

## Chapitre 9

# Approche bayésienne

L'approche bayésienne est une approche cohérente pour la résolution des problèmes inverses car elle permet de prendre en compte et de traiter de la même manière l'information *a priori* sur la solution et celle sur les données.

En effet, cette approche permet non seulement de prendre en compte le caractère incertain des erreurs au travers de l'attribution d'une loi de probabilité aux mesures, mais aussi l'information *a priori* au travers d'une loi de probabilité *a priori* que l'on attribue aux inconnues du problème.

La fusion de ces deux sources d'information s'effectue par la règle de Bayes. On peut alors dire que cette approche généralise la régularisation déterministe.

## 9.1 Introduction

L'approche bayésienne correspond formellement au cheminement suivant :

1. On explicite un ensemble d'hypothèses  $\mathcal{H}$  sur le problème concernant le modèle d'observation, les connaissances *a priori* sur les inconnues  $\mathbf{f}$  et sur le bruit  $\mathbf{b}$ .
2. On attribue une loi de probabilité *a priori*  $p(\mathbf{f}|\boldsymbol{\theta}_1, \mathcal{H})$  aux inconnues du problème pour traduire la connaissance initiale sur ces inconnues. Cette loi peut dépendre d'un certain nombre de paramètres  $\boldsymbol{\theta}_1$ .
3. On attribue une loi de probabilité  $p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_2; \mathcal{H})$  aux grandeurs mesurées pour traduire l'incertitude (due au bruit, aux erreurs de discrétisation et de quantification, au manque de précision de l'appareil de mesure, etc.) sur ces données. Cette loi aussi peut dépendre d'un certain nombre de paramètres  $\boldsymbol{\theta}_2$ . L'ensemble des paramètres  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  sont appelés les hyperparamètres du problème.
4. On utilise la règle de Bayes pour combiner l'information contenue dans les données et celle contenue dans la loi *a priori* et calculer ainsi la loi de probabilité *a posteriori*

$$p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) = \frac{p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_2; \mathcal{H}) p(\mathbf{f}|\boldsymbol{\theta}_1, \mathcal{H})}{p(\mathbf{g}|\boldsymbol{\theta}; \mathcal{H})}$$

où

$$p(\mathbf{g}|\boldsymbol{\theta}; \mathcal{H}) = \int p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_2; \mathcal{H}) p(\mathbf{f}|\boldsymbol{\theta}_1, \mathcal{H}) d\mathbf{f}.$$

Cette loi contient toute l'information disponible sur les inconnues  $\mathbf{f}$  après cette fusion.

5. On peut alors calculer n'importe quelle quantité. Par exemple, la moyenne *a posteriori* de  $\mathbf{f}$ :

$$E[\mathbf{f}] = \int \mathbf{f} p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) d\mathbf{f},$$

ou la moyenne de n'importe quelle autre fonction  $h(\mathbf{f})$ :

$$E[h(\mathbf{f})] = \int h(\mathbf{f}) p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) d\mathbf{f},$$

ou encore calculer la probabilité pour que  $\underline{\mathbf{f}} < \mathbf{f} \leq \overline{\mathbf{f}}$ :

$$P(\underline{\mathbf{f}} < \mathbf{f} \leq \overline{\mathbf{f}}) = \int_{\underline{\mathbf{f}}}^{\overline{\mathbf{f}}} p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) d\mathbf{f}$$

ou seulement s'intéresser à une des composantes  $f_j$  de  $\mathbf{f}$  et calculer

$$P(\underline{f}_j < f_j \leq \overline{f}_j) = \int_{\underline{f}_j}^{\overline{f}_j} p(f_j|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) df_j,$$

où

$$p(f_j|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) = \int \cdots \int p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) df_1 \cdots df_{j-1} df_{j+1} \cdots df_n$$

est la loi *a posteriori marginale*.

Lorsque  $\mathbf{f}$  est un scalaire ou un vecteur avec seulement deux composantes, on peut même tracer  $p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H})$  ou  $p(f_1, f_2|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H})$ , ce qui permettra d'avoir une idée plus précise sur la forme de ces distribution.

Mais en pratique, la manipulation d'une densité de probabilité n'est pas très comode dès lors qu'on a plus de deux variables. C'est pourquoi, on se contente souvent de fournir une description sommaire de la loi *a posteriori* en indiquant

– son mode :

$$M[\mathbf{f}] = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H})\},$$

– sa moyenne :

$$E[\mathbf{f}] = \int \mathbf{f} p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) d\mathbf{f},$$

– les modes de ses marginales :

$$M[f_j] = \arg \max_{f_j} \{p(f_j|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H})\},$$

– les  $\alpha$ -quantiles de ses marginales :

$$Q_\alpha(f_j) : P(f_j \leq Q_\alpha(f_j)) = \alpha,$$

– les médianes de ses marginales :

$$\text{Méd}[f_j] = Q_{\frac{1}{2}}(f_j) : P(f_j \leq \text{Méd}[f_j]) = \frac{1}{2},$$

– les régions de plus haute probabilité ou les supports  $\alpha$ -interquantiles :

$$[a, b] = [Q_{(1-\alpha)/2}(f_j), Q_{(1+\alpha)/2}(f_j)] : P(a \leq f_j \leq b) = 1 - \alpha,$$

– les variances :

$$\text{Var}[f_j] = \int (f_j - \bar{f}_j)^2 p(f_j|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}),$$

où

$$\bar{f}_j = E[f_j] = \int f_j p(f_j|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}),$$

– les covariances :

$$\text{Cov}[f_j, f_k] = \iint (f_j - \bar{f}_j)(f_k - \bar{f}_k) p(f_j, f_k|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) df_j df_k.$$

6. Une autre approche fondée sur la théorie de la décision consiste en la définition d'un estimateur  $\hat{\mathbf{f}}$  qui minimiserait la moyenne d'une fonction de coût  $C(\hat{\mathbf{f}}, \mathbf{f})$ . Nous reviendront sur cette approche et les choix de ces fonctions coût, ainsi que sur les liens qu'existe entre ces estimateurs et les caractéristiques de la loi *a posteriori*.
7. Finalement, une fois une solution choisie, il est impératif d'y attribuer un degré de confiance ainsi que de déterminer la sensibilité de cette solution vis-à-vis des erreurs de mesures et celles de la modélisation. Heureusement, en principe, cette approche bayésienne nous fourni tout ce qui est nécessaire pour pouvoir satisfaire cette exigence. Car disposant de la loi *a posteriori* on peut calculer, par exemple, la covariance de l'erreur de l'estimation, ce qui nous permet de mettre des barres d'erreurs sur les solutions estimées.

Voyons maintenant quelles sont les difficultés de l'utilisation de cette approche dans un problème réel et quelles sont les solutions existantes pour les résoudre. On peut classer ces difficultés dans les rubriques qui suivent :

1. choix de la forme de la loi *a priori*  $p(\mathbf{f}|\boldsymbol{\theta}_1; \mathcal{H})$  ;
2. choix ou estimation de ses paramètres  $\boldsymbol{\theta}_1$  ;
3. choix de la forme de la loi  $p(\mathbf{g}|\boldsymbol{\beta}, \boldsymbol{\theta}_2; \mathcal{H})$  ;
4. choix ou estimation de ses paramètres  $\boldsymbol{\theta}_2$  ;
5. calcul effectif de la loi *a posteriori*  $p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H})$  ;
6. choix d'une fonction de coût  $C(\hat{\mathbf{f}}, \mathbf{f})$  et calcul effectif de l'estimateur qu'en découle ;
7. détermination ou estimation *a posteriori* des hyperparamètres du problème, *i.e.* les paramètres  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  des lois de probabilités directes  $p(\mathbf{f}|\boldsymbol{\theta}_1; \mathcal{H})$  et  $p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_2; \mathcal{H})$ , à partir des données, c'est-à-dire en même temps que la solution ;

L'objectif principal de ce chapitre est de faire état de l'art concernant les solutions proposées pour résoudre ces problèmes et de fournir quelques exemples concrets (étude de cas) pour illustrer les difficultés et les performances des solutions proposées.

Mais, avant d'aller plus loin prenons un exemple, désormais classique, qui est le cas linéaire gaussien qui nous permettra de comprendre l'essentiel de ces différentes étapes.

## 9.2 Cas linéaire gaussien

S'il y a un cas où les calculs peuvent se faire facilement jusqu'au bout et où l'on peut facilement comprendre l'approche bayésienne est le cas où le modèle reliant les inconnues  $\mathbf{f}$  et les données  $\mathbf{g}$  est linéaire et où on peut attribuer des lois de probabilités gaussiennes à  $\mathbf{f}$  et aux mesures  $\mathbf{g}$ .

Nous allons suivre les étapes présentées en introduction pour la résolution du problème linéaire

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \mathbf{b} \quad (9.1)$$

– Information *a priori* sur  $\mathbf{f}$  :

Faisons l'hypothèse que nous puissions *a priori* connaître seulement la moyenne  $E[\mathbf{f}] = \mathbf{f}_0$  et la matrice de covariance  $E[(\mathbf{f} - \mathbf{f}_0)(\mathbf{f} - \mathbf{f}_0)^t] = \mathbf{R}_f = \sigma_f^2 \mathbf{P}_0$  de  $\mathbf{f}$ . Nous pouvons alors faire l'hypothèse que  $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}_0, \sigma_f^2 \mathbf{R}_f)$  :

$$p(\mathbf{f}) = A \exp \left[ \frac{-1}{2} (\mathbf{f} - \mathbf{f}_0)^t \mathbf{R}_f^{-1} (\mathbf{f} - \mathbf{f}_0) \right] = A \exp \left[ \frac{-1}{2\sigma_f^2} (\mathbf{f} - \mathbf{f}_0)^t \mathbf{P}_0^{-1} (\mathbf{f} - \mathbf{f}_0) \right] \quad (9.2)$$

avec  $A = (2\pi)^{-n/2} |\mathbf{R}_f|^{-1/2} = (2\pi\sigma_f^2)^{-n/2} |\mathbf{P}_0|^{-1/2}$ .

D'ailleurs, cette hypothèse peut être validée par le principe du maximum d'entropie.

– Information *a priori* sur  $\mathbf{b}$  :

Faisons l'hypothèse que  $E[\mathbf{b}] = \mathbf{0}$  et que sa matrice de covariance  $E[\mathbf{b}\mathbf{b}^t] = \mathbf{R}_b = \sigma_b^2 \mathbf{I}$ . Cette hypothèse est bien raisonnable. En effet, supposer que  $E[\mathbf{b}] = \mathbf{0}$  signifie que nous faisons l'hypothèse qu'il n'y a pas d'erreur systématique et l'hypothèse  $E[\mathbf{b}\mathbf{b}^t] = \sigma_b^2 \mathbf{I}$  signifie que le bruit est supposé blanc (non corrélé).

Avec les mêmes arguments que dans le cas précédent on fait alors l'hypothèse que  $p(\mathbf{b}) = \mathcal{N}(\mathbf{0}, \sigma_b^2 \mathbf{R}_b)$ . Maintenant, partant du modèle (9.1) et faisant l'hypothèse que le bruit  $\mathbf{b}$  est aussi indépendant du  $\mathbf{f}$  on peut facilement déduire

$$p(\mathbf{g}|\mathbf{f}) = B \exp \left[ \frac{-1}{2} (\mathbf{g} - \mathbf{H}\mathbf{f})^t \mathbf{R}_b^{-1} (\mathbf{g} - \mathbf{H}\mathbf{f}) \right] = B \exp \left[ \frac{-1}{2\sigma_b^2} (\mathbf{g} - \mathbf{H}\mathbf{f})^t (\mathbf{g} - \mathbf{H}\mathbf{f}) \right] \quad (9.3)$$

avec  $B = (2\pi)^{-m/2} |\mathbf{R}_b|^{-1/2} = (2\pi\sigma_b^2)^{-m/2}$ .

– Règle de Bayes :

Utilisant la règle de Bayes, il est facile de voir que

$$p(\mathbf{f}|\mathbf{g}) = C \exp \left[ \frac{-1}{2} J(\mathbf{f}) \right], \quad (9.4)$$

avec

$$\begin{aligned} J(\mathbf{f}) &= (\mathbf{g} - \mathbf{H}\mathbf{f})^t \mathbf{R}_b^{-1} (\mathbf{g} - \mathbf{H}\mathbf{f}) + (\mathbf{f} - \mathbf{f}_0)^t \mathbf{R}_f^{-1} (\mathbf{f} - \mathbf{f}_0) \\ &= \frac{1}{\sigma_b^2} \left[ (\mathbf{g} - \mathbf{H}\mathbf{f})^t (\mathbf{g} - \mathbf{H}\mathbf{f}) + \lambda (\mathbf{f} - \mathbf{f}_0)^t \mathbf{P}_0^{-1} (\mathbf{f} - \mathbf{f}_0) \right], \end{aligned}$$

où  $\lambda = \frac{\sigma_b^2}{\sigma_f^2}$ .

Avec un peu de calcul il n'est pas difficile de montrer que cette loi *a posteriori* est aussi gaussienne et on a

$$p(\mathbf{f}|\mathbf{g}) = \mathcal{N}(\hat{\mathbf{f}}, \hat{\mathbf{P}})$$

où

$$\begin{cases} \hat{\mathbf{f}} &= \mathbf{R}_f \mathbf{H}^t (\mathbf{H} \mathbf{R}_f \mathbf{H}^t + \mathbf{R}_b)^{-1} \mathbf{g} = \hat{\mathbf{P}} \mathbf{H}^t \mathbf{R}_b^{-1} \mathbf{g}, \\ \hat{\mathbf{P}} &= \mathbf{R}_f - \mathbf{R}_f \mathbf{H}^t (\mathbf{H} \mathbf{R}_f \mathbf{H}^t + \mathbf{R}_b)^{-1} \mathbf{H} \mathbf{R}_f = (\mathbf{R}_f^{-1} + \mathbf{H}^t \mathbf{R}_b^{-1} \mathbf{H})^{-1}, \end{cases}$$

ou encore

$$\begin{cases} \hat{\mathbf{f}} &= (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{P}_0^{-1})^{-1} \mathbf{H}^t \mathbf{g}, \\ \hat{\mathbf{P}} &= \sigma_b^2 (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{P}_0^{-1})^{-1}. \end{cases}$$

– Caractéristique de la solution :

Nous avons ainsi une expression analytique pour la loi *a posteriori* et nous pouvons facilement en déduire tout ce que l'on souhaite. Par exemple nous savons que la moyenne, la médiane et la mode sont identiques et égaux à  $\hat{\mathbf{f}}$ .

Ainsi tous ces estimateurs : le maximum *a posteriori* (MAP) ou la moyenne *a posteriori* (MP) sont identiques et égaux à  $\hat{\mathbf{f}}$  qui peut aussi être calculé par

$$\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g})\} = \arg \min_{\mathbf{f}} \{J(\mathbf{f}) = Q(\mathbf{f}) + \lambda \Omega(\mathbf{f})\}$$

avec

$$Q(\mathbf{f}) = (\mathbf{g} - \mathbf{H}\mathbf{f})^t (\mathbf{g} - \mathbf{H}\mathbf{f})$$

et

$$\Omega(\mathbf{f}) = (\mathbf{f} - \mathbf{f}_0)^t \mathbf{P}_0^{-1} (\mathbf{f} - \mathbf{f}_0),$$

ou encore, si on note par  $\mathbf{P}_0^{-1} = \mathbf{D}^t \mathbf{D}$  on a

$$\begin{aligned} J(\mathbf{f}) &= (\mathbf{g} - \mathbf{H}\mathbf{f})^t (\mathbf{g} - \mathbf{H}\mathbf{f}) + \lambda (\mathbf{f} - \mathbf{f}_0)^t \mathbf{D}^t \mathbf{D} (\mathbf{f} - \mathbf{f}_0) \\ &= \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \lambda \|\mathbf{D}(\mathbf{f} - \mathbf{f}_0)\|^2. \end{aligned}$$

On retrouve alors la notion de solution régularisée et la régularisation quadratique avec les avantages suivantes :

- Dans l'approche régularisation déterministe le choix des distances  $\Delta_1(\mathbf{g}, \mathbf{H}\mathbf{f}) = \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2$  et de  $\Delta_2(\mathbf{f}, \mathbf{f}_0) = \|\mathbf{D}(\mathbf{f} - \mathbf{f}_0)\|^2$  est plutôt descriptif et un peu arbitraire, alors que dans l'approche bayésienne, ces termes sont, respectivement, les conséquences des hypothèses faites sur la loi du bruit et sur la loi *a priori*.
- Dans régularisation déterministe, le coefficient de régularisation  $\lambda$  est un paramètre qui est choisi empiriquement, alors qu'ici,  $\lambda = \frac{\sigma_b^2}{\sigma_f^2}$ . Si l'on connaît les deux variances  $\sigma_b^2$  et  $\sigma_f^2$ , sa valeur est fixée. Mais, bien sûr, en pratique on ne connaît pas ces deux valeurs et on est encore amené à les fixer par expérience. Cependant, rien que l'expression de  $\lambda$  nous permet de mieux comprendre son comportement : plus le niveau du bruit est important, plus il faut choisir sa valeur grande pour obtenir un résultat satisfaisant. Nous verrons que l'approche bayésienne nous donne les outils nécessaires pour la déterminer (voir le chapitre sur l'estimation des hyperparamètres).

- Dans le cadre de la régularisation déterministe, si on se demande qu'elle est le degré de confiance que l'on peut avoir dans la solution  $\hat{\mathbf{f}}$  proposée, nous n'avons pas d'outil convenable pour répondre. Dans l'approche bayésienne, nous pouvons répondre à cette question, par exemple en calculant la matrice de covariance *a posteriori*  $\mathbf{P} = \mathbb{E}[(\mathbf{f} - \hat{\mathbf{f}})^2]$ . Sachant que les éléments diagonaux  $P_{jj} = \mathbb{E}[(f_j - \hat{f}_j)^2]$  sont les variances *a posteriori*, nous pouvons les utiliser pour mettre des barres d'erreur sur la solution. Les éléments non-diagonale  $P_{ij} = \mathbb{E}[(f_i - \hat{f}_i)(f_j - \hat{f}_j)]$  peuvent aussi nous renseigner sur le lien (corrélation) qui peut exister entre l'élément  $f_i$  et l'élément  $f_j$ .

## 9.3 Calcul de la loi *a posteriori*

### 9.3.1 Cas gaussien

Nous avons vu que dans le cas d'un modèle linéaire gaussien, la loi *a posteriori* est aussi gaussienne. Sachant qu'une loi gaussienne est entièrement définie par ses deux premiers moments, le problème devient assez simple. Nous disposons des expressions analytiques pour la moyenne et la matrice de covariance *a posteriori*, expressions dont les calculs nécessitent soit l'inversion des matrices soit l'optimisation d'un critère quadratique. Il reste cependant le coût de calcul. Souvent, on peut profiter de la structure des matrices à inverser ou de la forme quadratique du critère à optimiser pour fournir des algorithmes spécifiques pour effectuer ces calculs.

### 9.3.2 Cas général

Dans le cas général, le calcul complet de la loi *a posteriori* peut devenir plus délicat, et souvent, on se contente, soit de l'approximer par une loi gaussienne et de calculer la moyenne et la matrice de covariance, soit de définir un estimateur ponctuel à partir de cette loi, comme par exemple le *maximum a posteriori* (MAP) :

$$\hat{\mathbf{f}}_{\text{MAP}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g})\} = \arg \min_{\mathbf{f}} \{-\ln p(\mathbf{f}|\mathbf{g})\},$$

ou la *moyenne a posteriori* (en Anglais posterior mean PM) :

$$\hat{\mathbf{f}}_{\text{PM}} = \int \mathbf{f} p(\mathbf{f}|\mathbf{g}) \, d\mathbf{f},$$

ou encore le *maximum a posteriori marginal* (MAPmarginal) :

$$\hat{f}_j = \arg \max_{f_j} \{p(f_j|\mathbf{g})\},$$

où

$$p(f_j|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) = \int \cdots \int p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}; \mathcal{H}) \, df_1 \cdots df_{j-1} \, df_{j+1} \cdots df_n.$$

Notons aussi que le calcul des estimateurs PM et MAPmarginal nécessite le calcul des intégrales de dimension très élevée, alors que le calcul de l'estimateur MAP ne nécessite qu'une optimisation.

Notons aussi que dans le cas gaussien tous ces estimateurs sont identiques. En revanche, dans les autres cas, ils peuvent fournir des résultats très différents. Nous reviendrons sur les propriétés de ces estimateurs et les outils qui existent pour les calculer dans les chapitres suivants.



## 9.4 Choix des fonctions de coût

Dans l'approche théorie de la décision nous avons vu que, partant de la loi *a posteriori*, on peut définir une solution au problème, par l'intermédiaire d'une fonction de coût  $C(\hat{\mathbf{f}}, \mathbf{f})$ . En effet, dans cette approche on définit l'estimateur optimal  $\hat{\mathbf{f}}$  comme l'argument qui minimise le coût moyen :

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{z}} \{E[C(\mathbf{z}, \mathbf{f})]\} = \arg \min_{\mathbf{z}} \left\{ \int C(\mathbf{z}, \mathbf{f}) p(\mathbf{f}|\mathbf{g}) d\mathbf{f} \right\}.$$

Nous allons voir que suivant le choix de la fonction de coût on obtient des estimateurs différents, et en particulier, les estimateurs MAP, PM et MAPmarginal :

- Coût quadratique et estimateur PM :

$$C(\mathbf{f}, \hat{\mathbf{f}}) = (\mathbf{f} - \hat{\mathbf{f}})^t \mathbf{Q} (\mathbf{f} - \hat{\mathbf{f}})^t \longrightarrow \hat{\mathbf{f}} = \int \mathbf{f} p(\mathbf{f}|\mathbf{g}) d\mathbf{f}.$$

- Coût dirac et estimateur MAP :

$$C(\mathbf{f}, \hat{\mathbf{f}}) = 1 - \delta(\mathbf{f} - \hat{\mathbf{f}}) \longrightarrow \hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g})\}.$$

- Coût produit de diracs et estimateur MAPmarginal :

$$C(\mathbf{f}, \hat{\mathbf{f}}) = \prod_j 1 - \delta(x_j - \hat{x}_j) \longrightarrow \hat{f}_j = \arg \max_{x_j} \{p(x_j|\mathbf{g})\},$$

D'autres fonctions de coût peuvent aussi être utilisées, mais leur usage reste plus spécifique.

## 9.5 Choix de la loi *a priori*

Le problème de la conversion d'une information *a priori* en une loi de probabilité est un problème difficile et encore largement ouvert.

La principale difficulté réside dans le fait que rarement, en pratique, l'information *a priori* se présente directement sous une forme probabiliste. Par exemple, on nous dit que la grandeur recherchée est *positive ou bornée entre [0, 1]*, de variation *douce, croissante ou décroissante*, à *bande limitée ou à spectre limitée*, de variation *douce ou constante par morceaux ou par région*, de variation *impulsionnelle*, etc. La tâche, ensuite, est de construire des lois de probabilités qui puissent incorporer ces qualificatifs.

Les méthodes existantes peuvent être schématiquement regroupées en trois grandes classes.

Certaines reposent sur la théorie des *groupes de transformation* pour déterminer la mesure de référence "naturelle" pour la grandeur recherchée. On utilise ces méthodes surtout lorsqu'on sait peu de choses sur la grandeur recherchée, c'est à dire lorsque notre information *a priori* se résume à une connaissance qualitative sur la nature de la grandeur recherchée. Par exemple, nous savons que la grandeur recherchée représente un paramètre de *localisation* pouvant prendre une valeur quelconque sur l'axe des réels  $\mathbb{R}$  ou un paramètre *d'échelle* qui est par nature positif et pouvant prendre une valeur quelconque sur la demi-droite des réels positifs  $\mathbb{R}_+$ .

Mais en pratique, cette approche a permis de justifier après coup l'emploi de la mesure de Lebesgue pour les paramètres de localisation (fournissant ainsi une extension au cas continu de la distribution uniforme résultant de l'application du "Principe d'indifférence" de Bernoulli dans le cas discret) et de la mesure de Jeffreys dans le cas des paramètres d'échelle.

D'autres méthodes reposent sur des principes informationnel. Il s'agit principalement des méthodes dites "à maximum d'entropie" dans lesquelles on recherche une distribution qui soit la plus proche (au sens d'une distance de Kullback) d'une distribution de référence (souvent choisie par l'approche précédente) tout en vérifiant une information incomplète connue *a priori* sous la forme des moments (des contraintes linéaires) sur la loi recherchée. Mais là encore, cette approche a surtout permis de justifier après coup certains choix. De plus, elle n'est véritablement praticable que lorsque cette information *a priori* est faite de contraintes linéaires sur la distribution recherchée. On trouve alors la famille des *distributions exponentielles* comme nous l'avons vu dans l'étude du principe du maximum d'entropie.

Il existe enfin une dernière classe très importante, celle des constructions faites "à la main". C'est dans cette catégorie qu'entrent par exemple les modèles markoviens qui ont connu un développement spectaculaire depuis 1984 en traitement d'image, et qui permettent d'incorporer dans une distribution *a priori* des propriétés locales essentielles que doit posséder l'objet. La construction de ces modèles demande beaucoup de savoir-faire.

Bien entendu, il est hors de question de vouloir détailler toutes ces méthodes dans le cadre de ce document. Nous nous limiterons, dans la section suivante, à une description succincte des deux premières approches, plutôt à travers de quelques exemples que d'une manière formelle. Mais, en ce qui concerne la modélisation markovienne, nous lui réservons un chapitre tout entier.

### 9.5.1 Maximum d'entropie

Nous avons vu dans la section (8.3) que le principe du maximum d'entropie peut être utilisé pour attribuer une loi de probabilité à une grandeur traduisant une information *a priori* sous forme de moments. Nous allons illustrer ceci à travers quelques exemples.

Supposons que nous ayons un voltmètre et que nous cherchions à mesurer la tension de la ville  $X(t)$  à différents instants. Supposons que notre instrument soit un voltmètre idéal et ce que nous mesurons  $Y(t)$  peut être modélisé par

$$Y(t) = X(t) + B(t)$$

où  $B(t)$  représente à la fois les erreurs de lecture et le bruit externe. Nous voudrions appliquer l'approche bayésienne pour fournir une estimation  $\hat{X}(t)$  de  $X(t)$ . Considérons tout d'abord que nous avons fait juste une mesure  $y_1$  à partir de laquelle que nous voudrions estimer  $X$ . Nous avons vu que la première étape avant d'appliquer la règle de Bayes est d'attribuer les lois de probabilité  $p(x)$  et  $p(y|x)$ .

Concernant  $X$ , nous savons que sa moyenne est de 220 volts et que sa variation autour de cette moyenne est  $\pm 5$  volts (ceci dépend sûrement de pays et de la qualité des installations, mais supposons que nous soyons en France!). Comment traduire alors cette information en une loi de probabilité? C'est exactement dans ce cadre que le principe du ME peut fournir une réponse satisfaisante:

$$E[X] = x_0 = 220 \text{ volts}, \quad \text{Var}[X] = \sigma_x^2 = 25 \text{ volts}^2 \quad \xrightarrow{\text{ME}} \quad X \simeq \mathcal{N}(x_0, \sigma_x^2)$$

Concernant  $B$ , il est naturel de faire l'hypothèse que sa moyenne est nulle (il n'y a pas d'erreur systématique) et supposons que l'on puisse avoir une idée de sa puissance. En langage d'ingénieur: "le bruit est d'environ 10 dB". Là aussi, le principe du ME peut fournir une réponse satisfaisante:

$$E[B] = 0, \quad \text{Var}[B] = \sigma_b^2 = 1 \text{ volts}^2 \quad \xrightarrow{\text{ME}} \quad B \simeq \mathcal{N}(0, \sigma_b^2).$$

Nous avons maintenant tous les ingrédients pour calculer la loi *a posteriori*:

$$X|Y \simeq \mathcal{N}(\hat{x}, \hat{\sigma}^2),$$

et il n'est pas difficile de calculer  $\hat{x}$  et  $\hat{\sigma}^2$ :

$$\hat{x} = \frac{\sigma_x^2}{\sigma_x^2 + \sigma_b^2} y + \frac{\sigma_b^2}{\sigma_x^2 + \sigma_b^2} x_0$$

$$\hat{\sigma}^2 = \frac{\sigma_x^2 \sigma_b^2}{\sigma_x^2 + \sigma_b^2}.$$

Notons que

$$\begin{aligned} \text{si } \sigma_b^2 = 0 &\longrightarrow \hat{x} = y, \quad \hat{\sigma}^2 = 0 \quad \text{et} \\ \text{si } \sigma_x^2 = 0 &\longrightarrow \hat{x} = x_0, \quad \hat{\sigma}^2 = 0. \end{aligned}$$

Considérons maintenant que la grandeur  $X$  que nous voulons mesurer est l'énergie ou la puissance consommée dans un circuit et l'appareil de mesure est un wattmètre. Mais, avant même de choisir notre instrument nous savons que la grandeur à mesurer est positive et que son ordre de grandeur est d'environ 200 Watts:

$$X \geq 0, \quad E[X] = x_0 = 200 \text{ Watts} \quad \xrightarrow{\text{ME}} \quad X \simeq \mathcal{E}(x_0)$$

ce qui signifie :

$$p(x) = x_0 \exp[-x/x_0], \quad x > 0.$$

Si nous gardons les mêmes hypothèses sur  $B$  alors nous avons

$$p(y|x) \propto \exp\left[-\frac{1}{2\sigma_b^2}(y-x)^2\right]$$

ce qui nous donne pour la loi *a posteriori*:

$$p(x|y) \propto \exp\left[-\frac{1}{2\sigma_b^2}\left((y-x)^2 + \frac{2\sigma_b^2}{x_0}x\right)\right], \quad x > 0.$$

En réarrangeant le terme en exponentiel, on peut facilement voir que cette loi est une loi gaussienne tronquée:

$$p(x|y) \propto \exp\left[-\frac{1}{2\sigma_b^2}(x-\hat{x})^2\right], \quad x > 0, \quad \text{avec} \quad \hat{x} = y - \frac{\sigma_b^2}{x_0}$$

Revenons maintenant sur le premier exemple et considérons que nous avons fait  $M$  mesures aux intervalles réguliers  $\mathbf{Y} = \{y_1, \dots, y_M\}$  et que nous voulons estimer  $\mathbf{X} = \{X_1, \dots, X_M\}$ .

Une première hypothèse consiste à dire qu'*a priori*, il n'y a pas de raison de faire l'hypothèse que la loi de  $X_i$  soit différente de la loi de  $X_j$  et que  $X_i$  soit dépendante de  $X_j$  (hypothèse dite i.i.d.) et d'attribuer alors une loi de probabilité jointe:

$$p(\mathbf{x}) = \prod_{j=1}^M p(x_j) \longrightarrow \mathbf{X} \simeq \mathcal{N}(\mathbf{x}_0, \sigma_x^2 \mathbf{I}).$$

La même hypothèse concernant les éléments de  $\mathbf{B} = \{B_1, \dots, B_M\}$  nous permet d'écrire:

$$p(\mathbf{b}) = \prod_{j=1}^M p(b_j) \longrightarrow \mathbf{B} \simeq \mathcal{N}(\mathbf{0}, \sigma_b^2 \mathbf{I}).$$

Il n'est pas difficile de montrer que la loi *a posteriori* est:

$$p(\mathbf{x}|\mathbf{y}) = \prod_{j=1}^M p(x_j|y_j) \longrightarrow \mathbf{X}|\mathbf{Y} \simeq \mathcal{N}(\mathbf{Y}, \sigma^2 \mathbf{I})$$

Une deuxième hypothèse plus réaliste est de dire qu'*a priori*,  $X_j$  dépend de  $X_{j-1}$  (la valeur de tension à un instant ne peut pas être très différente de sa valeur à un pas d'échantillonnage plus loin). Comment peut-on alors traduire cette information? Supposons que nous puissions avoir une idée *a priori* sur le coefficient de corrélation entre  $X_j$  et  $X_{j-1}$ . Ainsi, si on résume tout ce que nous connaissons sur  $\mathbf{X}$ :

$$\mathbb{E}[X_j] = x_0, \quad \text{Var}[X_j] = \sigma_x^2, \quad \mathbb{E}[(X_j - x_0)(X_{j-1} - x_0)] = \beta \sigma_x^2.$$

Utilisant de nouveau le principe du ME, nous obtenons:

$$p(\mathbf{x}) \propto \exp\left[\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^t \mathbf{P}^{-1}(\mathbf{x} - \mathbf{x}_0)\right], \quad \text{avec} \quad \mathbf{P} = \sigma_x^2 \begin{pmatrix} 1 & \beta & & & \\ \beta & 1 & \beta & & \\ & \ddots & \ddots & \ddots & \\ & & & \beta & 1 & \beta \\ & & & & \beta & 1 \end{pmatrix}$$

Ainsi, d'une manière générale, l'utilisation du principe du maximum d'entropie nous conduit vers les lois exponentielles généralisées. En effet, si on suppose que notre information *a priori* est de la forme :

$$E[\phi_k(\mathbf{X})] = d_k, \quad k = 1, \dots, K,$$

nous avons vu que la loi à entropie maximale qui satisfait ces contraintes est de la forme :

$$p(\mathbf{x}) = \frac{1}{Z} \exp \left[ - \sum_{k=1}^K \lambda_k \phi_k(\mathbf{x}) \right],$$

où

$$Z(\lambda_1, \dots, \lambda_m) = \int p(\mathbf{x}) \exp \left[ - \sum_{k=1}^K \lambda_k \phi_k(\mathbf{x}) \right] d\mathbf{x}$$

et où les  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  sont déterminés par :

$$- \frac{\partial \ln Z(\lambda_1, \dots, \lambda_m)}{\partial \lambda_k} = d_k, \quad k = 1, \dots, K.$$

Un cas particulier intéressant est lorsque les fonctions  $\phi_k(x)$  sont séparables :

$$\phi_k(\mathbf{x}) = \sum_{j=1}^N \phi_{k_j}(x_j)$$

on a alors

$$p(\mathbf{x}) = \prod_{j=1}^N p(x_j), \quad \text{avec} \quad p(x_j) = \frac{1}{Z_j} \exp \left[ - \sum_{k=1}^K \lambda_k \phi_{k_j}(x_j) \right],$$

ce qui signifie que les  $x_j$  sont indépendants, et si de plus  $\phi_{k_j} = \phi_k$  on a

$$p(\mathbf{x}) = \prod_{j=1}^N p(x_j), \quad \text{avec} \quad p(x_j) = \frac{1}{z} \exp \left[ - \sum_{k=1}^K \lambda_k \phi_k(x_j) \right],$$

ce qui signifie que les  $x_j$  sont i.i.d (indépendants et de même loi).

On trouve parfois l'appellation *lois entropiques* pour les lois appartenant à ces deux derniers cas. Quelques cas particuliers de ces lois utilisées fréquemment sont :

– Gaussienne généralisée :

$$\phi(x_j) = |x_j|^r, \quad 1 < r < 2, \quad \longrightarrow \quad p(x_j) \propto \exp[-\alpha|x_j|^r].$$

La loi conjointe est alors de la forme :

$$p(\mathbf{x}) \propto \exp \left[ -\alpha \sum_{j=1}^N |x_j|^r \right] \quad \longrightarrow \quad -\ln p(\mathbf{x}) = cte + \alpha \sum_{j=1}^N |x_j|^r.$$

– Loi Gamma :

$$\phi(x_j) = -\alpha \ln\left(\frac{x_j}{m_j}\right) + \left(\frac{x_j}{m_j}\right), \quad \longrightarrow \quad p(x_j) \propto \left(\frac{x_j}{m_j}\right)^{-\alpha} \exp\left[\frac{x_j}{m_j}\right].$$

et

$$p(\mathbf{x}) = \prod_{j=1}^N p(x_j) \propto \left(\prod_{j=1}^N \left(\frac{x_j}{m_j}\right)^{-\alpha}\right) \exp\left[\sum_{j=1}^N \frac{x_j}{m_j}\right]$$

et par conséquent

$$-\ln p(\mathbf{x}) = cte - \alpha \sum_{j=1}^N \ln\left(\frac{x_j}{m_j}\right) + \sum_{j=1}^N \frac{x_j}{m_j}.$$

– Loi Béta :

$$\phi(x_j) = \alpha \ln(x_j) + \beta \ln(1 - x_j), \quad \longrightarrow \quad p(x_j) \propto x_j^\alpha (1 - x_j)^\beta.$$

et la loi conjointe est :

$$p(\mathbf{x}) \propto \prod_{j=1}^N x_j^\alpha (1 - x_j)^\beta \quad \longrightarrow \quad -\ln p(\mathbf{x}) = cte + \alpha \sum_{j=1}^N \ln(x_j) + \beta \sum_{j=1}^N \ln(1 - x_j).$$

Nous avons vu que dans le cas i.i.d., la forme générale de ces lois est :

$$p(\mathbf{x}) \propto \exp\left[-\sum_{k=1}^K \lambda_k \phi_k(\mathbf{x})\right], \quad \text{avec} \quad \phi_k(\mathbf{x}) = \sum_{j=1}^N \phi_k(x_j).$$

On peut se poser alors une question sur le choix des fonctions  $\phi_k(x_j)$ . Par exemple, on peut se demander s'il existe des familles de fonctions  $\phi_k$  de telle sorte que la lois  $p(x_j)$  et par conséquent la loi  $p(\mathbf{x})$  aient une propriété, par exemple d'invariance par changement d'échelle. Dans une étude récente, certains auteurs se sont posés de telles questions et ont essayé d'y répondre. Ici, nous résumons quelques uns de ces résultats.

L'idée de base se trouve dans la recherche des fonctions  $\phi_k$  de telle sorte que la famille de lois  $p(x_j)$  correspondantes soit fermée pour un groupe de transformations comme le changement d'échelle par exemple.

Prenons le cas du changement d'échelle pour illustrer cette idée. Si  $x_j$  représente une quantité et  $p(x_j)$  sa loi de probabilité, on voudrait que lors d'un changement d'échelle  $\alpha x_j$  avec  $\alpha > 0$ , la loi  $p(\alpha x_j)$  reste dans la même famille que  $p(x_j)$ . Autrement dit :

$$p(x_j) \propto \exp\left[-\sum_{k=1}^K \lambda_k \phi_k(x_j)\right] \longrightarrow p(\alpha x_j) \propto \exp\left[-\sum_{k=1}^K \lambda'_k \phi_k(x_j)\right]$$

où  $\lambda'_k$  ne doivent dépendre que des  $\lambda_k$  et de facteur d'échelle  $\alpha$ .

Il est alors facile de montrer que pour satisfaire à cette propriété il suffit d'avoir :

$$\forall x_j \ \& \ \forall \alpha > 0, \quad \sum_{k=1}^K \lambda_k \phi_k(x_j) = cte + \sum_{k=1}^K \lambda'_k \phi_k(x_j)$$

L'étude de cette propriété pour  $K = 1$  et  $K = 2$  donne les familles suivantes :

Lois à un paramètre  $K = 1$  :

$$\left\{ \phi(x) \right\} = \left\{ x^r, \ln x \right\}, \quad r > 0.$$

Lois à deux paramètres  $K = 2$  :

$$\left\{ (\phi_1(x), \phi_2(x)) \right\} = \left\{ (x^{r_1}, x^{r_2}), (x^{r_1}, \ln x), (x^{r_1}, x^{r_1} \ln x), (\ln x, \ln^2 x) \right\},$$

où  $r$ ,  $r_1$  et  $r_2$  sont des réels positifs.

On retrouve quelques cas particuliers intéressants en choisissant  $\phi_2(x) = x$  et pour différents choix de  $\phi_1(x)$  :

-  $\phi_1(x) = x^2$  ; on retrouve la loi gaussienne :

$$p(x) \propto \exp \left[ -\lambda x^2 - \mu x \right] \propto \exp \left[ -\lambda \left( x + \frac{\mu}{2\lambda} \right)^2 \right] = \mathcal{N} \left( m = \frac{-\mu}{\lambda}, \sigma^2 = \frac{1}{2\lambda} \right).$$

-  $\phi_1(x) = \ln x$  ; on retrouve la loi gamma :

$$p(x) \propto \exp \left[ -\lambda \ln x - \mu x \right] = x^{-\lambda} \exp \left[ -\mu x \right].$$

-  $\phi_1(x) = x \ln x$  ; on retrouve une loi dite de la forme entropie de Shannon :

$$p(x) \propto \exp \left[ -\lambda x \ln x - \mu x \right].$$

### 9.5.2 Modèles markoviens

La modélisation markovienne a pris une importance sans égale en traitement d'image depuis les travaux de Geman et Geman. L'objet de ce document n'est pas de donner une présentation complète de la modélisation markovienne, mais plutôt une présentation différente.

Nous avons vu dans la section précédente que, d'une manière générale, l'utilisation du principe du maximum d'entropie avec une information *a priori* de la forme :

$$E[\phi_k(\mathbf{x})] = d_k, \quad k = 1, \dots, K,$$

nous conduit aux lois de la forme :

$$p(\mathbf{x}) = \frac{1}{Z} \exp \left[ - \sum_{k=1}^K \lambda_k \phi_k(\mathbf{x}) \right].$$

Considérons maintenant le cas où les fonctions  $\phi_k(\mathbf{x})$  sont de la forme

$$\phi_k(\mathbf{x}) = \sum_j \sum_{i \in V_j} \phi_k(x_j, x_i)$$

où  $V_j$  signifie les (échantillons ou les pixels) voisins de  $j$ . Les fonctions  $\phi_k$  sont alors appelées *potentiels* et

$$U(\mathbf{x}) = \sum_j \sum_{i \in V_j} \phi_k(x_j, x_i)$$

*l'énergie*.

Un tel choix pour ces fonctions signifie que nous faisons l'hypothèse qu'il y a une corrélation locale entre les valeurs du signal aux instants successifs ou de l'image sur des pixels voisins.

Un cas particulier très courant est lorsque  $k = 1$  et  $\phi_k(x_j, x_i) = \phi_k(x_j - x_i)$ . On a alors

$$p(\mathbf{x}) \propto \exp \left[ -\lambda \sum_j \sum_{i \in V_j} \phi(x_j - x_i) \right]$$

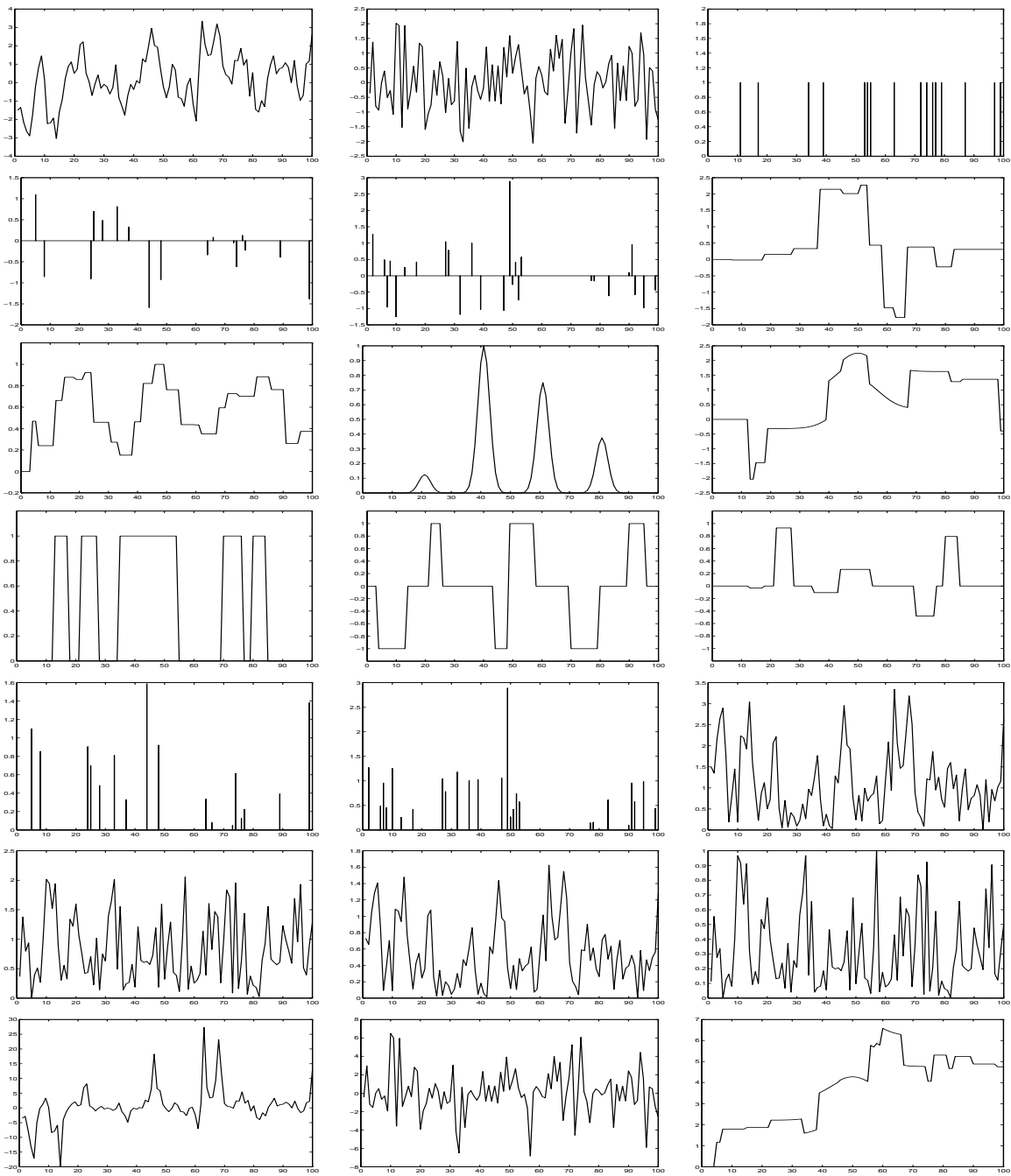
Le choix de la fonction potentiel  $\phi$  est crucial : les fonctions convexes permettront de modéliser des signaux et des champs continus et les fonctions non convexes permettront de modéliser des signaux et des champs qui contiennent des discontinuités. Une liste de ces fonctions classiquement utilisées en traitement du signal et des images est fournie en annexe D.



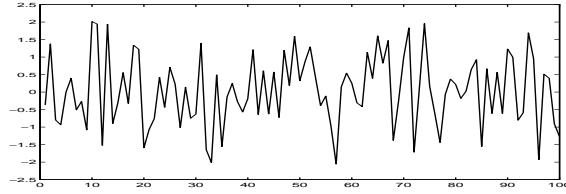
## 9.6 Quelques exemples de construction à la main

Dans cette section, nous allons montrer comment on peut choisir une loi de probabilité *a priori* à la main et avec du bon sens. La figure qui suit montre un ensemble de signaux très différents.

Nous allons, pour chacune de ces formes de signaux proposer une famille de lois convenable.

FIG. 9.1 - *Différents types de signaux.*

## 1. Modèle gaussien blanc pour des signaux continus à variation rapide



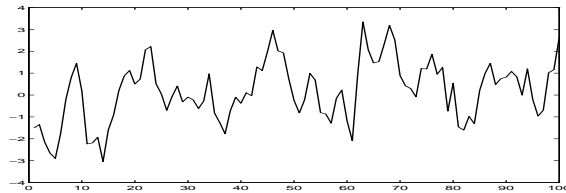
En regardant la forme de ce signal et en constatant que les échantillons du signal sont réparties d'une manière symétrique autour de zéro, on peut proposer une loi gaussienne centrée :

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_x = \sigma_x^2 \mathbf{I}) \longrightarrow p(\mathbf{x}) \propto \exp\left[-\frac{1}{2\sigma_x^2} \|\mathbf{x}\|^2\right]$$

ou encore

$$p(\mathbf{x}) \propto \exp\left[-\frac{1}{2\sigma_x^2} \sum_j x_j^2\right].$$

## 2. Modèle gaussien coloré pour des signaux continus à variation lente :



En regardant la forme de ce signal, et en la comparant au cas précédent, on peut proposer une loi centrée, gaussienne mais colorée :

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_x = \sigma_x^2 \mathbf{P}_0) \longrightarrow p(\mathbf{x}) \propto \exp\left[-\frac{1}{2\sigma_x^2} \mathbf{x}^t \mathbf{P}_0^{-1} \mathbf{x}\right],$$

avec  $\mathbf{P}_0$  une matrice symétrique et définie positive, ce qui permet d'écrire  $\mathbf{P}_0^{-1} = \mathbf{D}^t \mathbf{D}$  et on a alors

$$p(\mathbf{x}) \propto \exp\left[-\frac{1}{2\sigma_x^2} \|\mathbf{D}\mathbf{x}\|^2\right].$$

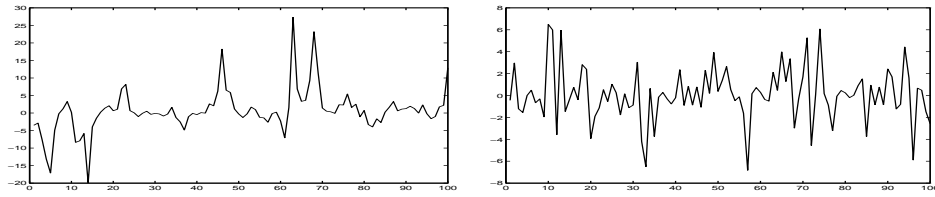
Un choix possible pour  $\mathbf{D}$  est la matrice de différences finies d'ordre un :

$$\mathbf{D}_1 = \begin{pmatrix} 1 & 0 & \cdots & \cdots & 0 \\ -1 & 1 & \ddots & & \vdots \\ 0 & -1 & 1 & \ddots & \vdots \\ \vdots & \ddots & & & 0 \\ 0 & \cdots & 0 & -1 & 1 \end{pmatrix}$$

ce qui permet d'écrire :

$$p(\mathbf{x}) \propto \exp\left[-\frac{1}{2\sigma_x^2} \sum_j (x_j - x_{j-1})^2\right]$$

### 3. Modèle gaussien généralisé pour des signaux continus plutôt impulsions



En regardant la forme de ce signal, et en la comparant au cas précédent, on peut constater que la proportion des échantillons ayant des valeurs proches de zéro est plus importante que celle des échantillons ayant des valeurs plus grandes et que cette proportion semble décroître très rapidement. On peut alors proposer une loi gaussienne généralisée :

$$p(\mathbf{x}) \propto \exp[-\beta \|\mathbf{D}_k \mathbf{x}\|^p], \quad 1 < p \leq 2,$$

avec  $\mathbf{D}_0 = \mathbf{I}$  et  $\mathbf{D}_k$  la matrice des différences finies d'ordre  $k$ .

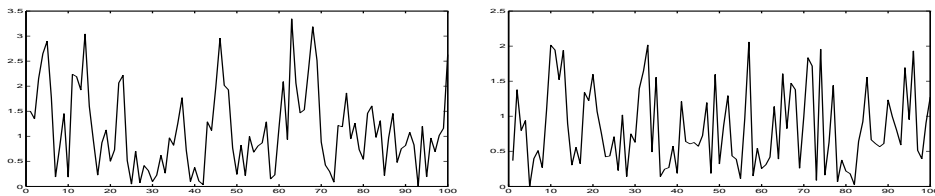
Dans le cas  $\mathbf{D}_0 = \mathbf{I}$  on a

$$p(\mathbf{x}) \propto \exp \left[ -\beta \sum_j |x_j|^p \right]$$

et dans le cas  $\mathbf{D}_1$  on a

$$p(\mathbf{x}) \propto \exp \left[ -\beta \sum_j |x_j - x_{j-1}|^p \right].$$

### 4. Modèle des signaux positifs



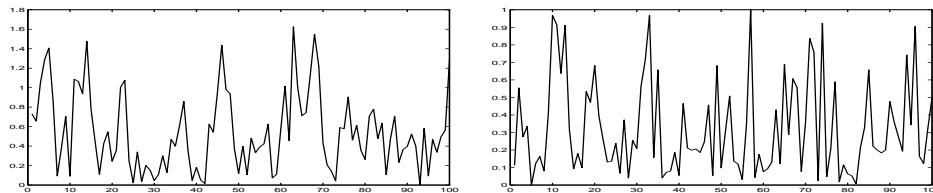
Loi gaussienne généralisée tronquée :

$$p(\mathbf{x}) \propto \exp \left[ -\beta \sum_j x_j^p \right], \quad x_j > 0.$$

Loi Gamma :

$$p(\mathbf{x}) \propto \exp \left[ -\alpha \sum_j \ln x_j - \beta \sum_j x_j \right], \quad x_j > 0.$$

## 5. Modèle des signaux continus mais bornés entre 0 et 1 : Loi Béta



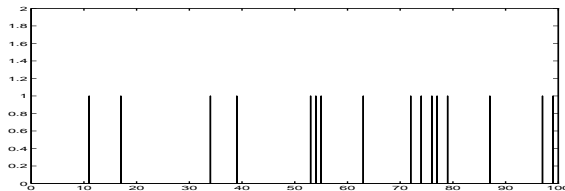
$$p(\mathbf{x}) = \prod_{j=1}^N p(x_j),$$

avec

$$\begin{aligned} p(x_j) &\propto x_j^{-\alpha} (1 - x_j)^{-\beta} \\ &\propto \exp[-\alpha \ln x_j - \beta \ln(1 - x_j)] \end{aligned}$$

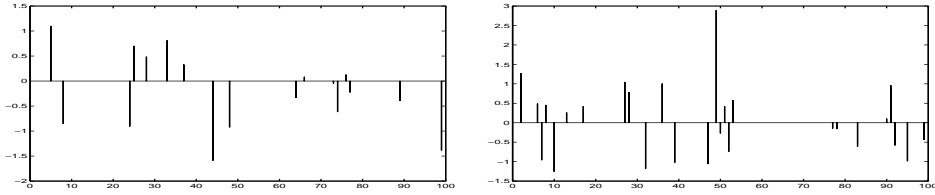
$$-\ln p(\mathbf{x}) = cte + \alpha \ln x_j + \beta \ln(1 - x_j)$$

## 6. Modèle de Bernoulli :



$$\begin{cases} P(X_j = 1) = \alpha, \\ P(X_j = 0) = 1 - \alpha \end{cases} \longrightarrow P(\mathbf{X} = \mathbf{x}) = \alpha^s (1 - \alpha)^{1-s} \text{ avec } s = \sum (x_j)$$

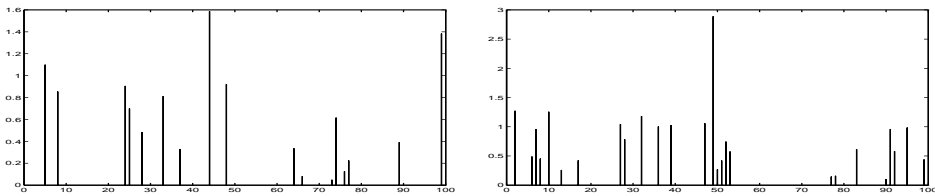
## 7. Modèle Bernoulli-Gaussien :



$$p(x_j|q_j) = \begin{cases} 0 & \text{si } q_j = 0, \\ \mathcal{N}(0, \sigma_x^2) & \text{si } q_j = 1 \end{cases} \quad \text{et} \quad \begin{cases} P(q_j = 1) = \lambda, \\ P(q_j = 0) = 1 - \lambda \end{cases}$$

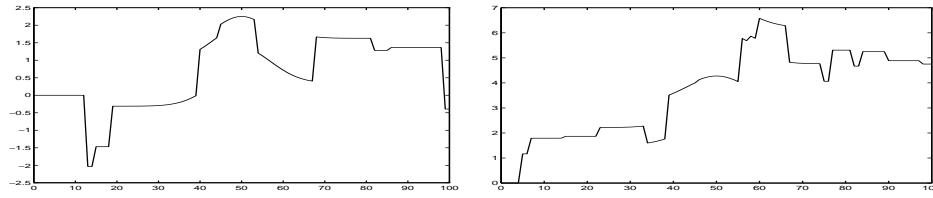
$$p(\mathbf{x}|\mathbf{q}) \propto \exp \left[ \frac{-1}{2q_j\sigma_x^2} \sum_j q_j x_j^2 \right]$$

## 8. Modèle Bernoulli-Gamma :



$$p(x_j|q_j) = \begin{cases} 0 & \text{si } q_j = 0, \\ \text{Gamma}(\alpha, \beta) & \text{si } q_j = 1 \end{cases} \quad \text{et} \quad \begin{cases} P(q_j = 1) = \lambda, \\ P(q_j = 0) = 1 - \lambda \end{cases}$$

## 9. Modèles pour les signaux continus par morceaux

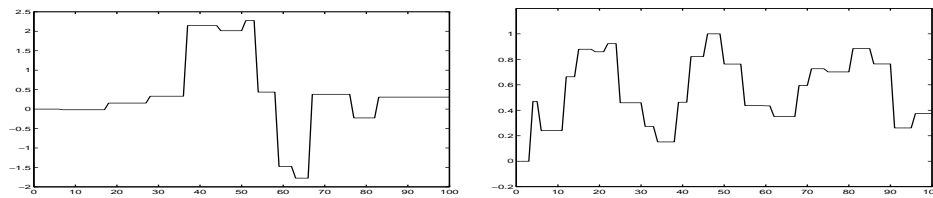


$$p(x_j|x_{j-1}, q_j) = \begin{cases} \mathcal{N}(x_{j-1}, \sigma_1^2) & \text{si } q_j = 0 \\ \mathcal{N}(0, \sigma_2^2) & \text{si } q_j = 1 \end{cases}$$

ou

$$p(\mathbf{x}|\mathbf{q}) \propto \exp \left[ -\frac{1}{\sigma_1^2} \sum_j q_j (x_j - x_{j-1}) \right]$$

## 10. Modèles pour les signaux constants par morceaux



$$p(x_j|x_{j-1}, q_j) = \begin{cases} \delta(x_j - x_{j-1}) & \text{si } q_j = 0 \\ \mathcal{N}(0, \sigma_2^2) & \text{si } q_j = 1 \end{cases}$$

## 9.7 Estimation des paramètres d'une loi à partir des observations directes

Nous avons vu que le PME peut nous permettre d'attribuer une loi de probabilité  $p(x)$  à une quantité  $X$  pour traduire une information incomplète sur  $X$  si cette information porte sur des moments  $E[\phi_k(X)] = d_k$ . Malheureusement, en pratique, rarement nous connaissons les valeurs numériques de  $d_k$ . C'est pourquoi, en général, l'attribution d'une loi de probabilité se fait en deux étapes :

1. Choix d'une famille de loi qui est équivalent au choix des fonctions  $\phi_k$  ;
2. Détermination des valeurs des paramètres à partir d'un certain nombre d'observations directes ou indirectes du  $X$ .

Dans cette section, nous allons nous limiter au cas où nous avons la possibilité d'observer directement  $X$ . Supposons maintenant que l'étape 1 est faite et que nous avons choisi une loi  $p(x; \theta)$  dépendant des paramètres  $\theta = [\theta_1, \dots, \theta_K]$ . Considérons alors le cas où nous avons  $N$  échantillons  $\{x_1, \dots, x_N\}$ . Nous cherchons à déterminer  $\theta$  à partir de ces échantillons. Deux méthodes classiquement utilisées sont : la méthode des moments (MM) et la méthode du maximum de vraisemblance (MV).

### 9.7.1 Méthode des moments

L'idée de base dans cette méthode est de déterminer  $\theta = (\theta_1, \dots, \theta_K)$  en utilisant un système d'équations reliant les moments théoriques  $E[X^k]$  et les moments empiriques  $\overline{X^k} = \frac{1}{N} \sum_{j=1}^N x_j^k$  :

$$E[X^k] = \int x^k p(x; \theta) dx = \frac{1}{N} \sum_{j=1}^N x_j^k, \quad k = 1, \dots, M$$

avec  $M \geq K$  et en espérant qu'il existe une solution unique à ce système d'équations (non linéaires en  $\theta$ ).

Pour illustrer cette méthode considérons les deux exemples suivants :

#### Exemple 1 :

Cas d'une loi gaussienne à deux paramètres  $\theta = [m, \sigma^2]$

$$p(x; m, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(x - m)^2\right].$$

Si on forme le système d'équations :

$$\begin{aligned} E[X] &= m = \frac{1}{N} \sum_{j=1}^N x_j \\ E[X^2] &= \sigma^2 + m^2 = \frac{1}{N} \sum_{j=1}^N x_j^2 \end{aligned}$$



il est alors facile de voir qu'il y a une solution qui est :

$$m = \bar{X} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$\sigma^2 = \overline{(X - m)^2} = \frac{1}{N} \sum_{j=1}^N (x_j - m)^2.$$

**Exemple 2 :**

Cas d'une loi gamma à deux paramètres  $\theta = [\alpha, \beta]$

$$p(x; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma \alpha} x^{\alpha-1} \exp[-\beta x].$$

Si on forme le système d'équations :

$$E[X] = \frac{\alpha}{\beta} = \frac{1}{N} \sum_{j=1}^N x_j,$$

$$E[X^2] = \frac{\alpha(1 + \alpha)}{\beta^2} = \frac{1}{N} \sum_{j=1}^N x_j^2,$$

et si on note par  $m = \bar{X} = \frac{1}{N} \sum_{j=1}^N x_j$  et par  $v = \overline{(X - m)^2} = \frac{1}{N} \sum_{j=1}^N (x_j - m)^2$ , il est alors facile de voir qu'il y a une solution qui est :

$$\alpha = \frac{2v - m^2}{v},$$

$$\beta = \frac{m}{v}.$$

Les inconvénients majeurs de cette méthode sont :

- l'absence de base théorique pour cette méthode;
- la sensibilité au nombre de données  $N$  de ces estimateurs.

### 9.7.2 Méthode du maximum de vraisemblance

Estimation au sens du MV consiste à définir la fonction de vraisemblance  $V(\theta) = p(x_1, \dots, x_n; \theta)$  ou plutôt la fonction log-vraisemblance  $L(\theta) = -\ln V(\theta)$  et à calculer l'argument  $\hat{\theta}$  qui la maximise :

$$\hat{\theta} = \arg \max_{\theta} \{V(\theta)\} = \arg \min_{\theta} \{L(\theta)\}$$

Le principal avantage de cet estimateur est qu'il est efficace.

Pour illustrer cette méthode et la comparer à la méthode des moments nous allons considérer les deux exemples de la section précédente.

**Exemple 1 :**

Cas d'une loi gaussienne à deux paramètres  $\theta = [m, \sigma^2]$

$$p(x; m, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{1}{2\sigma^2}(x - m)^2 \right]$$

La fonction du log-vraisemblance s'écrit :

$$L(m, \sigma^2) = \frac{N}{2} \ln(2\pi) + \frac{N}{2} \ln \sigma^2 + \frac{1}{2\sigma^2} \sum_{j=1}^N (x_j - m)^2.$$

Pour calculer l'argument qui la minimise, il faut annuler simultanément la dérivée par rapport à  $m$  et la dérivée par rapport à  $\sigma^2$ , ce qui permet d'obtenir un résultat bien connu :

$$\begin{aligned} m &= \frac{1}{N} \sum_{j=1}^N x_j, \\ \sigma^2 &= \frac{1}{N-1} \sum_{j=1}^N (x_j - m)^2. \end{aligned}$$

Notons la seule différence par rapport aux résultats de la méthode des moments apparaît dans l'estimation du  $\sigma^2$ .

**Exemple 2 :**

Cas d'une loi gamma à deux paramètres  $\theta = [\alpha, \beta]$

$$p(x; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp[-\beta x].$$

La fonction du log-vraisemblance s'écrit :

$$L(\alpha, \beta) = N[\alpha \ln \beta - \ln \Gamma(\alpha)] - (\alpha - 1) \sum_{j=1}^N \ln x_j - \beta \sum_{j=1}^N x_j.$$

Pour calculer l'argument qui la minimise, il faut annuler simultanément la dérivée par rapport à  $\alpha$  et la dérivée par rapport à  $\beta$  :

$$\begin{aligned} \frac{\partial L}{\partial \alpha} &= N \left[ \ln \beta - \frac{\partial \ln \Gamma(\alpha)}{\partial \alpha} \right] - \sum_{j=1}^N \ln x_j = 0 \\ \frac{\partial L}{\partial \beta} &= \frac{\alpha}{\beta} - \sum_{j=1}^N x_j = 0 \end{aligned}$$

Comme on peut le constater, il n'y a pas de solution analytique à ces deux équations, mais on peut la calculer numériquement.

Pour montrer cependant qu'il peut y avoir un lien entre ces deux méthodes, on peut réécrire les relations différemment.

Supposons que la loi  $p(x; \theta)$  soit de la forme générale :

$$p(x; \theta) \propto \exp \left[ - \sum_{k=1}^K \theta_k \phi_k(x) \right].$$

On a alors

$$L(\boldsymbol{\theta}) = -\ln p(x_1, \dots, x_N; \boldsymbol{\theta}) = \text{cte} + \sum_{j=1}^N \sum_{k=1}^K \theta_k \phi_k(x_j).$$

On peut alors montrer facilement que, dans la méthode des moments nous avons à résoudre

$$\int x^k p(x; \boldsymbol{\theta}) dx = \frac{1}{N} \sum_{j=1}^N x_j^k, \quad k = 0, \dots, K,$$

et dans la méthode du MV nous avons à résoudre

$$\int \phi_k(x) p(x; \boldsymbol{\theta}) dx = \frac{1}{N} \sum_{j=1}^N \phi_k(x_j), \quad k = 0, \dots, K,$$

où nous avons noté  $\phi_0(x) = 1$ .

On peut alors remarquer que les deux méthodes sont équivalentes pour la famille des lois exponentielles où  $\phi_k(x) = x^k$ .

## 9.8 Estimation des hyperparamètres

Il existe assez peu de méthodes de détermination des hyperparamètres. Dans le cas où ceux-ci se limitent au seul coefficient de régularisation, et où le régulariseur est quadratique, les méthodes de *validation croisée* fournissent des solutions acceptables. Mais il s'agit de méthodes déterministes par essence, qui minimisent un critère de risque dépendant de l'objet inconnu  $\mathbf{f}$ , ce qui ne peut être effectué qu'asymptotiquement puisque cet objet est évidemment inconnu.

Les méthodes bayésiennes, qui attribuent une distribution de probabilité *a priori* à l'objet, ne présentent pas cette limitation. Les hyperparamètres  $\theta$  constituent un second niveau de description du problème, indispensable pour “rigidifier” le premier niveau constitué par les paramètres eux-mêmes— c'est-à-dire l'objet  $\mathbf{f}$ . Dans un problème mal-posé, la valeur des hyperparamètres est importante pour obtenir une solution acceptable, mais ne présente pas d'intérêt en soi. Dans une approche bayésienne, on peut donc distinguer deux niveaux d'inférence. Le premier infère sur  $\mathbf{f}$ , pour une valeur donnée de  $\theta$ , au travers de la distribution *a posteriori*. Le second infère sur  $\theta$  grâce à une relation analogue :

$$p(\theta | \mathbf{g}, \mathbf{H}) = \frac{p(\theta | \mathbf{H}) p(\mathbf{g} | \theta, \mathbf{H})}{p(\mathbf{g} | \mathbf{H})}. \quad (9.5)$$

On retrouve là une caractéristique de l'utilisation de la règle de Bayes : la vraisemblance  $p(\mathbf{g} | \theta, \mathbf{H})$  attachée aux données dans le second niveau est le coefficient de normalisation dans le premier.

Si, comme cela est souvent le cas, ce terme est suffisamment “piqué”, l'influence de la distribution *a priori*  $p(\theta | \mathbf{H})$  est négligeable, et le second niveau d'inférence peut être résolu par maximisation de cette vraisemblance. Mais il faut pour cela résoudre un problème de marginalisation :

$$p(\mathbf{g} | \theta, \mathbf{H}) = \int p(\mathbf{f}, \mathbf{g} | \theta, \mathbf{H}) d\mathbf{f} = \int p(\mathbf{g} | \mathbf{f}, \theta, \mathbf{H}) p(\mathbf{f} | \theta) d\mathbf{f}. \quad (9.6)$$

Une telle intégrale conduit très rarement à une forme explicite.

Pour contourner cette difficulté, on peut introduire des “variables cachées”  $\mathbf{z}$  qui viennent compléter les observations  $\mathbf{g}$  de manière à ce que la nouvelle vraisemblance  $p(\mathbf{g}, \mathbf{z} | \theta, \mathbf{H})$  soit plus simple à calculer. On est alors conduit à maximiser des espérances conditionnelles par des techniques itératives, déterministes ou stochastiques (algorithmes EM et SEM).

On peut aussi remarquer que la *vraisemblance généralisée*

$$p(\mathbf{g}, \mathbf{f} | \theta, \mathbf{H}) = p(\mathbf{f} | \mathbf{g}, \theta, \mathbf{H}) p(\mathbf{g} | \theta, \mathbf{H}) = p(\mathbf{g} | \mathbf{f}, \theta, \mathbf{H}) p(\mathbf{f} | \theta) \quad (9.7)$$

résume toute l'information propre au premier niveau d'inférence, et vouloir en faire la maximisation conjointe par rapport à  $\mathbf{f}$  et  $\theta$ . Le problème d'intégration soulevé par (9.6) est évidemment évacué. A  $\theta$  fixé, le maximum de vraisemblance généralisée (MVG) coïncide avec le MAP. Par contre, à  $\mathbf{f}$  fixé, la situation est beaucoup moins favorable : la vraisemblance généralisée n'est pas, en général, majorée sur son domaine de définition, et il n'existe même pas toujours un maximum local.

A COMPLETER

## Chapitre 10

# Modélisation markovienne

Dans ce chapitre nous donnons un aperçu général de la modélisation markovienne des signaux et des images. Les modèles markoviens et particulièrement les champs de Markov et les champs de Gibbs sont devenus des outils indispensables pour la traduction d'une information *a priori* plus élaborée que la douceur ou la positivité en une loi de probabilité. En effet, la modélisation markovienne permet de construire des modèles qui peuvent prendre en compte les discontinuités dans un signal ou une image. Cette outil prend particulièrement son importance en traitement d'image où l'on peut construire les modèles composites (bords, contours, intensité) pour les images et les utiliser lors de la segmentation, restauration ou reconstruction d'image.