

Théorie de l'Information
Notes de Cours (part 5)
2006-2007

SIC-SICOM

Maria-João Rendas

19 novembre 2006

Table des matières

6	Identification de modèles :	
	le principe de longueur de description minimale (MDL)	95
6.1	Introduction	95
6.2	Apprentissage (de modèles) comme compression de données	98
6.3	Codes et lois de probabilité	100
6.4	MDL (codage en deux parties)	103
6.4.1	Complexité de Kolmogorov MDL "Idéal"	103
6.4.2	Codage en deux parties	104
6.5	Codeurs universels et MDL (codage en une partie)	107
6.5.1	Maximum de Vraisemblance Normalisé comme Modèle Uni- versel Optimal	110
6.5.2	MDL (codage en un partie) et complexité stochastique	112
6.6	Approximations de la complexité stochastique	113
6.6.1	Maximum de Vraisemblance Généralisé	113
6.6.2	MDL et Compression	114
6.6.3	Interprétation géométrique	114
6.6.4	Interprétation Bayésienne	117
6.6.5	Interprétation prédictive	118
6.7	MDL Général pour la sélection de modèles paramétriques	120
6.7.1	Complexité paramétrique infinie	120
6.7.2	Sommaire	122

Chapitre 6

Identification de modèles : le principe de longueur de description minimale (MDL)

6.1 Introduction

Le problème d'identification de modèles peut être formulé comme celui de choisir une parmi plusieurs explications alternatives pour un ensemble de données $x^{(n)}$ (en nombre n limité, en toute situation pratique). Dans ce Chapitre, nous exposons une méthode pour l'identification de modèles connue par le nom de "*principe de description minimale*" ou encore par l'acronyme **MDL** (de l'anglais *Minimum Description Length*). Ses premières versions ont été proposées par N. Rissanen dans les années 1970 [2], et bâtissent sur les travaux d'autres chercheurs sur la *théorie de la complexité* (notamment le chercheur russe Kolmogorov [3]) qui datent des années 60. Cette méthodologie est basée sur deux idées fondamentales :

1. *Apprendre un modèle* pour le système qui a engendré les données observées consiste à *déceler les régularités* présentes dans l'ensemble observé. L'identification de modèles est une tâche d'*apprentissage*, où l'objectif est de capturer tant que possible les *régularités des données*, et, donc, être capable de trouver sa description la plus simple : séparer ce qui est *structure* de ce qui est *accidentel* (bruit).
2. Les *régularités* des données (l'ensemble de règles, ou propriétés, qui cet ensemble satisfait) peuvent être exploitées pour les *compresser*, c'est à dire, pour les décrire (exactement) avec un nombre de symboles minimal (inférieur à celui qui décrit les observations). Un modèle correspond donc à un *langage* pour décrire les données. Dans le cas de modèles décrits par des *familles de distributions de probabilité* (hypothèses composées dans la terminologie des tests de décision statistique), un bon modèle doit correspondre à un *codeur universel* pour les données dans la famille considérée.

Ces deux observations restent valables **indépendamment de considérations sur l'existence d'un vrai modèle probabiliste qui aurait engendré les données observées**. L'approche du MDL se distingue d'une manière fondamentale, sur ce point particulier, des approches statistiques classiques où ce "vrai" modèle (probabiliste) est un ingrédient fondamental (comme l'approche Bayésienne, où la distribution *a priori* est censée traduire notre connaissance sur l'identité de ce *vrai* modèle). Les différentes distributions dans un même modèle correspondent à des différentes réalisations des mêmes régularités structurelles. Bien que n'excluant pas le type de modèles "signal + bruit", pour le MDL le bruit est le nombre de bits nécessaires pour identifier *la* séquence observée à l'aide d'un modèle qui contraint ses régularités. Un modèle avec un niveau de bruit élevé implique uniquement que les données ne sont pas compressibles avec le modèle considéré. Pour le MDL, l'inférence est basée *uniquement sur les données observées*.

Une procédure d'inférence statistique est *consistante* si elle identifie la *vraie* hypothèse avec probabilité 1 (asymptotiquement, pour un nombre de données n très grand). Bien que le MDL ne dépende pas, comme nous venons de le dire, de considérations sur l'existence de cette *vraie distribution*, il est important que, si jamais les données sont une réalisation d'une des distributions contenues dans les modèles considérés, cette distribution soit correctement identifiée. Comme nous le verrons, le *MDL est consistant*.

Comme nous le verrons avec plus de détail le long de ce Chapitre, la méthode MDL trouve des analogies et relations avec d'autres méthodes et problèmes d'analyse de données :

- **Rasoir d'Occam**. Le "rasoir" d'Occam (ou *principe de parcimonie*), énoncé par le philosophe médiéval William Occam (siècle XVI, Angleterre), sous la forme "**pluritas non est ponenda sine necessitate**," (une pluralité ne doit pas être invoquée sans besoin), conduit, dans le cas du problème d'identification de modèles, à choisir **le modèle le plus simple qui n'est pas contredit par les données**. Comme nous le verrons plus tard, le principe du MDL conduit lui aussi à la recherche d'un équilibre entre la *complexité* du modèle qui est ajusté aux données, et sa *capacité* pour les décrire exactement. En fait, quand deux modèles s'ajustent également bien aux données, MDL choisit le plus simple (dans le sens qu'il permet description plus courte), et dans ce sens, il réalise une sorte de "rasoir d'Occam". Nous remarquons cependant que cette préférence n'est pas un choix "philosophique", qui impliquerait une quelconque supposition sur la simplicité de la nature, mais uniquement la traduction du fait qu'il n'y a pas de sens à identifier des modèles complexes avec un volume de données réduit. En fait, comme nous le verrons dans une section de ce Chapitre, **pour une même source le modèle sélectionné par MDL peut être de plus en plus complexe, à mesure que plus de données sont observées**.
- Le problème de **sur-ajustement** (*overfitting*, en anglais). Ce problème est bien connu dans les problèmes d'identification de modèles paramétriques, où l'**ordre du modèle**, c'est à dire, le nombre de degrés de liberté du modèle qui est utilisé pour décrire les observations, est inconnu et doit être aussi déterminé à partir des données. Le MDL donne une réponse formellement justifiée (sur les proprié-



FIG. 6.1 – William Occam.

tés des codes universels) pour le choix de l'ordre du modèle. Il conduit ainsi à un équilibre entre les capacités de *description* des données et les propriétés de *généralisation* du modèle identifié.

- **Approche Bayésienne.** L'approche Bayésienne (du nom du révérend anglais Thomas Bayes, XVIII-ème siècle), fournit une méthode systématique pour *combiner (fusionner) des informations imprécises*. Dans le cadre du problème d'identification de modèles, l'approche Bayésienne conduit aussi à la recherche d'un *équilibre entre la plausabilité pour que les données soient engendrées par le modèle, et sa probabilité a priori*. Une des objections qui sont soulevées à l'approche Bayésienne est l'arbitrarité du choix de ces distributions *a priori*. Nous pouvons essayer de contre-carrer le problème de sur-ajustement en associant des probabilités *a priori* plus grandes aux modèles plus simples (ce que le "rasoir" d'Occam semble suggérer : "la nature préfère la simplicité"). Comme nous le verrons, le MDL possède des relations étroites avec l'approche Bayésienne, et peut être interprété comme une méthodologie pour choisir la distribution *a priori*.
- **Codage prédictif.** Plusieurs méthodes d'identification ont été proposées pour lesquelles le critère pour la sélection du modèle est la *capacité prédictive du modèle*, c'est à dire, pour prédire des valeurs qui n'ont pas encore été observées. Comme nous le discuterons plus tard, le MDL possède aussi des relations avec



FIG. 6.2 – Thomas Bayes.

avec k '1's :

$$x_3 \rightarrow (k, i_k(x_3))$$

Un exemple d'une nature différente (et plus près de beaucoup de cas pratiques de problèmes d'identification de modèle), et posé par le codage des valeurs du tableau suivant

x_i	y_i
-5	-286.4260
-4	-160.7825
-3	-85.7549
-2	-26.8838
-1	3.9952
0	16.7044
1	7.0395
2	7.0602
3	1.1894
4	-5.0388
5	8.2898

La figure 6.3 représente y_i en fonction de x_i . La forme de cette courbe suggère une dépendance polynomiale de y_i en x_i :

$$y_i = P_k(x_i) = a_k x_i^k + a_{k-1} x_i^{k-1} + \dots + a_1 x_i + a_0$$

Nous pouvons alors coder efficacement les données (y_i, x_i) en envoyant d'abord l'ordre k du polynome $P_k(\cdot)$ et ses coefficients $\{a_n\}_{n=0}^k$ (ou ses zéros) suivis de mots de code pour les valeurs de $\{x_i\}_{i=1}^n$ et pour les résidus $\{\epsilon_i\}_{i=1}^n$ du modèle polynomial (différence entre la valeur du polynome $P(x_i)$ et les valeurs y_i) :

$$\{(y_i, x_i)\}_{i=1}^n \rightarrow (k, \{a_n\}_{n=0}^k, \{x_i\}_{i=1}^n, \{\epsilon_i\}_{i=1}^n)$$

Si, $n \gg 1$ et $k \ll n$ cette description des données doit être plus compacte que la description directe des données : nous avons exploité la relation (polynomiale) sous-jacente à tous les paires pour les décrire plus efficacement

Code de préfixe pour les entiers

Dans l'exemple précédent nous avons eut besoin (pour coder les séquences x_1 et x_2 , de coder des entiers. Nous allons maintenant introduire un **code (de préfixe) pour les entiers** qui sera utilisé par la suite.

Pour toute séquence binaire $x \in \{0, 1\}^*$, (de taille $\ell(x) \geq 1$), nous pouvons construire un code de préfixe pour x , en le faisant précéder le code binaire de x par un mot $C_{sd}(m)$ qui indique la taille m de la séquence binaire ultérieurement utilisée pour coder la valeur de x :

$$x \rightarrow C(x) = C_{sd}(m)C_m(x)$$

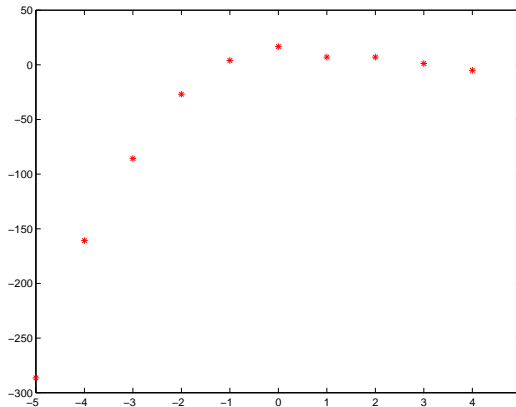


FIG. 6.3 – William Occam.

où $C_{sd}(m)$ est un code *auto-délimité* (*self-delimited*) pour l'entier m , et $C_m(x)$ est un code de longueur constante égale à m .

Le préfixe $C_{sd}(m)$ peut être, par exemple, une séquence de $m = \ell(x)$ zéros suivie d'un 1 :

$$x = 0100110; \quad \ell(x) = 7 \quad \Rightarrow \quad C(x) = 000000010100110,$$

où $m = 7$, $C_{sd}(7) = 00000001$ et $C_7(x) = 0100110$.

Le décodage de $C(x)$ est simple : nous identifions la fin du préfixe par l'occurrence du symbole '1'. Nous pouvons alors déterminer la taille m de la séquence qui code x , et donc, l'extraire du message.

Considérons maintenant le codage d'un entier $k \in \{1, 2, 3, \dots\}$. Comme $k \in \{1, 2^{\lceil \log k \rceil}\}$, nous pouvons coder k avec un code où le préfixe $C_{sd}(m)$ est une séquence de $m = \lceil \log k \rceil$ zéros suivie d'un 1.

Le **nombre total de bits** de ce code de préfixe pour les entiers est égal à $2^{\lceil \log k \rceil} + 1$. Par exemple, pour $k = 21$, le code sera

$$00000 1 10101$$

Le récepteur commence par compter le nombre de '0's : 5, il sait donc que les 5 bits après le '1' contiennent le code pour k .

6.3 Codes et lois de probabilité

Le lien entre le choix d'hypothèses/modèles et la compression des données n'a pas encore été établi. De l'étude du problème de codage source que nous avons menée dans les Chapitres précédents, nous savons que

1. à un code (de préfixe, complet) C sur un alphabet \mathcal{X} nous pouvons associer une loi de probabilité p_C sur \mathcal{X} :

$$C \rightarrow p_C(x) = 2^{-\ell(c(x))}, \forall x \in \mathcal{X}.$$

où $\ell(c(x))$ est la longueur du mot de code que C associe à $x \in \mathcal{X}$. L'inégalité de Kraft nous affirme que p_C est en général (pour des codes qui ne sont pas complets) une distribution non-normalisée (une *semi-mesure*), avec une masse totale inférieure ou égale à 1. Si le code est *complet*, alors, p_C est une loi de probabilité. Nous savons encore, du théorème de Shannon du codage sans pertes, que le code C est optimal (il conduit à une longueur moyenne de code minimale, égale à l'entropie de la source $H(X)$) :

$$L_C(X) = E[\ell(c(x))] = H(X).$$

2. à une loi de probabilité $p(x), x \in \mathcal{X}$, nous pouvons associer un code optimal C_p (de longueur moyenne minimale), avec des longueurs de mots de code

$$p(x) \rightarrow \ell(C_p(x)) = \lceil -\log p(x) \rceil.$$

Ce code (de Shannon-Fano) a une redondance inférieure à 1 bit. En fait, de la non-négativité de l'entropie relative, il découle que

$$E_p[-\log Q(x)] \geq E_p[-\log p(x)].$$

(*inégalité fondamentale de l'information.*)

Nous introduisons maintenant la terminologie précise qui sera utilisée dans tout ce Chapitre.

Définition 1 *Modèle probabiliste*

Un modèle probabiliste \mathcal{H} est un *ensemble* de sources décrites par lois de probabilité. Les *modèles paramétriques*, où les éléments de l'ensemble \mathcal{H} sont indexés par un vecteur (de dimension finie) de paramètres θ appartenant à un espace de paramètres Θ est particulièrement intéressant :

$$\mathcal{H} = \{p(X|\theta) : \theta \in \Theta\}, \quad \Theta \subset \mathbb{R}^k, k \geq 1.$$

△

Exemple 1 *Modèle Gaussien*

Pour des observations $x^{(n)} \in \mathbb{R}^n$, le modèle Gaussien

$$\mathcal{G} = \left\{ p(x^{(n)}|\theta) : \theta = (\mu, \Sigma) \right\}, \quad \mu \in \mathbb{R}^n, \quad \Sigma = \Sigma^T \geq 0 \in \mathbb{R}^{n \times n},$$

paramétré par la moyenne μ et la matrice de covariance Σ , où

$$p(x^{(n)}|\theta) = (2\pi \det(\Sigma))^{-1/2} \exp\left\{-\frac{1}{2}(x^{(n)} - \mu)^T \Sigma^{-1}(x^{(n)} - \mu)\right\},$$

est un modèle paramétrique, paramétré par $n + n(n+1)/2$ paramètres (où nous avons considéré les conditions de symétrie de la matrice de covariance. △

Définition 2 *Estimateur du Maximum de Vraisemblance*

Pour un modèle probabiliste \mathcal{H} et des données $x^{(n)} \in \mathcal{X}^{(n)}$, l'estimateur du Maximum de Vraisemblance est l'élément du modèle qui maximise la probabilité des données :

$$\hat{p}_{MV}(x^{(n)}) = \arg \max_{p \in \mathcal{H}} p(x^{(n)}).$$

Pour les modèles paramétriques, ce modèle est obtenue en prenant l'élément de \mathcal{H} qui est indexé par l'estimateur du Maximum de Vraisemblance des paramètres θ :

$$\hat{p}_{MV}(x^{(n)}) = p(\cdot | \hat{\theta}_{MV}(x^{(n)})), \quad \hat{\theta}_{MV}(x^{(n)}) = \arg \max_{\theta \in \Theta} p(x^{(n)} | \theta).$$

△

Définition 3 *Estimateur consistant*

Nous disons qu'un estimateur est consistant si, pour des observations

$$x^{(n)} \sim p_n \rightarrow_{n \rightarrow \infty} p^*, \quad x^{(n)} \in \mathcal{X}^n$$

nous avons

$$\lim_{n \rightarrow \infty} \hat{p}_{MV}(x^{(n)}) = p^*,$$

où la convergence est en probabilité. Pour les modèles paramétriques, où $p^* = p(\cdot | \theta^*)$,

$$\lim_{n \rightarrow \infty} \hat{\theta}_{MV}(x^{(n)}) = \theta^*.$$

△

Nottez que la consistance est une *propriété asymptotique*, elle concerne le comportement des estimateurs pour un grand nombre n de données : avec probabilité 1, l'estimateur tend vers la vraie mesure qui a engendré les données.

Exemple 2 *Modèle de Markov*

Nous avons défini dans un Chapitre précédant les chaînes de Markov d'ordre 1 (invariantes dans le temps) avec état $x_n \in \mathcal{X}$, où $|\mathcal{X}| = m$. Ce modèle est paramétré par la *distribution initiale* $p(x_0)$ ($m - 1$ paramètres) et par la *matrice de transition* $P_{ij} = p(x_{n+1} = i | x_n = j)$ ($m(m - 1)$ paramètres). La dimension de ce modèle est donc $(m + 1) \times (m - 1)$.

Nus introduisons maintenant les *modèles de Markov d'ordre k* , où la valeur à l'instant n dépend des k valeurs précédentes :

$$p(X_n = i | x_1^{n-1}) = p(X_n = i | x_{n-k}^{n-1}), \quad n > k.$$

Dans ce cas, le modèle est indexé par sa distribution initiale $p(x_1^k)$ avec $k \times (m - 1)$ paramètres, et par une matrice de transition de dimension $(m^k - 1) \times (m^k - 1)$.

Le modèle de Markov *binnaire* ($m = 2$) d'ordre k avec $\mathcal{X} = \{0, 1\}$, demande donc un

total de $k + (2^k - 1)^2$ paramètres.

Il peut être vérifié facilement que pour une chaîne de Markov d'ordre $k = 1$

$$\left[\hat{P}_{MV} \right]_{ij}(x^{(n)}) = \frac{n_{ij}}{n}, \quad (6.1)$$

où n_{ij} est le nombre de fois que le symbole i est précédé par le symbole j dans la séquence $x^{(n)}$. △

Définition 4 *Modèle de Bernoulli*

Le modèle de Bernoulli (valeurs binaires statistiquement indépendantes et identiquement distribuées), que nous noterons \mathcal{B} , peut être obtenu comme un modèle de Markov binaire d'ordre $k = 0$, avec une matrice de transition

$$P = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}.$$

Il est paramétré par un seul paramètre $\theta \in [0, 1]$, qui est la probabilité d'occurrence d'une des valeurs de \mathcal{X} :

$$\mathcal{B} = \left\{ p(x^{(n)}|\theta) : \theta \in [0, 1] \right\}, \quad p(x^{(n)}|\theta) = \prod_{i=1}^n p(x_i|\theta) = \theta^{n_1} (1 - \theta)^{n - n_1},$$

où n_1 est le nombre de '1's dans la séquence $x^{(n)}$. Il est aisé de vérifier que pour le modèle de Bernoulli

$$\hat{\theta}_{MV}(x^{(n)}) = \frac{n_1}{n}.$$

△

6.4 MDL (codage en deux parties)

6.4.1 Complexité de Kolmogorov MDL "Idéal"

Les exemples que nous avons présenté dans l'Introduction construisent des codes efficaces pour les messages considérés, mais ces codes sont trouvés, à chaque cas, par des méthodes *ad hoc*. Pour que l'approche MDL soit bien définie, il faut spécifier les types de codage qui sont des candidats pour coder les données.

Dans les années 60, **Kolmogorov** [3] a proposé la définition d'une **mesure de complexité intrinsèque d'un objet X** comme la *taille du plus petit programme qui serait capable de produire X* (une séquence, une image, une forme,...) *et s'arrêter*. Il a appelé cette mesure de complexité (que l'on peut démontrer qui est indépendante du langage de programmation choisi) **complexité algorithmique**. Dans le cadre de cette théorie (développée parallèlement par Solomonoff [4] et Chaitin [5]), les données sont décrites (codées) par un programme (qui doit s'exécuter sur une machine de Turing). La notion de complexité de Kolmogorov a depuis été reprise par plusieurs chercheurs, et est

d'une grande puissance formelle. Cependant, il est possible de démontrer qu'il n'existe pas un programme qui puisse la déterminer : **la complexité de Kolmogorov n'est pas calculable**, ce qui ne permet pas la considérer comme la base d'une méthode (pratique) d'identification.

Néanmoins, certains auteurs [6] ont proposé une version "idéale" du principe MDL, en considérant le choix du modèle dans lequel la complexité de Kolmogorov serait la plus petite.

Nous allons nous intéresser à des versions "pratiques" du principe du MDL, et donc à des méthodes de description qui soient *calculables*. Le prix à payer sera que nous devrons abandonner les langages de programmation génériques, pour passer à considérer la complexité d'un message par rapport à un ensemble de langages (une "encyclopédie"), et la complexité du message ne sera plus une propriété intrinsèque du message. Certaines séquences régulières ne seront pas comprimées.

Nous énonçons maintenant le principe du MDL d'une manière informelle, et présenterons une approche (codage en deux parties) qui, sous certaines conditions, en est une implémentation.

6.4.2 Codage en deux parties

Définition 5 *Principe du MDL (définition informelle)*

Soit \mathcal{M} un ensemble d'hypothèses (de modèles probabilistes \mathcal{H}), et $x^{(n)} \in \mathcal{X}^n$ des données. Le principe MDL nous dit de choisir le modèle $\mathcal{H}^* \in \mathcal{M}$ qui *comprime le plus les données*.

$$\mathcal{H}^* = \arg \min_{\mathcal{H} \in \mathcal{M}} \min_{H \in \mathcal{H}} L_H(x^{(n)}),$$

où $L_H(x)$ est la longueur du mot d'un *code optimal pour les données x correspondant au modèle H* . △

Les deux faits que nous avons établi dans la Section précédente justifient l'association de l'ensemble d'hypothèses probabilistes \mathcal{H} dans la définition 5 avec un ensemble de lois de probabilité (et donc de codes). La définition suivante précise une interprétation possible de la notion de "codage optimal avec un modèle" utilisée dans la définition informelle 5.

Définition 6 *MDL (codage en deux parties)*

Soient $\mathcal{H}^{(i)}$, $i = 1, 2, \dots$, un ensemble de modèles, et $x^{(n)}$ des données. Alors le *meilleur modèle pour les données*, \mathcal{H}^* est le plus petit modèle qui contient

$$H^* = \arg \min_{H \in \cup_{i=1,2,\dots} \mathcal{H}^{(i)}} L(H) + L(x^{(n)}|H), \quad (6.2)$$

où $L(H)$ est la taille nécessaire pour *décrire l'hypothèse*, et $L(x^{(n)}|H)$ le nombre de bits nécessaire pour *coder les données avec l'hypothèse (code) H* . △

Le premier terme dans l'équation (6.2) est la longueur de description de l'hypothèse H . Le deuxième terme est la longueur de description des données $x^{(n)}$ à l'aide de l'hypothèse H , et est une mesure de l'*ajuste du modèle aux données*. Remarquons que cette version en deux parties du MDL rend explicite l'idée de chercher un équilibre entre la **complexité du modèle utilisé** (modèles plus complexes demandront un $L(H)$ plus grand), et l'**ajuste des modèles aux données** $x^{(n)}$ (plus les données seront plausibles (probables) dans le modèle, plus petite sera la longueur de son code).

La définition précédente laisse sans réponse deux questions :

1. Comment détermine-t-on $L(H)$?
2. Comment calculer $L(x^{(n)}|H)$?

La réponse à la dernière question, le *choix de $L(x^{(n)}|H)$* , est relativement simple : si nous admettons qu'à chaque hypothèse H correspond une loi de probabilité $p(x^{(n)}|H)$ (car \mathcal{H} est un modèle probabiliste), alors nous devons prendre comme longueur de description des données la taille du code optimal :

$$L(x^{(n)}|H) = -\log p(x^{(n)}|H).$$

La longueur des mots de ce code est égale à moins la vraisemblance des données – ce qui correspond bien à une mesure d'ajustement. Nous savons que, si les données suivent la distribution $p(x^{(n)}|H^*)$, alors pour ce choix la taille $L(x^{(n)}|H)$ sera minimisée pour la bonne hypothèse $H = H^*$. Si (comme nous allons le voir) $L(H)$ est indépendante de n (ou si $L(H)$ croît sub-linéairement en n), alors, pour des grandes valeurs de n la somme dans l'équation (6.2) sera dominée par $L(x^{(n)}|H)$. La propriété d'équi-répartition asymptotique nous garanti alors que si les données sont une réalisation d'une des hypothèses contenues dans un des modèles, alors, la méthode MDL identifiera ce modèle avec probabilité 1 (elle est donc consistante).

La réponse à la première question, le *choix de $L(H)$* , est moins simple : si nous choisissons arbitrairement la façon dont les hypothèses sont codées, nous pouvons, avec des schémas différents, avoir des codes qui attribuent des tailles très différentes aux hypothèses, et nous n'aurions donc pas une façon objective de choisir le modèle. Comme nous le verrons, ce problème a conduit à la proposition de différentes versions alternatives du principe du MDL.

Souvent, l'ensemble de modèles a une *structure emboîtée*

$$\mathcal{M} = \cup_{k>0} \mathcal{H}^{(k)},$$

où la dimension de $\mathcal{H}^{(k)}$ croît avec k , et

$$\mathcal{H}^{(k)} \subset \mathcal{H}^{(k+1)}, k = 1, 2, \dots$$

C'est le cas, par exemple, quand chaque $\mathcal{H}^{(k)}$ est le modèle d'observations bruitées de l'ensemble de polynômes d'ordre k ou pour l'ensemble modèles de Markov d'ordre k . Dans ces cas, pour chaque valeur de k nous pouvons utiliser un code fixe, en deux parties, pour coder les différentes hypothèses contenues dans $\mathcal{H}^{(k)}$, comme c'est illustré par l'exemple suivant.

Exemple 3 *Un test du MDL pour les modèles de Markov*

Nous pouvons décrire un modèle de Markov (binaire) d'ordre k en indiquant en premier lieu la valeur de k , suivie de la valeur du paramètre $\theta \in [0, 1]^{2^k}$:

$$H \in \mathcal{H}^{(k)} \rightarrow (k, \theta).$$

L'entier k peut être codé avec le code présenté page 100, qui demande $2\lceil \log k \rceil + 1$ bits :

$$L(k) = 2\lceil \log k \rceil + 1.$$

Pour coder le paramètre θ , nous devons indiquer les valeurs des 2^k probabilités p_{1i} , pour que le symbole 1 soit précédé par la i -ème séquence s_i de k symboles binaires, $s_i \in \{0, 1\}^k$, (cela détermine la probabilité pour que le symbole '0' soit précédé par la même séquence : $p_{0i} = 1 - p_{1i}$).

Nous remarquons maintenant que nous pouvons restreindre l'ensemble $\mathcal{H}^{(k)}$ de tous les modèles de Markov d'ordre k à un sous-ensemble discret fini $\mathcal{H}^{(k)'}$. En fait, uniquement les modèles qui peuvent être identifiés à partir d'une séquence de taille n sont pertinents, ce qui implique que les probabilités p_{1i} soient de la forme de l'équation (6.1), page 103, où, nécessairement

$$n_{ij} \in \{0, \dots, n\}.$$

Construire un code qui associerait des mots de code à des hypothèses qui ne correspondent pas à des estimés possibles aurait conduit à une valeur supérieure pour $L(H)$ sans que les valeurs correspondantes de $L(x^{(n)}|H)$ puissent être plus petites.

Il n'y a donc que $n + 1$ valeurs possibles pour chaque probabilité p_{1i} (fréquence n_{1j}), qui peuvent être codés avec un nombre de bits non supérieur à $\lceil \log(n + 1) \rceil \simeq \log(n + 1)$. Comme le paramètre θ du modèle de Markov est défini par 2^k de ces valeurs, pour identifier un modèle particulier nous avons besoin d'un nombre de bits total égal à

$$L(H) = L(k) + L(\theta) = 2 \log k + 1 + 2^k \log(n + 1).$$

Le critère pour choisir le modèle de Markov (l'ordre et les paramètres) est donc la solution de

$$\min_k \min_{\theta} -\log p(x^{(n)}|\theta) + 2 \log k + 1 + 2^k \log(k + 1),$$

ou encore, en faisant appel à la définition d'estimateur de Maximum de Vraisemblance

$$\min_k -\log p(x^{(n)}|\hat{\theta}_{MV}) + 2 \log k + 1 + 2^k \log(k + 1).$$

△

6.5 Codeurs universels et MDL (codage en une partie)

La Section précédente a illustré l'application du principe du MDL avec une approche de codage en deux parties : codage de l'hypothèse (avec $L(H)$ bits) suivie du codage des données dans l'hypothèse (avec $-\log \hat{p}_{\mathcal{H}}(x^{(n)})$ bits). Cependant, le choix du codage effectué est heuristique, et n'identifie pas d'une façon formelle le code qui doit être associée à un modèle probabiliste \mathcal{H} . Nous allons introduire la notion *codeur universel dans une classe de modèles*, qui permet de dépasser ce problème.

Définition 7 *Codeur universel (dans une classe de modèles)*

Soit \mathcal{H} un modèle probabiliste (et donc un ensemble de longueurs de codes optimaux) pour des données $x^{(n)}$. Un code \bar{L} est universel pour \mathcal{H} s'il est capable (asymptotiquement) de coder la séquence $x^{(n)}$ avec un nombre de bits égal au code optimal pour $x^{(n)}$ dans le modèle \mathcal{H} :

$$\forall x^{(n)} \quad \lim_{n \rightarrow \infty} \frac{1}{n} \bar{L}(x^{(n)}) = \lim_{n \rightarrow \infty} \min_{H \in \mathcal{H}} \frac{1}{n} L(x^{(n)}|H). \quad (6.3)$$

Notons que si la séquence $x^{(n)}$ est une réalisation d'une des hypothèses dans \mathcal{H} , alors ce limite est égal au taux d'entropie $\bar{H}(X)$. \triangle

Nous remarquons que cette définition de code universel correspond à une notion de code universel *ponctuel*, c'est à dire, la propriété (6.3) est vérifiée *pour toutes les séquences* $x^{(n)}$. Dans le Chapitre précédent nous avons présenté, dans le contexte de l'étude de l'algorithme de Lempel-Ziv, une notion de code universel qui est basée sur la longueur *moyenne* des mots de code.

Les deux exemples suivants montrent que *les codes en deux parties sont des codes universels*.

Exemple 4 *Nombre fini d'hypothèses*

Soit $x^{(n)}$ une séquence binaire, et considérons le modèle suivant :

$$\mathcal{H} = \{\mathcal{B}(\theta_i), i = 1, \dots, 9\}, \quad \theta_i = \frac{i}{10},$$

de façon que

$$-\log p(x^{(n)}|\mathcal{B}(\theta_i)) = -n_1 \log(i/10) - (n - n_1) \log\left(\frac{10-i}{10}\right).$$

Nous pouvons construire un code universel (en deux parties) pour ce modèle de la façon suivante. Nous codons d'abord i^* ,

$$i^* = \arg \max_{i=1, \dots, 9} p(x^{(n)}|\mathcal{B}(\theta_i))$$

en utilisant un code uniforme en $\{1, \dots, 9\}$, avec donc $\lceil \log 9 \rceil$ bits. Le message $x^{(n)}$ est ensuite codé avec le code $\mathcal{B}_{\theta_{i^*}}$, ce qui demande

$$L_{i^*} = -\log p\left(x^{(n)}|\mathcal{B}(\theta_{i^*})\right) = L_{\mathcal{B}(\theta_{i^*})}(x^{(n)})$$

bits. Ce code a une taille totale

$$\bar{L}_{\mathcal{H}}(x^{(n)}) = \lceil \log 9 \rceil - \log p\left(x^{(n)}|\mathcal{B}(\theta_{i^*})\right) = \lceil \log 9 \rceil + \min_{H \in \mathcal{H}} L_H(x^{(n)}).$$

Nous pouvons facilement constater que

$$\lim_{n \rightarrow \infty} \frac{1}{n} \bar{L}(x^{(n)}) = \lim_{n \rightarrow \infty} \frac{1}{n} \min_{H \in \mathcal{H}} L_H(x^{(n)}),$$

qui est la limite atteinte par le meilleur code pour $x^{(n)}$ dans le modèle \mathcal{H} , et que le code présenté est donc universel pour le modèle considéré. \triangle

Dans cet exemple nous avons utilisé un *code uniforme* pour coder les éléments $H \in \mathcal{H}$ (l'index i^* de l'exemple précédent), ce qui correspond à admettre une distribution uniforme pour toutes les hypothèses H du modèle \mathcal{H} . Autres choix seraient possibles. Cependant, ce choix minimise la redondance du code dans le *pire cas* (le cas où la vraie hypothèse reçoit le mot de code le plus long pour la distribution admise).

Exemple 5 *Modèle avec un ensemble dénombrable d'hypothèses*

Considérons maintenant le cas d'un modèle avec un nombre infini d'hypothèses :

$$\mathcal{H} = \{H_1, H_2, \dots\}.$$

Un codage en deux parties pour ce modèle peut être obtenu, d'une façon analogue à l'exemple précédent, en considérant le codage de l'index i^* avec le code pour les entiers que nous avons introduit page 100, qui demande $2\lceil \log k \rceil + 1$ bits pour coder l'entier k . Nous devons donc maintenant choisir

$$i^* = \arg \inf_{i=1,2,\dots} \left\{ 2\lceil \log i \rceil + 1 + L_{H_i}(x^{(n)}) \right\},$$

et donc

$$\bar{L}(x^{(n)}) = \inf_{i=1,2,\dots} \left\{ 2\lceil \log i \rceil + 1 + L_{H_i}(x^{(n)}) \right\}.$$

Contrairement au cas précédent, nous ne pouvons plus borner par une constante la redondance de ce code par rapport au code qui atteint la longueur de code minimale dans le modèle \mathcal{H} . Par contre,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \bar{L}(x^{(n)}) = \lim_{n \rightarrow \infty} \inf_{i=1,2,\dots} \left\{ \frac{2 \log i + 1}{n} - \frac{1}{n} \log p(x^{(n)}|H_i) \right\},$$

et, si les observations $x^{(n)}$ sont une réalisation d'une des lois de probabilité $p(\cdot|H_{i_0})$ ¹ dans \mathcal{H} ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \bar{L}(x^{(n)}) = \bar{H}(p_{i_0})$$

¹et si le modèle correspondant est ergodique.

où $\overline{H}(p_{i^0})$ est le taux d'entropie associé à la loi de probabilité correspondante à la vraie distribution des données. Pour cet exemple, le taux avec lequel le nombre de bits par symbole approche le taux d'entropie de la source n'est plus constant comme pour l'exemple précédent. \triangle

L'association entre codes et lois de probabilité nous permet d'associer aux codes utilisés pour décrire le modèle (avec des longueurs de code $L(H)$) des lois de probabilité définies dans l'ensemble des hypothèses contenues dans \mathcal{H}

$$p(H) = 2^{-L(H)}, \quad H \in \mathcal{H},$$

établissant ainsi un pont entre le codage MDL en deux parties et l'approche Bayésienne. L'exemple suivant illustre cette relation, montrant que les marginales de Bayes sont aussi des modèles (codes) universels.

Exemple 6 *Modèle Bayésien universel*

Soit \mathcal{H} un ensemble fini ou dénombrable d'hypothèses (lois de probabilité), paramétrées par un vecteur de paramètres $\theta \in \Theta$:

$$\mathcal{H} = \left\{ p(x^{(n)}|\theta) : \theta \in \Theta \right\}.$$

Soit W une distribution de probabilité dans Θ . À chaque distribution W nous pouvons associer un modèle de mélange pour les observations :

$$p(x^{(n)}|\mathcal{H}, W) = \sum_{\theta \in \Theta} W(\theta)p(x^{(n)}|\theta). \quad (6.4)$$

Il est immédiat que $p(x^{(n)}|\mathcal{H}, W)$ est un modèle universel pour \mathcal{H} . La longueur du code correspondant est :

$$\overline{L}_{\mathcal{H}, W}(x^{(n)}) = -\log p(x^{(n)}|\mathcal{H}, W).$$

Et donc

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \overline{L}_{\mathcal{H}, W}(x^{(n)}) &\stackrel{(a)}{\leq} - \lim_{n \rightarrow \infty} \frac{1}{n} \log W(\theta)p(x^{(n)}|\theta) \\ &\leq - \lim_{n \rightarrow \infty} \frac{1}{n} \max_{\theta} \left(\log W(\theta) + \log p(x^{(n)}|\theta) \right) \end{aligned} \quad (6.5)$$

où (a) est vraie pour tout choix de θ dans le membre droit. Ceci découle du fait que tous les termes de la somme sont non-négatifs, et de la monotonie du logarithme. Si nous notons θ_0 la valeur de θ correspondante à ce maximum,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \overline{L}_{\mathcal{H}, W}(x^{(n)}) \leq - \lim_{n \rightarrow \infty} \frac{1}{n} \log p(x^{(n)}|\theta_0) = \overline{H}(p(\cdot|\theta_0)),$$

car, parce que $W(\theta_0)$ est une constante indépendante de n ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log W(\theta_0) = 0.$$

Ceci démontre que la distribution de mélange Bayésienne de l'équation (6.4) conduit à un code universel pour les données. Elle est donc un exemple de *modèle universel* (pour le modèle \mathcal{H}). Une analyse de l'équation (6.5) nous montre que **le code universel de Bayes est supérieur au codage en deux parties** avec un code $L(H)$ dérivé de la distribution *a priori* W . Dans ce dernier cas, nous serions conduits à une longueur de code

$$-\min_{\theta} \log p(x^{(n)}|\theta) + \log W(\theta),$$

qui, sauf pour le cas où $p(x^{(n)}|\theta) = 0, \forall \theta \neq \theta_0$, est strictement supérieur à la valeur dans le membre gauche de (6.5), correspondante au modèle de mélange. \triangle

6.5.1 Maximum de Vraisemblance Normalisé comme Modèle Universel Optimal

Définition 8 *Pénalité* $\mathcal{P}_{\bar{P}, \mathcal{H}, x^{(n)}}$

Soit \mathcal{H} un modèle probabiliste, et soit \bar{P} une loi de probabilité définie en \mathcal{X}^n . Soient $x^{(n)}$ les observations, $x^{(n)} \in \mathcal{X}^n$.

La *pénalité de la loi \bar{P} par rapport au modèle \mathcal{H} pour les observations $x^{(n)}$* est, par définition

$$\mathcal{P}_{\bar{P}, \mathcal{H}, x^{(n)}} = -\log \bar{P}(x^{(n)}) - \min_{p \in \mathcal{H}} \left(-\log p(x^{(n)}) \right).$$

Cette pénalité est la différence entre (i) la taille du code associé à \bar{P} , et le (ii) nombre de bits nécessaire pour coder les observations *avec le meilleur code dans le modèle \mathcal{H}* . Pour le cas de modèles paramétriques, $\mathcal{H} = \{p(\cdot|\theta) : \theta \in \Theta\}$,²

$$\mathcal{P}_{\bar{P}, \mathcal{H}, x^{(n)}} = -\log \bar{P}(x^{(n)}) + \log p \left(x^{(n)} | \hat{\theta}_{MV}(x^{(n)}) \right). \quad (6.6)$$

Avec cette définition, nous pouvons déjà constater qu'une loi \bar{P} sera un modèle (code) universel pour le modèle \mathcal{H} si pour toute séquence $x^{(n)}$ la pénalité croît plus lentement que n

$$\forall x^{(n)} \quad \lim_{n \rightarrow \infty} \frac{1}{n} \mathcal{P}_{\bar{P}, \mathcal{H}, x^{(n)}} = 0.$$

La notion de pénalité, qui dépend de la séquence observée $x^{(n)}$, ne doit pas être confondue avec celle de *redondance*, introduite dans les Chapitres précédents, et qui concerne le comportement *en moyenne* de la longueur des mots du code.

Il est intéressant d'étudier la pénalité pour un modèle (code) universel pour \mathcal{H} . Cependant, la valeur de la pénalité, telle qu'elle est proposée dans la définition précédente, dépend de la séquence observée $x^{(n)}$. Pour certaines séquences elle peut même être négative, et en même temps qu'elle prend des valeurs grandes pour d'autres séquences. La définition suivante considère la plus grande valeur (sur toutes les observations $x^{(n)}$)

²Nous admettons ici que l'estimée du Maximum de Vraisemblance est bien définie.

de la pénalité pour un code \bar{P} et un modèle \mathcal{H} .

Définition 9 *Pénalité au pire cas* $\mathcal{R}_{\bar{P},\mathcal{H}}$

Soit \mathcal{H} un modèle probabiliste, et \bar{P} une loi de probabilité en \mathcal{X}^n . La *pénalité au pire cas* de \bar{P} par rapport à \mathcal{H} est, par définition

$$\mathcal{R}_{\bar{P},\mathcal{H}} = \max_{x^{(n)} \in \mathcal{X}^n} \mathcal{P}_{\bar{P},\mathcal{H},x^{(n)}}. \quad (6.7)$$

Nous sommes en ce moment en position de définir un **modèle universel optimal** \bar{P}^* comme celui qui *minimise la pénalité au pire cas* :

$$\begin{aligned} \bar{P}^* &= \arg \min_{\bar{P}} \mathcal{R}_{\bar{P},\mathcal{H}} \\ &= \arg \min_{\bar{P}} \max_{x^{(n)} \in \mathcal{X}^n} \left\{ -\log \bar{P}(x^{(n)}) + \log p \left(x^{(n)} | \hat{\theta}_{MV}(x^{(n)}) \right) \right\}, \end{aligned}$$

où nous avons considéré le cas d'un modèle probabiliste paramétrique. Ce problème d'optimisation a une solution qui est égale à

$$\bar{P}_{m\bar{v}n}(x^{(n)}) = \frac{p \left(x^{(n)} | \hat{\theta}_{MV}(x^{(n)}) \right)}{\sum_{y^{(n)} \in \mathcal{X}^n} p \left(y^{(n)} | \hat{\theta}_{MV}(y^{(n)}) \right)}, \quad (6.8)$$

quand le dénominateur de cette expression est fini.

La distribution (6.8) est connue par le nom de *distribution de Shtarkov*, du nom du chercheur qui a établi son optimalité. Elle associe à chaque possible séquence observée $x^{(n)}$ une probabilité qui est proportionnelle à celle qui lui est assignée par le modèle correspondant à l'estimateur du Maximum de Vraisemblance pour $x^{(n)}$ (la distribution qui lui attribue la probabilité la plus grande). Ceci explique la notation "**m\bar{v}n**" dans l'expression précédente, indiquant "*Maximum de Vraisemblance Normalisé*".

La distribution définie en (6.8) est toujours bien définie quand \mathcal{H} est un ensemble *fini*. Autrement, la somme dans le dénominateur peut être infinie, et donc $\bar{P}_{m\bar{v}n}$ n'est pas défini.

Définition 10 *Complexité paramétrique* $\mathcal{C}(\mathcal{H})$

Soit \mathcal{H} un modèle probabiliste paramétrique. Sa *complexité* est, par définition

$$\mathcal{C}(\mathcal{H}) = \log \sum_{x^{(n)} \in \mathcal{X}^n} p \left(x^{(n)} | \hat{\theta}_{MV}(x^{(n)}) \right). \quad (6.9)$$

Avec cette définition, la distribution de Shtarkov s'écrit

$$\log \bar{P}_{m\bar{v}n}(x^{(n)}) = \log p \left(x^{(n)} | \hat{\theta}_{MV}(x^{(n)}) \right) - \mathcal{C}(\mathcal{H}).$$

De par sa définition, nous pouvons constater que la complexité paramétrique d'un modèle, $\mathcal{C}(\mathcal{H})$ est d'autant plus grande que le modèle \mathcal{H} peut décrire (avec probabilité élevée) un grand nombre de séquences $x^{(n)}$. Comme nous le verrons, la complexité stochastique $\mathcal{C}(\mathcal{M})$ est liée au nombre de degrés de liberté du modèle \mathcal{M} .

La démonstration de l'optimalité de $\bar{P}_{m\,v\,n}$ est immédiate. Si nous utilisons (6.8) dans la définition (6.6), nous obtenons

$$\mathcal{R}_{\bar{P}_{m\,v\,n}, \mathcal{H}} = \log \sum_{x^{(n)} \in \mathcal{X}^n} p \left(x^{(n)} | \hat{\theta}_{MV}(x^{(n)}) \right) = \mathcal{C}(\mathcal{H}),$$

indépendamment de $x^{(n)}$. La pénalité devient donc indépendante de la séquence $x^{(n)}$. Comme $\forall p \neq \bar{P}_{m\,v\,n}$ il doit exister au moins un $z^{(n)} \in \mathcal{X}^n$ tel que $p(z^{(n)}) < \bar{P}_{m\,v\,n}(z^{(n)})$,

$$\begin{aligned} \mathcal{R}_{p, \mathcal{H}} &= \max_{x^{(n)} \in \mathcal{X}^n} \mathcal{P}_{p, \mathcal{H}, x^{(n)}} \\ &\geq \mathcal{P}_{p, \mathcal{H}, z^{(n)}} = -\log p(z^{(n)}) + \log p \left(z^{(n)} | \hat{\theta}_{MV}(z^{(n)}) \right) \\ &> -\log \bar{P}_{m\,v\,n}(z^{(n)}) + \log p \left(z^{(n)} | \hat{\theta}_{MV}(z^{(n)}) \right) \\ &= \mathcal{R}_{\bar{P}_{m\,v\,n}, \mathcal{H}, z^{(n)}} = \mathcal{R}_{\bar{P}_{m\,v\,n}, \mathcal{H}}. \end{aligned}$$

Ceci montre que la pénalité au pire cas de toutes les autres distributions p dans \mathcal{H} doit être supérieure à celle de $\bar{P}_{m\,v\,n}$.

6.5.2 MDL (codage en un partie) et complexité stochastique

Nous présentons maintenant une formulation du principe de la longueur de description minimale, qui fait appel à la notion de modèle universel optimal présentée dans la section précédente, et qui permet de résoudre quelques problèmes associés aux codes en deux parties étudiés dans la Section 6.4.

Nous conduisons cette présentation dans le cadre simple de [choix entre deux modèles \$\mathcal{H}_1\$ et \$\mathcal{H}_2\$](#) , et en admettant que la *complexité de ces deux modèles est finie*, de façon que les *distributions de Shtarkov correspondantes*, que nous notons $\bar{P}_{m\,v\,n}(\cdot | \mathcal{H}_i)$, $i = 1, 2$, sont bien définies. Nous reviendrons dans une section postérieure sur le problème de choisir entre un nombre infini de modèles.

Définition 11 *Principe MDL (codage en une partie)*

Soient $\mathcal{H}_1, \mathcal{H}_2$ des modèles alternatifs pour les observations $x^{(n)}$, et $\bar{P}_{m\,v\,n}(\cdot | \mathcal{H}_i)$, $i = 1, 2$, les distributions qui minimisent la pénalité au pire cas par rapport aux modèles \mathcal{H}_i , $i = 1, 2$, respectivement. Alors, le principe du MDL nous dit de choisir le modèle \mathcal{H}_{j^*} pour lequel la longueur du mot de code qui est associé aux observations par le modèle universel optimal est minimale :

$$j^* = \arg \min_{j=1,2} -\log \bar{P}_{m\,v\,n} \left(x^{(n)} | \mathcal{H}_j \right) \quad (6.10)$$

$$= \arg \min_{j=1,2} \left\{ -\log p \left(x^{(n)} | \hat{\theta}_j(x^{(n)}) \right) + \mathcal{C}(\mathcal{H}_j) \right\}. \quad (6.11)$$

△

Cette dernière expression nous donne une interprétation en deux parties de cette version reformulée du MDL : le modèle optimal \mathcal{H}_{j^*} réalise un **équilibre entre la capacité de décrire exactement les données** (le terme $p(x^{(n)}|\hat{\theta}_j(x^{(n)}))$) et sa **complexité paramétrique** (mesurée par le terme $\mathcal{C}(\mathcal{H}_j)$).

La définition précédente du principe du MDL implique donc que le modèle choisit pour les données est **celui par lequel le code universel a le meilleur comportement (pénalité minimale) dans le pire cas**.

Définition 12 *Complexité stochastique des données* $\bar{L}(x^{(n)}, \mathcal{H})$

Soit \mathcal{H} un modèle probabiliste (paramétrique). La complexité des observations $x^{(n)}$ par rapport au modèle \mathcal{H} est, par définition, la probabilité qui lui est attribuée par la distribution de Shtarkov associée à \mathcal{H} :

$$\bar{L}(x^{(n)}, \mathcal{H}) = -\log \bar{P}_{m\bar{v}n}(x^{(n)}|\mathcal{H}).$$

△

Le principe du MDL que nous venons de formuler, nous indique donc de prendre *le modèle pour lequel la complexité stochastique des observations est minimale*.

En général, la complexité paramétrique d'un modèle (et donc la complexité stochastique des observations dans ce modèle) ne peut pas être calculée, ni même numériquement, la seule exception connue étant le cas Gaussien. Cependant, des approximations de la complexité stochastique des données $\bar{L}(x^{(n)}, \mathcal{H})$, valables pour n grand, peuvent être déterminées, comme c'est présenté dans la section suivante.

6.6 Approximations de la complexité stochastique

6.6.1 Maximum de Vraisemblance Généralisé

Pour les problèmes de décision entre plusieurs hypothèses composées \mathcal{H}_i , (où chaque hypothèse est associée à une *famille* de lois de probabilité), le principe du **Maximum de Vraisemblance Généralisé** nous dit de choisir l'hypothèse \mathcal{H}_{i^*} qui maximise

$$p_{\hat{p}_i}(x^{(n)}) = \max_{p \in \mathcal{H}_i} p(x^{(n)}) + c,$$

où c est une constante qui détermine la performance du test (la probabilité des différents types d'erreur). Souvent, une version simplifiée de ces tests est appliquée, qui consiste à ignorer la constante c et à choisir simplement l'hypothèse qui maximise la probabilité des données. Ceci correspond, dans la perspective du MDL que nous venons d'énoncer dans la section précédente, à négliger la complexité paramétrique de chaque modèle, $\mathcal{C}(\mathcal{H}_i)$, et pour des modèles avec un nombre de degrés de liberté différents, en général

conduit à un choix systématique du modèle le *plus complexe* (celui qui a le plus grand nombre de degrés de liberté, et qui peut donc décrire plus précisément les observations $x^{(n)}$).

Le MDL essaie de contrarier cette tendance pour choisir un modèle de complexité élevée en considérant la distribution de Shtarkov, qui utilise, comme nous avons vu, une version *normalisée* du Maximum de Vraisemblance. Le plus grand sera le nombre d'observations qui peuvent être bien décrites par les distributions contenues dans un modèle, plus grande sera sa complexité paramétrique $\mathcal{C}(\mathcal{H})$, et donc plus le modèle doit finement décrire les observations $x^{(n)}$ pour qu'il puisse être choisi.

6.6.2 MDL et Compression

L'idée originale de Rissanen qui justifie le principe du MDL est de considérer le problème d'identification de modèles comme celui d'apprendre le modèle qui le mieux exprime les régularités présentes dans les données. Nous venons de voir que cela correspond à choisir le modèle pour lequel la distribution de Shtarkov associée conduit à un code de longueur minimale. L'utilisation de $\bar{P}_{m\bar{v}n}$ semble justifiée par les deux remarques suivantes :

1. le mieux la meilleure distribution dans \mathcal{H}_j décrira les données, le plus petite sera la longueur de code.
2. aucune distribution dans chaque \mathcal{H}_j n'est donnée une préférence : en effet, la pénalité (qui a le sens d'une redondance *pour chaque séquence*) est la même pour toutes les séquences $x^{(n)}$, indépendamment de la loi qui a générée la séquence observée. Ce code est **le seul code de préfixe avec cette propriété** : $\bar{P}_{m\bar{v}n}$ traite toutes les distributions de la même manière.

La complexité stochastique des données $x^{(n)}$ dans un modèle \mathcal{H} peut être interprétée comme la *quantité d'information sur le modèle contenue dans les données*. Elle est la somme de deux termes. La complexité paramétrique du modèle, qui mesure la quantité de *structure* dans les données (dans le contexte du modèle \mathcal{H}), et le terme $-\log p(x^{(n)}|\hat{\theta}(x^{(n)}))$, qui mesure la quantité de *bruit* dans les données.

6.6.3 Interprétation géométrique

La complexité paramétrique d'un modèle peut être interprétée comme le **nombre de distributions discernibles dans le modèle**. Intuitivement, comme nous l'avons déjà remarqué, plus de distributions différentes un modèle contient plus grand est le nombre de messages qu'il peut décrire, et donc plus grand est le risque de over-fitting. Cependant, si ces distributions sont "similaires", dans le sens qu'elles décrivent essentiellement le même ensemble d'observations, elles ne doivent pas contribuer séparément à la complexité stochastique du modèle. Dans cette perspective, la complexité devrait être une mesure du nombre de distributions "différentes" qu'un modèle donné contient. L'analyse suivante montre dans un cas simple que $\mathcal{C}(\mathcal{H})$ mesure exactement ceci.

Soit $\mathcal{H} = \{p(\cdot|\theta_i), i \in \{1, \dots, M\}\}$ un modèle avec un nombre fini de distribu-

tions. Alors

$$\begin{aligned}
\mathcal{C}(\mathcal{H}) &= \log \sum_{x^{(n)} \in \mathcal{X}^n} p(x^{(n)} | \hat{\theta}_{MV}(x^{(n)})) \\
&= \log \sum_{j=1, \dots, M} \sum_{x^{(n)} : \hat{\theta}_{MV}(x^{(n)}) = \theta_j} p(x^{(n)} | \theta_j) \\
&= \log \sum_{j=1, \dots, M} \left(1 - \sum_{x^{(n)} : \hat{\theta}_{MV}(x^{(n)}) \neq \theta_j} p(x^{(n)} | \theta_j) \right) \\
&= \log \left(M - \sum_{j=1, \dots, M} \Pr \{ \hat{\theta}_{MV}(x^{(n)}) \neq \theta_j | \theta_j \} \right) \\
&\leq \log M .
\end{aligned}$$

Nous voyons donc que la complexité paramétrique est dans ce cas simple égale à la différence entre le nombre de modèles dans \mathcal{H} (M) et la probabilité pour que les modèles soient confondus. L'existence de distributions qui peuvent être confondues avec une probabilité élevée conduit à une diminution de la complexité paramétrique (ces distributions ne sont pas comptabilisées d'une façon indépendante). Pour des grandes valeurs de n , cette probabilité d'erreur tend vers zéro (sauf dans des cas pathologiques), et la complexité paramétrique tend vers $\log M$, le (logarithme du) nombre de modèles différentes dans \mathcal{H} .

Ce fait a été formellement établi en [1], à l'aide d'arguments sur la géométrie de la variété Riemannienne associée à des modèles probabilités paramétriques $\{p(x^{(n)} | \theta) : \theta \in \Theta\}$ par la définition d'une métrique qui est donnée par la matrice de Fisher $I(\theta)$. (On admet que Θ est compact.) Nous présentons ici très sommairement l'argument utilisé. Considérons les ellipsoïdes centrés sur une grille discrète de points $\theta_i \in \Theta$:

$$(\theta - \theta_i)^T I(\theta_i) (\theta - \theta_i) \leq d(n)$$

où la grille $\{\theta_i\}$ est telle que les rectangles maximales contenus dans ces ellipsoïdes déterminent un partitionnement non-uniforme de l'espace de paramètres Θ . Considérons maintenant que $d(n)$ dépend du nombre de données n de telle façon que le volume de ces ellipsoïdes peut être approximé par

$$V_i(n) = |\det(I(\theta_i))|^{-1/2} \left(\frac{2\pi}{n} \right)^{k/2},$$

de façon que si θ_i et θ_j sont des éléments voisins dans cette grille, alors, la probabilité d'erreur

$$\Pr\{\hat{\theta}(x^{(n)}) = \theta_i | \theta_j\} \rightarrow_{n \rightarrow \infty} 0, \quad i \neq j.$$

Le taux de diminution du volume $V_i(n)$ des ellipsoïdes (et donc d'augmentation du nombre de points dans la grille discrète θ_i) est critique, dans le sens que pour une variation de $d(n)$ qui induirait une diminution plus rapide de $V_i(n)$, la probabilité de confondre deux points voisins de la grille ne convergerait pas vers zéro. L'auteur montre

que le nombre de points dans la grille ainsi construite est égal à la complexité paramétrique $\mathcal{C}(\mathcal{H})$, qu'il désigne par *complexité géométrique*. Nous pouvons donc comprendre le terme correspondant à la complexité paramétrique dans le code universel optimal de Shtarkov (maximum de vraisemblance normalisé) comme la longueur des mots de code pour une distribution uniforme $p(\theta_i) = 1/\mathcal{C}(\mathcal{H})$ dans cet ensemble de modèles discernibles en \mathcal{H} pour les observations $x^{(n)}$. Le terme $-\log p(x^{(n)}|\hat{\theta}(x^{(n)}))$, code les données comme une des séquences qui correspondent à l'élément de la partition associée au modèle choisit, c'est à dire, les "détails" des données.

Remarque 1

L'analyse que nous venons de présenter montre que pour des modèles *finis* les codes en deux parties avec une distribution a priori uniforme sont asymptotiquement optimales. Comme nous l'avons vu (exemple 4) la pénalité de ces codes est $\log M$, et nous venons de voir que pour le code universel optimal la pénalité prend asymptotiquement la valeur de $\log M$. Cependant, pour des valeurs de n petits, et pour certaines séquences $x^{(n)}$, la probabilité d'erreur est non-nulle, et donc la complexité paramétrique est inférieure à cette valeur asymptotique. Le code en deux parties considère donc une pénalité qui sera pour certaines séquences, supérieure à celle du code optimal en une seule partie.

△

Pour le cas plus intéressant de modèles qui contiennent un nombre *infini* de distributions, la complexité paramétrique a une interprétation comme un *quotient* de deux volumes. Cette interprétation est basée dans l'expansion asymptotique de la distribution $\bar{P}_{m\nu n}$ présentée dans le lemme suivant.

Lemme 1

Soit \mathcal{H} un modèle paramétrique de dimension m . Sous certaines conditions de régularité

$$\mathcal{C}(\mathcal{H}) = \frac{k}{2} \log \frac{n}{2\pi} + \log \int_{\theta \in \Theta} \sqrt{\det(I(\theta))} d\theta + o(1), \quad (6.12)$$

où n est la longueur de la séquence observée, $I(\theta)$ est la **matrice de Fisher** pour le vecteur de paramètres θ , d'élément générique

$$[I(\theta)]_{ij} = \lim_{n \rightarrow \infty} \frac{1}{n} E_{\theta} \left\{ -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p(x^{(n)}|\theta) \right\},$$

et $o(1) \rightarrow 0$ quand $n \rightarrow \infty$.

△

Le dernier terme en (6.12) ne dépend pas de n et donc sa contribution devient négligeable pour des valeurs de n grandes, impliquant que pour n très grand,

$$\mathcal{C}(\mathcal{H}) \simeq \frac{k}{2} \log \frac{n}{2\pi}, \quad n \gg 1. \quad (6.13)$$

Le membre droit de cette expression coïncide (à part un facteur $k/2 \log 2\pi$) avec le terme de pénalité utilisé par la méthode **BIC (Bayesian Information Criterion)** pour le

problème d'identification de modèles. Le fait que MDL et BIC, pour des valeurs suffisamment grandes de n , conduisent à la même pénalité a conduit à l'affirmation (erronée) que MDL et BIC sont équivalents.

En particulier, le Lemme 1 est vrai si \mathcal{H} est une famille exponentielle :

$$p(x^{(n)}|\theta) = \exp(\beta^T t(x^{(n)})) f(x^{(n)}) g(\beta).$$

Cette famille contient un grand nombre de densités usuelles, comme les Bernoulli et multinomiale Gaussienne, Gamma, etc.

Le premier terme de (6.12) mesure directement le nombre de degrés de liberté du modèle, k . Le deuxième terme est une correction qui dépend de la forme fonctionnelle du modèle. C'est une correction qui ne dépend pas de n , et qui donc peut être négligée pour des valeurs de n très grandes.

Exemple 7 Complexité du modèle de Bernoulli

Pour le modèle de Bernoulli introduit dans la page 103 la dimension est $k = 1$. La matrice de Fisher pour ce modèle est

$$I(\theta) = \frac{1}{\theta(1-\theta)}.$$

L'utilisation de ce résultat dans l'expansion asymptotique (6.12) conduit à

$$\mathcal{C}(\mathcal{B}) = \frac{1}{2} \log n + \frac{1}{2} \log \frac{\pi}{2} - 3 + o(1).$$

△

6.6.4 Interprétation Bayésienne

Les tests de décision Bayésiens ont un lien fort avec la méthode MDL. Considérons le [test entre deux hypothèses composées](#) suivant

$$H_1 : x^{(n)} \sim p_1(x) \in \mathcal{H}_\infty = \left\{ p(x^{(n)}|\theta^1), \theta^1 \in \Theta^1 \right\}, \theta^1 \sim w^1(\theta^1)$$

$$H_2 : x^{(n)} \sim p_2(x) \in \mathcal{H}_\infty = \left\{ p(x^{(n)}|\theta^2), \theta^2 \in \Theta^2 \right\}, \theta^2 \sim w^2(\theta^2)$$

où les hypothèses sont équiprobables : $\Pr\{\mathcal{H}_1\} = \Pr\{\mathcal{H}_2\}$. Le test optimal au sens de Bayes (pour le critère de probabilité d'erreur minimale) choisit l'hypothèse i^* qui a la *plus grande probabilité a posteriori* :

$$i^* = \arg \max_{i \in \{1,2\}} \bar{p}(x^{(n)}|\mathcal{H}_i), \quad \bar{p}(x^{(n)}|\mathcal{H}_i) = \int_{\Theta^i} p(x^{(n)}|\theta^i) w^i(\theta^i) d\theta^i. \quad (6.14)$$

Quand les \mathcal{H}_i sont des **familles exponentielles**, et sous conditions de régularité, une expansion de Laplace (développement en série de la fonction intégrée dans l'expression

précédante) conduit à l'approximation suivante

$$\begin{aligned} -\log \bar{p}(x^{(n)}|\mathcal{H}_i) &= -\log p(x^{(n)}|\hat{\theta}^i(x^{(n)})) + \frac{k}{2} \log \frac{n}{2\pi} - \log w(\hat{\theta}^i(x^{(n)})) \\ &+ \log \sqrt{\det(I(\theta^i(x^{(n)})))} + o(1). \end{aligned}$$

Si nous comparons cette expression à l'équation (6.12), nous pouvons constater que la longueur de code atteinte par ce code "de Bayes" *diffère par une constante* de la longueur optimale $-\log \bar{P}_{mnn}$. Pour des grandes valeurs de n , les deux approches conduiront donc au choix du même modèle.

Si nous considérons le cas particulier où les distributions *a priori* sont les **distributions de Jeffrey** ("least informative prior", introduite par Jeffrey en 1946) :

$$w(\theta) = \frac{\sqrt{\det(I(\theta))}}{\int_{\theta' \in \Theta} \sqrt{\det(I(\theta'))} d\theta'}, \quad (6.15)$$

nous pouvons facilement constater que (6.15) coïncide exactement avec (6.12) : **pour des familles exponentielles, pour n grand, l'approche Bayésienne avec une distribution *a priori* de Jeffrey est équivalente au principe du MDL.**

Quand les modèles \mathcal{H}_i n'appartiennent pas à la famille exponentielle, l'expression suivante est valable sous des conditions de régularité :

$$\begin{aligned} -\log \bar{p}(x^{(n)}|\mathcal{H}_i) &= -\log p(x^{(n)}|\hat{\theta}^i(x^{(n)})) + \frac{k}{2} \log \frac{n}{2\pi} \\ &- \log w(|\hat{\theta}^i(x^{(n)})|) + \log \sqrt{\det(\hat{I}(x^{(n)}))} + o(1), \end{aligned} \quad (6.16)$$

où $\hat{I}(x^{(n)})$ est l'information *observée*.

Pour les familles exponentielles, l'information observée coïncide avec l'information de Fisher, et nous sommes conduits à l'expression précédente. En dehors des modèles exponentiels, si la **vraie** distribution des données appartient à un des modèles \mathcal{H}_i , l'information observée converge encore vers l'information de Fisher, et l'approche Bayésienne est encore asymptotiquement optimale.

6.6.5 Interprétation prédictive

Soit p une distribution en \mathcal{X}^n . Alors

$$p(x^{(n)}) = \prod_{i=1}^n \frac{p(x_1^i)}{p(x_1^{i-1})} = \prod_{i=1}^n p(x_i|x_1^{i-1}), \quad (6.17)$$

et donc

$$-\log p(x^{(n)}) = -\sum_{i=1}^n \log p(x_i|x_1^{i-1}). \quad (6.18)$$

Nous pouvons interpréter le terme $-\log p(x_i|x_1^{i-1})$ comme la pénalité associée à l'observation $X_i = x_i$ quand nous essayons de prédire la valeur de X_i avec la distribution

$p(\cdot \cdot \cdot | x_1^{i-1})$ construite avec les observations précédentes. Cette pénalité sera d'autant plus petite que la valeur observée aura une probabilité élevée pour ce modèle conditionnel. L'expression (6.18) nous dit donc que la longueur des mots du code associé à la distribution p est la somme des pénalités pour la prédiction de chaque valeur X_i de la séquence basée sur toutes les valeurs précédentes x_1^{i-1} (valeurs observées).

L'équation (6.17) établit une relation entre les *modèles de probabilité* pour des séquences $x^{(n)}$ et des *stratégies de prédiction*, qui associent à chaque possible séquence passée x_1^{i-1} une loi de probabilité pour la valeur future X_i . De la même façon, cette équation nous permet d'associer une loi de probabilité définie en \mathcal{X}^n à des stratégies de prédiction.

Soit maintenant \mathcal{H} un modèle paramétrique, et \bar{P} un *code universel* par rapport au modèle \mathcal{H} . Dans le cadre de l'estimation Bayésienne, et pour des observation i.i.d., il est bien connu que la distribution prédictive converge (quand $n \rightarrow \infty$), vers la distribution du Maximum de Vraisemblance $p(\cdot | \hat{\theta}_1^{i-1})$. Le même comportement est obtenu pour tous les codes universels, de façon que nous pouvons approximer les distributions conditionnelles $\bar{P}(\cdot | x_1^{i-1})$ par $p(\cdot | \hat{\theta}_1^{i-1})$:

$$\bar{P}(X_i | x_1^{i-1} - 1) \simeq p(\cdot | \hat{\theta}_1^{i-1}), \quad (6.19)$$

et donc

$$-\log \bar{P}(x^{(n)}) \simeq -\sum_{i=1}^n \log p(\cdot | \hat{\theta}_1^{i-1}). \quad (6.20)$$

Sous des conditions de régularité, il peut être démontré que le membre droit de cette équation peut être écrit comme

$$-\sum_{i=1}^n \log p(\cdot | \hat{\theta}_1^{i-1}) = -\log p(\cdot | \hat{\theta}(x^{(n)})) + \frac{k}{2} \log n + o(1),$$

qui diffère par une constante de l'expression (6.12), démontrant que cette approche prédictive conduit asymptotiquement au même choix que le principe du MDL. Cette analyse nous fournit une interprétation alternative du MDL comme choisissant le **modèle pour lequel l'erreur de prédiction accumulée est minimale**.

Nous remarquons finalement que l'estimée $\hat{\theta}$ en (6.20) peut ne pas être l'estimée du Maximum de Vraisemblance, pouvant être remplacée par un autre estimateur qui converge vers l'estimateur MV. Dans certains cas, comme le montre l'exemple suivant, le comportement peut même être supérieur.

Exemple 8 *Modèle de Bernoulli*

Considérons le modèle de Bernoulli, décrit dans la page 103, où $\theta \in [0, 1]$ est la probabilité d'observer un 1. Admettons que

$$x^{(n)} = 0 \ 0 \ 1 \ \dots,$$

et donc l'estimateur du Maximum de Vraisemblance de θ est

$$\hat{\theta}_{MV}(x_1^2) = 0,$$

impliquant que

$$p(x_3|\hat{\theta}_{MV}(x_1^2)) = 0 \Rightarrow -\log p(x_3|\hat{\theta}_{MV}(x_1^2)) = -\log p(1|\hat{\theta}_{MV}(x_1^2)) = \infty,$$

démontrant que ce modèle n'est pas universel. Cependant, la distribution prédictive basée sur l'estimateur modifié suivant (proposé par Laplace)

$$\hat{\theta}_L(x_1^2) = \frac{n_1 + \lambda}{n + 2\lambda},$$

conduit à un code universel. En effet, il peut être démontré que pour $\lambda = 1$, la distribution prédictive avec cet estimateur coïncide avec la distribution marginale de Bayes relativement à une distribution *a priori* uniforme pour θ . Pour $\lambda = 2$, nous obtenons l'estimateur Baysien pour la distribution *a priori* de Jeffrey pour le modèle de Bernoulli, qui atteint (asymptotiquement) le même comportement que le code universel optimal (maximum de vraisemblance normalisé).

Cette approche a des ressemblances avec la méthode du Maximum de Vraisemblance Généralisé, décrite dans la section 6.6.1. La différence fondamentale entre les deux approches, et qui explique que l'approche prédictive qui vient d'être décrite ne souffre pas des problèmes de *overfitting* associés aux tests de Vraisemblance Généralisée, est le fait qu'ici chaque observation est prédite (codée) avec le modèle déterminé par les observations **passées**, ce qui n'est pas le cas pour le test Maximum de Vraisemblance Généralisé. Celui-ci utilise dans un premier pas *toutes les observations* pour estimer le paramètre θ , et utilise ensuite cette estimée pour coder $x^{(n)}$.

6.7 MDL Général pour la sélection de modèles paramétriques

Le matériel présenté dans le Chapitre précédent considère le cas où le nombre de modèles est fini, et la solution avancée (la distribution de Shtarkov) requiert que la complexité paramétrique des modèles soit finie. Dans ce Chapitre nous présentons une version plus générale du MDL qui peut traiter des modèles avec une complexité infinie.

6.7.1 Complexité paramétrique infinie

La complexité paramétrique des modèles les plus communs est infinie. Un exemple important est celui des modèles Gaussiens, comme le montre l'exemple suivant.

Exemple 9 Complexité du modèle Gaussien

Soit \mathcal{H} la famille de distributions Gaussiennes avec variance σ^2 fixée :

$$\mathcal{H}_{\sigma^2} = \left\{ p_{\sigma^2}(x|\mu) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} : \mu \in \mathfrak{R} \right\},$$

étendue à $x^{(n)} \in \mathfrak{R}^n$ par hypothèse d'indépendance (produit des densités marginales). L'estimateur du Maximum de Vraisemblance de la moyenne μ est

$$\hat{\mu}(x_1^i) = \frac{1}{i} \sum_{j=1}^i x_j.$$

Alors,

$$\mathcal{C}(\mathcal{H}) = \log \int_{x^{(n)} \in \mathfrak{R}^+} p_{\sigma^2}(x^{(n)} | \hat{\mu}(x^{(n)})) dx^{(n)} = \infty,$$

et donc le code universel optimal n'est pas défini pour ce modèle.

La matrice de Fisher pour ce modèle est

$$I(\mu) = \frac{1}{\sigma^2},$$

et donc nous avons également

$$\int_{\mathfrak{R}} \sqrt{\det(I(\mu))} d\mu = \infty,$$

et donc le mélange Bayésien pour la distribution *a priori* de Jeffrey n'est pas défini non plus.

Cependant, si nous considérons que $\mu \in [a, b]$, avec $a, b < \infty$, la complexité pour ce modèle "limité" est finie :

$$\int_{x^{(n)} : \hat{\mu}(x^{(n)}) \in [a, b]} p_{\sigma^2}(x^{(n)} | \hat{\mu}(x^{(n)})) dx^{(n)} = \frac{b-a}{\sqrt{2\pi}\sigma} \sqrt{n}.$$

Soient alors les modèles emboîtés suivants :

$$\mathcal{H}_K = \left\{ p_{\sigma^2, K}(x | \mu) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, |\mu| \leq K \right\}, K \geq 0, \quad (6.21)$$

de façon que

$$\mathcal{H} = \cup_{K \geq 0} \mathcal{H}_K.$$

Maintenant, pour chaque K , $\mathcal{C}(\mathcal{H}_K) < \infty$, et donc les codes universels optimaux correspondants, $\bar{P}_{mvn}(x^{(n)} | \mathcal{H}_K)$ existent.

Nous pouvons maintenant coder les données avec un *code en deux parties*, $\bar{P}_{meta}(x^{(n)} | \mathcal{H})$, qui code dans un premier temps la constante K , et utilise ensuite le code optimal correspondant pour coder les données. Ce code a une longueur

$$-\log \bar{P}_{meta}(x^{(n)} | \mathcal{H}) = \min_K \left\{ -\log \bar{P}_{mvn}(x^{(n)} | \mathcal{H}_K) + L(K) \right\}, \quad (6.22)$$

où $L(K)$ est la longueur du mot de code utilisé pour K .

Cette approche, basée sur un codage en deux parties, est sous-optimale. La cause de la non-optimalité est liée au fait que le code réserve plusieurs mots de code pour

la même séquence, une pour chaque valeur possible de K . Une alternative à cette approche, qui est basée dans la restriction de l'espace d'observations, consiste à limiter l'espace des paramètres. Dans [7] Rissanen propose une approche alternative, basée sur l'utilisation d'une version re-normalisée du code universel optimal :

$$\bar{P}_{rmvn}(x^{(n)})|\mathcal{H} = \frac{\bar{P}_{mvm}(x^{(n)})|\mathcal{H}_{|\hat{\mu}(x^{(n)})}}}{\int_{y^{(n)} \in \mathcal{X}^n} \bar{P}_{mvm}(y^{(n)})|\mathcal{H}_{|\hat{\mu}(x^{(n)})} dy^{(n)}}$$

6.7.2 Sommaire

Le matériel présenté indique que si nous souhaitons appliquer le principe du MDL pour choisir entre différents modèles, nous devons chercher à définir un **modèle universel** pour l'ensemble de modèles, qui soit capable de coder toutes les séquences d'une taille donnée n .

Si l'ensemble de modèles est **fini**, nous utilisons une distributions *a priori* **uniforme** pour les modèles (longueur de code constante). Dans le cas contraire, la distribution uniforme n'existe plus, et nous sommes forcés à donner une préférence à certains modèles sur les autres.

Quand la complexité paramétrique des modèles est **infinie**, et donc le modèle universel optimal de Shtarkov n'existe pas, nous devons décomposer \mathcal{H} comme l'union de sous-modèles emboîtés \mathcal{H}_k , de complexité croissante en k . Un code universel pour l'ensemble de modèles est ensuite construit, avec une pénalité qui est proche de la pénalité associée au sous-modèle de complexité minimale qui contient l'estimée du Maximum de Vraisemblance.

Bibliographie

- [1] Vijay Balasubramanian, "A geometric framework for Occam's razor for inference of parametric distributions," Princeton Physics Preprint PUPT-1588, Princeton, NJ, USA, 1996. (<http://arxiv.org/pdf/adap-org/9601001>).
- [2] J. Rissanen, "Modelling by shortest data description," *Automatica*, 14 :465 :471, 1978.
- [3] A. Kolmogorov, "Three approaches to the quantitative definition of information," *Problems Inform. Transmission*, (1), 1 :7, 1965.
- [4] Ray Solomonoff, "A Formal Theory of Inductive Inference, Part I (II)," *Information and Control*, Part I : Vol 7, No. 1(2), pp. 1 :22 (224 :254), March(June) 1964. (<http://world.std.com/rjs/1964pt1.pdf>/<http://world.std.com/rjs/1964pt2.pdf>)
- [5] Gregory Chaitin, "On the length of programs for computing finite binary sequences," *Journal of the ACM* 13 (1966), pp. 547-569. (<http://www.cs.auckland.ac.nz/CDMTCS/chaitin/acm66.pdf>)
- [6] Paul Vitanyi, Ming Li, "Minimum description Length Induction, Bayesianism and Kolmogorov Complexity," *IEEE Trans. Inf. the.*, 46 :2, 446 :464, 2000. (<http://www.cwi.nl/paulv/papers/mdlindbayeskolmcompl.pdf>)
- [7] J. Rissanen, "Strong Optimality of the Normalized ML Models as Universal Codes and Information in Data," *IEEE Trans. Information Theory*, vol. 47(5), pp. 1712-1717, 2001. (<http://www.mdl-research.org/pub/bound2.ps>)