0, 1, \dots$	$\frac{rq}{p}$	$\frac{rq}{p^2}$	$\left(\frac{p}{1 - qz} \right)^r$
Pascal	$\mathbb{P}(X = k) = C_{k-1}^{r-1} p^r q^{k-r}$ $q = 1 - p, \quad k = r, r + 1, \dots$	$\frac{r}{p}$	$\frac{rq}{p^2}$	$\left(\frac{pz}{1 - qz} \right)^r$

$$\text{Fonction hypergéométrique : } F(a, b; c; z) = \sum_{n=0}^{+\infty} \frac{a(a+1) \dots (a+n-1) b(b+1) \dots (b+n-1) z^n}{c(c+1) \dots (c+n-1) n!}$$

- La somme de n v.a. indépendantes suivant la loi de Bernoulli de paramètre p suit une loi binomiale $\mathcal{B}(n, p)$.
- La somme de deux v.a. indépendantes suivant les lois binomiales $\mathcal{B}(m, p)$ et $\mathcal{B}(n, p)$ suit la loi binomiale $\mathcal{B}(m+n, p)$.
- La somme de deux v.a. indépendantes suivant les lois de Poisson $\mathcal{P}(\lambda)$ et $\mathcal{P}(\mu)$ suit la loi de Poisson $\mathcal{P}(\lambda + \mu)$.
- La somme de deux v.a. indépendantes suivant les lois binomiales négatives de paramètres (r, p) et (s, p) suit la loi binomiale négative de paramètres $(r + s, p)$.
- La somme de r v.a. indépendantes suivant la loi géométrique $\mathcal{G}(p)$ suit la loi de Pascal de paramètres (r, p) .

B.2 Lois continues

distribution	loi de probabilité	$\mathbb{E}(X)$	$\text{var}(X)$	fonction caract. $\mathbb{E}(e^{itX})$
Uniforme $\mathcal{U}(a, b)$	$\frac{1}{b-a} \mathbb{1}_{[a,b]}(x)$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\frac{e^{ibt} - e^{iat}}{i(b-a)t}$
Exponentielle $\mathcal{E}(\lambda)$	$\lambda e^{-\lambda x} \mathbb{1}_{\mathbb{R}^+}(x)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	$\frac{\lambda}{\lambda - it}$
Normale $\mathcal{N}(m, \sigma^2)$	$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-m)^2}{2\sigma^2}\right)$	m	σ^2	$e^{imt - \frac{1}{2}\sigma^2 t^2}$
Weibull $\mathcal{W}(\lambda, a)$	$\lambda a x^{a-1} e^{-\lambda x^a} \mathbb{1}_{]0, +\infty[}(x)$	$\lambda^{-\frac{1}{a}} \Gamma\left(\frac{1}{a} + 1\right)$	$\lambda^{-\frac{2}{a}} [\Gamma\left(\frac{2}{a} + 1\right) - \Gamma\left(\frac{1}{a} + 1\right)^2]$	
Cauchy $\mathcal{C}(a, b)$	$\frac{a}{\pi(a^2 + (x-b)^2)}$	non définie	non définie	$e^{ibt - a t }$
Gamma $\Gamma(a, \lambda)$	$\frac{\lambda^a}{\Gamma(a)} x^{a-1} e^{-\lambda x} \mathbb{1}_{]0, +\infty[}(x)$	$\frac{a}{\lambda}$	$\frac{a}{\lambda^2}$	$\left(\frac{\lambda}{\lambda - it}\right)^a$
Bêta $B(a, b)$	$\frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1} \mathbb{1}_{]0,1[}(x)$	$\frac{a}{a+b}$	$\frac{ab}{(a+b)^2(a+b+1)}$	$M(a, a+b; it)$
Khi-Deux $\chi^2(n)$	$\frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} \mathbb{1}_{]0, +\infty[}(x)$	n	$2n$	$(1-2it)^{-n/2}$
Student $\mathcal{T}(n)$	$\frac{\Gamma(\frac{n+1}{2})}{\sqrt{\pi n} \Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}$	0 si $n > 1$	$\frac{n}{n-2}$ si $n > 2$	$\frac{2}{\Gamma(\frac{n}{2})} \left(\frac{ t \sqrt{n}}{2}\right)^{\frac{n}{2}} K_{\frac{n}{2}}(t \sqrt{n})$
Fisher $\mathcal{F}(m, n)$	$\frac{m^{\frac{m}{2}} n^{\frac{n}{2}}}{B(\frac{m}{2}, \frac{n}{2})} \frac{x^{\frac{m}{2}-1}}{(mx+n)^{\frac{m+n}{2}}} \mathbb{1}_{]0, +\infty[}(x)$	$\frac{n}{n-2}$ si $n > 2$	$\frac{2n^2(m+n-2)}{m(n-4)(n-2)^2}$ si $n > 4$	$M\left(\frac{m}{2}; -\frac{n}{2}; -\frac{n}{m}it\right)$

Fonction Gamma : $\Gamma(a) = \int_0^{+\infty} x^{a-1} e^{-x} dx$

Fonction Bêta : $B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx$

Fonction de Kummer : $M(a; b; z) = \sum_{n=0}^{+\infty} \frac{a(a+1)\dots(a+n-1)}{b(b+1)\dots(b+n-1)} \frac{z^n}{n!}$

Fonction de Bessel modifiée : $K_\nu(z) = \frac{\pi}{2} \frac{I_{-\nu}(z) - I_\nu(z)}{\sin \pi \nu}$ où $I_\nu(z) = \left(\frac{z}{2}\right)^\nu \sum_{n=0}^{+\infty} \frac{1}{n! \Gamma(n+\nu+1)} \left(\frac{z^2}{4}\right)^n$

- La somme de n v.a. indépendantes suivant la loi exponentielle $\mathcal{E}(\lambda)$ suit la loi Gamma $\Gamma(n, \lambda)$.
- La somme de deux v.a. indépendantes suivant les lois Gamma $\Gamma(a, \lambda)$ et $\Gamma(b, \lambda)$ suit la loi Gamma $\Gamma(a+b, \lambda)$.
- Si les v.a. indépendantes X et Y suivent les lois Gamma $\Gamma(a, \lambda)$ et $\Gamma(b, \lambda)$, alors $\frac{X}{X+Y}$ suit la loi Bêta $B(a, b)$.
- La somme de deux v.a. indépendantes suivant les lois normales $\mathcal{N}(m_1, \sigma_1^2)$ et $\mathcal{N}(m_2, \sigma_2^2)$ suit la loi normale $\mathcal{N}(m_1 + m_2, \sigma_1^2 + \sigma_2^2)$.
- Le quotient de deux variables indépendantes suivant la loi normale $\mathcal{N}(0, 1)$ suit la loi de Cauchy $\mathcal{C}(1, 0) = \mathcal{T}(1)$.
- La somme des carrés de n v.a. indépendantes suivant la loi normale $\mathcal{N}(0, 1)$ suit la loi du Khi-Deux $\chi^2(n) = \Gamma(\frac{n}{2}, \frac{1}{2})$.
- Si les v.a. indépendantes X et Y suivent les lois normale $\mathcal{N}(0, 1)$ et du Khi-Deux $\chi^2(n)$, alors $\frac{X}{\sqrt{Y/n}}$ suit la loi de Student $\mathcal{T}(n)$.
- Si les v.a. indépendantes X et Y suivent les lois du Khi-Deux $\chi^2(m)$ et $\chi^2(n)$, alors $\frac{mX}{nY}$ suit la loi de Fisher $\mathcal{F}(m, n)$.

C Éléments sur la simulation pseudo-aléatoire

Les générateurs pseudo-aléatoires sont des éléments central des méthodes de simulation : elles reposent toutes sur la transformation de variables uniformes $\mathcal{U}(0, 1)$.

Définition 56 Générateur pseudo-aléatoire. *Un générateur pseudo-aléatoire est une transformation déterministe Ψ de $]0, 1[$ dans $]0, 1[$ telle que, pour toute valeur initiale u_0 et tout n , la suite*

$$\{u_0, \Psi(u_0), \Psi(\Psi(u_0)), \dots, \Psi^n(u_0)\}$$

a le même comportement statistique qu'une suite iid $\mathcal{U}(0, 1)$.

Sans appel au "hasard", la suite déterministe $(u_0, u_1 = \Psi(u_0), \dots, u_n = \Psi(u_{n-1}))$ doit ressembler à une suite aléatoire.

- En **Python**, il faut faire appel à la procédure `random.seed()`

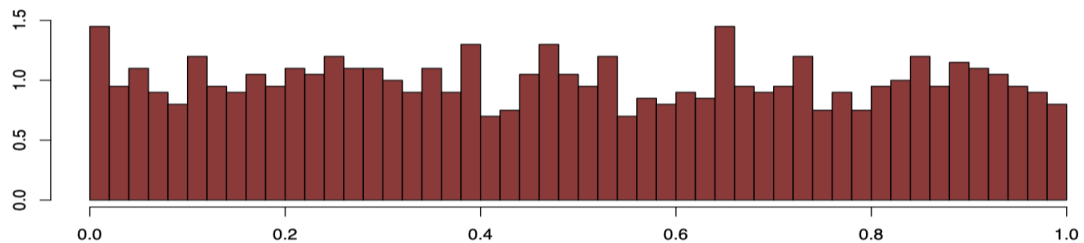
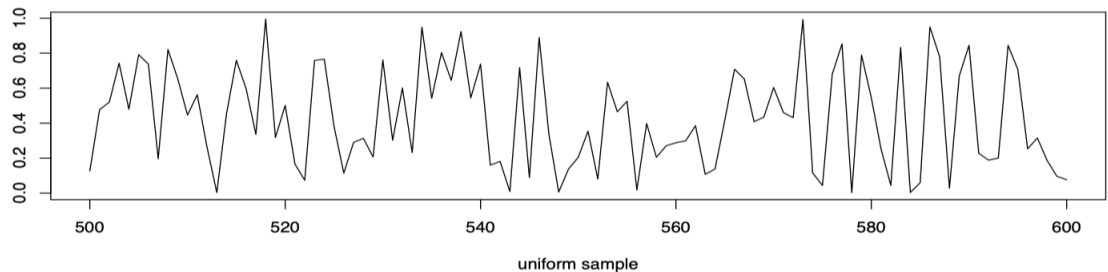
Description :

'random.seed(a=None, version=2)' generates pseudo-random values with seed 'a'.

Exemple :

`u = random.seed(20)`

Generally, the seed value is the previous number generated by the generator. However, When the first time you use the random generator, there is no previous value. So by-default current system time is used as a seed value.



- En **C**, il faut faire appel à la procédure `rand() / random()`

SYNOPSIS

```
# include <stdlib.h>
long int random(void);
```

DESCRIPTION

The random() function uses a non-linear additive feedback random number generator employing a default table of size 31 long integers to return successive pseudo-random numbers in the range from 0 to RAND MAX. The period of this random generator is very large, approximately $16 \cdot ((2^{31}) - 1)$.

RETURN VALUE

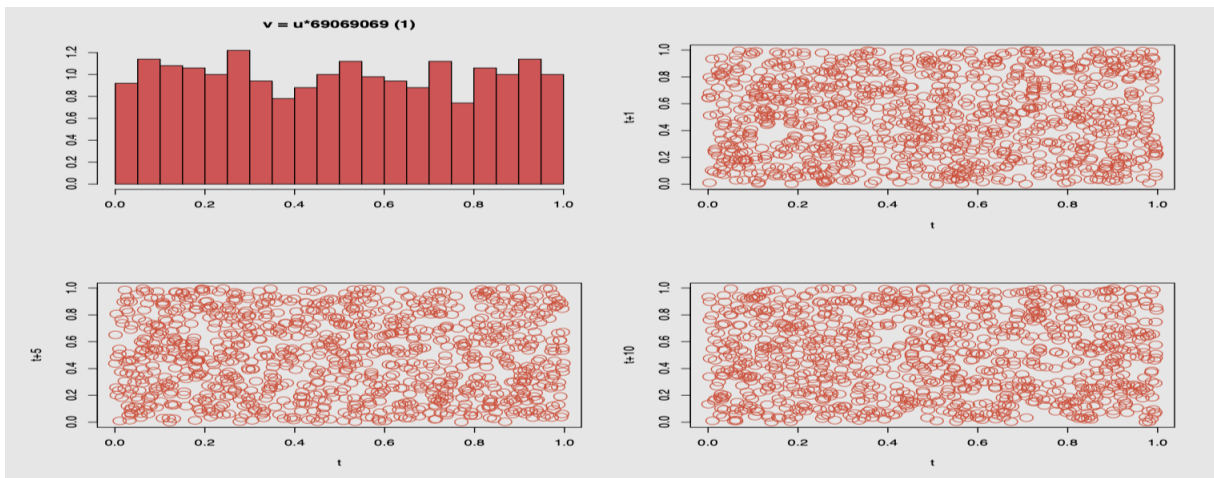
random() returns a value between 0 and RAND MAX.

Un générateur usuel est le suivant :

Définition 57 *Générateur congruenciel* Le générateur congruenciel

$$D(x) = (ax + b) \bmod (M + 1).$$

est de période M pour les bons choix de (a, b) et se transforme en générateur sur $]0, 1[$ par division par $M + 2$.



Il faut toujours utiliser la fonction appropriée sur l'ordinateur ou le logiciel en service plutôt que de construire un générateur aléatoire de mauvaise qualité.

D Rappels sur les chaînes de Markov

Rappelons les deux définitions exprimées dans le cours.

Définition 58 Noyau de transition. Une chaîne de Markov homogène est déterminée par un noyau de transition, défini sur $\Theta \times \mathcal{B}(\Theta)$ à l'itération i par

$$\mathcal{K}(\theta|A) = P(\theta^{(i)} \in A | \theta^{(i-1)} = \theta) = \int_A \underbrace{\kappa(\theta, \tilde{\theta})}_{\text{densité de transition sur } \tilde{\theta}} d\tilde{\theta},$$

telle que $\mathcal{K}(\cdot|A)$ est mesurable $\forall A \in \mathcal{B}(\Theta)$. Cette notion de noyau généralise au cadre continu celle de matrice de transition d'un état à un autre dans un cadre discret.

Toute la structure d'une chaîne de Markov, que l'on considèrera toujours d'ordre 1 dans ce cours, dépend seulement du choix d'un noyau de transition et de l'état initial (ou la distribution initiale) de la chaîne, comme l'exprime la définition suivante.

Définition 59 Chaîne de Markov. Sachant un noyau de transition \mathcal{K} , une suite $\theta_0, \dots, \theta_n, \dots$ de variables aléatoires est une chaîne de Markov d'ordre 1 si, $\forall n \geq 0$, la distribution de θ_n conditionnelle à la σ -algèbre (filtration) générée par $\theta_{n-1}, \theta_{n-2}, \dots, \theta_0$ est la même que celle de $\theta_n | \theta_{n-1}$:

$$\pi(\theta_n \in \mathcal{A} | \theta_{n-1}, \theta_{n-2}, \dots, \theta_0) = \pi(\theta_n \in \mathcal{A} | \theta_{n-1}), = \mathcal{K}(\theta_{n-1} | \mathcal{A}).$$

EXEMPLE 36. Marche aléatoire. Soit $\Theta \in \mathbb{R}^p$. La marche aléatoire (random walk) gaussienne est une chaîne de Markov de noyau $\mathcal{K}(\theta|\cdot)$ associé à la distribution $\mathcal{N}_p(\theta, \tau^2 I_p)$:

$$\theta_{n+1} = \theta_n + \tau \epsilon_n \text{ avec } \epsilon_n \sim \mathcal{N}(0, 1).$$

L'irréductibilité est une mesure de la sensibilité de la chaîne de Markov aux conditions initiales, qui fournit une garantie de convergence de cette chaîne : tout ensemble de Θ a une chance d'être visité par la chaîne de Markov.

Proposition 15 Irréductibilité.

- Si Θ est discret, la chaîne est irréductible si tous les états communiquent :

$$P_\theta(\tau_{\theta'} < \infty) > 0 \quad \forall (\theta, \theta') \in \Theta^2$$

où $\tau_{\theta'}$ est le premier temps (> 0) de visite de θ' .

- Si Θ est continu, la chaîne est irréductible pour une mesure ψ si, $\forall \theta \in \Theta$ et pour presque tout $\mathcal{A} \in \mathcal{B}(\Theta)$ avec $\psi(\mathcal{A}) > 0$, $\forall n < \infty$

$$\mathcal{K}^n(\theta | \mathcal{A}) > 0.$$

L'irréductibilité est une condition trop faible pour être sûr que $(\theta_n)_n$ visite suffisamment de fois n'importe quel sous-ensemble $\mathcal{A} \in \mathcal{B}(\Theta)$. Il faut également vérifier des conditions de stabilité pour garantir une approximation acceptable de la distribution-cible : la notion de *réccurrence* formalise de telles conditions.

Pour un espace d'état Θ discret, la récurrence d'un état est équivalent à une probabilité 1 de retour certain en cet état. Elle est dite **récurrente positive** lorsque le temps moyen de retour est fini (récurrente nulle sinon). Lorsque les chaînes sont irréductibles et si Θ est fini (borné), l'irréductibilité implique la récurrence positive. De manière plus générale, on oppose la *Harris-réccurrence* à la *transcience* :

Définition 60 Harris-réccurrence et transcience. Un ensemble \mathcal{A} est Harris-récurrent si

$$P_\theta(\eta_{\mathcal{A}} = \infty) = 1 \quad \forall \theta \in \Theta \tag{35}$$

où $\eta_{\mathcal{A}}$ désigne le nombre de visites dans un ensemble \mathcal{A} . La propriété (35) implique que

$$\mathbb{E}[\eta_{\mathcal{A}}] = \infty.$$

La transience correspond à la propriété contraposée :

$$\mathbb{E}[\eta_{\mathcal{A}}] < \infty.$$

L'étude d'une chaîne de Markov est aussi l'étude de son éventuelle *mesure invariante* π :

$$\theta_{n+1} \sim \pi \quad \text{si} \quad \theta_n \sim \pi.$$

Définition 61 Mesure invariante. π est invariante par $\mathcal{K}(\theta|\mathcal{A})$ si

$$\pi(\mathcal{A}) = \int_{\Theta} \mathcal{K}(\theta, \mathcal{A}) d\Pi(\theta) \quad \forall \mathcal{A} \in \mathcal{B}(\Theta).$$

On peut comprendre l'invariance de π de la façon suivante : soit $\pi_{\mu}(\theta_n \in \cdot)$ la loi de θ à l'étape n de la chaîne, μ désignant la loi de départ de cette chaîne. Si une mesure limite γ_{μ} existe telle que, $\forall \mathcal{A} \in \mathcal{B}(\Theta)$,

$$\pi_{\mu}(\theta_n \in \mathcal{A}) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \gamma_{\mu}(\mathcal{A}),$$

alors

$$\begin{aligned} \gamma_{\mu}(\mathcal{A}) &= \lim_{n \rightarrow \infty} \int_{\Theta} \mu(d\theta) \mathcal{K}^n(\theta, \mathcal{A}), \\ &= \lim_{n \rightarrow \infty} \int_{\Theta} \int_{\Theta} \mathcal{K}^{n-1}(\theta, d\theta) \mathcal{K}(\theta, \mathcal{A}), \\ &= \int_{\Theta} \gamma_{\mu}(d\theta) \mathcal{K}(\theta, \mathcal{A}) \end{aligned}$$

car la convergence de $\int_{\Theta} \mu(d\theta) \mathcal{K}^n(\theta, \cdot)$ implique la convergence des intégrales de fonctions mesurables bornées. Ainsi, si la chaîne de Markov converge vers une distribution limite, il s'agit d'une mesure invariante.

Cette mesure-limite, invariante (dite aussi *stationnaire*), peut présenter une propriété tout à fait intéressante : elle peut être **ergodique**, c'est-à-dire indépendante de la loi initiale μ . Cette propriété se traduit de façon générale par une **convergence en norme en variation totale** entre la mesure $\mathcal{K}^n(\theta, \cdot)$ et la loi stationnaire $\pi(\theta|\dots)$, qui peut être raffinée dans le cas où Θ est discret. Dans le cas discret et fini, les équations de Chapman-Kolmogorov permettent d'obtenir la mesure-limite.

Définition 62 Norme en variation totale entre mesures. Soit (μ_1, μ_2) deux mesures sur \mathcal{A} . Alors la norme en variation totale entre μ_1 et μ_2 est

$$\|\mu_1 - \mu_2\|_{TV} = \sup_{\mathcal{A}} |\mu_1(\mathcal{A}) - \mu_2(\mathcal{A})|.$$

Théorème 28 Ergodicité d'une chaînes de Markov. Si $(\theta_n)_n$ est Harris récurrente positive et apériodique, alors, pour presque toute distribution initiale μ , la loi stationnaire $\pi(\cdot)$ est approchée par la convergence suivante :

$$\lim_{n \rightarrow \infty} \left\| \int_{\Theta} \mathcal{K}^n(\theta, \cdot) \mu(d\theta) - \pi(\cdot) \right\|_{TV} = 0$$

Cette convergence en variation totale implique que pour presque toute fonction bornée $h : \Theta \rightarrow \mathbb{R}$,

$$\lim_{n \rightarrow \infty} |\mathbb{E}_{\mu}[h(\theta_n)] - \mathbb{E}_{\pi}[h(\theta)]| = 0.$$

On déduit de ce résultat le théorème ergodique (théorème 22), qui pour permettre l'obtention d'un théorème de limite centrale nécessite que la chaîne de Markov soit *réversible* :

$$\theta_{n+1}|\theta_{n+2} = x \sim \theta_{n+1}|\theta_n = x.$$

E Calcul bayésien avec OpenBUGS et JAGS

E.1 Contexte de développement

WinBUGS et son successeur OpenBUGS font partie du projet BUGS (*Bayesian inference Using Gibbs Sampler*) qui vise à rendre simple la pratique des méthodes MCMC aux statisticiens. Il a été développé par l'université de Cambridge. **Seul OpenBUGS est actuellement maintenu.**

WinBUGS et OpenBUGS peuvent être utilisés de différentes manières :

- Via une interface « clique-bouton » qui permet de contrôler l'analyse,
- En utilisant des modèles définis par des interfaces graphiques, appelés `DoddleBUGS`,
- Via d'autres logiciels tels que R (en particulier via le package `R2WinBUGS`).

WinBUGS et OpenBUGS sont des logiciels libres et gratuits, Cependant, afin d'accéder à la version non restreinte de WinBUGS, il est nécessaire d'obtenir la clé d'utilisation. Le site internet pour WinBUGS et OpenBUGS, **The BUGS Project** présente les deux logiciels et fournit de la documentation. Le site spécialisé pour OpenBUGS est <http://www.openbugs.net/>

JAGS est une version "rapide" de BUGS développée par Martyn Plummer. Il repose sur le même langage, à quelques différences subtiles près. Il a la particularité de ne pas présenter d'interface graphique pour la gestion des chaînes (mais il peut aussi être appelé par R). Traditionnellement, il possède moins de distributions de probabilité que sous OpenBUGS.

D'autres outils logiciels. Régulièrement, d'autres outils logiciels sont élaborés permettant de s'attaquer à des modèles à état latents, ou/et de grande dimension, mettant en oeuvre des MCMC accélérées... Citons le package Python `ELFI`²², le langage `STAN` ou encore le langage `Birch`²³.

E.2 Un exemple "fil rouge" : le modèle bêta-binomial

Dans une parcelle, on compte le nombre d'arbres de l'espèce A. On répète l'opération en J parcelles. On veut modéliser ce processus d'échantillonnage.

On se place dans une parcelle i donnée. Si on suppose que l'espèce de chaque arbre est indépendante de l'espèce de ses voisins, on peut modéliser le nombre d'arbre de l'espèce A par une loi binomiale de paramètres p et N .

$$Y_i \sim \mathcal{B}(N_i, p_i) \quad (1.1)$$

avec N_i le nombre d'arbres de la parcelle, p_i la proportion inconnue d'arbre de type A.

Si on suppose en prime que toutes les parcelles sont équivalentes et que la proportion d'arbres d'espèce A est la même dans toutes les parcelles alors on ajoute l'hypothèse

$$p_i = p \quad \text{pour tous les } i \quad (1.2)$$

Enfin si on suppose que les parcelles sont indépendantes les unes des autres, le modèle s'écrit

$$Y_i \stackrel{i.i.d}{\sim} \mathcal{B}(N_i, p) \quad (1.3)$$

22. <https://elfi.readthedocs.io/en/latest/usage/tutorial.html>

23. <https://birch-lang.org/talks/automated-learning-slides.pdf>

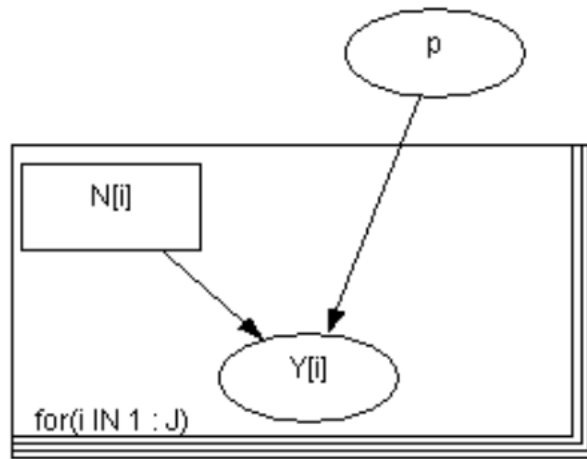


FIG. 1.1 – DAG du modèle binomial

Que sait-on sur le paramètre p qui régit la proportion des arbres de type A dans nos parcelles ?

1. A peu près rien. C'est une proportion donc ça prend des valeurs continues entre 0 et 1. on peut encoder ça au travers d'une loi de probabilité uniforme sur $[0; 1]$.
2. on a une expérience d'autres étude set on se souvient qu'en moyenne p tourne autour de 10% mais avec une certaine variabilité. On décide d'encoder cette connaissance par une loi Beta.

$$p \sim \beta(1, 9)$$

3. On sait d'expérience que p ne dépasse "jamais" 0.2, et qu'en moyenne il vaut 0.1. Ca signifie que $\mathbb{E}(p) = 0.1$, or si $p \sim \beta(a, b)$ alors $\mathbb{E}(p) = a/(a + b)$, donc $b = 9a$. On "résoud" ensuite l'inégalité qui dit que p est très rarement plus grand que 0.2, i.e $\mathbb{P}(p > 0.2) < 0.001$, par exemple par dichotomie dans R.

On peut mener un calcul *a posteriori* par conjugaison. Avec

$$p \sim \beta(a, b) \quad \text{et} \quad y|p \sim \mathcal{B}(N, p) \quad \text{alors}$$

$$p|y \sim \beta(a + y, b + N - y).$$

E.3 Fonctionnement résumé d'OpenBUGS

1. Ouverture de la fenêtre OpenBUGS
2. Création d'un fichier "modele.txt" contenant l'écriture formelle de la *vraisemblance* et de la *distribution a priori*
 - nécessité d'utiliser une boucle sur les données pour la vraisemblance

- langage BUGS différent de R (mais plutôt compréhensible)
3. Création d'un fichier "data.txt" avec des données entrées en vectoriel (possibilité de tableaux)
 - les "données" regroupent aussi les constantes du problème : taille des données n , etc.
 4. Éventuellement création d'un fichier d'initialisation pour les paramètres (chaînes MCMC)

Testons ici l'implémentation du cas d'étude bêta-binomial. Ouvrons un fichier `model-beta-binomial.txt` (à ouvrir comme "Text") :

```
#-----
# A simple binomial model
#-----|

model {

  # Likelihood of the binomial distribution
  for (k in 1:n)
  {
    x[k] ~ dbin(p, N)
  }

  # The prior of the unknown parameter
  p ~ dbeta(a,b)
}
```

1. **Vérification du modèle** via "Model/Model Specification" puis "Check Model "

- \Rightarrow *model is syntactically correct*

2. **Enregistrement des données** via :

- (a) sélection du fichier `data-beta-binomial.txt`

```
list(N=10,
x=c(7, 7, 5, 8, 3, 4, 6, 4, 5, 4, 4, 4, 6, 3, 5, 1, 5, 7, 7, 3),
n=20,
a=12.5,
b=112.5
)
```

- (b) "Load data "

- \Rightarrow *data loaded*

3. **Sélection du nombre de chaînes MCMC puis compilation du modèle** via "compile "

- \Rightarrow *model compiled*

4. **Initialisation des chaînes MCMC** : 2 façons possibles

- (a) Sélection dans le fichier `init-beta-binomial.txt` puis "load inits", chaîne par chaîne
- (b) Génération automatique via "gen inits" pour toutes les chaînes

- \Rightarrow *model initialized*

5. **Ouverture de la fenêtre de monitoring des chaînes** via "Inference/Sample Monitor Tool "

- Écrire "p" dans la fenêtre "node", valider avec "set "

6. Ouverture de la fenêtre de lancement des chaînes via "Model/Update"

- cliquer sur "update" pour lancer une première fois les chaînes
- \Rightarrow *model is updating*

7. Monitorer les chaînes via la fenêtre consacrée :

- Aller chercher "p" dans la fenêtre "node", puis cliquer sur :
- "trace" pour tracer l'évolution des chaînes associées à p
- "trace" pour tracer la densité *a posteriori* courante (approximative)
- "coda" pour récupérer les chaînes
- "stats" pour obtenir un résumé statistique de la loi *a posteriori* courante
- etc.

E.4 Quelques détails supplémentaires concernant OpenBUGS

Menu "Inference".

1. Sous-menu "Correlation"

- "Correlation Tool"
 - scatter \Rightarrow trace un "scatterplot" entre 2 dimensions
 - matrix \Rightarrow dessine la matrice de corrélation (par niveaux de gris)
 - print \Rightarrow calcule le coefficient de corrélation linéaire

2. Sous-menu "Compare"

- "Comparison Tool"
 - boxplot \Rightarrow trace une "boîte à moustaches" d'une dimension sélectionnée

Menu "Model".

- Commande "*latex*" : fournit le code latex du fichier sélectionné (utile pour le fichier de modèle !)

E.5 Liste des distributions de probabilités disponibles

Discrete Univariate

[Bernoulli](#)
[Binomial](#)
[Categorical](#)
[Negative Binomial](#)
[Poisson](#)
[Geometric](#)
[Geometric \(alternative\)](#)
[Non-central Hypergeometric](#)

Continuous Univariate

[Beta](#)
[Chi-squared](#)
[Double Exponential](#)
[Exponential](#)
[Flat](#)
[Gamma](#)
[Generalized Extreme Value](#)
[Generalized F](#)
[Generalized Gamma](#)
[Generalized Pareto](#)
[Generic LogLikelihood Distribution](#)
[Log-normal](#)
[Logistic](#)
[Normal](#)
[Pareto](#)
[Student-t](#)
[Uniform](#)
[Weibull](#)

Discrete Multivariate

[Multinomial](#)

Continuous Multivariate

[Dirichlet](#)
[Multivariate Normal](#)
[Multivariate Student-t](#)
[Wishart](#)

E.6 Noeuds logiques et indexation

Les noeuds logiques sont définis par une flèche et sont **toujours indexés** :

```
mu[i] <- beta0 + beta1 * z1[i] + beta2 * z2[i] + b[i]
```

On peut utiliser une fonction de lien (log, logit, probit) :

```
logit(mu[i]) <- beta0 + beta1 * z1[i] + beta2 * z2[i] + b[i]
```

On peut définir des tableaux :

```
Y[(i + j) * k, 1]
```

Par ailleurs, toute variable (noeud logique ou stochastique ~) ne peut apparaître qu'une fois dans la partie gauche d'une expression (sauf dans le cas d'une transformation de données) du type :

```
for (i in 1:N) {  
  z[i] <- sqrt(y[i])  
  z[i] ~ dnorm(mu, tau)  
}
```

Enfin, on peut créer des noeuds multiparités : soient μ et τ deux vecteurs de taille K . On peut alors définir la boucle suivante :

```
for (i in 1 : I) {  
  x[i, 1 : K] ~ dnmnorm(mu[, ], tau[, ])  
}
```

L'aide sur les fonctions utiles est disponible ici : <http://www.openbugs.net/Manuals/ModelSpecification.html>

E.7 Pièges à éviter

Deux pièges peuvent fréquemment survenir :

- **Se tromper de paramétrisation.** Il faut toujours aller vérifier de quelle façon les distributions de probabilité sont définies ! C'est piégeant en particulier pour les lois normale et log-normale :

Log-normal

$x \sim \text{dlnorm}(\mu, \tau)$ $\sqrt{\frac{\tau}{2\pi}} \frac{1}{x} \exp\left(-\frac{\tau}{2}(\log x - \mu)^2\right); \quad x > 0$

Normal

$x \sim \text{dnorm}(\mu, \tau)$ $\sqrt{\frac{\tau}{2\pi}} \exp\left(-\frac{\tau}{2}(x - \mu)^2\right); \quad -\infty < x < \infty$

- **Confondre censure et troncature.** La *censure* est possible en utilisant la notation suivante :

```
x ~ dnorm(mu, tau)C(lower, upper)
```

Il s'agit d'une censure par intervalle. On laisse un blanc à gauche (*resp.* à droite) si la donnée est censurée à droite (*resp.* à gauche). La *truncation* est possible en utilisant la notation suivante :

```
x ~ dnorm(mu, tau)T(lower, upper)
```

Notons enfin que Les priors impropres (non intégrables) ne sont pas utilisables en BUGS. Il faut les approcher avec des distributions propres mais de variance très large (ce qui est dangereux hélas). Une "règle du pouce" pour les paramètres de variance inverse est la suivante :

$$\tau \sim \text{dgamma}(0.001, 0.001)$$

mais l'usage d'OpenBUGS, JAGS, etc. nécessite donc de toujours mener des études de sensibilité *a posteriori*. Enfin, l'implémentation d'une nouvelle distribution est possible, mais il faut utiliser des subterfuges, typiquement par une transformation de variables (ex : Box-Müller pour simuler des gaussiennes).

E.8 Utilisation d'OpenBUGS avec R

Plusieurs packages sont disponibles pour appeler un code OpenBUGS depuis R : BRugs, R2WinBUGS et R2OpenBUGS. Nous recommandons l'usage du dernier, qui nécessite d'installer également le package CODA (gestion des MCMC) :

1. Installer le package R2OpenBUGS (nécessite CODA).
2. Le charger dans un programme R avec `library(R2OpenBUGS)`.
3. Définir le répertoire où se trouve les fichiers BUGS comme répertoire de travail courant, par exemple :

```
> setwd("E:/RepertoireTravail Courant/Docs Perso")
> getwd()
[1] "E:/RepertoireTravail Courant/Docs Perso"
```

`getwd()` sert à vérifier que le changement de répertoire a bien été pris en compte et doit retourner le chemin du répertoire de travail courant.

4. Spécification du modèle :

```
# Définir le nom du fichier contenant le modèle BUGS
filename <- "model-beta-binomial.txt"

# Définir les données
donnees <- list(N=10,
x=c(7, 7, 5, 8, 3, 4, 6, 4, 5, 4, 4, 4, 6, 3, 5, 1, 5, 7, 7, 3),
n=20,
a=12.5,
b=112.5
)

# Définir les paramètres
params <- c("p")

# Initialisation des chaînes
inits <- function()
{
  list(p=0.5)
  list(p=0.1)
  list(p=0.9)
}

# Simulation des chaînes de Markov
out <- bugs(donnees,inits,params,filename,n.chains=3,debug=T,n.iter=1000,working.directory=getwd())

print(out)
plot(out)
```

5. Obtention des résultats : avec `debug=T`, fermer la fenêtre OpenBUGS qui s'est ouverte pour achever le traitement numérique. Le répertoire courant doit se présenter ainsi :

CODAchain1.txt	03/03/2016 04:07	Document texte	12 Ko
CODAchain2.txt	03/03/2016 04:07	Document texte	12 Ko
CODAchain3.txt	03/03/2016 04:07	Document texte	12 Ko
CODAindex.txt	03/03/2016 04:07	Document texte	1 Ko
data.txt	03/03/2016 04:07	Document texte	1 Ko
data-beta-binomial.txt	02/03/2016 17:30	Document texte	1 Ko
exemple-beta-binomial.r	03/03/2016 04:10	Fichier R	1 Ko
init-beta-binomial.txt	02/03/2016 17:31	Document texte	1 Ko
inits1.txt	03/03/2016 04:07	Document texte	1 Ko
inits2.txt	03/03/2016 04:07	Document texte	1 Ko
inits3.txt	03/03/2016 04:07	Document texte	1 Ko
log.odc	03/03/2016 04:07	Microsoft Office D...	15 Ko
log.txt	03/03/2016 04:07	Document texte	1 Ko
model-beta-binomial.txt	02/03/2016 17:22	Document texte	1 Ko
script.txt	03/03/2016 04:07	Document texte	1 Ko

E.9 Détails sur JAGS : exemple de script

```
% Script for the bash execution of the posterior computation

model in model.txt
data in SKJ-jags.r
compile, nchains(3)
initialize

update 10
update 10

monitor M
monitor logP0
monitor Ft
monitor Ff
monitor F
monitor selectivity.at.age
monitor recovery.rate
monitor tau2
monitor tau.selec
monitor mu.selec
monitor pR
monitor sigma

update 100

coda *
```

E.10 D'autres outils (R/Python)

Nimble est un outil permettant de lancer des tâches de calcul bayésien à partir de R. Pour Python 3, on pourra explorer et utiliser Pymc.

Références

- [1] H. Akaike. On entropy maximization principle. In : Krishnaiah, P.R. (Editor). *Applications of Statistics, North-Holland, Amsterdam*, pages 27–41, 1977.
- [2] Anonyme. *Measuring River Eischarge in High Flow (Flood) or High Sediment Concentration Conditions*. Application Note : R&E Instruments – Acoustic Eoppler Current Profilers. Communication Technology Technical Report, 1999.
- [3] D.J. Benjamin, J.O. Berger, and V.E. Johnson. Redefine statistical significance. *Nature Human Behavior*, pages DOI :10.1038/s41562–017–0189–z, 2017.
- [4] N. Bouleau. *Probabilités de l'ingénieur*. Hermann, 1986.
- [5] S.C.Y. Chan, Y. Niv, and K.A. Norman. A probability distribution over latent causes, in the orbitofrontal cortex. *The Journal of Neuroscience*, 36 :7817–7828, 2016.
- [6] E. Cœur and M. Lang. L'information historique des inondations : l'histoire ne donne-t-elle que des leçons ? *La Houille Blanche*, 2 :79–84, 2000.
- [7] M.K. Cowles and B.P. Carlin. Markov Chain Monte Carlo Convergence Diagnostics : A Comparative Review. *Journal of the American Statistical Association*, 91 :883–904, 1996.
- [8] S. Dehaene. *Consciousness and the Brain : Deciphering How the Brain Codes our Thoughts*. Viking Press, 2014.
- [9] M. Evans. Measuring statistical evidence using relative belief. *Computational Structural Biotechnology Journal*, 14 :91–96, 2016.
- [10] M. Evans and H. Moshonov. Checking for prior-data conflict. *Bayesian Analysis*, 1 :893–914, 2006.
- [11] R.A. Fisher. *Statistical Methods for Research Workers*. Oliver and Boyd, Edinburgh, 1926.
- [12] V. Fortuin. Priors in Bayesian Deep Learning : A Review. *International Statistical Review*, 2022.
- [13] M. Galevski. La corrélation entre les pluies torrentielles and l'intensité de l'érosion (avant-propos de P. Reneuve). *Annales de l'École Nationale des Eaux and Foêts and de la station de recherches and expériences*, 14 :379–428, 1955.
- [14] Jakob Gawlikowski, Cedrique Rovile Njiteucheu Tassi, Mohsin Ali, Jongseo Lee, Matthias Humt, Jianxiang Feng, Anna M. Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, M. Shahzad, Wen Yang, Richard Bamler, and Xiaoxiang Zhu. A survey of uncertainty in deep neural networks. *ArXiv*, abs/2107.03342, 2021.
- [15] J.I. Gold and H.R. Heekeren. *Neural Mechanisms for Perceptual Decision Making*. In : Neuroeconomics (chapter 19), P.W. Glimcher and R. Fehr (eds), Second Edition, 2013.
- [16] C. Gourerious and A. Monfort. *Statistique et modèles économétriques*. Economica, Paris, 1996.
- [17] S. Greenland, S.J. Senn, K.J. Rothman, J.B. Carlin, C. Poole, S.N. Goodman, and D. Altman. Statistical tests, p values, confidence interval, and power : a guide to misinterpretations. *European Journal of Epidemiology*, 31 :227–350, 2016.
- [18] V.E. Johnson. Revised standards for statistical evidence. *Proceedings of the National Academy of Science*, 110 :19313–19317, 2013.
- [19] M. Keller, A. Pasanisi, and E. Parent. Réflexions sur l'analyse d'incertitudes dans un contexte industriel : information disponible et enjeux décisionnels. *Journal de la Société Française de Statistiques*, 2012.
- [20] M.Y. Kim and X. Xue. The analysis of multivariate interval-censored survival data. *Statistics in Medicine*, 21 :3715–3726, 2002.

- [21] A.N. Kolmogorov. *Foundations of the Theory of Probability*. Chelsea Publishing Co., Oxford, 1950.
- [22] J.F. Le Gall. *Intégration, Probabilités and Processus Aléatoires*. Cours de l'École Normale Supérieure, 2006.
- [23] H. Lebesgue. *Oeuvres scientifiques (en cinq volumes)*. Institut de Mathématiques de l'Université de Genève, 1972.
- [24] E.L. Lehman. *Fisher, Neyman, and the creation of classical statistics*. New York : Springer, 2011.
- [25] N.R. Mann, R.E. Schafer, and N.D. Singpurwalla. *Methods for Statistical Analysis of Reliability and Life Data*. Wiley Series in Probability and Statistics, 1974.
- [26] J.-M. Marin and C.P. Robert. *Bayesian Core : A Practical Approach to Computational Bayesian Statistics*. Springer, 2007.
- [27] C. Neves and M. Isabel Fraga Alves. Testing extreme value conditions – an overview and recent approaches. *REVSTAT*, 6 :83–100, 2008.
- [28] R. Nuzzo. Scientific method : Statistical errors. *Nature*, 506 :150–152, 2014.
- [29] G.W. Oehlert. A Note on the Delta Method. *The American Statistician*, 46 :27–29, 1992.
- [30] B. O'Neill. Exchangeability, correlation and Bayes' Effect. *International Statistical Review*, 77 :241–250, 2011.
- [31] E. Parent and J. Bernier. *Le raisonnement bayésien. Modélisation and inférence*. Springer, 2007.
- [32] A. Pasanisi, M. Keller, and E. Parent. Estimation of a quantity of interest in uncertainty analysis : some help from Bayesian Decision Theory. *Reliability Engineering and System Safety*, 100 :93–101, 2012.
- [33] O. Payastre. Utilité de l'information historique pour l'étude du risque de crues. *14ième Journées Scientifiques de l'Environnement : l'Eau, la Ville, la Vie, 12-13 mai*, 2003.
- [34] J. Planzalg and R. Hamböcker. *Parametric statistical theory*. Walter de Gruyter, Berlin, 1994.
- [35] A. Pouget, J.M. Beck, W.J. Ma, and P.E. Latham. Probabilistic brains : knowns and unknowns. *Nature Neuroscience*, 16 :1170–1178, 2016.
- [36] S.J. Press. *Subjective and Objective Bayesian Statistics (second edition)*. Wiley : New York, 2003.
- [37] D.S. Reis and J.R. Stedinger. Bayesian mcmc flood frequency analysis with historical information. *Journal of Hydrology*, 313 :97–116, 2005.
- [38] C.P. Robert. *The Bayesian Choice : From Decision-Theoretic Foundations to Computational Implementation (2nd edition)*. Springer, 2007.
- [39] C.P. Robert and G. Casella. *Monte Carlo Statistical Methods (second edition)*. Springer, 2004.
- [40] A. Sabourin. Semi-parametric modeling of excesses above high multivariate thresholds with censored data. *Journal of Multivariate Analysis*, 136 :126–146, 2015.
- [41] E. Salinas. Prior and prejudice. *Nature Neuroscience*, 14 :943–945, 2011.
- [42] L. Sanders. The probabilistic mind. *Science News*, 180 :18, 2011.
- [43] G. Saporta. *Probabilités, analyses des données and statistiques*. Technip, 2006.
- [44] Gideon E. Schwarz, H. Estimating the dimension of a model. *Annals of Statistics*, 6 :461–464, 1978.
- [45] J. Sprenger. *Bayésianisme versus fréquentisme en inférence statistique*. I. Drouet (ed.). Éditions Matériologiques, Paris, 2017.

- [46] M.A. Stephens. Edf statistics for goodness of fit and some comparisons. *Journal of the American Statistical Association*, 69 :730–737, 1974.
- [47] R. van de Schoot, S. Depaoli, R. King, B. Kramer, K. Martens, M.G. Tadesse, M. Vannucci, A. Gelman, D. Veen, J. Willemsen, and C. Yau. Bayesian statistics and modelling. *Nature Reviews. Methods Primer*, 2021.
- [48] A.W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.
- [49] Cochran W.G. The χ^2 test of goodness of fit. *Annals of Mathematical Statistics*, 23 :315–345, 1952.
- [50] W. Xie and Barton R.B. Nelson, B.L. Multivariate input uncertainty in output analysis for stochastic simulation. *soumis*, 2016.