

Université Paris-Sud
Laboratoire de Mathématiques

Mémoire de Stage de Master

**Redondance bayésienne et minimax,
Sources stationnaires sans mémoire en
alphabet infini**

Dominique Bontemps

14 septembre 2007

Sous la direction de Mme Elisabeth Gassiat

Remerciements

Pour l'écriture de ce document je suis particulièrement redevable à Elisabeth Gassiat, mon maître de stage, qui m'a guidé dans la compréhension des différents articles que j'ai étudiés, et dans l'obtention des applications nouvelles qui sont présentées ici.

Merci aussi à mes camarades qui ont écoutées mes questions techniques et m'ont aidé à les résoudre, en particulier à Laurent Veysseire et à Fabrice Mathurin.

Résumé

On dispose d'une famille de lois de probabilité (P_θ) appelées *sources*, qui produisent des signaux indépendants identiquement distribués sur un alphabet A dénombrable, et on désire trouver un code qui compresse sans perte ces signaux de manière optimale. La qualité de la compression est mesurée par la *redondance*, ou entropie relative.

Le but de ce mémoire est d'explorer dans l'optique du codage / compression de données les applications possibles des encadrements de la *redondance bayésienne* et de la *redondance minimax* fournis par l'article « Mutual Information, Metric Entropy and Cumulative Relative Entropy Risk » de Haussler et Opper ([HO97]). Ceux-ci exploitent des transformées de Laplace de la *distance de Hellinger*, et peuvent être reliées à l'*entropie métrique* de la famille des sources pour cette même distance.

En combinant ces résultats avec des techniques issues des travaux de Xie et Barron sur le mélange de Dirichlet en alphabet fini, nous serons en mesure de fournir un résultat asymptotique nouveau, en exprimant un équivalent de la redondance minimax de certaines classes de lois en alphabet infini : les *classes enveloppe à décroissance exponentielle*.

Table des matières

1	Introduction	5
1.1	Longueur de code	6
1.2	Redondance d'un code	7
2	Encadrer la Redondance bayésienne	10
2.1	α -affinité, I-divergence et D-divergence	10
2.1.1	α -affinité	10
2.1.2	I-divergence et D-divergence	13
2.2	Les résultats	14
2.3	Exemple d'application : Les familles dénombrables	22
3	Encadrer la Redondance minimax	24
3.1	Entropie métrique	24
3.2	Les théorèmes sur la redondance minimax	27
4	Classes enveloppe	34
4.1	État des lieux	34
4.2	Application des résultats précédents aux classes enveloppes	35
4.3	Minoration de la redondance de la classe à décroissance exponentielle	41
5	Et au-delà ?	47

1 Introduction

Ce mémoire est motivé principalement par le codage (compression) de sources sans mémoire en alphabet infini, mais le modèle utilisé est plus général.

Supposons que nous disposons d'une famille paramétrée de sources $\Lambda = (\mathbf{P}_\theta)_{\theta \in \Theta}$ qui produisent un signal $\mathbf{X} = (X_n)_{n \geq 1} \in A^{\mathbb{N}^*}$ à valeur dans un alphabet A fini ou dénombrable¹. En pratique, si A est fini de taille k il sera identifié à $\{1, 2, \dots, k\}$, et si A est infini il sera identifié à \mathbb{N}^* .

Sous la loi \mathbf{P}_θ la loi de X_1 est notée P_θ , et la loi de $X_{1:n} = (X_1, \dots, X_n)$ est P_θ^n . Nous nous plaçons dans l'hypothèse de sources stationnaires sans mémoire : le processus \mathbf{X} est i.i.d. sous \mathbf{P}_θ , et $P_\theta^n = P_\theta^{\otimes n}$.

D'autre part nous notons ν la mesure de comptage sur A ²; pour tout $\theta \in \Theta$, P_θ est absolument continue devant ν , et nous noterons dP_θ sa densité par rapport à ν ³.

Enfin, nous supposons que si $\theta \neq \theta^*$ sont deux éléments de Θ , \mathbf{P}_θ et \mathbf{P}_{θ^*} sont distincts. Nous munissons Θ de la topologie de la convergence faible des lois \mathbf{P}_θ , et de la tribu borélienne associée; par la suite les priors définis sur Θ seront des mesures par rapport à cette tribu.

Dans ce contexte, nous sommes intéressés par les problématiques suivantes : premièrement, trouver un code qui atteigne un « bon » taux de compression quel que soit $\theta \in \Theta$ (on ne connaît pas la « vraie » loi \mathbf{P}_{θ^*} , on sait seulement que $\theta^* \in \Theta$); deuxièmement, estimer quel est le meilleur taux de compression que l'on puisse espérer sur toute la famille Λ .

Les codes concernés sont les codes binaires (le message est codé par une suite de bits 0 ou 1) uniquement décodables : chaque message codé correspond à un unique message non codé.

Donner un sens mathématique à ces questions nécessite de définir la notion de redondance, et de rappeler la correspondance entre codage et loi de probabilité.

Avant de rentrer plus avant dans la problématique du codage de données, il nous faut remarquer que le modèle développé ici rentre dans une problématique statistique plus large d'estimation avec perte logarithmique (ou d'entropie relative), dont les applications

¹Dans un modèle plus général, on demande que le processus \mathbf{X} soit à valeurs dans un espace métrique \mathcal{Y} séparable et complet. Voir [HO97] pour plus de précisions.

²Dans le cas où \mathbf{X} est à valeurs dans un espace métrique \mathcal{Y} , il faut supposer qu'il existe une mesure σ -finie ν qui domine toute la classe $\Lambda^1 = (P_\theta)_{\theta \in \Theta}$ des premières marginales. Toutes les mesures sur \mathcal{Y} considérées par la suite sont alors supposées absolument continues par rapport à ν . Le choix d'un tel ν n'a pas d'influence sur les résultats. [HO97] montre que cette condition est en fait plus faible que la simple existence de la redondance minimax, qui est la quantité que nous voulons étudier dans ce travail.

³Ceci afin de préserver la notation intégrale, valable lorsque \mathbf{X} est à valeurs dans un ensemble \mathcal{Y} non nécessairement dénombrable.

s'étendent de la théorie de l'information (capacité d'un canal informatique) aux mathématiques financières (gestion de portefeuilles), aux théories de l'apprentissage (estimation de la loi d'une source) et des jeux, en incluant bien sûr la compression de données y compris de sources avec mémoire. Les outils présentés et utilisés dans ce mémoire sont utilisables et sont même souvent apparus dans ces autres domaines. Nous invitons à se reporter à [HO97], qui a fourni la base de notre travail et qui est rédigé dans cette optique plus large.

1.1 Longueur de code

Nous nous intéressons maintenant aux liens entre longueur de code et loi de probabilité. Ces liens sont illustrés par les théorèmes ci-après. Il s'agit de résultats classiques dont la démonstration n'est pas donnée ici — les personnes intéressées peuvent se reporter par exemple à [CT91].

Nous utilisons les notations suivantes : Si \mathcal{X} est un ensemble fini ou dénombrable, $\mathcal{X}^* = \bigcup_{m \in \mathbb{N}} \mathcal{X}^m$ est l'ensemble des mots formés à partir de l'alphabet \mathcal{X} . Un code est une application $B \subset A^* \rightarrow \{0, 1\}^*$, et la longueur de code $l(\cdot)$ est l'application qui à un mot de code associe le nombre de bits qui le composent.

Théorème 1 (association d'une sous-probabilité à un code). *Si f est un code uniquement décodable de $B \subset A^*$ dans $\{0, 1\}^*$, alors*

$$\sum_{x \in B} 2^{-l(f(x))} \leq 1.$$

Ceci permet d'associer à un code f une sous probabilité $Q_f(x) = 2^{-l(f(x))}$.

Théorème 2 (code de Shannon). *Soit Q une sous-probabilité sur l'ensemble $B \subset A^*$. Alors il existe un code f uniquement décodable tel que*

$$\forall x \in B, \quad l(f(x)) = \lceil -\log_2 Q(x) \rceil.$$

En réalité on connaît bien un code qui réalise ce théorème : c'est le code de Shannon qui est en outre un code préfixe⁴.

Pour la suite, l'optique d'un codage binaire nous amène à utiliser le logarithme en base 2 ; nous le noterons « log », et nous utiliserons « ln » pour le logarithme népérien⁵.

Ces deux théorèmes permettent de faire le lien entre longueur de code et probabilité d'un message. En particulier ils permettent d'associer un code à une loi de probabilité. Remarquons qu'un code associé à une sous-probabilité peut être amélioré en un code associé à une probabilité.

Lorsqu'on s'intéresse à la longueur moyenne de code on voit apparaître l'entropie de la source.

⁴pour des précisions sur les codes préfixes, voir [CT91].

⁵Si cette manière de faire a clairement son intérêt lorsqu'on parle de codage, elle a aussi son inconvénient lorsqu'on veut utiliser des résultats portant sur l'entropie relative qui ne sont pas spécifiques au codage de données, comme ce sera le cas dans nos sections 2 et 3.

1 Introduction

Définition 1 (entropie binaire). Soit P une loi de probabilité sur l'alphabet A . L'entropie de P , ou entropie binaire, est définie par

$$H(P) = - \sum_{x \in A} P(x) \log P(x) = \mathbb{E}_P[-\log P(X)]$$

Théorème 3 (Shannon-Kraft-Mc Millan). Soit P^n est loi de probabilité sur A^n . Si f est un code uniquement décodable sur A^n , alors

$$\mathbb{E}_{P^n} l(f(X_{1:n})) \geq H(P^n).$$

Réciproquement, il existe un code uniquement décodable f (en fait préfixe) tel que

$$\mathbb{E}_{P^n} l(f(X_{1:n})) \leq H(P^n) + 1.$$

On sait construire de tels codes via le théorème 2 (codes de Shannon ou de Huffman) : il s'agit du code associé à la « vraie » loi de probabilité P^n . Cependant ils présentent plusieurs difficultés :

- la connaissance de la vraie loi de codage P^n est requise pour construire un code qui atteint la borne du théorème 3 ;
- ces codes peuvent être trop lourds à implémenter en terme de temps de calculs, et en pratique un soin particulier doit être apporté à chercher des codes raisonnablement implémentables ;
- si la suite des probabilités $(P^n)_{n \in \mathbb{N}}$ à partir desquels on construit les codes n'est pas compatible, il ne peut qu'on ne puisse pas réaliser d'implémentation en ligne qui réutilise le code de $X_{1:n}$ pour coder $X_{1:n+1}$.

La nécessité de fournir des codes qui se comportent bien même lorsque la loi sous-jacente de la source est inconnue nous pousse à introduire la notion de redondance d'un code.

1.2 Redondance d'un code

Désormais nous utilisons la correspondance code/probabilité de la manière suivante : nous recherchons une loi de probabilité Q^n sur A^n telle que $-\log Q^n(X_{1:n})$ soit minimale en un certain sens.

Plus précisément, nous voulons comparer la longueur de code obtenue avec Q^n avec la longueur de code optimale obtenue avec la « vraie » probabilité comme dans le théorème 3. Selon que veut minimiser cet écart en moyenne ou en norme infinie, cela donne lieu à deux mesures de la qualité d'un code : redondance ou risque individuel (regret).

Définition 2. Soit P^n une loi de probabilité (source) sur A^n . Soit Q^n une autre loi de probabilité sur A^n .

- (i) Le regret, ou risque individuel, de la loi Q^n par rapport à la loi P^n est

$$\sup_{x_{1:n} \in A^n} \log P^n(x_{1:n}) - \log Q^n(x_{1:n}).$$

1 Introduction

(ii) La redondance, ou redondance moyenne⁶, de la loi Q^n par rapport à la loi P^n est définie par

$$R_n(Q^n, P^n) = \mathbb{E}_{P^n}[-\log Q^n(X_{1:n}) + \log P^n(X_{1:n})] = D(P^n \parallel Q^n)$$

où $D(P \parallel Q) = (\log e) D_{KL}(P \parallel Q)$ désigne la version de la divergence de Kullback-Liebler, ou entropie relative, qui utilise le logarithme en base 2⁷.

Le regret comme la redondance admettent une écriture sous forme de regret cumulatif (règle de chaînage) ; Dans le cas de la redondance, qui nous intéresse plus particulièrement dans ce mémoire, cela s'écrit :

$$R_n(Q^n, P^n) = \sum_{k=1}^n \mathbb{E}_{P^k} D(P(X_k = \cdot | X_{1:k-1}) \parallel Q(X_k = \cdot | X_{1:k-1}))$$

où P^k désigne la loi marginale $P_{X_{1:k}}$.

Si la source P^n est une source sans mémoire $P^{\otimes n}$, la formule se simplifie en

$$R_n(Q^n, P^n) = \sum_{k=1}^n \mathbb{E}_{P^k} D(P \parallel Q(X_k = \cdot | X_{1:k-1})).$$

Nous désirons mesurer le comportement d'un code par rapport à toute la famille $\Lambda = (P_\theta)_{\theta \in \Theta}$ dans laquelle se trouve la vraie loi inconnue de la source. Deux manières de faire : une approche bayésienne et une approche minimax.

Définition 3. (i) Le regret minimax de la classe Λ est définie par

$$R_n^*(\Lambda) = \inf_{Q^n} \sup_{P \in \Lambda} \sup_{x_{1:n} \in A^n} \log P^n(x_{1:n}) - \log Q^n(x_{1:n}).$$

(ii) Soit μ une loi de probabilité sur Θ . La redondance bayésienne de la classe Λ par rapport au prior μ est définie par

$$R_{n,\mu}^{\text{Bayes}}(\Lambda) = \inf_{Q^n} \int_{\Theta} D(P_\theta^n \parallel Q^n) d\mu(\theta).⁸$$

(iii) La redondance minimax de la classe Λ est définie par

$$R_n^{\text{minimax}}(\Lambda) = \inf_{Q^n} \sup_{\theta \in \Theta} D(P_\theta^n \parallel Q^n).$$

à chaque fois l'infimum en Q^n est pris sur toutes les distributions de probabilité sur A^n .

⁶Selon un autre terminologie, ce que nous avons appelé *regret* est appelé *redondance individuelle*, et ce que nous avons appelé *redondance* est appelée *redondance moyenne*. Nous avons choisi cette manière de dire pour alléger dans cet exposé le terme *redondance moyenne minimax* en *redondance minimax*.

⁷Rappelons la définition de la divergence de Kullback-Liebler (entropie relative) dans sa version « habituelle » :

$$D_{KL}(P \parallel Q) = - \int dP \ln \frac{dQ}{dP}.$$

⁸Ailleurs on utilise le terme *redondance bayésienne* pour désigner la *redondance (moyenne) maximin*, qui vaut $\sup_{\mu} R_{n,\mu}^{\text{Bayes}}$.

1 Introduction

Nous avons cité le regret minimax par soucis de complétude, mais dans ce mémoire nous nous focalisons sur la redondance bayésienne et la redondance minimax. Le lien entre les deux est illustré dans le résultat suivant :

Théorème 4. *Soit $(\Theta^*, X_{1:n})$ un vecteur aléatoire à valeurs dans $\Theta \times A^n$, tel que la loi de Θ^* est μ et la loi de $X_{1:n}$ sachant $\Theta^* = \theta^*$ est $P_{\theta^*}^n$. Alors*

- (i) *L'unique distribution qui réalise l'infimum dans la définition de $R_{n,\mu}^{\text{Bayes}}$ est le mélange de Bayes $M_{n,\mu}$ défini par*

$$M_{n,\mu}(x_{1:n}) = \int_{\Theta} P_{\theta}^n(x_{1:n}) d\mu(\theta),$$

et

$$R_{n,\mu}^{\text{Bayes}} = \int_{\Theta} D(P_{\theta}^n \| M_{n,\mu}) d\mu(\theta) = I(\Theta^*, X_{1:n}),$$

où $I(X, Y)$ désigne l'information mutuelle entre les variables X et Y ⁹.

(ii)

$$R_n^{\text{minimax}} = \sup_{\mu} R_{n,\mu}^{\text{Bayes}},$$

où le supremum est pris sur tous les mesures de probabilités (priors) sur Θ . De plus

$$R_n^{\text{minimax}} = \inf_{\mu} \sup_{\theta \in \Theta} D(P_{\theta}^n \| M_{n,\mu}).$$

La deuxième partie de ce théorème — redondance minimax et redondance maximin sont égales — est la base de la suite de ce mémoire. On pourra se référer à [Hau96] pour sa démonstration dans un cadre très général.

Nous allons dans la section 2 présenter des encadrements de la redondance bayésienne, d'où sont tirées dans la section 3 des encadrements de la redondance minimax. Enfin, dans notre section 4 nous en déduirons diverses applications nouvelles à un type particulier de classes Λ — les classes enveloppe — en nous attachant plus particulièrement aux classes enveloppe à décroissance exponentielle.

⁹L'information mutuelle est logiquement définie ici par $I(X, Y) = D(P_{(X,Y)} \| P_X \otimes P_Y)$, et elle hérite du facteur $(\log e)$ qui vient de l'utilisation du logarithme en base 2 dans la définition que nous prenons de $D(P \| Q)$.

2 Encadrer la Redondance bayésienne

Cette section, ainsi que la suivante, est essentiellement basée sur l'article [HO97] qui utilise la métrique de Hellinger et certaines intégrales de Laplace pour fournir des encadrements originaux de la redondance bayésienne et de la redondance minimax. Nous présenterons les résultats eux-mêmes mais laisserons de côté les preuves que l'on peut trouver dans l'article cité¹.

Nous allons dans un premier temps définir les objets mathématiques utilisés : I-divergence, D-divergence et condition de α -affinité ; puis nous exposerons les résultats ; enfin nous en donnerons un exemple d'application aux familles dénombrables.

2.1 α -affinité, I-divergence et D-divergence

Nous allons maintenant définir divers quantités utilisées dans le courant de ce mémoire. Les notations suivront celles de [HO97].

2.1.1 α -affinité

Une condition qui est amenée à jouer un rôle central dans les résultats les plus intéressants de ce mémoire est la condition de α -affinité.

Définition 4 (α -affinité). *Soit $\alpha > 1$, et P et Q deux distributions sur A . L' α -affinité entre P et Q est définie par*

$$\rho_\alpha(P, Q) = \int_A (dP(x))^\alpha (dQ(x))^{1-\alpha} d\nu(x).$$

Heuristiquement, l'utilité de faire appel à l' α -affinité peut s'expliquer ainsi : la qualité d'un code dépend de la divergence de Kullbach entre la « vraie » loi de la source et la distribution associée au code. Nous remplacerons la divergence de Kullbach par la distance de Hellinger, mais en faisant cela on perd des informations : il se peut que la divergence de Kullbach soit infinie alors que la distance de Hellinger est finie. Des conditions sur l' α -affinité permettent de palier à ce problème².

Le paramètre α étant appelé à jouer un autre rôle dans notre exposé, il est remplacé dans la pratique par $1 + \lambda$.

¹Si nous avons choisi de ne pas inclure — à quelques exceptions près — les démonstrations de ces résultats c'est pour deux raisons : d'une part elles alourdiraient notre texte, et d'autre part nous pensons que la démonstration originale est très propre. Cependant nous effectuerons ponctuellement certains calculs laissés par les auteurs à la vérification du lecteur.

²Ceci est visible en particulier dans le lemme 5 de [HO97].

2 Encadrer la Redondance bayésienne

Les conditions concernant la $(1 + \lambda)$ -affinité qui apparaîtront par la suite concernent l'existence d'un $\lambda > 0$ pour lequel l'une des deux quantités suivantes est finie :

Définition 5. Soit $\lambda > 0$, et μ un prior.

(i)

$$R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} = \inf_Q \int_{\Theta} \rho_{1+\lambda}(P_\theta, Q) d\mu(\theta)$$

(ii)

$$R_{1,\rho_{1+\lambda}}^{\text{minimax}} = \inf_Q \sup_{\theta \in \Theta} \rho_{1+\lambda}(P_\theta, Q)$$

dans les deux cas l'infimum en Q est pris sur toutes les distributions de probabilité sur A .

Les deux quantités sont reliées par la formule $R_{1,\rho_{1+\lambda}}^{\text{minimax}} \geq \sup_{\mu} R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}}$ qui découle de la relation toujours vraie maximin \leq minimax. Mais ici, à la différence de la redondance, il n'y a pas de résultat d'égalité.

Laquelle des deux quantités interviendra dans les théorèmes, cela dépendra de l'objet que l'on veut encadrer : la redondance bayésienne ou la redondance minimax.

Ces quantités méritent que l'on s'y arrête un peu plus avant. Commençons par la version bayésienne.

Condition portant sur risque bayésien

Proposition 1. Soit μ un prior, et $\lambda > 0$.

(i) Supposons que $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} < \infty$. L'infimum dans la définition de $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}}$ est alors réalisé par l'unique mesure $Q = Q_\mu$ définie par

$$dQ_\mu(x) = \frac{\left(\int_{\Theta} (dP_\theta(x))^{1+\lambda} d\mu(\theta)\right)^{1/(1+\lambda)}}{C_{\lambda,\mu}},$$

où

$$C_{\lambda,\mu} = \int_A \left(\int_{\Theta} (dP_\theta(x))^{1+\lambda} d\mu(\theta)\right)^{1/(1+\lambda)} d\nu(x);$$

en outre

$$R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} = C_{\lambda,\mu}^{1+\lambda}.$$

(ii) Réciproquement, si $C_{\lambda,\mu} < \infty$, alors $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} < \infty$.

Démonstration. Supposons que $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} < \infty$, et posons $\alpha = 1 + \lambda$. La définition de $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}}$ s'écrit

$$R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} = \inf_Q \int_{\Theta} \int_A (dP_\theta(x))^\alpha (dQ(x))^{1-\alpha} d\nu(x) d\mu(\theta),$$

2 Encadrer la Redondance bayésienne

avec Q variant dans l'ensemble C des mesures de probabilités sur A . Dans le cas particulier d'un alphabet dénombrable, ν est la mesure de comptage et il n'y a pas à justifier que Q est absolument continue devant ν ³.

Sur C nous pouvons définir la fonction f à valeurs dans $[0, \infty]$ par

$$f(Q) = \int_A \left(\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) \right) \frac{1}{(dQ(x))^{\alpha-1}} d\nu(x),$$

et $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}}$ s'écrit plus simplement $\inf_{Q \in C} f(Q)$. Comme $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} < \infty$, il existe $Q \in C$ tel que $f(Q) < \infty$, et pour tout $x \in A$

$$\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) < \infty.$$

Quitte à restreindre A , on peut supposer que

$$\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) > 0.$$

Pour Q tel que $f(Q) < \infty$, nous posons

$$h_Q(x) = \frac{dQ(x)}{\left(\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) \right)^{1/\alpha}}.$$

Alors

$$\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) = \left(\frac{dQ(x)}{h_Q(x)} \right)^{\alpha},$$

et

$$f(Q) = \int_A \left(\frac{1}{h_Q(x)} \right)^{\alpha} dQ(x) d\nu(x) < \infty.$$

La fonction $J : x \mapsto x^{\alpha}$ est convexe et nous pouvons appliquer le théorème de Jensen avec la mesure $dQ(x) d\nu(x)$:

$$\begin{aligned} f(Q) &\geq J \left(\int_A \frac{dQ(x)}{h_Q(x)} d\nu(x) \right) \\ &= \left(\int_A \left(\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) \right)^{1/\alpha} d\nu(x) \right)^{\alpha} \\ &= C_{\lambda,\mu}^{1+\lambda}. \end{aligned}$$

En conséquence $C_{\lambda,\mu} < \infty$ et Q_{μ} est bien définie. En outre

$$\begin{aligned} f(Q_{\mu}) &= \int_A \left(\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) \right) C_{\lambda,\mu}^{\alpha-1} \left(\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) \right)^{\frac{1-\alpha}{\alpha}} d\nu(x) \\ &= C_{\lambda,\mu}^{\lambda} \int_A \left(\int_{\Theta} (dP_{\theta}(x))^{\alpha} d\mu(\theta) \right)^{1/\alpha} d\nu(x) \\ &= C_{\lambda,\mu}^{1+\lambda}. \end{aligned}$$

³Si ce n'était pas le cas, il faudrait en outre changer les « pour tout x » en « pour presque tout x ».

2 Encadrer la Redondance bayésienne

En conséquence $f(Q) \geq f(Q_\mu)$ pour tout $Q \in C$, Q_μ réalise le minimum de f et $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} = C_{\lambda,\mu}^{1+\lambda}$.

Réciproquement $C_{\lambda,\mu} < \infty$ entraîne $f(Q_\mu) < \infty$ et $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} < \infty$. □

Condition portant sur risque minimax

Si nous nous intéressons maintenant à la condition « il existe $\lambda > 0$ tel que $R_{1,\rho_{1+\lambda}}^{\text{minimax}} < \infty$ », il faut remarquer qu'elle est moins forte qu'une autre condition que l'on rencontre dans différents travaux, et qui est en particulier vérifiée pour les classes enveloppe traitées dans la section 4 : la condition d'enveloppe intégrable $\int_A (\sup_{\theta \in \Theta} dP_\theta(x)) d\nu(x) < \infty$ ⁴. Voyons cela de plus près.

Supposons qu'il existe une fonction g sur A d'intégrale finie par rapport à ν , et telle que pour tout $\theta \in \Theta$ et tout $x \in A$, $dP_\theta(x) \leq g(x)$ ⁵. Cette condition est en particulier vérifiée si A est finie⁶. On note $S(\Theta) = \int_A g(x) d\nu(x)$. Soit Q la distribution sur A qui admet pour densité par rapport à ν la fonction $g/S(\Theta)$, et $\lambda > 0$.

On a alors, pour tout $\theta \in \Theta$,

$$\begin{aligned} \rho_{1+\lambda}(P_\theta, Q) &\leq \int (dP_\theta) S(\Theta)^\lambda d\nu \\ &\leq S(\Theta)^\lambda, \end{aligned}$$

d'où $R_{1,\rho_{1+\lambda}}^{\text{minimax}} < \infty$.

On pourra trouver dans [HO97] une discussion plus complète au sujet des différentes conditions qui apparaissent chez les auteurs traitant des sujets liés à l'entropie relative.

2.1.2 I-divergence et D-divergence

L'objectif final est déterminer la redondance de la classe Λ en fonction de ses propriétés géométriques par rapport à la distance de Hellinger. Pour passer de la divergence de Kullback — qui dans la définition de la redondance mesure la « distance » entre deux probabilités — à la distance de Hellinger, nous avons besoin d'introduire deux familles de divergences liées à l' α -affinité : l'I-divergence et la D-divergence.

Définition 6. Soit P et Q deux mesures de probabilité sur A .

(i) Soit $\alpha \neq 1$. L'I-divergence d'ordre α entre P et Q est définie par

$$I_\alpha(P \parallel Q) = \frac{1}{\alpha - 1} \ln \rho_\alpha(P, Q) = \frac{1}{\alpha - 1} \ln \int (dP)^\alpha (dQ)^{1-\alpha}.$$

⁴A fortiori la condition d'enveloppe intégrable entraîne que pour tout $\lambda > 0$ et tout prior μ , $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} < \infty$.

⁵Cette façon de faire résout la question de la mesurabilité de $\sup_{\theta \in \Theta} dP_\theta$ dans le cas où \mathbf{X} est à valeurs dans un ensemble \mathcal{Y} non nécessairement dénombrable. Pour obtenir une quantité qui ne dépend pas du choix arbitraire de g ni de ν , on peut aussi prendre l'infimum sur les fonctions g mesurables qui dominant la famille $(dP_\theta)_{\theta \in \Theta}$.

⁶Dans le cas où le processus est à valeurs dans un espace complet \mathcal{Y} (voir la note 1 page 5), l'hypothèse correspondante est que \mathcal{Y} est une partie bornée de \mathbb{R}^k .

(ii) Soit $0 < \alpha < 1$. La D -divergence d'ordre α entre P et Q est définie par

$$D_\alpha(P, Q) = \frac{1}{1-\alpha} (1 - \rho_\alpha(P, Q)) = \frac{1}{1-\alpha} \left(1 - \int (dP)^\alpha (dQ)^{1-\alpha} \right).$$

(iii) Pour $\alpha = 1$, on définit

$$D_1(P, Q) = I_1(P \| Q) = (\ln 2) D(P \| Q) \text{ } ^7 = \int \left(dQ - dP - dP \ln \frac{dQ}{dP} \right).$$

En écrivant

$$D_\alpha(P, Q) = \frac{1}{1-\alpha} \int (\alpha dP + (1-\alpha)dQ - (dP)^\alpha (dQ)^{1-\alpha}),$$

on obtient facilement $D_\alpha(P, Q) \geq 0$ à partir de l'inégalité $\alpha x + (1-\alpha)y - x^\alpha y^{1-\alpha} \geq 0$ valable pour tout $x, y \geq 0$ et $0 < \alpha < 1$.

A partir de $-\ln x \geq 1 - x$ on a $I_\alpha(P \| Q) \geq D_\alpha(P, Q)$, d'où $I_\alpha(P \| Q) \geq 0$.

Remarquons que les deux quantités sont proches lorsque $\rho_\alpha(P, Q)$ est proche de 1.

Un cas particulier important est le carré de la distance de Hellinger, obtenu lorsque $\alpha = 1/2$:

$$D_{1/2}(P, Q) = D_{HL}^2(P, Q) = \int \left(\sqrt{dP} - \sqrt{dQ} \right)^2 d\nu.$$

A la différence des autres distances ou divergences dont on a parlé jusqu'ici, la distance de Hellinger D_{HL} est une distance au sens propre (une métrique) : elle est symétrique et vérifie l'inégalité triangulaire. De nombreux auteurs⁸ l'ont utilisée pour encadrer le risque de procédures d'estimation en statistique.

L'intérêt de cette valeur particulière de α vis à vis des autres possibles sera discuté plus en détail à la fin de la partie 2.2 (page 21).

Les résultats que nous présenterons dans la section 3 feront le lien entre la redondance minimax de la classe Λ et ses propriétés métriques pour la distance de Hellinger.

2.2 Les résultats

Nous allons donner ici des résultats concernant la redondance bayésienne $R_{n,\mu}^{\text{Bayes}}(\Lambda)$ en fonctions du logarithme de la transformée de Laplace de l'I-divergence, puis de la D -divergence. Commençons donc par fixer un prior μ sur Θ et un α dans $(0, 1)$.

Les résultats originaux de [HO97] donnent non seulement des encadrements de $R_{n,\mu}^{\text{Bayes}}$, mais aussi des encadrements de la redondance ponctuelle $R_{n,P_\mu^{\text{Bayes}}}(\theta)$ de la stratégie de

⁷La présence du facteur $\ln 2$ est rendue nécessaire par la définition de la divergence de Kullback donnée dans l'introduction (définition 2), qui utilise le logarithme en base 2. Le même problème surviendra dans l'énoncé des différents résultats des sections 2 et 3. On remarquera au passage que l'article [HO97] a été rédigé en utilisant le logarithme népérien et évite donc ce problème.

⁸[HO97] cite en particulier des travaux de Le Cam, Birgé, Hasminskii et Ibragimov, et van der Geer.

2 Encadrer la Redondance bayésienne

Bayes $P_\mu^{\text{Bayes}} = M_{n,\mu}$, valables pour θ variant dans des ensembles de grande probabilité pour μ . Rappelons que

$$R_{n,P_\mu^{\text{Bayes}}}(\theta) = D(P_\theta^n \| M_{n,\mu}),$$

où P_θ est la vraie loi de la source.

Cependant nous ne présenterons ici que les points qui concernent la redondance bayésienne, qui nous intéresse davantage en vue des sections suivantes.

Théorème 5. *Soit μ une mesure a priori sur Θ , et soit $0 < \alpha < 1$. Pour tout $\theta \in \Theta$, soit Q_θ une distribution de probabilité arbitraire sur A ⁹ sachant θ . Pour tout $n \geq 1$,*

$$\begin{aligned} & - \int_{\Theta} d\mu(\theta) \log \int_{\Theta} d\mu(\tilde{\theta}) \exp[-n(1-\alpha)I_\alpha(P_\theta \| P_{\tilde{\theta}})] \\ & \leq R_{n,\mu}^{\text{Bayes}}(\Lambda) \\ & \leq - \int_{\Theta} d\mu(\theta) \log \int_{\Theta} d\mu(\tilde{\theta}) \exp[-nI_1(P_\theta \| Q_{\tilde{\theta}})]. \end{aligned}$$

En utilisant la relation $I_\alpha(P \| Q) \geq D_\alpha(P, Q)$, on peut remplacer dans la minoration I_α par D_α , et ainsi obtenir la forme de la qui nous intéressera en fin de compte. En revanche dans la majoration, si on voit bien apparaître la même structure (intégrale du logarithme d'une transformée de Laplace) que dans la minoration, la divergence de Kullbach est encore présente. L'objet des résultats qui viennent sera de la transformer en une divergence D_α . L'introduction d'une nouvelle famille Q_θ permettra justement de faire cette transformation; bien sûr on peut aussi obtenir un cas particulier important en choisissant $Q_\theta = P_\theta$.

Une première variante du théorème utilisant la divergence D_α est obtenue de la manière suivante. Pour $0 < \alpha < 1$ et $x > 0$, on définit

$$b_\alpha(x) = \frac{(1-\alpha)(x - \ln x - 1)}{\alpha + (1-\alpha)x - x^{1-\alpha}},$$

ainsi que $b_\alpha(0) = \infty$. On peut montrer¹⁰ que $b_\alpha(x)$ décroît strictement en x , tend vers 1 à l'infini, et vers l'infini en 0. Soit

$$B_\alpha(\Theta) = \sup_{x \in A, \theta^*, \theta \in \Theta} b_\alpha \left(\frac{dP_{\theta^*}(x)}{dP_\theta(x)} \right).$$

Corollaire 1. *Pour tout $0 < \alpha < 1$ et $n \geq 1$,*

$$\begin{aligned} & - \int_{\Theta} d\mu(\theta) \log \int_{\Theta} d\mu(\tilde{\theta}) \exp[-n(1-\alpha)D_\alpha(P_\theta, P_{\tilde{\theta}})] \\ & \leq R_{n,\mu}^{\text{Bayes}}(\Lambda) \\ & \leq - \int_{\Theta} d\mu(\theta) \log \int_{\Theta} d\mu(\tilde{\theta}) \exp[-nB_\alpha(\Theta)D_\alpha(P_\theta, P_{\tilde{\theta}})]. \end{aligned}$$

⁹Le théorème original est rédigé avec $A = \mathcal{Y}$ un espace métrique séparable et complet, comme signalé dans l'introduction.

¹⁰Nous le ferons à l'occasion du lemme 1

2 Encadrer la Redondance bayésienne

Ce corollaire découle du lemme suivant, dont nous donnons la démonstration ¹¹.

Lemme 1. *Pour toutes distributions P et Q sur A et pour tout $0 < \alpha < 1$,*

$$I_1(P \parallel Q) \leq \left(\sup_{x \in A} b_\alpha \left(\frac{dQ(x)}{dP(x)} \right) \right) D_\alpha(P, Q).$$

Démonstration. Avant de démontrer le lemme lui-même, nous commençons par montrer que $b_\alpha(x)$ décroît strictement en x , tend vers 1 à l'infini, et vers l'infini en 0.

Les limites s'obtiennent facilement en raisonnant par équivalents :

$$b_\alpha(x) \underset{x \rightarrow 0^+}{\sim} \frac{-(1-\alpha) \ln x}{\alpha} \xrightarrow{x \rightarrow 0^+} \infty ;$$

$$b_\alpha(x) \underset{x \rightarrow \infty}{\sim} \frac{(1-\alpha)x}{(1-\alpha)x} = 1.$$

En revanche pour la décroissance stricte de b_α nous raisonnons par dérivées successives :

$$b'_\alpha(x) = \frac{1-\alpha}{(\alpha + (1-\alpha)x - x^{1-\alpha})^2} f(x),$$

avec $f(x) = \alpha - \alpha x^{-1} - \alpha x^{1-\alpha} + \alpha x^{-\alpha} + (1-\alpha) \ln x - (1-\alpha)x^{-\alpha} \ln x$ et $b'_\alpha(1) = f(1) = 0$.

$$f'(x) = x^{-1-\alpha} g(x),$$

avec $g(x) = -\alpha(1-\alpha)x + (1-\alpha)x^\alpha + (\alpha - \alpha^2 - 1) + \alpha x^{\alpha-1} + \alpha(1-\alpha) \ln x$ et $f'(1) = g(1) = 0$.

$$g'(x) = \alpha(1-\alpha)x^{\alpha-2} h(x),$$

avec $h(x) = -x^{2-\alpha} + x - 1 + x^{1-\alpha}$. On a $g'(1) = h(1) = 0$ et

$$h'(x) = x^{-\alpha} i(x),$$

avec $i(x) = -(2-\alpha)x + x^\alpha + (1-\alpha)$. On a $h'(1) = i(1) = 0$, $\lim_{x \rightarrow 0^+} i(x) = 1 - \alpha > 0$ et

$$i'(x) = x^{\alpha-1} j(x),$$

avec $j(x) = \alpha - (2-\alpha)x^{1-\alpha}$. j est croissante et s'annule seulement au point $\lambda_\alpha = \left(\frac{\alpha}{2-\alpha}\right)^{1/(1-\alpha)}$. En remarquant que les termes mis en facteur dans les calculs de dérivées ci-dessus sont tous strictement positifs sur $(0, \infty)$, on peut alors résumer ce qu'on obtient dans le tableau de variation suivant, dans lequel les signes + et - désignent des valeurs strictement positives ou strictement négatives :

¹¹Cette démonstration est simple, et c'est l'occasion de corriger une erreur mineure qui s'est glissée dans le texte original. Nous en profitons pour vérifier les propriétés de b_α laissées dans [HO97] au soin du lecteur. Cependant ces preuves n'ont rien d'essentiel à la compréhension du sujet et peuvent être laissées de côté lors de la lecture.

2 Encadrer la Redondance bayésienne

x	0	λ_α	1	∞
$j(x)$	+	0	-	
$i(x)$		+	0	-
$h(x)$		-	0	-
$g(x)$		+	0	-
$f(x)$			0	-
$b'_\alpha(x)$		-	0	-

Ceci démontre bien la décroissance stricte de b_α . Passons maintenant à la démonstration du lemme lui-même.

Si $dP = dQ$ presque partout, alors $I_1(P \| Q) = 0$ et la relation est vraie. On peut donc se restreindre au cas où $D_\alpha(P, Q) > 0$. Soit $S = \{x \in A : dP(x) = 0\}$.

$$\begin{aligned}
 \frac{I_1(P \| Q)}{D_\alpha(P, Q)} &= \frac{(1 - \alpha) \int_{A-S} dP \left(\frac{dQ(x)}{dP(x)} - \ln \frac{dQ(x)}{dP(x)} - 1 \right) + (1 - \alpha) \int_S dQ}{\int_{A-S} dP \left(\alpha + (1 - \alpha) \frac{dQ(x)}{dP(x)} - \left(\frac{dQ(x)}{dP(x)} \right)^{1-\alpha} \right) + (1 - \alpha) \int_S dQ} \\
 &\leq \frac{\left(\sup_{x \in A} b_\alpha \left(\frac{dQ(x)}{dP(x)} \right) \right) \int_{A-S} dP \left(\alpha + (1 - \alpha) \frac{dQ(x)}{dP(x)} - \left(\frac{dQ(x)}{dP(x)} \right)^{1-\alpha} \right) + (1 - \alpha) \int_S dQ}{\int_{A-S} dP \left(\alpha + (1 - \alpha) \frac{dQ(x)}{dP(x)} - \left(\frac{dQ(x)}{dP(x)} \right)^{1-\alpha} \right) + (1 - \alpha) \int_S dQ} \\
 &\leq \sup_{x \in A} b_\alpha \left(\frac{dQ(x)}{dP(x)} \right),
 \end{aligned}$$

puisque $b_\alpha \geq 1$. □

Nous allons maintenant présenter un autre résultat utilisant de manière plus fine les divergences D_α grâce à un recours à l' α -affinité.

Théorème 6. *Soit μ une mesure a priori sur Θ , et soit $0 < \alpha < 1$. On suppose que $R_{1, \mu, \rho_{1+\lambda}}^{\text{Bayes}} < \infty$. Pour tout $n \geq 1$, on a alors*

$$\begin{aligned}
 & - \int_{\Theta} d\mu(\theta) \ln \int_{\Theta} d\mu(\tilde{\theta}) \exp[-n(1 - \alpha)D_\alpha(P_\theta, P_{\tilde{\theta}})] \\
 & \leq (\ln 2) R_{n, \mu}^{\text{Bayes}}(\Lambda) \\
 & \leq - \int_{\Theta} d\mu(\theta) \ln \int_{\Theta} d\mu(\tilde{\theta}) \exp \left[-(n \ln n) \frac{4(1 - \alpha)(1 + o(1))}{\alpha \lambda} D_\alpha(P_\theta, P_{\tilde{\theta}}) \right] \\
 & \quad + R_{1, \mu, \rho_{1+\lambda}}^{\text{Bayes}} + o(1),
 \end{aligned}$$

où, pour α et λ fixés, $o(1)$ est à chaque fois une fonction $f(n)$ qui tend vers 0 à l'infini. En outre le même résultat reste valable si on remplace $D_\alpha(P_\theta, P_{\tilde{\theta}})$ par $I_\alpha(P_\theta \| P_{\tilde{\theta}})$.

Vis-à-vis du corollaire 1, l'intérêt du théorème 6 est de supprimer le recours à $B_\alpha(\Theta)$ qui est infini si les densités dP_θ ne sont pas bornées et maintenues loin de 0 en tout

2 Encadrer la Redondance bayésienne

point $x \in A$. Cependant le théorème ne prétend pas fournir les meilleures constantes ni le résultat le plus fin. Nous verrons dans la section suivante que cela suffit déjà pour obtenir des résultats intéressants concernant l'ordre de grandeur de la redondance minimax.

Concernant la condition $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} < \infty$, nous avons déjà remarqué plus haut ¹² qu'elle était vérifiée si la famille Λ admettait une enveloppe intégrable ; ce sera en particulier le cas dans notre section 4. Cependant le théorème s'applique à de nombreux autres cas. Pour caractériser les couples (Θ, μ) couverts ou non par le théorème 6, définissons la fonction $f_{\Theta,\mu}$ à valeurs dans $[0, \infty]$ par

$$f_{\Theta,\mu}(\lambda) = \frac{1}{\lambda} \ln R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}}$$

pour $\lambda > 0$ et

$$f_{\Theta,\mu}(0) = (\ln 2) R_{1,\mu}^{\text{Bayes}}.$$

Nous montrons ci-dessous que pour tous Θ et μ , si $f_{\Theta,\mu}(\lambda) < \infty$ pour un $\lambda > 0$, alors pour tout $0 < \lambda' < \lambda$, $f_{\Theta,\mu}(\lambda') < \infty$ et

$$\lim_{\lambda \rightarrow 0} f_{\Theta,\mu}(\lambda) = f_{\Theta,\mu}(0).$$

En conséquence il y a trois cas possibles :

1. $f_{\Theta,\mu}(\lambda) < \infty$ pour un $\lambda > 0$. Dans ce cas $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} < \infty$ et donc le théorème 6 s'applique et peut être utilisé pour obtenir des bornes sur $R_{n,\mu}^{\text{Bayes}}$.
2. $f_{\Theta,\mu}(\lambda) = \infty$ pour tout $\lambda > 0$ et $f_{\Theta,\mu}(0) = \infty$ ¹³. Dans ce cas $R_{n,\mu}^{\text{Bayes}} = \infty$ pour tout n . le théorème 6 ne s'applique pas mais encadrer $R_{n,\mu}^{\text{Bayes}}$ est trivial.
3. $f_{\Theta,\mu}(\lambda) = \infty$ pour tout $\lambda > 0$ mais $f_{\Theta,\mu}(0) < \infty$. C'est le seul cas non trivial où le théorème 6 ne s'applique pas. Selon la terminologie de [HO97], le couple (Θ, μ) est alors dit irrégulier ; un exemple est proposé plus bas.

Démonstration. Nous donnons une preuve partielle, n'étant pas encore parvenu à justifier une étape des calculs qui suivent.

Rappelons que $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} = C_{\lambda,\mu}^{1+\lambda}$, avec

$$C_{\lambda,\mu} = \int_A \left(\int_{\Theta} (dP_{\theta}(x))^{1+\lambda} d\mu(\theta) \right)^{1/(1+\lambda)} d\nu(x).$$

Soit $\lambda > 0$ tel que $f_{\Theta,\mu}(\lambda) < \infty$. Alors pour ν -presque tout $x \in A$,

$$\left(\int_{\Theta} (dP_{\theta}(x))^{1+\lambda} d\mu(\theta) \right)^{1/(1+\lambda)} = \|\theta \mapsto dP_{\theta}(x)\|_{1+\lambda} < \infty$$

¹²Lors de la définition de $R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}}$, page 13.

¹³Selon [HO97] la fonction $f_{\Theta,\mu}$ est en outre croissante — et en conséquence $f_{\Theta,\mu}(0) = \infty$ implique $f_{\Theta,\mu}(\lambda) = \infty$ pour tout $\lambda > 0$ — mais nous n'avons pas pu le vérifier.

2 Encadrer la Redondance bayésienne

en utilisant la norme de $L^{1+\lambda}(\mu)$; pour tout $0 < \lambda' < \lambda$,

$$\|\theta \mapsto dP_\theta(x)\|_{1+\lambda'} \leq \|\theta \mapsto dP_\theta(x)\|_{1+\lambda}$$

puisque μ est une mesure de probabilités. On en tire que $C_{\lambda',\mu} \leq C_{\lambda,\mu}$ et $f_{\Theta,\mu}(\lambda') < \infty$.

Montrons que $\lim_{\lambda' \rightarrow 0} C_{\lambda',\mu} = 1$. Soit $x \in A$ tel que $\|\theta \mapsto dP_\theta(x)\|_{1+\lambda} < \infty$. En dominant par la fonction $1 \wedge (dP_\theta(x))^{1+\lambda}$ on a

$$\lim_{\lambda' \rightarrow 0} \int_{\Theta} (dP_\theta(x))^{1+\lambda'} d\mu(\theta) = \int_{\Theta} dP_\theta(x) d\mu(\theta).$$

En dominant l'intégrale sur A par $\|\theta \mapsto dP_\theta(x)\|_{1+\lambda}$, on obtient

$$\lim_{\lambda' \rightarrow 0} C_{\lambda',\mu} = \int_A \left(\int_{\Theta} dP_\theta(x) d\mu(\theta) \right) d\nu(x) = 1$$

par le théorème de Fubini.

En conséquence $\lim_{\lambda \rightarrow 0} R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} = 1$, et la règle de l'Hôpital donne

$$\lim_{\lambda \rightarrow 0} f_{\Theta,\mu}(\lambda) = \left. \frac{d}{d\lambda} \left(R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} \right) \right|_{\lambda=0},$$

si cette dernière dérivée existe. Vérifions que c'est le cas et calculons-la.

Soit $\lambda' < \lambda$, et $0 < \epsilon < \lambda - \lambda'$. Soit $x \in A$ tel que $\|\theta \mapsto dP_\theta(x)\|_{1+\lambda} < \infty$. Posons $M = \sup_{t \geq 1} t^{-\epsilon} \ln t$. En séparant les cas supérieurs et inférieurs à 1, on vérifie que

$$(dP_\theta(x))^{1+s} |\ln(dP_\theta(x))| \leq e^{-1} + M(dP_\theta(x))^{1+\lambda}$$

pour tout $s \in (0, \lambda - \epsilon)$, grâce à la relation $\min_{0 < t < 1} t \ln t = -e^{-1}$. Grâce à cette domination nous pouvons dériver en λ' :

$$\frac{d}{d\lambda'} \left(\int_{\Theta} (dP_\theta(x))^{1+\lambda'} d\mu(\theta) \right) = \int_{\Theta} (dP_\theta(x))^{1+\lambda'} \ln(dP_\theta(x)) d\mu(\theta).$$

Notons maintenant B la quantité

$$\begin{aligned} B(\lambda', x) &= \frac{d}{d\lambda'} \left(\int_{\Theta} (dP_\theta(x))^{1+\lambda'} d\mu(\theta) \right)^{1/(1+\lambda')} \\ &= \frac{1}{(1+\lambda')^2} \left(\int_{\Theta} (dP_\theta(x))^{1+\lambda'} d\mu(\theta) \right)^{1/(1+\lambda')} \left[-\ln \left(\int_{\Theta} (dP_\theta(x))^{1+\lambda'} d\mu(\theta) \right) \right. \\ &\quad \left. + (1+\lambda') \left(\int_{\Theta} (dP_\theta(x))^{1+\lambda'} d\mu(\theta) \right)^{-1} \int_{\Theta} (dP_\theta(x))^{1+\lambda'} \ln(dP_\theta(x)) d\mu(\theta) \right]. \end{aligned}$$

En utilisant le théorème de Jensen avec la fonction convexe $x \ln x$ on obtient que $B(\lambda', x) \geq 0$.

A ce stade la logique voudrait que l'on majore B par une fonction ν -intégrable indépendante de s , au moins pour s variant dans un voisinage de l'origine. *Nous n'y sommes*

2 Encadrer la Redondance bayésienne

pas parvenus. Supposons cependant que nous l'ayons fait. Alors nous pouvons dériver sous le signe intégrale, et

$$\frac{d}{d\lambda} (C_{\lambda,\mu}) = \int_A B(\lambda, x) d\nu(x).$$

Comme d'autre part

$$\frac{d}{d\lambda} \left(R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} \right) = C_{\lambda,\mu}^{1+\lambda} \ln C_{\lambda,\mu} + (1 + \lambda) C_{\lambda,\mu}^\lambda \frac{d}{d\lambda} (C_{\lambda,\mu}),$$

nous pouvons calculer, en utilisant la relation $\int_A dP_\theta(x) d\nu(x) = 1$,

$$\begin{aligned} \left. \frac{d}{d\lambda} \left(R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} \right) \right|_{\lambda=0} &= 0 + \left. \frac{d}{d\lambda} (C_{\lambda,\mu}) \right|_{\lambda=0} \\ &= \int_A B(0, x) d\nu(x) \\ &= \int_A \left(\int_{\Theta} dP_\theta(x) \ln dP_\theta(x) d\mu(\theta) \right) \\ &\quad - \left(\int_{\Theta} dP_\theta(x) d\mu(\theta) \right) \ln \left(\int_{\Theta} dP_\theta(x) d\mu(\theta) \right) \nu(x) \\ &= \int_{A \times \Theta} dP_\theta(x) \ln \frac{dP_\theta(x)}{\int_A dP_\theta(x) d\nu(x) \cdot \int_{\Theta} dP_\theta(x) d\mu(\theta)} d(\nu \otimes \mu)(x, \theta) \\ &= (\ln 2) I(\Theta^*, X_1) \\ &= (\ln 2) R_{1,\mu}^{\text{Bayes}}. \end{aligned}$$

□

Exemple 1. Soit $A = \{1, 2, 3, \dots\}$, $\Theta = \{3, 4, 5, \dots\}$ et pour tout $\theta \in \Theta$ et tout $x \in A$, $P_\theta(X = x)$ vaut $1 - (1/\ln \theta)$ si $x = 1$, $1/\ln \theta$ si $x = \theta$ et 0 dans les autres cas. On choisit $\mu(\theta) = c/(\theta \ln^2 \theta)$, avec $c = \sum_{i=3}^{\infty} 1/(i \ln^2 i)$. Alors (Θ, μ) est irrégulier.

Démonstration. Soit $\lambda > 0$. Calculons $C_{\lambda,\mu}$:

$$\begin{aligned} C_{\lambda,\mu} &= \sum_{x \geq 1} \left(\sum_{\theta \geq 3} (dP_\theta(x))^{1+\lambda} \right)^{1/(1+\lambda)} \\ &= \left(c \sum_{\theta \geq 3} \frac{(1 - (1/\ln \theta))^{1+\lambda}}{\theta \ln^2 \theta} \right)^{1/(1+\lambda)} + \sum_{x \geq 3} \frac{c^{1/(1+\lambda)}}{x^{1/(1+\lambda)} (\ln x)^{3/(1+\lambda)}} \\ &= \infty \quad \text{quelque soit } \lambda > 0. \end{aligned}$$

La dernière égalité se justifie ainsi : dans la deuxième ligne le premier terme est sommable mais pas le second.

2 Encadrer la Redondance bayésienne

Vérifions maintenant que la redondance bayésienne est finie ; Soit P_μ le mélange bayésien.

$$\begin{aligned} P_\mu(1) &= c \sum_{\theta \geq 3} \frac{1 - (1/\ln \theta)}{\theta \ln^2 \theta} \\ P_\mu(2) &= 0 \\ P_\mu(x) &= \frac{c}{x \ln^2 x} \quad \text{pour tout } x \geq 3. \end{aligned}$$

$$\begin{aligned} (\ln 2) D(P_\theta \| P_\mu) &= \left(1 - \frac{1}{\ln \theta}\right) \ln \frac{1 - (1/\ln \theta)}{P_\mu(1)} + \frac{1}{\ln \theta} \ln \frac{\theta \ln^2 \theta}{c} \\ &= \left(1 - \frac{1}{\ln \theta}\right) \ln \left(1 - \frac{1}{\ln \theta}\right) - \left(1 - \frac{1}{\ln \theta}\right) \ln P_\mu(1) + 1 + \frac{2 \ln \ln \theta}{\ln \theta} - \frac{1}{\ln \theta} \ln c. \\ R_{1,\mu}^{\text{Bayes}} &= c \sum_{\theta \geq 3} \frac{1}{\theta \ln^2 \theta} D(P_\theta \| P_\mu) \\ (\ln 2) R_{1,\mu}^{\text{Bayes}} &= c \sum_{\theta \geq 3} \left[\left(\left(1 - \frac{1}{\ln \theta}\right) \ln \left(1 - \frac{1}{\ln \theta}\right) - \frac{1}{\ln \theta} \ln c \right) \frac{1}{\theta \ln^2 \theta} \right. \\ &\quad \left. + \left(1 - \left(1 - \frac{1}{\ln \theta}\right) \ln P_\mu(1)\right) \frac{1}{\theta \ln^2 \theta} + 2 \frac{\ln \ln \theta}{\theta \ln^3 \theta} \right] \\ &< \infty. \end{aligned}$$

En effet chacun des termes de la deuxième ligne se distribue en un terme sommable : le premier est équivalent à $-(1 + \ln c) \frac{1}{\theta \ln^3 \theta}$, le deuxième est équivalent à $(1 - \ln P_\mu(1)) \frac{1}{\theta \ln^2 \theta}$ et le dernier est négligeable devant $\frac{1}{\theta \ln^2 \theta}$. \square

Choix $\alpha = 1/2$.

Nous donnons ici un argument heuristique pour justifier l'intérêt du choix $\alpha = 1/2$. Si l'on considère les intégrales qui apparaissent dans le théorème 5, la méthode de Laplace nous enseigne que ce qui compte est le comportement près du maximum de la fonction dont on prend la transformée de Laplace, c'est-à-dire lorsque $I_\alpha(P_\theta \| P_{\hat{\theta}})$ est proche de 0. Ce phénomène est particulièrement visible dans le premier exemple d'application que donne [HO97]. Mais il se trouve que lorsque P et Q deviennent proches, au sens où $dP/dQ \rightarrow 1$ uniformément, alors

$$\frac{I_1(P \| Q)}{(1 - \alpha)I_\alpha(P \| Q)} \rightarrow \frac{1}{\alpha(1 - \alpha)}.$$

Pour diminuer l'écart entre la minoration et la majoration il faut donc choisir α de manière à minimiser $\frac{1}{\alpha(1-\alpha)}$, c'est-à-dire prendre $\alpha = 1/2$. A cela s'ajoute l'avantage non négligeable déjà mentionné : $D_{1/2}$ est le carré de la distance de Hellinger, et Θ muni de cette distance devient un espace métrique.

2.3 Exemple d'application : Les familles dénombrables

Dans le cas où la famille de lois Λ est paramétrique et suffisamment régulière on dispose déjà d'estimations de la redondance bayésienne en fonction de l'information de Fisher, grâce notamment aux résultats présentés dans [CB94]. Sous ce rapport les théorèmes 5 et 6 n'apportent pas d'éléments nouveaux¹⁴. En réalité le vrai intérêt de ces théorèmes est qu'ils s'appliquent aussi à des familles Λ paramétriques mais non régulières, et à des familles non paramétriques. En plus de fournir des encadrements de la redondance bayésienne pour ces familles, cela permettra aussi de déduire les résultats présentés dans la prochaine section.

Nous nous sommes intéressés dans ce mémoire à deux exemples d'application des théorèmes 5 et 6 : les familles dénombrables et les classes enveloppe. L'objet de ce paragraphe est d'étudier les classes dénombrables, tandis que les classes enveloppe, qui sont un exemple de classes « régulières » mais non paramétriques, seront l'objet de la section 4.

Les familles dénombrables ont un comportement particulier qui mérite qu'on s'y arrête. En effet, à condition que la loi μ soit d'entropie finie, la proposition 2¹⁵ affirme que la redondance bayésienne reste bornée quand n varie. Les différents résultats qui viennent sont intéressants sans être réellement nouveaux¹⁶.

Proposition 2. *On suppose que $\Theta = \{\theta_i\}$ est dénombrable. Soit $(\Theta^*, X_{1:n})$ un vecteur aléatoire à valeurs dans $\Theta \times A^n$, tel que la loi de Θ^* est μ et la loi de $X_{1:n}$ sachant $\Theta^* = \theta_i$ est $P_{\theta_i}^n$. Soit $H(\Theta^*) = -\sum_i \mu(\theta_i) \log \mu(\theta_i)$ l'entropie de la loi μ , à valeurs dans $[0, \infty]$. Pour tout $n \in \mathbb{N}$,*

$$R_{n,\mu}^{\text{Bayes}} = I(\Theta^*, X_{1:n}) \leq H(\Theta^*),$$

et

$$\lim_{n \rightarrow \infty} R_{n,\mu}^{\text{Bayes}} = H(\Theta^*).$$

Ce comportement est à rapprocher du celui, déjà connu par ailleurs, de la redondance minimax dans le cas où Θ est fini. Le lien est la relation $\sup_{\mu} H(\Theta^*) = \log |\Theta|$.

Corollaire 2. *On suppose que Θ est dénombrable. Si Θ est infini, on pose $|\Theta| = \infty$. Alors, pour tout $n \in \mathbb{N}$,*

$$R_n^{\text{minimax}} \leq \log |\Theta|,$$

et

$$\lim_{n \rightarrow \infty} R_n^{\text{minimax}} = \log |\Theta|.$$

¹⁴[HO97] préfère se concentrer sur la redondance ponctuelle de la stratégie de Bayes $R_{n,P_{\mu}^{\text{Bayes}}}(\theta)$: le théorème 3 — qui correspond à notre théorème 6 — permet d'obtenir un ordre de grandeur de cette quantité déjà connu grâce à [CB90] et [CB94], offrant ainsi un résultat intéressant sans être novateur. Il nous faut mentionner qu'à cette occasion [HO97] illustre l'intérêt de choisir $\frac{1}{2}$ pour valeur de la constante α : cette valeur optimise la minoration de $R_{n,P_{\mu}^{\text{Bayes}}}(\theta)$ obtenue.

¹⁵Il s'agit du corollaire 3 de [HO97].

¹⁶[HO97] cite divers ouvrages contenant des résultats semblables.

2 Encadrer la Redondance bayésienne

Concernant la vitesse de convergence, on peut tirer du théorème 5 le corollaire suivant, valable lorsque Θ est de cardinal fini :

Proposition 3. *On reprend les hypothèses et les notations de la proposition 2. Pour tout $n \in \mathbb{N}$,*

$$H(\Theta^*) - I(\Theta^*, X_{1:n}) \leq (|\Theta| - 1) \left(\max_{1 \leq i < j \leq |\Theta|} \sum_{x \in A} \sqrt{P_{\theta_i}(x)P_{\theta_j}(x)} \right)^n .$$

Comme P_{θ_i} et P_{θ_j} sont supposés différents pour $i \neq j$, $\sum_{x \in A} \sqrt{P_{\theta_i}(x)P_{\theta_j}(x)} < 1$ et la proposition 3 garantit une vitesse de convergence exponentielle.

3 Encadrer la Redondance minimax

Dans cette section nous allons utiliser les résultats de la section précédente sur la redondance bayésienne pour obtenir des encadrements de la redondance minimax. Nous avons déjà exprimé la redondance en fonction d'intégrales de Laplace d'une D_α -divergence, avec le cas particulier $\alpha = 1/2$ qui correspond à la distance de Hellinger ; ici nous irons plus loin dans cette direction, en reliant la redondance minimax aux propriétés métriques de la famille de lois envisagée, plus précisément à l'entropie métrique introduite par Kolmogorov et Tikhomirov dans [KT61].

3.1 Entropie métrique

La distance de Hellinger entre les lois marginales d'ordre 1 P_θ issues de la classe Λ nous fournit une métrique sur l'ensemble Θ , définie par

$$h(\theta^*, \theta) = D_{HL}(P_{\theta^*}, P_\theta),$$

sous l'hypothèse déjà exprimée que P_θ et P_{θ^*} sont distincts si $\theta \neq \theta^*$. Les caractéristiques métriques de l'espace (Θ, h) qui nous intéressent sont les nombres « packing » et « covering »¹ et l'entropie métrique associée.

Définition 7. Soit (S, ρ) un espace métrique séparable et complet. Soit $\epsilon > 0$.

- (i) Nous notons $\mathcal{D}_\epsilon(S, \rho)$ le cardinal de la plus petite partition finie de S de diamètre au plus ϵ , ou nous posons $\mathcal{D}_\epsilon(S, \rho) = \infty$ si une telle partition finie n'existe pas. S est dit totalement borné si $\mathcal{D}_\epsilon(S, \rho) < \infty$ pour tout $\epsilon > 0$.
- (ii) L'entropie métrique de (S, ρ) est définie par

$$\mathcal{H}_\epsilon(S, \rho) = \ln \mathcal{D}_\epsilon(S, \rho).$$

- (iii) Une ϵ -couverture de S est une partie A de S telle que pour tout $x \in S$ il existe $y \in A$ vérifiant $\rho(x, y) < \epsilon$. Le nombre covering $\mathcal{N}_\epsilon(S, \rho)$ est le cardinal de la plus petite ϵ -couverture finie de S , ou nous posons $\mathcal{N}_\epsilon(S, \rho) = \infty$ si une telle ϵ -couverture finie n'existe pas.
- (iv) Une partie ϵ -séparée de S est un sous-ensemble A de S tel que pour tous $x, y \in A$ distincts, $\rho(x, y) > \epsilon$. Le nombre packing $\mathcal{M}_\epsilon(S, \rho)$ est le cardinal de la plus grande partie ϵ -séparée de S , ou nous posons $\mathcal{M}_\epsilon(S, \rho) = \infty$ s'il existe des parties ϵ -séparées de cardinal arbitrairement grand.

¹Nous avons choisi de ne pas traduire les termes anglais « packing number » et « covering number ».

3 Encadrer la Redondance minimax

L'intérêt de ces différentes définitions est donné par le lemme suivant, qui permet de choisir la quantité la plus convenable pour estimer l'entropie métrique :

Lemme 2. *Pour tout $\epsilon > 0$,*

$$\mathcal{M}_{2\epsilon}(S, \rho) \leq \mathcal{D}_{2\epsilon}(S, \rho) \leq \mathcal{N}_\epsilon(S, \rho) \leq \mathcal{M}_\epsilon(S, \rho).$$

Ce résultat est très classique ; nous n'en donnons pas la démonstration, qui peut se trouver par exemple dans [vdVW96].

Liée à l'entropie métrique on trouve la notion de dimension d'un espace métrique.

Définition 8. *Les dimensions métriques supérieure et inférieure de S sont définies respectivement par :*

$$\begin{aligned} \overline{\dim}(S) &= \limsup_{\epsilon \rightarrow 0} \frac{\mathcal{H}_\epsilon(S, \rho)}{\ln(1/\epsilon)}, \\ \underline{\dim}(S) &= \liminf_{\epsilon \rightarrow 0} \frac{\mathcal{H}_\epsilon(S, \rho)}{\ln(1/\epsilon)}. \end{aligned}$$

Lorsque $\overline{\dim}(S) = \underline{\dim}(S)$, alors cette quantité est notée $\dim(S)$ et appelée dimension métrique de S . En conséquence

$$\dim(S) = \lim_{\epsilon \rightarrow 0} \frac{\mathcal{H}_\epsilon(S, \rho)}{\ln(1/\epsilon)}.$$

Nous désirons donc utiliser les résultats de la section 2, avec $\alpha = \frac{1}{2}$. Le lien entre la redondance bayésienne et la redondance minimax est fournie par le théorème 4 au point (ii). Pour passer à l'entropie métrique, commençons par définir

$$b(\epsilon) = \sup \left\{ \frac{I_1(P_\theta \| P_{\theta^*})}{D_{HL}^2(P_\theta, P_{\theta^*})} : \theta, \theta^* \in \Theta \text{ et } D_{HL}^2(P_\theta, P_{\theta^*}) \leq \epsilon \right\}.$$

Cette quantité est à comparer avec $B_{1/2}(\Theta)$ utilisé dans le corollaire 1. De fait c'est là le cœur de l'utilisation de l'entropie métrique : utiliser des petites boules de diamètre ϵ dans lesquelles le passage de D_{KL} à D_{HL} ne soit pas trop pénalisant.

Les outils techniques qui vont nous permettre d'encadrer la redondance minimax peuvent maintenant être exprimés : ce sont les deux lemmes qui suivent.

Lemme 3. *On suppose que (Θ, h) est totalement borné. Alors, pour tout $n \geq 1$,*

(i)

$$\begin{aligned} (\ln 2) R_n^{\text{minimax}} &\geq \sup_{\epsilon > 0} \left\{ -\ln \left(\frac{1}{\mathcal{M}_\epsilon(\Theta, h)} + \exp \left(-\frac{n\epsilon^2}{2} \right) \right) \right\} \\ &\geq \sup_{\epsilon > 0} \min \left\{ \mathcal{H}_\epsilon(\Theta, h), \frac{n\epsilon^2}{8} \right\} - \ln 2. \end{aligned}$$

3 Encadrer la Redondance minimax

(ii)

$$\begin{aligned} (\ln 2) R_n^{\text{minimax}} &\leq \inf_{\epsilon > 0} \{ \mathcal{H}_\epsilon(\Theta, h) + b(\epsilon)n\epsilon^2 \} \\ &\leq \inf_{\epsilon > 0} \{ \mathcal{H}_\epsilon(\Theta, h) + B_{1/2}(\Theta)n\epsilon^2 \}. \end{aligned}$$

Si de plus il existe $\lambda > 0$ tel que $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$,

$$(\ln 2) R_n^{\text{minimax}} \leq \inf_{\epsilon > 0} \left\{ \mathcal{H}_\epsilon(\Theta, h) + \frac{4(1 + o(1))n \ln n \epsilon^2}{\lambda} \right\} + R_{1, \rho_{1+\lambda}}^{\text{minimax}} + o(1),$$

où $o(1)$ est à chaque fois une fonction $f(n)$ qui tend vers 0 à l'infini.

Si ce lemme montre bien l'intérêt de faire recours à l'entropie métrique, il n'est pas toujours utilisable en l'état. En effet la dépendance en ϵ de $\mathcal{H}_\epsilon(\Theta, h)$ peut être non continue, et même si elle est continue il n'est pas toujours facile de voir quelle sorte de bornes est obtenue via le lemme 3. Il est alors utile de poser les définitions suivantes : Soit Θ fixé et soit $f_l(x)$ et $f_u(x)$ des fonctions continues et croissantes de $(0, \infty)$ dans $(0, \infty)$, vérifiant

$$\liminf_{\epsilon \rightarrow 0} \frac{\mathcal{H}_\epsilon(\Theta, h)}{f_l(1/\epsilon)} \geq 1 \text{ et } \limsup_{\epsilon \rightarrow 0} \frac{\mathcal{H}_\epsilon(\Theta, h)}{f_u(1/\epsilon)} \leq 1.$$

Pour tout $s > 0$, soit $\epsilon_l(s)$ l'unique solution de l'équation $f_l(1/\epsilon) = s\epsilon^2$, et $\epsilon_u(s)$ l'unique solution de l'équation $f_u(1/\epsilon) = s\epsilon^2$. Soit

$$(3.1) \quad F_l(s) = f_l\left(\frac{1}{\epsilon_l(s)}\right) = s\epsilon_l^2(s),$$

$$(3.2) \quad F_u(s) = f_u\left(\frac{1}{\epsilon_u(s)}\right) = s\epsilon_u^2(s).$$

Ceci nous permet d'énoncer le lemme suivant

Lemme 4. Soit F_l et F_u les fonctions définies ci-dessus.

(i)

$$\liminf_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{(\log e) F_l(n/8)} \geq 1$$

(ii) Si $\lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty$, alors, pour toute fonction $g(n)$ vérifiant $g(n) \xrightarrow[n \rightarrow \infty]{} \infty$,

$$\limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{(\log e) F_u(ng(n))} \leq 1.$$

De même, s'il existe $\lambda > 0$ tel que $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$, alors, pour toute fonction $g(n)$ comme ci-dessus,

$$\limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{(\log e) F_u(ng(n) \ln n)} \leq 1.$$

Dans la pratique il est courant que $F_l(n)$ et $F_u(n \ln n)$ soient proches asymptotiquement. Dans ce cas le taux de croissance asymptotique de R_n^{minimax} est essentiellement obtenue en résolvant l'équation $\mathcal{H}_\epsilon(\Theta, h) = n\epsilon^2$. Nous pouvons utiliser cette approche pour obtenir les résultats présentés dans les paragraphes qui vient.

3.2 Les théorèmes sur la redondance minimax

Les résultats présentés ici sont deux théorèmes qui utilisent les lemmes du paragraphe précédent : le premier donne de manière globale le comportement de la redondance minimax en fonction de la dimension métrique de Θ , tandis que le second entre plus en détail dans l'encadrement de la redondance minimax. Nous donnerons en corollaires des versions légèrement modifiées des mêmes résultats pour traiter le cas où on n'a pas à sa disposition l'entropie métrique de Θ , mais seulement des minoration ou des majoration ; ce sont en effet ces versions que nous utiliserons en pratique par la suite.

Nous utilisons la notation $a(x) \underset{x \rightarrow l}{\asymp} b(x)$ — avec $l \in \overline{\mathbb{R}}$ — pour signifier qu'il existe $0 < c, d < \infty$ et un voisinage V de l tels que $c \cdot b(x) \leq a(x) \leq d \cdot b(x)$ pour tout $x \in V$.

Théorème 7. *On suppose qu'il existe $\lambda > 0$ tel que $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$.*

(i) *Si Θ est fini, alors*

$$R_n^{\text{minimax}} \xrightarrow{n \rightarrow \infty} \log |\Theta|.$$

(ii) *Si $\dim(\Theta, h) = 0$, alors*

$$R_n^{\text{minimax}} \in o(\log n).$$

(iii) *Si $\dim(\Theta, h) = D$, avec $0 < D < \infty$, alors*

$$R_n^{\text{minimax}} \sim \frac{D}{2} \log n.$$

(iv) *Si $\dim(\Theta, h) = \infty$ mais (Θ, h) est totalement borné, alors*

$$R_n^{\text{minimax}} \in O(n) \quad \text{mais} \quad R_n^{\text{minimax}} \notin O(\log n).$$

(v) *Si (Θ, h) n'est pas totalement borné, alors*

$$\text{si } R_1^{\text{minimax}} < \infty, \text{ alors } R_n^{\text{minimax}} \underset{x \rightarrow l}{\asymp} n, \text{ sinon } R_n^{\text{minimax}} = \infty \text{ pour tout } n.$$

La condition $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$ ne peut pas être enlevée. L'exemple suivant, tiré de [Gar06], en donne confirmation :

Exemple 2. Soit f une fonction positive strictement décroissante sur \mathbb{N}^* , et telle que $f(1) < 1$. Pour $k \in \mathbb{N}^*$, Soit p_k la probabilité sur \mathbb{N} définie par

$$p_k(l) = \begin{cases} 1 - f(k) & \text{si } l = 0 ; \\ f(k) & \text{si } l = k ; \\ 0 & \text{sinon.} \end{cases}$$

Soit $\Lambda^1 = \{p_1, p_2, \dots\}$, et soit Λ la famille des lois stationnaires sans mémoire dont les premières marginales sont dans Λ^1 .

Le premier résultat est la proposition 13 de [Gar06]² :

$$\limsup_{k \rightarrow \infty} f(k) \ln k = \infty \Leftrightarrow \forall n \in \mathbb{N}, R_n^{\text{minimax}}(\Lambda) = \infty.$$

²Nous avons légèrement corrigé l'énoncé en remplaçant la limite simple par la limite supérieure. La démonstration reste inchangée.

3 Encadrer la Redondance minimax

Nous voulons maintenant comparer le comportement de cette famille avec l'énoncé du théorème 7. Nous affirmons tout d'abord que Λ^1 est totalement bornée pour la distance de Hellinger si et seulement si $\lim_{k \rightarrow \infty} f(k) = 0$. Cela semble contredire le théorème 7 si on choisit f qui tend vers 0 moins vite que $1/(\ln k)$, par exemple $f(k) = 1/\sqrt{\ln k}$. En réalité, c'est la condition $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$ qui n'est pas vérifiée. En effet, nous montrons que $R_{1, \rho_{1+\lambda}}^{\text{minimax}} = \infty$ dès que $\limsup_{k \rightarrow \infty} k^{\frac{\lambda}{1+\lambda}} f(k) = \infty$, ce qui est en particulier vérifié pour tout $\lambda > 0$ dès que $\limsup_{k \rightarrow \infty} f(k) \ln k = \infty$.

Démonstration. Commençons par calculer $D_{HL}(p_k, p_l)$ pour $k < l$.

$$\begin{aligned} D_{HL}(p_k, p_l)^2 &= \sum_{m \geq 0} \left(\sqrt{p_k(m)} - \sqrt{p_l(m)} \right)^2 \\ &= \left(\sqrt{1 - f(k)} - \sqrt{1 - f(l)} \right)^2 + f(k) + f(l) \\ &= 2 \left(1 - \sqrt{(1 - f(k)) \cdot (1 - f(l))} \right) \end{aligned}$$

Comme $f(l) < f(k)$, on obtient

$$\sqrt{2f(l)} < D_{HL}(p_k, p_l) < \sqrt{2f(k)}.$$

Soit maintenant $\epsilon > 0$. Si $f(k) \rightarrow 0$, alors il existe l_ϵ tel que $f(l_\epsilon) < \frac{\epsilon^2}{2}$. En conséquence la boule $B_{D_{HL}}(P_{l_\epsilon}, \epsilon)$ contient tous les P_k pour $k \geq l_\epsilon$, et $\mathcal{N}_\epsilon(\Lambda^1, D_{HL}) \leq l_\epsilon$. Comme c'est vrai pour tout $\epsilon > 0$, Λ^1 est totalement bornée.

Réciproquement, si Λ^1 est totalement bornée, alors pour tout $\epsilon > 0$, \mathcal{D}_ϵ est fini et il existe $k \neq l$ tels que $D_{HL}(p_k, p_l) < \epsilon$. En conséquence $\lim f(k) \leq \frac{1}{2} D_{HL}(p_k, p_l)^2 < \frac{\epsilon^2}{2}$, et $f(k) \rightarrow 0$.

Pour montrer que $R_{1, \rho_{1+\lambda}}^{\text{minimax}} = \infty$, il suffit de prouver que $R_{1, \mu, \rho_{1+\lambda}}^{\text{Bayes}}$ n'est pas borné lorsque le prior μ varie. En effet $R_{1, \rho_{1+\lambda}}^{\text{minimax}} \geq \sup_\mu R_{1, \mu, \rho_{1+\lambda}}^{\text{Bayes}}$. Par la proposition 1, nous savons en outre que cela revient à montrer que $C_{\lambda, \mu}$ n'est pas borné en μ . Soit $m \geq 2$ et μ la distribution uniforme sur $\{1, 2, \dots, m\}$.

$$\begin{aligned} C_{\lambda, \mu} &= \sum_{l \geq 0} \left(\frac{1}{m} \sum_{k=1}^m p_k(l)^{1+\lambda} \right)^{1/(1+\lambda)} \\ &= \left(\frac{1}{m} \sum_{k=1}^m (1 - f(k))^{1+\lambda} \right)^{1/(1+\lambda)} + \sum_{l=1}^m \left(\frac{1}{m} f(l)^{1+\lambda} \right)^{1/(1+\lambda)} \\ &\geq \frac{1}{m^{1/(1+\lambda)}} \sum_{l=1}^m f(l) \\ &\geq m^{\frac{\lambda}{1+\lambda}} f(m). \end{aligned}$$

Si $\limsup_{k \rightarrow \infty} k^{\frac{\lambda}{1+\lambda}} f(k) = \infty$, $C_{\lambda, \mu}$ n'est donc pas borné en μ . □

Nous donnons maintenant le corollaire annoncé :

3 Encadrer la Redondance minimax

Corollaire 3. *On suppose qu'il existe $\lambda > 0$ tel que $R_{1,\rho_{1+\lambda}}^{\text{minimax}} < \infty$.*

(i) *Si $\overline{\dim}(\Theta, h) \leq D$, avec $0 < D < \infty$, alors*

$$R_n^{\text{minimax}} \leq (1 + o(1)) \frac{D}{2} \log n.$$

(ii) *Si $\underline{\dim}(\Theta, h) \geq D$, avec $0 < D < \infty$, alors*

$$R_n^{\text{minimax}} \geq (1 + o(1)) \frac{D}{2} \log n.$$

Démonstration. La démonstration est celle que donne [HO97]. Nous la reproduisons parce qu'elle permet de voir la démarche utilisée en se basant sur les lemmes 3 et 4. Nous traiterons uniquement la majoration, la minoration étant tout à fait semblable.

Puisque $\overline{\dim}(\Theta, h) \leq D$, nous pouvons choisir

$$f_u(x) = D \ln x$$

En résolvant l'équation $D \ln(1/\epsilon) = n\epsilon^2$ on trouve, comme nous le vérifions plus bas,

$$(3.3) \quad \epsilon_u(n) \sim \sqrt{\frac{D}{2n} \ln n},$$

et en conséquence l'équation (3.2) entraîne

$$F_u(n) \sim \frac{D}{2} \ln n.$$

On choisit, pour appliquer le lemme 4, $g(n) = \ln n$, d'où il vient

$$\begin{aligned} (\log e) F_u(n \ln^2 n) &\sim \frac{D}{2} \log n ; \\ \limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{(D/2) \log n} &\leq 1. \end{aligned}$$

Il nous reste maintenant à justifier la relation (3.3). En posant $x_n = \epsilon_u(n)^{-2}$, l'équation à résoudre devient $x_n \ln x_n = 2n/D$, soit encore $f(x_n) = 2n/D$ en posant $f(x) = x \ln x$. Comme $2n/D > 0$, $x_n > 1$. f est strictement croissante sur $[1, \infty)$, donc l'équation admet une unique solution.

Si n est assez grand, $2n/D > f(e) = e$ et $x_n > e$, d'où $x_n < f(x_n)$. D'autre part

$$\begin{aligned} f\left(\frac{2n}{D \ln n}\right) &= \frac{2n}{D \ln n} (\ln n + \ln 2 - \ln D - \ln \ln n) \\ &= \frac{2n}{D} - \frac{2n \ln \ln n + \ln D - \ln 2}{D \ln n} \\ &< f(x_n) \text{ si } n \text{ est assez grand,} \end{aligned}$$

d'où $x_n > \frac{2n}{D \ln n}$.

3 Encadrer la Redondance minimax

Par ailleurs

$$\frac{f(x)}{f(y)} = \frac{x \ln(x/y) + \ln y}{y \ln y} = \frac{x}{y} + \frac{1}{\ln y} f\left(\frac{x}{y}\right),$$

d'où

$$\frac{x}{y} = \frac{f(x)}{f(y)} - \frac{1}{\ln y} f\left(\frac{x}{y}\right).$$

En prenant $x = \frac{2n}{D \ln n}$ et $y = x_n$, on obtient

$$\frac{1}{\ln n} < \frac{x_n D \ln n}{2n} < 1,$$

d'où

$$-e^{-1} \leq f\left(\frac{x_n D \ln n}{2n}\right) \leq 0$$

et

$$-\frac{1}{\ln x_n} f\left(\frac{x_n D \ln n}{2n}\right) \xrightarrow{n \rightarrow \infty} 0.$$

Comme

$$f\left(\frac{2n}{D \ln n}\right) f(x_n)^{-1} = 1 - \frac{\ln \ln n + \ln D - \ln 2}{\ln n} \xrightarrow{n \rightarrow \infty} 1,$$

on a

$$\frac{x_n D \ln n}{2n} \xrightarrow{n \rightarrow \infty} 1 \quad \text{et} \quad x_n \sim \frac{2n}{D \ln n}.$$

□

Si le théorème 7 est assez précis lorsque la dimension métrique de Θ pour la distance de Hellinger est finie, en particulier lorsque Θ est une partie de \mathbb{R}^n , ce n'est plus le cas lorsque la dimension de Θ est infinie. Cependant les lemmes 3 et 4 permettent de caractériser le taux d'accroissement asymptotique de la redondance minimax en fonction de l'entropie métrique de Θ dans de nombreux cas. En pratique l'entropie métrique croît habituellement en $(1/\epsilon)^\alpha \ln(1/\epsilon)^\beta$, et le théorème 8 permet de traiter ce cas.

Théorème 8. *On suppose qu'il existe $\lambda > 0$ tel que $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$. Soit $l(x)$ une fonction continue et croissante, définie sur $(0, \infty)$, et telle que pour tous $\gamma \geq 0$ et $C > 0$,*

(i)

$$\lim_{x \rightarrow \infty} \frac{l(Cx(l(x))^\gamma)}{l(x)} = 1$$

(ii)

$$\lim_{x \rightarrow \infty} \frac{l(Cx(\ln x)^\gamma)}{l(x)} = 1.$$

Alors

1.

$$\text{Si } \mathcal{H}_\epsilon(\Theta, h) \sim l\left(\frac{1}{\epsilon}\right), \text{ alors } R_n^{\text{minimax}} \sim (\log e) l(\sqrt{n}).$$

3 Encadrer la Redondance minimax

2. Si, pour un certain $\alpha \geq 0$,

$$\mathcal{H}_\epsilon(\Theta, h) \asymp \left(\frac{1}{\epsilon}\right)^\alpha l\left(\frac{1}{\epsilon}\right),$$

alors

(i) Si $\lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty$, alors

$$R_n^{\text{minimax}} \asymp n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)}) \right]^{2/(\alpha+2)}$$

(ii) sinon

$$\liminf_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)}) \right]^{2/(\alpha+2)}} > 0$$

et

$$\limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)}) \right]^{2/(\alpha+2)} (\log n)^{\alpha/(\alpha+2)}} < \infty.$$

Les conditions sur la fonction l portent sur sa vitesse de croissance, qui ne doit pas être trop grande. Une fonction logarithme convient, et permet de traiter les entropies de la forme $\mathcal{H}_\epsilon(\Theta, h) \sim C \cdot (1/\epsilon)^\alpha \ln(1/\epsilon)^\beta$. Cependant il faut faire attention à la condition $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$, qui ne peut pas plus être supprimée ici que dans le théorème 7. En revanche les bornes inférieures fournies par le théorème restent valables sans cette condition mais rien ne garantit alors qu'elles sont fines.

Ici aussi il nous faut énoncer une version qui sépare les minoration et les majorations, pour traiter le cas de classes Λ dont on ne connaît pas l'entropie métrique à un équivalent près.

Corollaire 4. Soit $l(x)$ une fonction comme dans le théorème 8.

1. Soit $g(\epsilon)$ une fonction et $\gamma > 0$ un réel tels que

$$g(\epsilon) \underset{\epsilon \rightarrow 0}{\sim} \gamma l\left(\frac{1}{\epsilon}\right).$$

(i) Si $\mathcal{H}_\epsilon(\Theta, h) \geq g(\epsilon)$ pour tout $\epsilon > 0$, alors

$$\liminf_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{\gamma (\log e) l(\sqrt{n})} \geq 1.$$

(ii) Si $\mathcal{H}_\epsilon(\Theta, h) \leq g(\epsilon)$ pour tout $\epsilon > 0$, et s'il existe en outre $\lambda > 0$ tel que $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$, ou bien si $\lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty$, alors

$$\limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{\gamma (\log e) l(\sqrt{n})} \leq 1.$$

3 Encadrer la Redondance minimax

2. (i) Si pour $\alpha > 0$ il existe $c > 0$ tel que

$$\liminf_{\epsilon \rightarrow 0} \mathcal{H}_\epsilon(\Theta, h) \left(\frac{1}{\epsilon}\right)^{-\alpha} l \left(\frac{1}{\epsilon}\right)^{-1} \geq c,$$

alors

$$\liminf_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{\alpha/(\alpha+2)} [l(n^{1/(\alpha+2)})]^{2/(\alpha+2)}} > 0.$$

(ii) Si pour $\alpha > 0$ il existe $C < \infty$ tel que

$$\limsup_{\epsilon \rightarrow 0} \mathcal{H}_\epsilon(\Theta, h) \left(\frac{1}{\epsilon}\right)^{-\alpha} l \left(\frac{1}{\epsilon}\right)^{-1} \leq C,$$

alors

(a) Si $\lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty$, alors

$$\limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{\alpha/(\alpha+2)} [l(n^{1/(\alpha+2)})]^{2/(\alpha+2)}} < \infty,$$

(b) S'il existe $\lambda > 0$ tel que $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$, alors

$$\limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{\alpha/(\alpha+2)} [l(n^{1/(\alpha+2)})]^{2/(\alpha+2)} (\ln n)^{\alpha/(\alpha+2)}} < \infty.$$

Démonstration. Ici encore il ne s'agit que d'une adaptation de la preuve fournie par [HO97]. Seule la partie 1 apporte des aménagements, aussi nous ne nous occupons que d'elle. Comme plus haut, nous ne traiterons que la majoration, la minoration étant semblable (remarquons seulement que la condition $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$ n'est plus nécessaire).

Nous supposons donc que $\mathcal{H}_\epsilon(\Theta, h) \leq g(\epsilon)$, avec $g(\epsilon) \sim \gamma l(1/\epsilon)$. Nous pouvons donc poser

$$f_u(x) = \gamma l(x),$$

ce qui entraîne, comme nous le montrons plus bas,

$$(3.4) \quad \epsilon_u(n) \sim \sqrt{\frac{\gamma l(\sqrt{n})}{n}},$$

d'où

$$F_u(n) \sim \gamma l(\sqrt{n}).$$

En choisissant $g(n) = \ln n$ dans le lemme 4, il vient

$$F_u(n \ln^2 n) \sim \gamma l(\sqrt{n} \ln n) \sim \gamma l(\sqrt{n}),$$

ce qui permet de conclure.

3 Encadrer la Redondance minimax

Vérifions la relation (3.4). En posant $x_n = \epsilon_u(n)^{-1}$, l'équation à résoudre devient

$$x_n^2 l(x_n) = n/\gamma.$$

Comme n/γ est strictement positif et tend vers l'infini, nous observons que $l(x_n) > 0$ et $x_n \rightarrow \infty$. En partant de l'équation $x_n \sqrt{l(x_n)} = \sqrt{n/\gamma}$, nous pouvons écrire

$$\frac{x_n \sqrt{l(x_n)}}{\sqrt{l(x_n \sqrt{l(x_n)})}} = \frac{\sqrt{n/\gamma}}{\sqrt{l(\sqrt{n/\gamma})}},$$

d'où

$$x_n \sqrt{\frac{l(\sqrt{n/\gamma})}{n/\gamma}} = \sqrt{\frac{l(x_n \sqrt{l(x_n)})}{l(x_n)}}.$$

Comme $x_n \rightarrow \infty$, le terme de droite tend vers 1 par propriété de l , et donc

$$\epsilon_u(n) \sim \sqrt{\frac{l(\sqrt{n/\gamma})}{n/\gamma}} \sim \sqrt{\frac{\gamma l(\sqrt{n})}{n}}.$$

□

Maintenant que nous avons énoncé les outils dont nous avons besoin, nous allons pouvoir dans la section qui vient les appliquer au sujet qui nous intéresse : les classes enveloppe.

4 Classes enveloppe

Nous allons dans un premier temps définir ce que sont les classes enveloppe et citer divers résultats déjà acquis les concernant. Ensuite nous verrons quelles applications aux classes enveloppe on peut tirer des théorèmes énoncés dans les sections précédentes, et nous obtiendrons en particulier une majoration originale de la redondance des classes à décroissance exponentielle. Enfin, nous ferons appel à d'autres outils pour obtenir la minoration correspondante, complétant ainsi notre résultat en un équivalent de la redondance minimax des classes à décroissance exponentielle.

4.1 État des lieux

C'est sur les classes enveloppe que se concentre notre intérêt pour les résultats de [HO97] appliqués au codage de données. Ici l'alphabet est infini dénombrable et identifié à \mathbb{N}^* .

Définition 9. Soit f une fonction de \mathbb{N}^* dans $[0, 1]$. La classe de sources stationnaires sans mémoire d'enveloppe f est définie par

$$\Lambda_f = \{\mathbf{P} = P^{\otimes \infty} : \forall k \in \mathbb{N}^*, P(k) \leq f(k)\}.$$

Bien entendu il faut que $\sum_{k \geq 1} f(k) \geq 1$ pour que Λ_f soit non vide. On peut aussi accepter des valeurs de f supérieures à 1 mais elles ne changent rien à Λ_f . Le résultat suivant donne un premier résultat important sur le comportement des classes enveloppe en fonction de leur enveloppe. Il regroupe la proposition 5 et le théorème 3 de [BGG06].

Théorème 9. Soit f une fonction de \mathbb{N}^* dans $[0, 1]$, et Λ_f la classe enveloppe associée. Alors, pour tout $n \in \mathbb{N}$,

$$R_n^{\text{minimax}}(\Lambda_f) < \infty \Leftrightarrow R_n^*(\Lambda_f) < \infty \Leftrightarrow \sum_{k \geq 1} f(k) < \infty.$$

La première équivalence est propre aux classes enveloppe alors que la seconde est valable pour toutes les classes de sources stationnaires sans mémoire sur un alphabet dénombrable. On remarque que le comportement est le même pour toutes les valeurs de n (mais ce n'est pas une nouveauté).

Les classes enveloppes intéressantes seront donc celles pour lesquelles f décroît assez vite. Heuristiquement, on peut dire que la redondance de la classe enveloppe dépend essentiellement de la vitesse de décroissance de l'enveloppe. [BGG06] fournit un résultat général pour encadrer la redondance des classes enveloppe¹ et l'applique à deux types de

¹C'est leur théorème 5.

4 Classes enveloppe

classes enveloppe qui nous intéresseront plus particulièrement : les classes à décroissance polynomiale et les classes à décroissance exponentielle.

Définition 10. Nous utilisons la notation classique $\zeta(\alpha) = \sum_{k \geq 1} k^{-\alpha}$, pour $\alpha > 1$.

- (i) Soit $\alpha > 1$ et C tel que $C > 1 \wedge \frac{2^\alpha}{\zeta(\alpha)}$. La classe à décroissance polynomiale $\Lambda_{C \cdot -\alpha}$ est la classe associée à la fonction enveloppe $f_{\alpha, C} : x \mapsto 1 \wedge Cx^{-\alpha}$.
- (ii) Soit $\alpha > 1$ et C tel que $C > e^{2\alpha}$. La classe à décroissance exponentielle $\Lambda_{Ce^{-\alpha}}$ est la classe associée à la fonction enveloppe $g_{\alpha, C} : x \mapsto 1 \wedge Ce^{-\alpha x}$.

Théorème 10. Dans les deux cas ci-dessous on suppose que α et C vérifient les conditions données dans la définition 10.

(i)

$$\begin{aligned} n^{1/\alpha} A(\alpha) \log[C\zeta(\alpha)] &\leq R_n^{\text{minimax}}(\Lambda_{C \cdot -\alpha}) \\ &\leq R_n^*(\Lambda_{C \cdot -\alpha}) \leq \left(\frac{2Cn}{\alpha - 1}\right)^{1/\alpha} (\log n)^{1-1/\alpha} + O(1), \end{aligned}$$

où

$$A(\alpha) = \frac{1}{\alpha} \int_1^\infty \frac{1 - e^{-1/(\zeta(\alpha)u)}}{u^{1-1/\alpha}} du.$$

(ii)

$$\frac{1}{8\alpha} \log^2 n (1 - o(1)) \leq R_n^{\text{minimax}}(\Lambda_{Ce^{-\alpha}}) \leq R_n^*(\Lambda_{Ce^{-\alpha}}) \leq \frac{1}{2\alpha} \log^2 n + O(1).$$

Ces deux résultats nous serviront de base de comparaison pour les applications aux classes enveloppe des théorèmes d'encadrement de la redondance minimax fournis dans la section 3.

4.2 Application des résultats précédents aux classes enveloppes

Comme l'indique le théorème 9, les classes enveloppe qui présentent un intérêt sont celles pour lesquelles l'enveloppe est intégrable. En conséquence la condition $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$ est réalisée pour tout $\lambda > 0$, et nous pourrons bien appliquer les différents théorèmes qui l'exigent.

Étudions de plus près le comportement des classes enveloppe. Pour cela, fixons f une enveloppe sommable. Une quantité importante qui permet de caractériser le comportement de Λ_f est la queue de la série, que nous notons \bar{f} :

$$(4.1) \quad \bar{f}(n) = \sum_{k \geq n+1} f(k).$$

4 Classes enveloppe

Considérons deux lois P_1 et P_2 correspondant aux premières marginales de deux éléments de Λ_f . Alors,

$$\begin{aligned} D_{HL}(P_1, P_2) &= \sqrt{\sum_{k \geq 1} \left(\sqrt{P_1(k)} - \sqrt{P_2(k)} \right)^2} \\ &\leq \sqrt{\sum_{k \geq 1} \left(2\sqrt{f(k)} \right)^2} \\ &= 2\bar{f}(1). \end{aligned}$$

Nous pouvons faire deux remarques concernant l'expression de $D_{HL}(P_1, P_2)$:

- L'espace (Λ_f^1, D_{HL}) est en isométrie avec la partie $A_f \cap \{\|x\| = 1\}$ de ℓ^2 , où A_f est défini par

$$(4.2) \quad A_f = \{(x_k)_{k \in \mathbb{N}^*} \in \ell^2 : \forall k \in \mathbb{N}^*, 0 \leq x_k \leq \sqrt{f(k)}\},$$

grâce à la transformation $P \mapsto (\sqrt{P(k)})_{k \in \mathbb{N}^*}$. L'entropie métrique \mathcal{H}_ϵ sera calculée sur A_f .

- Il suffit qu'il existe deux entiers k_1 et k_2 dans le support de f tels que $\sum_{k \neq k_1, k_2} f(k) \geq 1$ pour que l'on puisse construire deux lois P_1 et P_2 telles que $D_{HL}(P_1, P_2)$ soit aussi petit que l'on veut mais $D_{KL}(P_1 \| P_2) = \infty$. En conséquence $b(\epsilon) = \infty$ quelque soit $\epsilon > 0$, et on ne peut utiliser les parties du lemme 3 et du théorème 8 qui supposent que $b(\epsilon)$ est fini. Cette situation est en fait vérifiée pour la plupart des enveloppes f , et en particulier pour les enveloppes $f_{\alpha, C}$ et $g_{\alpha, C}$ des classes à décroissance polynomiales ou exponentielles dès que C est assez grand.

Les idées que nous utilisons pour calculer l'entropie métrique des classes enveloppes sont principalement deux. La première est que les coefficients correspondant à k grand varient peu et contribuent peu à la distance de Hellinger, et que nous pourrions donc tronquer la suite des coefficients de manière à se ramener à un rectangle en dimension finie. La seconde est que le nombre de points ϵ -écartés dans une partie de \mathbb{R}^n est plus petit que le nombre de boules de rayon $\epsilon/2$ que l'on peut mettre dans le même ensemble élargi de $\epsilon/2$. Voyons comment exploiter ces idées pour majorer l'entropie métrique.

Soit $\epsilon > 0$. Commençons par définir le point où nous voulons tronquer les coefficients. Pour cela posons

$$(4.3) \quad N_\epsilon = \inf \left\{ n \geq 1 : \bar{f}(n) \leq \frac{\epsilon^2}{16} \right\}.$$

Vérifions que notre troncature « fonctionne ». Soit $x = (x_n)_{n \geq 1}$ un élément de A_f et

4 Classes enveloppe

$\tilde{x} = (x_n \mathbf{1}_{n \leq N_\epsilon})_{n \geq 1}$ sa version tronquée. Alors

$$\begin{aligned} \|x - \tilde{x}\| &= \sqrt{\sum_{n \geq N_\epsilon+1} x_n^2} \\ &\leq \sqrt{\sum_{n \geq N_\epsilon+1} \sqrt{f(n)}^2} \\ &= \sqrt{\bar{f}(N_\epsilon)} \\ &\leq \frac{\epsilon}{4}. \end{aligned}$$

Soit maintenant S une $\epsilon/4$ -couverture de $\{x \in A_f : \forall n \geq N_\epsilon, x_n = 0\}$. Soit $x \in A_f$. Alors il existe $y \in S$ tel que $\|\tilde{x} - y\| \leq \epsilon/4$, d'où $\|x - y\| \leq \epsilon/2$, et S est une $\epsilon/2$ -couverture de A_f .

Rassemblons toutes ces points :

$$\begin{aligned} \mathcal{D}_\epsilon(\Lambda_f^1, D_{HL}) &= \mathcal{D}_\epsilon(A_f \cap \{\|x\| = 1\}, \|\cdot\|_{\ell^2}) \\ &\leq \mathcal{D}_\epsilon(A_f, \|\cdot\|_{\ell^2}) \\ &\leq \mathcal{N}_{\epsilon/2}(A_f, \|\cdot\|_{\ell^2}) \\ &\leq \mathcal{N}_{\epsilon/4} \left(\prod_{1 \leq k \leq N_\epsilon} [0, \sqrt{f(k)}], \|\cdot\|_{\mathbb{R}^{N_\epsilon}} \right) \\ &\leq \mathcal{M}_{\epsilon/4} \left(\prod_{1 \leq k \leq N_\epsilon} [0, \sqrt{f(k)}], \|\cdot\|_{\mathbb{R}^{N_\epsilon}} \right) \\ &\leq \frac{\text{Vol} \left(\prod_{1 \leq k \leq N_\epsilon} \left[-\frac{\epsilon}{8}, \sqrt{f(k)} + \frac{\epsilon}{8}\right] \right)}{\text{Vol} \left(B_{\mathbb{R}^{N_\epsilon}} \left(0, \frac{\epsilon}{8}\right) \right)} \\ &\leq \left(\frac{\epsilon}{8}\right)^{-N_\epsilon} \frac{1}{\text{Vol} \left(B_{\mathbb{R}^{N_\epsilon}}(0, 1) \right)} \prod_{k=1}^{N_\epsilon} \left(\sqrt{f(k)} + \frac{\epsilon}{4} \right). \end{aligned}$$

Remarquons au passage que la classe Λ_f^1 est totalement bornée.

Nous pouvons alors écrire le résultat suivant :

Proposition 4.

$$\mathcal{H}_\epsilon \leq N_\epsilon \ln(1/\epsilon) + 3N_\epsilon \ln 2 + A(\epsilon) + B(\epsilon),$$

où

$$A(\epsilon) = -\ln \text{Vol} \left(B_{\mathbb{R}^{N_\epsilon}}(0, 1) \right) = \ln \frac{\Gamma\left(\frac{N_\epsilon}{2} + 1\right)}{\pi^{\frac{N_\epsilon}{2}}}$$

et

$$B(\epsilon) = \sum_{k=1}^{N_\epsilon} \ln \left(\sqrt{f(k)} + \frac{\epsilon}{4} \right).$$

4 Classes enveloppe

En outre

$$A(\epsilon) \sim \frac{N_\epsilon}{2} \ln N_\epsilon,$$

et

$$-N_\epsilon \ln(1/\epsilon) - 2N_\epsilon \ln 2 \leq B(\epsilon) \leq \frac{\epsilon}{4} N_\epsilon.$$

L'encadrement concernant $B(\epsilon)$ nous montre que $B(\epsilon)$ est un terme qui diminue plutôt l'entropie, à la différence de $A(\epsilon)$ qui y contribue; on peut remarquer que si $\ln N_\epsilon \asymp \ln(1/\epsilon)$, alors l'ordre de grandeur du majorant est $N_\epsilon \ln(1/\epsilon)$, ce qui correspond au point 2 du théorème 8.

Démonstration. Il reste à justifier nos affirmations sur $A(\epsilon)$ et $B(\epsilon)$. Concernant le premier, nous utilisons la formule de Stirling (dans sa version modifiée par [WW96, chapitre XII]) :

$$(4.4) \quad \Gamma(x) = \sqrt{2\pi} x^{x-1/2} e^{-x} e^{\frac{\beta}{12x}}, \quad \text{avec } \beta \in [0, 1].$$

Nous pouvons alors écrire

$$\begin{aligned} A(\epsilon) &= -\frac{N_\epsilon}{2} \ln \pi + \ln \Gamma\left(\frac{N_\epsilon}{2} + 1\right) \\ &= -\frac{N_\epsilon}{2} \ln \pi + \left(\frac{N_\epsilon + 1}{2}\right) \ln\left(\frac{N_\epsilon}{2} + 1\right) - \left(\frac{N_\epsilon}{2} + 1\right) + \frac{\ln(2\pi)}{2} + \frac{\beta}{12\left(\frac{N_\epsilon}{2} + 1\right)}. \end{aligned}$$

Comme $N_\epsilon \rightarrow \infty$, l'équivalence annoncée est vraie.

L'encadrement de $B(\epsilon)$, quant à lui, ne pose pas de difficultés. □

Voyons maintenant ce que ce résultat donne lorsqu'il est appliqué aux classes à décroissance polynomiale ou exponentielle.

Corollaire 5. *On choisit $\Lambda_f = \Lambda_{C, -\alpha}$ une classe enveloppe à décroissance polynomiale. Alors il existe une constante $K > 0$ telle que*

$$R_n^{\text{minimax}} \leq K n^{1/\alpha} \ln n.$$

Démonstration. Commençons par calculer N_ϵ :

$$N_\epsilon \leq \inf \left\{ n \geq 1 : \sum_{k \geq n+1} C k^{-\alpha} \leq \frac{\epsilon^2}{16} \right\};$$

or

$$\sum_{k \geq n+1} k^{-\alpha} \leq \int_n^\infty x^{-\alpha} dx = \frac{1}{(\alpha - 1)n^{\alpha-1}},$$

4 Classes enveloppe

et en conséquence

$$\begin{aligned} N_\epsilon &\leq \inf \left\{ n \geq 1 : n^{\alpha-1} \geq \frac{16C}{(\alpha-1)\epsilon^2} \right\} \\ &\leq \left\lceil \left(\frac{16C}{(\alpha-1)\epsilon^2} \right)^{\frac{1}{\alpha-1}} \right\rceil. \end{aligned}$$

Nous pouvons donc écrire que $N_\epsilon \leq g(\epsilon)$, avec

$$g(\epsilon) \sim \left(\frac{16C}{(\alpha-1)\epsilon^2} \right)^{\frac{1}{\alpha-1}} = \left(\frac{16C}{(\alpha-1)} \right)^{\frac{1}{\alpha-1}} (1/\epsilon)^{\frac{2}{\alpha-1}}.$$

En utilisant la remarque après la proposition 4 il vient alors

$$\mathcal{H}_\epsilon \leq K(1/\epsilon)^{\frac{2}{\alpha-1}} \ln(1/\epsilon)$$

pour un $K > 0$ et pour ϵ au voisinage de 0.

Nous pouvons alors utiliser le corollaire 4 avec $l = \ln$ pour obtenir

$$\limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{1/\alpha} \ln n} < \infty,$$

puisque

$$\frac{2}{\alpha-1} \left(\frac{2}{\alpha-1} + 2 \right)^{-1} = \frac{1}{\alpha}.$$

□

Ce résultat est moins bon que celui du théorème 10 issu de [BG06]. On peut penser que la majoration de \mathcal{H}_ϵ peut être améliorée, cependant nous remarquons que la puissance $\frac{2}{\alpha-1}$ qui porte sur $(1/\epsilon)$ est correcte, puisque c'est celle qui est demandée pour produire le terme principal $n^{1/\alpha}$. Malheureusement elle entraîne nécessairement un facteur $(\ln n)^{1/\alpha}$ entre la majoration et la minoration du théorème 8, ce qui est moins bon, dès que $\alpha < 2$, que l'écart constaté avec le théorème 10. Cependant nous ne savons pas à ce stade si le théorème 8 permet d'améliorer la minoration fournie par le théorème 10 : si notre majoration de l'entropie n'est pas trop brutale nous pouvons espérer un ordre de grandeur en $n^{1/\alpha}(\ln n)^{1-1/\alpha}$ pour la minoration de R_n^{minimax} , ce qui serait très bon !

Corollaire 6. *On choisit $\Lambda_f = \Lambda_{C e^{-\alpha}}$ une classe enveloppe à décroissance exponentielle. Alors*

$$R_n^{\text{minimax}} \leq (1 + o(1)) \frac{1}{4\alpha \log e} \log^2 n.$$

4 Classes enveloppe

Démonstration. Cette fois-ci nous avons

$$\begin{aligned}
N_\epsilon &\leq \inf \left\{ n \geq 1 : \sum_{k \geq n+1} C e^{-\alpha k} \leq \frac{\epsilon^2}{16} \right\} \\
&\leq \inf \left\{ n \geq 1 : \frac{C e^{-\alpha(n+1)}}{1 - e^{-\alpha}} \leq \frac{\epsilon^2}{16} \right\} \\
&\leq \inf \left\{ n \geq 1 : -\alpha(n+1) \leq -\ln \frac{16C}{(1 - e^{-\alpha})\epsilon^2} \right\} \\
&\leq \left\lceil \frac{2}{\alpha} \ln(1/\epsilon) + \frac{1}{\alpha} \ln \frac{16C}{1 - e^{-\alpha}} - 1 \right\rceil,
\end{aligned}$$

et nous pouvons écrire

$$N_\epsilon \leq \frac{2}{\alpha} \ln(1/\epsilon) + \frac{1}{\alpha} \ln \frac{16C}{1 - e^{-\alpha}} \sim \frac{2}{\alpha} \ln(1/\epsilon).$$

Cela nous permet de majorer $B(\epsilon)$:

$$\begin{aligned}
B(\epsilon) &\leq \sum_{k=1}^{N_\epsilon} \ln \left(\frac{\epsilon}{4} + \sqrt{C} e^{-\frac{\alpha}{2}k} \right) \\
&\leq \int_0^{N_\epsilon} \ln \left(\frac{\epsilon}{4} + \sqrt{C} e^{-\frac{\alpha}{2}x} \right) dx \\
&= \int_0^{N_\epsilon} \frac{\ln C}{2} - \frac{\alpha}{2}x + \ln \left(1 + \frac{\epsilon \cdot e^{\frac{\alpha}{2}x}}{4\sqrt{C}} \right) dx.
\end{aligned}$$

D'autre part

$$\begin{aligned}
N_\epsilon &\leq \frac{2}{\alpha} \ln(1/\epsilon) + \frac{1}{\alpha} \ln \frac{16C}{1 - e^{-\alpha}} \\
e^{\frac{\alpha}{2}N_\epsilon} &\leq \frac{1}{\epsilon} \cdot \sqrt{\frac{16C}{1 - e^{-\alpha}}} \\
\frac{\epsilon \cdot e^{\frac{\alpha}{2}N_\epsilon}}{4\sqrt{C}} &\leq \frac{1}{\sqrt{1 - e^{-\alpha}}},
\end{aligned}$$

et donc

$$B(\epsilon) \leq -\frac{\alpha}{4} N_\epsilon^2 + \left(\frac{\ln C}{2} + \frac{1}{\sqrt{1 - e^{-\alpha}}} \right) N_\epsilon \sim -\frac{\alpha}{4} N_\epsilon^2 \sim -\frac{1}{\alpha} \ln^2(1/\epsilon).$$

Par ailleurs

$$A(\epsilon) \sim \frac{2}{\alpha} \ln(1/\epsilon) \cdot \ln \ln(1/\epsilon) = o(\ln^2(1/\epsilon)),$$

et, en rassemblant les différents termes de la proposition 4, on obtient

$$\mathcal{H}_\epsilon \leq g(\epsilon)$$

avec $g(\epsilon) \sim \frac{1}{\alpha} \ln^2(1/\epsilon) = \frac{1}{\alpha \log^2 e} \log^2(1/\epsilon)$.

Il ne reste plus qu'à appliquer le corollaire 4 avec $l = \log^2$ pour conclure. □

Cette fois-ci la comparaison avec le théorème 10 montre que cette majoration est très intéressante. Mieux, nous serons en mesure d'obtenir dans le paragraphe 4.3 une minoration semblable qui complétera notre résultat en une équivalence. Ce bon comportement des classes à décroissance exponentielle peut s'analyser ainsi : la décroissance de f est suffisamment rapide pour garantir une entropie qui croît plus lentement qu'un polynôme, ce qui permet d'utiliser la partie 1 du corollaire 4 qui donne une vraie équivalence.

Et la minoration ?

Concernant la minoration de l'entropie métrique nous pouvons utiliser une démarche analogue à celle suivie pour la majoration. Malheureusement nous ne pouvons pas nous débarrasser aussi facilement de la condition $\|x\| = 1$. Nous avons essayé de résoudre le problème dans le cas des classes à décroissance exponentielle, pour lesquelles nous avons de bons espoirs, mais nous avons été bloqués par le calcul de volumes un peu compliqués (l'intersection d'une couronne et d'un rectangle) en grande dimension.

Aussi nous avons eu recours à une toute autre méthode, dérivée du travail de [XB97], pour minorer la redondance des classes à décroissance exponentielle.

4.3 Minoration de la redondance de la classe à décroissance exponentielle

L'approche que nous développons dans ce paragraphe est basée sur la recherche d'un prior μ sur la classe Λ_f^1 , où $\Lambda_f = \Lambda_{Ce^{-\alpha \cdot}}$ est une classe enveloppe à décroissance exponentielle, tel que $R_{n,\mu}^{\text{Bayes}}(\Lambda_f)$ soit arbitrairement proche de la majoration $\frac{1}{4\alpha \log e} \log^2 n$ obtenue dans le paragraphe précédent. Le prior que nous choisissons est dérivé du prior de Dirichlet, et nous aurons besoin d'une version modifiée des résultats de [XB97] qui utilise ce prior pour les sources stationnaires sans mémoire en alphabet fini. Voyons de plus près ce dont il s'agit.

Considérons un alphabet $A = \{a_1, a_2, \dots, a_k\}$ fini de cardinal $k \in \mathbb{N}^*$. Alors un processus \mathbf{P} stationnaire sans mémoire à valeurs dans A est caractérisé par les valeurs $P(a_i)$, $1 \leq i \leq k$, prises par sa première marginale sur les éléments de A , avec la contrainte $\sum_{1 \leq i \leq k} P(a_i) = 1$. Nous pouvons alors paramétrer la classe des sources sans mémoire sur l'alphabet A sous la forme (\mathbf{P}_θ) , avec $\theta \in S_k$ le simplexe de \mathbb{R}^k

$$S_k = \{(\theta_1, \theta_2, \dots, \theta_k) \in [0, 1]^k : \sum_{1 \leq i \leq k} \theta_i = 1\},$$

en posant $P_\theta(a_i) = \theta_i$.

Une autre notation du simplexe est obtenue en posant $\theta_1 = 1 - \sum_{2 \leq i \leq k} \theta_i$, avec $(\theta_2, \dots, \theta_k)$ variant dans l'ensemble

$$S'_k = \{(\theta_2, \dots, \theta_k) \in [0, 1]^{k-1} : \sum_{2 \leq i \leq k} \theta_i \leq 1\}.$$

4 Classes enveloppe

Cela permet de définir facilement la mesure de Lebesgue $d\theta$ sur le simplexe de \mathbb{R}^k par restriction de la mesure de Lebesgue sur $[0, 1]^{k-1}$.

Si maintenant nous considérons $X_{1:n}$ la suite des n premiers symboles produits par la source \mathbf{P}_θ , nous pouvons poser

$$T_i = \sum_{j=1}^n \mathbf{1}_{X_j=a_i} \quad \text{pour } 1 \leq i \leq k,$$

et la probabilité de la séquence $X_{1:n}$ sous \mathbf{P}_θ est alors

$$P_\theta^n(X_{1:n}) = \theta_1^{T_1} \theta_2^{T_2} \dots \theta_k^{T_k}.$$

Le mélange de Dirichlet² consiste à choisir le prior μ de densité proportionnelle à $\theta_1^{-1/2} \theta_2^{-1/2} \dots \theta_k^{-1/2}$ par rapport à la mesure de Lebesgue. La loi de Bayes sur A^n associée vérifie donc

$$P_\mu^n(X_{1:n}) = \frac{\int_{S'_k} \theta_1^{T_1-1/2} \theta_2^{T_2-1/2} \dots \theta_k^{T_k-1/2} d\theta}{\int_{S'_k} \theta_1^{-1/2} \theta_2^{-1/2} \dots \theta_k^{-1/2} d\theta}.$$

Cette expression peut être simplifiée grâce aux expressions des intégrales de Dirichlet utilisant la fonction Γ :

$$\begin{aligned} D_k(\lambda_1, \dots, \lambda_k) &= \int_{S'_k} \theta_1^{\lambda_1-1} \theta_2^{\lambda_2-1} \dots \theta_k^{\lambda_k-1} d\theta \\ &= \frac{\Gamma(\lambda_1) \Gamma(\lambda_2) \dots \Gamma(\lambda_k)}{\Gamma\left(\sum_{i=1}^k \lambda_i\right)}. \end{aligned}$$

Proposition 5. Soit θ un élément du simplexe S_k , et P_μ^n le mélange de Bayes sur A^n défini ci-dessus. Alors

$$D(P_\theta^n \parallel P_\mu^n) \geq \frac{k-1}{2} \log \frac{n}{2\pi} + \log \frac{\Gamma(1/2)^k}{\Gamma(k/2)} - \frac{5k}{3} \log e.$$

Démonstration. La démonstration, comme la proposition elle-même, est une adaptation de la proposition 1 de [XB97].

La formule est invariante lorsqu'on change la base du logarithme, aussi nous choisirons le logarithme népérien pour cette démonstration. Nous notons \mathbb{E}_θ l'espérance sous la loi \mathbf{P}_θ . Alors

$$\begin{aligned} (4.5) \quad D(P_\theta^n \parallel P_\mu) &= \mathbb{E}_\theta \ln \frac{P_\theta^n(X_{1:n})}{P_\mu^n(X_{1:n})} \\ &= \mathbb{E}_\theta \sum_{i=1}^k T_i \ln \theta_i - \mathbb{E}_\theta \ln \frac{D_k(T_1 + \frac{1}{2}, \dots, T_k + \frac{1}{2})}{D_k(\frac{1}{2}, \dots, \frac{1}{2})} \\ &= \ln \frac{\Gamma(1/2)^k}{\Gamma(k/2)} + \overbrace{\sum_{i=1}^k n \theta_i \ln \theta_i - \mathbb{E}_\theta \ln \frac{\prod_{i=1}^k \Gamma(T_i + \frac{1}{2})}{\Gamma(n + \frac{k}{2})}}^{(A)}. \end{aligned}$$

²Dans ce cas il s'agit aussi du mélange de Jeffrey.

4 Classes enveloppe

Nous utilisons maintenant la formule de Stirling (4.4), en notant β_0 le coefficient correspondant à l'équation de $\Gamma(n + \frac{k}{2})$ et β_i celui correspondant à $\Gamma(T_i + \frac{1}{2})$.

$$\begin{aligned} \text{(A)} &= \sum_{i=1}^k n\theta_i \ln \theta_i - \mathbb{E}_\theta \frac{\prod_{i=1}^k (\sqrt{2\pi}(T_i + \frac{1}{2})^{T_i})}{\sqrt{2\pi}(n + \frac{k}{2})^{n+(k-1)/2}} - \sum_{i=1}^k \mathbb{E}_\theta \frac{\beta_i}{12(T_i + \frac{1}{2})} + \frac{\beta_0}{12(n + \frac{k}{2})} \\ &\geq -\frac{k-1}{2} \ln 2\pi - \frac{k}{6} + \underbrace{\sum_{i=1}^k \left(n\theta_i \ln \theta_i - \mathbb{E}_{\theta_i} T_i \ln \left(T_i + \frac{1}{2} \right) \right)}_{\text{(B)}} + \underbrace{\left(n + \frac{k-1}{2} \right) \ln \left(n + \frac{k}{2} \right)}_{\text{(C)}}. \end{aligned}$$

Nous minorons maintenant les termes (B) et (C) séparément. Pour le second nous avons

$$\begin{aligned} \text{(C)} &= \left(n + \frac{k-1}{2} \right) \ln n + \left(n + \frac{k-1}{2} \right) \ln \left(1 + \frac{k}{2n} \right) \\ &\geq n \ln n + \frac{k-1}{2} \ln n. \end{aligned}$$

Nous pouvons écrire le terme (B) sous la forme

$$\text{(B)} = -n \ln n + \sum_{i=1}^k \underbrace{\left(n\theta_i \ln n\theta_i - \mathbb{E}_{\theta_i} T_i \ln T_i \right)}_{\text{(B}_1)} - \sum_{i=1}^k \overbrace{\mathbb{E}_{\theta_i} T_i \ln \left(1 + \frac{1}{2T_i} \right)}^{\leq 1/2}.$$

En partant de l'équation $t \ln t \leq t(t-1)$, et en posant $t = \frac{T_i}{n\theta_i}$, nous avons

$$\begin{aligned} \frac{T_i}{n\theta_i} (\ln T_i - \ln n\theta_i) &\leq \frac{T_i}{n\theta_i} \frac{T_i - n\theta_i}{n\theta_i} \\ T_i \ln T_i - T_i \ln n\theta_i &\leq \frac{(T_i - n\theta_i)^2}{n\theta_i} + (T_i - n\theta_i) \\ \mathbb{E}_{\theta_i} T_i \ln T_i - n\theta_i \ln n\theta_i &\leq \frac{\text{Var } T_i}{n\theta_i} = 1 - \theta_i, \end{aligned}$$

grâce à la relation $\mathbb{E}_{\theta_i} T_i = n\theta_i$. En conséquence

$$\text{(B}_1) \geq -(1 - \theta_i) \geq -1,$$

$$\text{(B)} \geq -n \ln n - \frac{3k}{2},$$

et

$$\text{(A)} \geq \frac{k-1}{2} \ln \frac{n}{2\pi} - \frac{5k}{3}.$$

Il ne reste plus qu'à rassembler les différents éléments de l'équation (4.5) pour obtenir le résultat annoncé. \square

4 Classes enveloppe

Une fois acquis ce résultat en dimension finie, nous voudrions l'exploiter pour traiter le cas de la classe $\Lambda_{Ce^{-\alpha}}$ en dimension infinie. Soit $m \in \mathbb{N}^*$ fixé tel que

$$\sum_{i \geq m} Ce^{-\alpha i} < 1,$$

et soit

$$\Theta_k = \left\{ \theta = (\theta_1, 0, \dots, 0, \theta_m, \dots, \theta_k, 0, \dots) : \theta_1 = 1 - \sum_{i=m}^k \theta_i \text{ et } \forall m \leq i \leq k, 0 \leq \theta_i \leq Ce^{-\alpha i} \right\}.$$

Nous choisissons ensuite le prior μ_k sur Θ_k proportionnel au prior de Dirichlet μ sur S_{k-m+2} portant sur les coordonnées $(\theta_1, \theta_m, \dots, \theta_k)$:

$$\begin{aligned} d\mu_k(\theta_1, 0, \dots, 0, \theta_m, \dots, \theta_k, 0, \dots) &= \frac{\theta_1^{-1/2} \theta_m^{-1/2} \dots \theta_k^{-1/2} d\theta}{\int_{\Theta_k} \theta_1^{-1/2} \theta_m^{-1/2} \dots \theta_k^{-1/2} d\theta} \\ &= \frac{\int_{S_{k-m+2}} \theta_1^{-1/2} \theta_m^{-1/2} \dots \theta_k^{-1/2} d\theta}{\int_{\Theta_k} \theta_1^{-1/2} \theta_m^{-1/2} \dots \theta_k^{-1/2} d\theta} d\mu(\theta_1, \theta_m, \dots, \theta_k), \end{aligned}$$

où $d\theta$ est la mesure de Lebesgue sur le simplexe S_{k-m+2} indexé par $(\theta_1, \theta_m, \dots, \theta_k)$, et Θ_k est identifié à sa projection sur ce simplexe. Par la suite nous poserons

$$\begin{aligned} C(k) &= \int_{\Theta_k} \theta_1^{-1/2} \theta_m^{-1/2} \dots \theta_k^{-1/2} d\theta, \\ D(k) &= \int_{S_{k-m+2}} \theta_1^{-1/2} \theta_m^{-1/2} \dots \theta_k^{-1/2} d\theta \\ &= \frac{\Gamma(1/2)^{k-m+2}}{\Gamma(\frac{k-m+2}{2})}. \end{aligned}$$

En pratique cela revient à se restreindre à l'alphabet $A_k = \{a_1, a_m, \dots, a_k\} \subset A$. Soit donc $x_{1:n}$ un élément de A_k^n . Alors

$$\begin{aligned} P_{\mu_k}^n(x_{1:n}) &= \int_{\Theta_k} P_\theta(x_{1:n}) d\mu_k(\theta) \\ &= \frac{D(k)}{C(k)} \int_{\Theta_k} P_\theta(x_{1:n}) d\mu(\theta) \\ &\leq \frac{D(k)}{C(k)} \int_{S_{k-m+2}} P_\theta(x_{1:n}) d\mu(\theta) \\ &= \frac{D(k)}{C(k)} P_\mu^n(x_{1:n}), \end{aligned}$$

4 Classes enveloppe

et si θ est un élément de Θ_k ,

$$\begin{aligned}
 D(P_\theta^n \| P_{\mu_k}^n) &= \int_{A_k^n} \log \frac{P_\theta^n(x_{1:n})}{P_{\mu_k}^n(x_{1:n})} dP_\theta^n(x_{1:n}) \\
 &\geq \int_{A_k^n} \log \frac{P_\theta^n(x_{1:n})}{P_\mu^n(x_{1:n})} dP_\theta^n(x_{1:n}) + \log C(k) - \log D(k) \\
 &= D(P_\theta^n \| P_\mu^n) + \log C(k) - \log D(k) \\
 &\geq \frac{k-m+1}{2} \log \frac{n}{2\pi} + \log C(k) - \frac{5(k-m+2)}{3} \log e,
 \end{aligned}$$

en vertu de la proposition 5. En conséquence

$$\begin{aligned}
 R_{n,\mu_k}^{\text{Bayes}} &= \int_{\Theta_k} D(P_\theta^n \| P_{\mu_k}^n) d\mu_k(\theta) \\
 (4.6) \quad &\geq \log C(k) + \frac{k}{2} \log n - \frac{10 \log e + 3 \log 2\pi}{6} k - \frac{m-1}{2} \log n \\
 &\quad + \left(\frac{m-1}{2} \log 2\pi + \frac{5(m-2)}{3} \log e \right).
 \end{aligned}$$

Calculons maintenant $C(k)$. Remarquons tout d'abord que nous avons choisi m de telle sorte que tous les choix de $(\theta_m, \dots, \theta_k)$ dans le rectangle $[0, Ce^{-\alpha m}] \times \dots \times [0, Ce^{-\alpha k}]$ sont possibles, ce qui permet d'écrire les intégrales sur Θ_k comme intégrales sur ce rectangle. Nous avons alors

$$\begin{aligned}
 C(k) &= \int_{\Theta_k} \frac{d\theta_m \cdots d\theta_k}{\sqrt{\theta_1 \theta_m \cdots \theta_k}} \\
 &\geq \prod_{i=m}^k \int_0^{Ce^{-\alpha i}} \frac{d\theta_i}{\sqrt{\theta_i}} \\
 &= \prod_{i=m}^k 2\sqrt{C} e^{-\frac{\alpha}{2} i},
 \end{aligned}$$

d'où

$$\begin{aligned}
 \log C(k) &\geq (k-m+1) \left(1 + \frac{\log C}{2} \right) - \frac{\alpha \log e}{2} \sum_{i=m}^k i \\
 &= (k-m+1) \left(1 + \frac{\log C}{2} \right) - \frac{\alpha \log e}{4} (k(k+1) - m(m+1)).
 \end{aligned}$$

Nous pouvons alors injecter ce résultat dans la relation 4.6, et obtenir $R_{n,\mu_k}^{\text{Bayes}} \geq g(n, k)$,

4 Classes enveloppe

avec

$$g(n, k) = -\frac{\alpha \log e}{4} k^2 + \frac{k}{2} \log n + \left[1 + \frac{\log C}{2} - \frac{\alpha \log e}{4} - \frac{5 \log e}{3} - \frac{\log 2\pi}{2} \right] k \\ - \frac{m-1}{2} \log n + \left[\frac{\alpha(m^2 + m - 1) \log e}{4} - (m-1) \left(1 + \frac{\log C}{2} \right) \right. \\ \left. + \frac{m-1}{2} \log 2\pi + \frac{5(m-2) \log e}{3} \right].$$

En choisissant $k_n = \frac{1}{\alpha \log e} \log n$, seuls les deux premiers termes comptent, et nous obtenons

$$g(n, k_n) \sim \frac{1}{4\alpha \log e} \log^2 n.$$

Il ne nous reste plus qu'à énoncer le résultat obtenu :

Théorème 11. *On considère $\Lambda_{Ce^{-\alpha}}$ une classe enveloppe à décroissance exponentielle. Alors, lorsque $n \rightarrow \infty$,*

$$R_n^{\text{minimax}}(\Lambda_{Ce^{-\alpha}}) \geq (1 + o(1)) \frac{1}{4\alpha \log e} \log^2 n,$$

et en conséquence

$$R_n^{\text{minimax}}(\Lambda_{Ce^{-\alpha}}) \sim \frac{1}{4\alpha \log e} \log^2 n.$$

5 Et au-delà ?

Avec notre théorème 11 nous avons illustré l'intérêt du recours à l'entropie métrique et aux résultats de [HO97] pour calculer la redondance de classes de sources sans mémoire.

Dans le cas des classes enveloppe à décroissance rapide, de type exponentielle, nous avons constaté que l'entropie métrique restait « petite » et assurait des résultats fins sur l'encadrement de la redondance minimax. En revanche, des classes enveloppe à décroissance plus lente, de type polynomiale, entraîne un écart entre l'ordre de grandeur de la majoration et celui de la minoration. On peut cependant espérer que la minoration fournisse le bon ordre de grandeur, mais il reste un travail à accomplir pour minorer efficacement l'entropie métrique et ainsi tester cette hypothèse.

Par ailleurs dans ce travail a été développée, en collaboration étroite avec Elisabeth Gassiat, une technique pour minorer la redondance minimax, issue des travaux de Xie et Barron sur le mélange de Dirichlet. Peut-être pourrions-nous réutiliser cela pour d'autres classes de sources ?

Enfin il reste à traduire ce travail théorique sur la redondance en la découverte éventuelle d'un nouveau code adapté aux classes enveloppe à décroissance exponentielle.

Index

- α -affinité, 10, 17
- alphabet
 - dénombrable, 5, 12, 34, 44
 - fini, 13, 41
- Bayes
 - mélange de Bayes, 9, 15
 - prior, 5, 9, 14, 41
- classe enveloppe, 34, 35
 - à décroissance exponentielle, 35, 39, 41, 46
 - à décroissance polynomiale, 35, 38
- code, 5
 - binaire, 5
 - de Shannon, 6
 - longueur de code, 6
- covering number, 24
- D-divergence, 13, 14, 24
- dimension métrique, 25, 27, 30
- Dirichlet
 - mélange, 42
 - prior, 41, 44
- entropie
 - binaire, 7
 - finie, 22
 - métrique, 24, 25, 27, 30, 31, 36
 - relative, 6, 8, 13
- enveloppe intégrable, 13, 18, 35
- Fisher (information de), 22
- Hellinger (distance de), 10, 14, 21, 24, 36
- Huffman (code de), 7
- I-divergence, 13, 14
- information mutuelle, 9
- irrégulière (famille), 18, 20
- Kullback (divergence de), 8, 10, 15
- Laplace
 - méthode de Laplace, 21
 - transformée de Laplace, 10, 14, 15, 21
- packing number, 24
- redondance, 8
 - bayésienne, 8, 14, 15, 17, 41, 45
 - minimax, 8, 25, 27, 30, 34, 35, 39
- regret, 7
 - minimax, 8, 35
- Shannon, 6
- Shannon-Kraft-Mc Millan, 7
- Stirling (formule de), 38, 43
- totalemment borné, 24, 25, 27, 37

Bibliographie

- [BGG06] Stéphane Boucheron, Aurélien Garivier et Elisabeth Gassiat: *Coding on countably infinite alphabets*. article soumis, dec. 2006.
- [CB90] Bertrand S. Clarke et Andrew R. Barron: *Information-theoretic asymptotics of Bayes methods*. IEEE Transactions on Information Theory, 36(3) :453–471, 1990.
- [CB94] Bertrand S. Clarke et Andrew R. Barron: *Jeffrey’s prior is asymptotically least favorable under entropy risk*. Journal of Statistical Planning and Inference, 41 :37–60, 1994.
- [CT91] Thomas M. Cover et Joy A. Thomas: *Elements of Information Theory*. Wiley Series in Telecommunications. John Wiley & sons, 1991, ISBN 0-471-06259-6.
- [Gar06] Aurélien Garivier: *Modèles contextuels et alphabets infinis en théorie de l’information*. Thèse de doctorat, Université Paris-Sud 11, nov. 2006.
- [Hau96] David Haussler: *A general minimax result for relative entropy*. Rapport technique UCSC-CRL-96-26, 1996. <http://cite-seer.ist.psu.edu/haussler96general.html>.
- [HO97] David Haussler et Manfred Opper: *Mutual information, metric entropy and cumulative relative entropy risk*. The Annals of Statistics, 25(6) :2451–1492, 1997.
- [KT61] A. N. Kolmogorov et V. M. Tikhomirov: *ϵ -entropy and ϵ -capacity of sets in functional spaces*. Amer. Math. Soc. Trans. Ser. 2, 17 :277–364, 1961.
- [vdVW96] Aad W. van der Vaart et Jon A. Wellner: *Weak Convergence and Empirical Processes*. Springer Series in Statistics. Springer-Verlag, 1996, ISBN 0-387-94640-3.
- [WW96] E. T. Whittaker et G. N. Watson: *A Course of Modern Analysis*. Cambridge Mathematical Library. Cambridge University Press, Cambridge, reprint of the fourth (1927) edition édition, 1996, ISBN 0-521-58807-3. An Introduction to the general theory of infinite processes and of analytic functions; with an account of the principal transcendental functions.
- [XB97] Qun Xie et Andrew R. Barron: *Minimax redundancy for the class of memoryless sources*. IEEE Transactions on Information Theory, 43(2) :646–657, 1997.