

## Chapitre 8

# Approches probabilistes

Nous avons distingué quatre approches probabilistes différentes pour la résolution des problèmes inverses. L'objet de ce chapitre est de détailler un peu plus ces approches.

Dans ce chapitre, nous aborderons les trois premières approches, en montrant par des exemples simples leur limites. Nous réservons une place privilégiée à l'approche bayésienne qui sera détaillée dans un chapitre séparé.

## 8.1 Approche n'utilisant que les moments

Dans cette approche, on considère que les grandeurs observées et inconnues sont bien des fonctions aléatoires, on se limite à caractériser les lois de probabilités de ces grandeurs que par leurs moments, souvent jusqu'à l'ordre deux. Ensuite on établit un lien entre les moments des grandeurs observées et ceux des grandeurs inconnues puis on cherche à l'inverser.

Mais, un point qui est souvent laissé pour compte est : l'estimation des moments des grandeurs observées. Souvent on fait l'hypothèse que l'on peut estimer ces moments par les moyennes empiriques. Ceci peut être raisonnable lorsqu'on se limite aux moments d'ordre un et deux, mais cela devient inefficace pour les moments d'ordre supérieur. C'est d'ailleurs dans le premier cas que ces méthodes trouvent particulièrement leur intérêt : c'est l'estimation en moyenne quadratique.

Pour illustrer cette approche, nous allons considérer le cas du filtrage optimal de Wiener pour la résolution du problème de déconvolution de signaux.

### Filtrage optimal de WIENER :

Considérons le problème de la déconvolution des signaux:

$$g(t) = h(t) * f(t) + b(t)$$

supposons que  $f(t)$ ,  $b(t)$  et  $g(t)$  sont des fonctions aléatoires et cherchons à lier les moments d'ordre un :

$$E[g(t)], \quad E[b(t)] \quad \text{et} \quad E[f(t)]$$

et les moments d'ordre deux :

$$\begin{aligned} R_{gg}(\tau) &= E[g(t)g(t+\tau)], \\ R_{ff}(\tau) &= E[f(t)f(t+\tau)], \\ R_{bf}(\tau) &= R_{fb}(-\tau) = E[b(t)f(t+\tau)] \\ R_{gf}(\tau) &= R_{fg}(-\tau) = E[g(t)f(t+\tau)] \end{aligned}$$

de ces grandeurs. Faisant l'hypothèse que  $f(t)$  et  $b(t)$  sont indépendantes. On obtient :

$$\begin{aligned} E[g(t)] &= h(t) * E[f(t)] + E[b(t)] \\ R_{gg}(\tau) &= h(t) * h(t) * R_{ff}(\tau) + R_{bb}(\tau) \\ R_{gf}(\tau) &= h(t) * R_{ff}(\tau) \end{aligned}$$

Passant dans le domaine de FOURIER on a :

$$\begin{aligned} S_{gg}(\omega) &= |H(\omega)|^2 S_{ff}(\omega) + R_{bb}(\omega) \\ S_{gf}(\omega) &= H(\omega) S_{ff}(\omega), \\ S_{fg}(\omega) &= H^*(\omega) S_{ff}(\omega), \end{aligned}$$

L'objectif du filtrage optimal est de fournir une solution  $\hat{f}(t)$  qui s'obtient par un filtrage linéaire de  $g(t)$  :

$$\hat{f}(t) = w(t) * g(t)$$

et la réponse impulsionnelle du filtre  $w(t)$  est telle que l'erreur quadratique moyenne

$$EQM = E \left[ \left( f(t) - \hat{f}(t) \right)^2 \right]$$

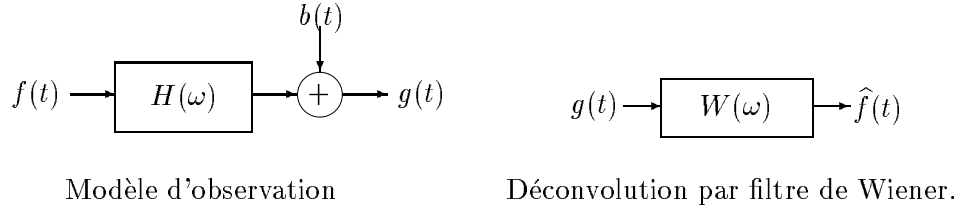


FIG. 8.1 - Filtrage de WIENER.

soit minimale.

En utilisant le principe d'orthogonalité:

$$E[(f(t) - w(t) * g(t)) g(t - \tau)] = 0 \quad \forall t, \tau \longrightarrow R_{fg}(\tau) = w(t) * R_{gg}(\tau).$$

Passant dans le domaine de FOURIER on a:

$$W(\omega) = \frac{S_{fg}(\omega)}{S_{gg}(\omega)} = \frac{H^*(\omega)S_{ff}(\omega)}{|H(\omega)|^2 S_{ff}(\omega) + S_{bb}(\omega)}.$$

Remplaçant pour les expressions de  $S_{fg}(\omega)$  et  $S_{gg}(\omega)$  on obtient :

$$W(\omega) = \frac{H^*(\omega)S_{ff}(\omega)}{|H(\omega)|^2 S_{ff}(\omega) + S_{bb}(\omega)} = \frac{1}{H(\omega)} \frac{|H(\omega)|^2}{|H(\omega)|^2 + \frac{S_{bb}(\omega)}{S_{ff}(\omega)}}$$

La difficulté dans cette approche commence lors de la mise en œuvre réelle de la méthode. En effet, il faut pouvoir estimer  $S_{gg}(\omega)$  et  $S_{fg}(\omega)$  ce qui nécessite la connaissance *a priori* de  $S_{ff}(\omega)$  et de  $S_{bb}(\omega)$ .

Souvent, en pratique on fait l'hypothèse que

$$\frac{S_{bb}(\omega)}{S_{ff}(\omega)} = r,$$

avec  $r$  une constante représentant l'inverse du rapport signal à bruit. Avec cette hypothèse on a

$$W(\omega) = \frac{1}{H(\omega)} \frac{|H(\omega)|^2}{|H(\omega)|^2 + r}.$$

Cette hypothèse n'est cependant pas très réaliste car elle considère que le signal et le bruit ont la même forme de spectre, ce qui est rarement le cas. En effet, on peut souvent considérer le bruit blanc, mais le signal que l'on cherche n'a pas le même spectre.

Notons que lorsque  $r \mapsto 0$  on retrouve le filtrage inverse.

Une autre hypothèse plus réaliste consiste à supposer

$$\frac{S_{bb}(\omega)}{S_{ff}(\omega)} = |D(\omega)|^2,$$

où  $D(\omega)$  représente la fonction de transfert d'un filtre passe-haut. Avec cette hypothèse on a

$$W(\omega) = \frac{1}{H(\omega)} \frac{|H(\omega)|^2}{|H(\omega)|^2 + |D(\omega)|^2}.$$

Nous verrons dans le chapitre suivant qu'il existe un lien entre cette approche et celle de l'estimation bayésienne dans le cas linéaire avec des lois gaussiennes.

## 8.2 Estimation au sens du maximum de vraisemblance

Dans cette approche on prend en compte explicitement le caractère incertain des mesures en les caractérisant par une loi de probabilité  $p(\mathbf{g}|\mathbf{f})$ .

L'idée de base ensuite est de considérer  $p(\mathbf{g}|\mathbf{f})$  comme une fonction  $V(\mathbf{f}) = p(\mathbf{g}|\mathbf{f})$  appelée fonction *vraisemblance*. On définit alors la solution  $\hat{\mathbf{f}}$  du problème inverse comme l'argument qui maximise cette fonction :

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

### Exemple 1: Loi Gaussienne :

Considérons le modèle  $\mathbf{g} = \mathbf{H}\mathbf{f} + \mathbf{b}$  et supposons que  $\mathbf{b}$  puisse être modélisé par un vecteur aléatoire centrée, blanc et gaussien :

$$b_i \sim \mathcal{N}(0, \sigma_b^2) \longrightarrow \mathbf{b} \sim \mathcal{N}(\mathbf{0}, \sigma_b^2 \mathbf{I}) \longrightarrow \mathbf{g}|\mathbf{f} \sim \mathcal{N}(\mathbf{H}\mathbf{f}, \sigma_b^2 \mathbf{I}).$$

On en déduit alors :

$$p(\mathbf{g}|\mathbf{f}) = K \exp \left[ \frac{-1}{2\sigma_b^2} [\mathbf{g} - \mathbf{H}\mathbf{f}]^t [\mathbf{g} - \mathbf{H}\mathbf{f}] \right].$$

L'estimé au sens du MV devient alors :

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{-\ln p(\mathbf{g}|\mathbf{f})\} = \arg \min_{\mathbf{f}} \{\|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2\}$$

et on retrouve l'estimation au sens des moindres carrés.

Si nous connaissons un peu plus sur les caractéristiques des  $b_i$ , par exemple si les  $b_i$  ont des variances différentes, alors

$$b_i \sim \mathcal{N}(0, \sigma_i^2)$$

et il n'est pas difficile de montrer que

$$p(\mathbf{g}|\mathbf{f}) = K \exp \left[ \frac{-1}{2} (\mathbf{g} - \mathbf{H}\mathbf{f})^t \mathbf{W}^{-1} (\mathbf{g} - \mathbf{H}\mathbf{f}) \right]$$

avec  $\mathbf{W} = \text{diag} \{\sigma_1^2, \dots, \sigma_m^2\}$ . Alors,

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{-\ln p(\mathbf{g}|\mathbf{f})\} = \arg \min_{\mathbf{f}} \{\|\mathbf{g} - \mathbf{H}\mathbf{f}\|_{\mathbf{W}}^2\}$$

et on retrouve l'estimation au sens des moindres carrés généralisés.

L'extension au cas d'un modèle non linéaire se fait aussi très facilement.

### Exemple 2: Loi Gaussienne généralisée :

Supposons maintenant que les  $b_i$  suivent une loi gaussienne généralisée :

$$b_i \sim p(b_i) \propto \exp[-\beta |b_i|^p], \quad \beta > 0, \quad 1 < p \leq 2.$$

Si on fait l'hypothèse que  $\mathbf{b}$  est blanc alors on a :

$$\mathbf{b} \sim p(\mathbf{b}) \propto \exp \left[ -\beta \sum_{i=1}^M |b_i|^p \right] \longrightarrow -\ln p(\mathbf{b}) = K + \beta \sum_{i=1}^M |b_i|^p$$

On en déduit alors :

$$-\ln p(\mathbf{g}|\mathbf{f}) = K + \beta \sum_{i=1}^M |g_i - [\mathbf{H}\mathbf{f}]_i|^p.$$

L'estimé au sens du MV devient alors :

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{-\ln p(\mathbf{g}|\mathbf{f})\} = \arg \min_{\mathbf{f}} \left\{ \sum_{i=1}^M |g_i - [\mathbf{H}\mathbf{f}]_i|^p \right\} = \arg \min_{\mathbf{f}} \{\|\mathbf{g} - \mathbf{H}\mathbf{f}\|^p\}$$

et on retrouve l'estimation au sens des moindres carrés.

### Exemple 3: Loi Gamma :

Supposons maintenant que les  $b_i$  suivent une loi de gamma :

$$b_i \sim \mathcal{G}(\alpha, \beta) \longrightarrow p(b_i) = K_i b_i^{-\alpha} \exp[-\beta b_i]$$

et si on fait l'hypothèse que  $\mathbf{b}$  est blanc ( $b_i$  et  $b_j$ ) indépendantes  $\forall i, j, i \neq j$ , on a

$$p(\mathbf{b}) = \prod_{i=1}^M p(b_i) \longrightarrow \ln p(\mathbf{b}) = K + \sum_{i=1}^M \alpha \log(b_i) + \beta b_i.$$

On en déduit :

$$\ln p(\mathbf{g}|\mathbf{f}) = K + \sum_{i=1}^M \alpha \log(g_i - [\mathbf{H}\mathbf{f}]_i) + \beta(g_i - [\mathbf{H}\mathbf{f}]_i)$$

et

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{-\ln p(\mathbf{g}|\mathbf{f})\}.$$

Il n'y a pas d'expression analytique pour cette solution, mais elle peut être calculé numériquement par un algorithme itératif.

En conclusion :

- L'approche MV permet de mieux prendre en compte la nature du bruit en lui attribuant une loi de probabilité  $p(\mathbf{b})$ , ou, d'une manière plus générale, l'incertaine sur les mesures en leur attribuant une loi de probabilité  $p(\mathbf{g}|\mathbf{f})$  ;
- Dans le cas gaussien on retrouve l'estimation au sens des moindres carrés ;
- Cette approche peut fournir des résultats satisfaisants lorsqu'elle est utilisée en combinaison avec une modélisation paramétrique des inconnues avec un nombre de paramètres très inférieure au nombre de données indépendantes. Malheureusement, cette approche donne rarement des résultats satisfaisants pour la résolution des problèmes inverses dans le cadre algébrique où le nombre d'inconnues est du même ordre de grandeur voir plus grand que le nombre des mesures. C'est pourquoi, certains auteurs ont proposé une approche dite *Maximum de vraisemblance pénalisée* qui consiste à définir la solution par :

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{-\ln p(\mathbf{g}|\mathbf{f}) + \phi(\mathbf{f})\}$$

où  $\phi(\mathbf{f})$  est une fonction de pénalisation choisie d'une manière adhoc pour permettre d'obtenir une solution satisfaisante. Nous verrons alors que cette approche peut mieux être interprétée dans le cadre de l'inférence bayésienne où  $\phi(\mathbf{f}) = -\ln p(\mathbf{f})$  avec  $p(\mathbf{f})$  la loi *a priori*.

### 8.3 Approche du maximum d'entropie

Avant de voir comment le principe du maximum d'entropie peut être utilisé pour la résolution des problèmes inverses, nous allons rappeler un certain nombre de définitions et de notions qui seront essentielles pour la compréhension de cette approche.

#### 8.3.1 Définition de l'entropie

Le principe du ME peut être approché de différentes manières. L'approche de la théorie de l'information [1] est sans doute la mieux adaptée à notre problème. JAYNES ([2, 3, 4]) est parmi les premiers auteurs modernes à avoir introduit le formalisme du ME par l'approche de la théorie de l'information. Cette notion d'entropie est introduite de la manière suivante : considérons une variable aléatoire discrète  $X$  produisant des réalisations  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$  et attribuons les probabilités  $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$  à ces réalisations pour représenter notre information partielle sur cette variable. On définit la quantité  $I_i = \ln(1/p_i)$  comme la quantité d'information apportée par la réalisation  $x_i$ .

Le raisonnement intuitif conduisant à cette expression est le suivant : plus un événement est rare, plus le gain d'information obtenu par sa réalisation est grand. L'utilisation du logarithme rend additif le gain total d'information obtenu par la réalisation de plusieurs événements indépendants. On définit alors l'entropie d'un processus par la somme pondérée des informations individuelles de chaque réalisation. C'est la définition de l'entropie donnée par SHANNON :

$$S(\mathbf{p}) = \sum_{i=1}^n p_i \ln \frac{1}{p_i} = - \sum_{i=1}^n p_i \ln p_i. \quad (8.1)$$

L'entropie  $S$  est une mesure d'incertitude de la distribution  $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$ , déterminée uniquement par certaines règles élémentaires de cohérence logique et d'additivité ([3, 4, 5]).

Généralisant ce concept, on définit  $-\ln(q_i/p_i)$  comme le gain d'information, sur une probabilité *a priori*  $p_i$ , apporté par la connaissance de la probabilité  $q_i$  de réalisation de l'événement  $x_i$ . On définit alors :

$$S(\mathbf{q}, \mathbf{p}) = - \sum_{i=1}^n q_i \ln \frac{p_i}{q_i} = \sum_{i=1}^n q_i \ln \frac{q_i}{p_i}, \quad (8.2)$$

appelée entropie croisée ou entropie relative de la distribution  $q_i$  par rapport à la distribution  $p_i$ . Notons ici qu'il y a un changement de signe pour des raisons historiques et il est clair que la minimisation de l'entropie croisée  $S(\mathbf{q}, \mathbf{p})$  se réduit à la maximisation de l'entropie  $S(\mathbf{q})$  si l'*a priori*  $\mathbf{p}$  est uniforme.

Sous réserve de quelques précautions, on peut généraliser ce qui précède au cas de distributions continues, et définir l'entropie [1] par :

$$S(p) = - \int p(x) \ln p(x) dx, \quad (8.3)$$

et l'entropie croisée par :

$$S(q, p) = - \int q(x) \ln \frac{p(x)}{q(x)} dx. \quad (8.4)$$

### 8.3.2 Lois à maximum d'entropie

Voyons maintenant ce que signifie le choix d'une distribution de probabilité à maximum d'entropie contenant une information *a priori*, ou qui soit compatible avec des contraintes connues sur cette distribution. Pour cela, précisons tout d'abord la signification de l'information contenue dans une distribution de probabilité  $\{p_1, p_2, \dots, p_n\}$ . À l'évidence il faut que l'on puisse extraire cette information de cette distribution. Considérons une variable aléatoire discrète  $X$  qui peut prendre des valeurs  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$  avec une distribution de probabilité  $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$ . Maintenant si on nous demande quelle est la meilleure estimée  $\hat{\phi}$  d'une fonction  $\phi(X)$  au sens du minimum de l'erreur quadratique moyenne. La solution est immédiate :

$$\hat{\phi} = E[\phi] = \sum_{i=1}^n p_i \phi(x_i).$$

C'est un problème direct et bien-posé. Inversement, si l'on se pose la question d'ajuster une distribution  $\mathbf{p}$  pour incorporer une information donnée sur la fonction  $\phi(X)$ , il faut se poser la question suivante :

connaissant une règle d'estimation précise, par exemple la minimisation de l'erreur quadratique moyenne, comment choisir  $\mathbf{p}$  pour que l'on ait  $\hat{\phi} = E[\phi]$ ?

Le problème est qu'il existe, en général, beaucoup de distributions qui satisfont cette contrainte. Il s'agit d'un problème inverse mal posé au sens où la solution n'est pas unique. Le principe du maximum d'entropie nous permet alors de choisir une solution.

Soit maintenant le cas plus général où nous considérons les  $m$  fonctions  $\{\phi_1(X), \phi_2(X), \dots, \phi_m(X)\}$  pour lesquelles nous avons un ensemble de données  $\{d_1, d_2, \dots, d_m\}$  pouvant s'exprimer sous la forme des  $m$  contraintes simultanées suivantes :

$$E[\phi_k(X)] = \sum_{i=1}^n p_i \phi_k(x_i) = d_k, \quad k = 1, \dots, m.$$

Dans chaque problème, ces données  $\{d_1, d_2, \dots, d_m\}$  peuvent avoir des interprétations physiques différentes et la difficulté consiste à incorporer ces données (contraintes) dans notre distribution de probabilité. Nous voulons ajuster la distribution de probabilité  $\{p_1, p_2, \dots, p_n\}$  à nos données. En général, le nombre des contraintes  $m$  est inférieur à  $n$  et il y a une infinité de solutions à ce problème. En d'autres termes, on peut trouver une infinité de distributions de probabilité  $\{p_1, p_2, \dots, p_n\}$  qui satisfont ces contraintes. C'est ici que le principe du maximum d'entropie est mis en œuvre pour choisir, parmi ces solutions possibles, celle qui a l'entropie maximale, c'est-à-dire la loi qui satisfait toutes les contraintes (toute l'information connue) et qui est la moins compromettante vis-à-vis de toute autre information inconnue. Le mot *information* devant être pris au sens de la définition de l'information moyenne ou de l'entropie de Shannon (équation 3).

Le problème se formule ainsi :

**Problème P1 :**

$$\begin{aligned} &\text{maximiser} \quad S(\mathbf{p}) = - \sum_{i=1}^n p_i \ln p_i, \\ &\text{sous les contraintes} \quad \sum_{i=1}^n p_i \phi_k(x_i) = d_k, \quad k = 1, \dots, m. \end{aligned}$$



La solution est obtenue par une technique variationnelle de multiplicateurs de Lagrange et est donnée par :

$$p_i = \frac{1}{Z(\lambda_1, \dots, \lambda_m)} \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x_i) \right], \quad i = 1, \dots, n, \quad (8.5)$$

où

$$Z(\lambda_1, \dots, \lambda_m) = \sum_{i=1}^n \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x_i) \right] \quad (8.6)$$

est la fonction de partition, et les  $\{\lambda_1, \lambda_2, \dots, \lambda_m\}$  sont déterminés par le système d'équations :

$$- \frac{\partial \ln Z(\lambda_1, \dots, \lambda_m)}{\partial \lambda_k} = d_k, \quad k = 1, \dots, m, \quad (8.7)$$

avec  $m$  données  $d_k$  et  $m$  inconnues  $\lambda_k$ . Notons que ce système d'équations n'est autre que le système d'équations des contraintes constituées par les données.

La valeur maximale de l'entropie est :

$$S_{\max} = \ln Z + \sum_{k=1}^m \lambda_k d_k.$$

Si on note  $\lambda_0 = \ln Z$  on peut vérifier facilement les propriétés suivantes :

$$- \frac{\partial \lambda_0}{\partial \lambda_k} = E[\phi_k(x_i)] = d_k, \quad (8.8)$$

$$- \frac{\partial^2 \lambda_0}{\partial \lambda_k^2} = E[\phi_k^2(x_i)] - d_k^2 = \text{Var}[\phi_k(x_i)], \quad (8.9)$$

$$- \frac{\partial^2 \lambda_0}{\partial \lambda_k \partial \lambda_l} = E[\phi_k(x_i) \phi_l(x_i)] - d_k d_l = \text{Cov}\{\phi_k(x_i), \phi_l(x_i)\}. \quad (8.10)$$

### 8.3.3 Lois à minimum d'entropie relative

L'extension de ce qui précède au cas de l'entropie croisée se formule ainsi :

#### Problème P2 :

Étant donné une distribution de probabilité *a priori*  $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$ , déterminer la distribution de probabilité *a posteriori*  $\mathbf{q} = \{q_1, q_2, \dots, q_n\}$  qui minimise

$$S(\mathbf{q}, \mathbf{p}) = - \sum_i q_i \ln \frac{p_i}{q_i} = \sum_i q_i \ln \frac{q_i}{p_i},$$

et qui satisfait les  $m$  contraintes

$$\sum_{i=1}^n q_i \phi_k(x_i) = d_k, \quad k = 1, \dots, m.$$

La solution est obtenue, là aussi, par une technique variationnelle de multiplicateurs de Lagrange et donnée par :

$$q_i = \frac{1}{Z} p_i \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x_i) \right], \quad i = 1, \dots, n, \quad (8.11)$$

où

$$Z(\lambda_1, \dots, \lambda_m) = \sum_{i=1}^n p_i \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x_i) \right] \quad (8.12)$$

est la fonction de partition, et les  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  sont déterminés par le système d'équations (9).

### 8.3.4 Extension au cas continu

Dans le cas continu, les définitions (8.3) et (8.4) conduisent à :

**Problème P3 :**

$$\begin{aligned} &\text{maximiser } S(p) = - \int p(x) \ln p(x) dx, \\ &\text{sous les contraintes } \int \phi_k(x) p(x) dx = d_k, \quad k = 1, \dots, m. \end{aligned}$$

La solution est

$$p(x) = \frac{1}{Z} \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x) \right], \quad (8.13)$$

où

$$Z(\lambda_1, \dots, \lambda_m) = \int \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x) \right] dx \quad (8.14)$$

est la fonction de partition, et les  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  sont déterminés par le système d'équations (9).

**Problème P4 :**

Étant donné une densité de probabilité *a priori*  $p(x)$ , déterminer la densité de probabilité *a posteriori*  $q(x)$  qui minimise l'entropie croisée

$$S(q, p) = - \int q(x) \ln \frac{p(x)}{q(x)} dx,$$

et qui satisfait les  $m$  contraintes

$$\int \phi_k(x) q(x) dx = d_k, \quad k = 1, \dots, m.$$

La solution est

$$q(x) = \frac{1}{Z} p(x) \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x) \right], \quad (8.15)$$

où

$$Z(\lambda_1, \dots, \lambda_m) = \int p(x) \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x) \right] dx \quad (8.16)$$

est la fonction de partition, et les  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  sont déterminés par (9).

Notons également que le vecteur  $\hat{\lambda} = (\lambda_1, \dots, \lambda_m)$  dans tous ces problèmes peut être considéré comme la solution du problème d'optimisation suivant :

**Problème dual des problèmes P1-P4 :**

$$\hat{\lambda} = \arg \min_{\lambda} \left\{ D(\lambda) = \lambda^t \mathbf{d} + \ln Z(\lambda) \right\} \quad (8.17)$$

où  $\mathbf{d} = (d_1, \dots, d_m)$ . Ce problème d'optimisation est appelé *problème dual* des problèmes P1-P4. Notons que l'expression de  $Z(\boldsymbol{\lambda})$  est différente suivant le problème.

Une présentation légèrement différente de ces relations peut être obtenue si on définit

$$\phi_0(x) = 1, \quad \text{et} \quad d_0 = 1,$$

ce qui permet d'inclure la contrainte de normalisation de  $q(x)$  et de formuler le problème précédent de la manière suivante :

$$\text{minimiser} \quad S(q, p) = - \int q(x) \ln \frac{p(x)}{q(x)} dx,$$

$$\text{sous les contraintes} \quad \int \phi_k(x) q(x) dx = d_k, \quad k = 0, \dots, m.$$

La solution peut être écrite sous la forme

$$q(x) = p(x) \exp \left[ - \sum_{k=0}^m \lambda_k \phi_k(x) \right], \quad (8.18)$$

et les  $\{\lambda_0, \lambda_1, \lambda_2, \dots, \lambda_n\}$  sont déterminés par :

$$G(\lambda_0, \lambda_1, \lambda_2, \dots, \lambda_n) = - \int \phi_k(x) p(x) \exp \left[ - \sum_{l=0}^m \lambda_l \phi_l(x) \right] dx = d_k, \quad k = 0, \dots, m. \quad (8.19)$$

On remarque que  $\lambda_0$  est relié à la fonction de partition  $Z$  par  $Z = - \exp[\lambda_0]$ , et on a :

$$\lambda_0 = - \ln Z = - \ln \int p(x) \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(x) \right] dx. \quad (8.20)$$

Les autres coefficients  $\lambda_k$  sont déterminés par le système d'équations :

$$- \frac{\partial}{\partial \lambda_k} \lambda_0(\lambda_1, \dots, \lambda_n) = d_k, \quad k = 1, \dots, m. \quad (8.21)$$

Notons aussi que la valeur minimale  $S_{\min}(q, p)$  peut être exprimée en fonction des  $\lambda_k$  par :

$$S_{\min}(q, p) = -\lambda_0 - \sum_{k=1}^m \lambda_k \phi_k(x). \quad (8.22)$$

Il n'est en général pas possible d'obtenir une relation explicite pour les coefficients  $\lambda_k$ . On résoud alors le système d'équations (8.19) numériquement par des méthodes itératives. Cependant, dans certaines situations simples on peut résoudre le problème d'une façon analytique [3, 6]. Le tableau 1 montre quelques exemples de lois, que l'on obtient sous forme analytique, en résolvant le problème P3 :

$$\text{maximiser} \quad S(p) = - \int p(x) \ln p(x) dx,$$

$$\text{sous les contraintes} \quad \int \phi_k(x) p(x) dx = d_k, \quad k = 1, \dots, m.$$

Tableau 1: Exemples de lois à maximum d'entropie.		
contraintes	domaine de $x$	loi de probabilité
$\phi_1(x) = x$	$x \in \mathbb{R}_+$	$p(x) = \frac{1}{\mu} \exp[-\mu x]$ loi exponentielle
$\phi_1(x) =  x $	$x \in \mathbb{R}$	$p(x) = \frac{1}{2\mu} \exp[-\mu x ]$ loi de LAPLACE
$\begin{cases} \phi_1(x) = x \\ \phi_2(x) = x^2 \end{cases}$	$x \in \mathbb{R}$	$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-m)^2}{2\sigma^2}\right]$ loi de Gausse
$\begin{cases} \phi_1(x) = x \\ \phi_2(x) = \ln x \end{cases}$	$x \in \mathbb{R}_+$	$p(x) \propto \exp[-\lambda \ln x - \mu x]$ $\propto x^{-\lambda} \exp[-\mu x]$ loi Gamma
$\begin{cases} \phi_1(x) = \ln x \\ \phi_2(x) = \ln(1-x) \end{cases}$	$x \in ]0, 1[$	$p(x) \propto \exp[-\lambda \ln x - \mu \ln(1-x)]$ $\propto x^{-\lambda} (1-x)^{-\mu}$ loi Béta
$\begin{cases} \phi_1(x) = \ln x \\ \phi_2(x) = x^2 \end{cases}$	$x > 0$	$p(x) \propto \exp[-\lambda \ln x - \mu x^2]$ $\propto x^{-\lambda} \exp[-\mu x^2]$ loi de Rayleigh

### 8.3.5 Extension au cas multivariable

Toutes ces relations peuvent être généralisées au cas multivariable. Par exemple le problème P4 devient :

#### Problème P5 :

Étant donné une densité de probabilité *a priori*  $p(\mathbf{x})$ , déterminer la densité de probabilité *a posteriori*  $q(\mathbf{x})$  qui minimise l'entropie croisée

$$S(q, p) = - \int q(\mathbf{x}) \ln \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x},$$

et qui satisfait les  $m$  contraintes

$$\int \phi_k(\mathbf{x}) q(\mathbf{x}) d\mathbf{x} = d_k, \quad k = 1, \dots, m.$$

La solution est

$$q(\mathbf{x}) = \frac{1}{Z} p(\mathbf{x}) \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(\mathbf{x}) \right], \quad (8.23)$$

où

$$Z(\lambda_1, \dots, \lambda_m) = \int p(\mathbf{x}) \exp \left[ - \sum_{k=1}^m \lambda_k \phi_k(\mathbf{x}) \right] d\mathbf{x} \quad (8.24)$$

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

$$-\frac{\partial \ln Z(\lambda_1, \dots, \lambda_m)}{\partial \lambda_k} = d_k, \quad k = 1, \dots, m. \quad (8.25)$$

Ici aussi, dans certains cas on peut obtenir une relation explicite pour la solution. Par exemple, si  $p(\mathbf{x})$  est une distribution exponentielle multivariée de la forme séparable

$$p(\mathbf{x}) = \prod_{i=1}^n \frac{1}{\alpha_i} \exp \left[ -\frac{x_i}{\alpha_i} \right],$$

et si les contraintes sont des contraintes sur les moments d'ordre un

$$E[X_i] = \int x_i q(x_i) dx_i = m_i, \quad i = 1, \dots, n,$$

alors la solution  $q(\mathbf{x})$  reste une fonction exponentielle multivariée séparable :

$$q(\mathbf{x}) = \prod_{i=1}^n \frac{1}{m_i} \exp \left[ -\frac{x_i}{m_i} \right].$$

De même si  $p(\mathbf{x})$  est une fonction gaussienne multivariée séparable :

$$p(\mathbf{x}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\alpha_i} \exp \left[ -\frac{1}{2} \left( \frac{x_i - m_i}{\alpha_i} \right)^2 \right],$$

et si les contraintes sont des contraintes du second ordre :

$$E[(X_i - m_i)^2] = \int (x_i - m_i)^2 q(x_i) dx_i = \sigma_i^2,$$

alors la solution  $q(\mathbf{x})$  reste une fonction gaussienne multivariée séparable :

$$q(\mathbf{x}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp \left[ -\frac{1}{2} \left( \frac{x_i - m_i}{\sigma_i} \right)^2 \right].$$

Ce dernier résultat peut être généralisé au cas où  $p(\mathbf{x})$  est une fonction gaussienne multivariée non dégénérée :

$$p(\mathbf{x}) = \frac{|\Sigma|^{-1/2}}{(2\pi)^{n/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mathbf{m})^t \Sigma^{-1} (\mathbf{x} - \mathbf{m}) \right],$$

et si les contraintes sont des contraintes sur les moments d'ordre un et deux non singulières ( $\Sigma'$  inversible) :

$$\begin{aligned} E[\mathbf{x}] &= \int \mathbf{x} q(\mathbf{x}) d\mathbf{x} = \mathbf{m}', \\ E[(\mathbf{x} - \mathbf{m}')(\mathbf{x} - \mathbf{m}')^t] &= \int (\mathbf{x} - \mathbf{m}')(\mathbf{x} - \mathbf{m}')^t q(\mathbf{x}) d\mathbf{x} = \Sigma', \end{aligned}$$

alors la solution  $q(\mathbf{x})$  reste une fonction gaussienne multivariée non dégénérée :

$$q(\mathbf{x}) = \frac{|\Sigma'|^{-1/2}}{(2\pi)^{n/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mathbf{m}')^t \Sigma'^{-1} (\mathbf{x} - \mathbf{m}') \right].$$

## 8.4 Utilisation pour les problèmes inverses

Arrivés à ce stade, nous avons les outils nécessaires pour comprendre comment le principe du maximum d'entropie peut être utilisé dans la résolution des problèmes inverses.

Il existe au moins deux catégories de méthodes qui utilisent le principe du maximum d'entropie pour la résolution des problèmes inverses :

### 8.4.1 Maximum d'entropie classique

Dans cette approche qui ne s'applique qu'aux objets positifs, l'idée de base consiste à assimiler l'objet  $\mathbf{f}$  à une distribution de probabilité et les données  $\mathbf{g}$  à des contraintes linéaires sur celle-ci. Sachant qu'en général, ces contraintes ne sont pas suffisantes pour définir une unique solution au problème, on utilise le PME pour choisir une solution parmi l'ensemble des solutions admissibles qui peuvent être définies soit par

$$\{\mathbf{f} : \mathbf{H}\mathbf{f} = \mathbf{g}\},$$

ou par

$$\{\mathbf{f} : \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 \leq c\}.$$

On choisit alors dans cet ensemble la solution qui maximise l'entropie

$$-\sum_{j=1}^N f_j \log f_j,$$

ou d'une manière plus générale celle qui minimise

$$-\sum_{j=1}^N \left[ f_j \log \left( \frac{f_j}{m_j} \right) + (f_j - m_j) \right],$$

où  $\mathbf{m}$  est une solution par défaut (*a priori*).

Notons que dans cette approche on assimile  $\mathbf{f}$  à une distribution de probabilité et les données  $\mathbf{g}$  à des contraintes linéaires sur celle-ci.

### 8.4.2 Maximum d'entropie sur la moyenne

Dans cette approche

- On fait l'hypothèse que  $\mathbf{f} \in \mathcal{C}$  où  $\mathcal{C}$  est un ensemble convexe muni d'une mesure de référence  $\mu(\mathbf{f})$  (ou d'une densité de probabilité *a priori*  $d\mu(\mathbf{f}) = p_0(\mathbf{f}) d\mathbf{f}$  et on suppose que

$$\begin{aligned} \mathbf{m} &= \int_{\mathcal{C}} \mathbf{f} d\mu(\mathbf{f}) \\ &= \int_{\mathcal{C}} \mathbf{f} p_0(\mathbf{f}) d\mathbf{f}. \end{aligned}$$

- On considère  $\mathbf{f}$  comme la moyenne d'un vecteur aléatoire  $\mathbf{F}$  pour lequel on définit une loi de probabilité  $P(\mathbf{f})$  (ou une densité de probabilité  $dP(\mathbf{f}) = p(\mathbf{f}) d\mathbf{f}$ :

$$\begin{aligned} \mathbf{f} = \mathbf{E}[\mathbf{F}] &= \int_{\mathcal{C}} \mathbf{f} dP(\mathbf{f}) \\ &= \int_{\mathcal{C}} \mathbf{f} p(\mathbf{f}) d\mathbf{f}, \end{aligned}$$

et les données  $\mathbf{g} = \mathbf{H}\mathbf{f}$  comme des contraintes linéaires sur la loi de probabilité  $p(\mathbf{f})$  :

$$\begin{aligned}\mathbf{g} = \mathbf{H}\mathbf{f} = \mathbf{H}\mathbf{E}[\mathbf{F}] = \mathbf{E}[\mathbf{H}\mathbf{F}] &= \int_{\mathcal{C}} \mathbf{H}\mathbf{f} \, dP(\mathbf{f}) \\ &= \int_{\mathcal{C}} \mathbf{H}\mathbf{f} p(\mathbf{f}) \, d\mathbf{f}.\end{aligned}$$

- On cherche alors parmi l'ensemble des lois de probabilités satisfaisant ces contraintes celle qui minimise la distance de Kulback entre  $dP(\mathbf{f})$  et  $d\mu(\mathbf{f})$  :

$$- \int \log \frac{dP(\mathbf{f})}{d\mu(\mathbf{f})} dP(\mathbf{f}),$$

ou, d'une manière équivalente, celle qui minimise l'entropie relative

$$- \int p(\mathbf{f}) \log \left( \frac{p(\mathbf{f})}{p_0(\mathbf{f})} \right) d\mathbf{f}.$$

La solution, nous l'avons vu, peut être obtenue via le Lagrangien :

$$\begin{aligned}\mathcal{L}(\mathbf{f}, \boldsymbol{\lambda}) &= \int_{\mathcal{C}} \left[ \log \frac{dP(\mathbf{f})}{d\mu(\mathbf{f})} - \sum_{i=1}^M \lambda_i (g_i - [\mathbf{H}\mathbf{f}]_i) \right] dP(\mathbf{f}) \\ &= \int_{\mathcal{C}} \left[ \log \frac{dP(\mathbf{f})}{d\mu(\mathbf{f})} - \boldsymbol{\lambda}^t (\mathbf{g} - \mathbf{H}\mathbf{f}) \right] dP(\mathbf{f}) \\ &= \int_{\mathcal{C}} \left[ \log \frac{p(\mathbf{f})}{p_0(\mathbf{f})} - \boldsymbol{\lambda}^t (\mathbf{g} - \mathbf{H}\mathbf{f}) \right] p(\mathbf{f}) \, d\mathbf{f}\end{aligned}$$

et donnée par

$$dP(\mathbf{f}, \boldsymbol{\lambda}) = \exp \left[ \boldsymbol{\lambda}^t [\mathbf{H}\mathbf{f}] - \log Z(\boldsymbol{\lambda}) \right] d\mu(\mathbf{f}),$$

ou d'une manière équivalente

$$p(\mathbf{f}, \boldsymbol{\lambda}) = \exp \left[ \boldsymbol{\lambda}^t [\mathbf{H}\mathbf{f}] - \log Z(\boldsymbol{\lambda}) \right] p_0(\mathbf{f}) \, d\mathbf{f},$$

où

$$Z(\boldsymbol{\lambda}) = \int_{\mathcal{C}} \exp \left[ \boldsymbol{\lambda}^t [\mathbf{H}\mathbf{f}] \right] d\mu(\mathbf{f}) = \int_{\mathcal{C}} \exp \left[ \boldsymbol{\lambda}^t [\mathbf{H}\mathbf{f}] \right] p_0(\mathbf{f}) \, d\mathbf{f}.$$

Les paramètres de Lagrange sont obtenus en cherchant la solution du système d'équations non linéaires suivant :

$$\frac{\partial \log Z(\boldsymbol{\lambda})}{\partial \lambda_i} = g_i, \quad i = 1, \dots, M.$$

- Une fois  $p(\mathbf{f})$  déterminée on définit la solution du problème par

$$\hat{\mathbf{f}} = \mathbf{E}[\mathbf{f}] = \int \mathbf{f} \, dP(\mathbf{f}) = \int \mathbf{f} p(\mathbf{f}) \, d\mathbf{f}.$$

On peut montrer que, cette solution  $\hat{\mathbf{f}}$ , qui dépend ainsi de  $\boldsymbol{\lambda}$ , peut être obtenue directement en minimisant un critère  $\Omega(\mathbf{f}, \mathbf{m})$  sous les contraintes  $\mathbf{g} = \mathbf{H}\mathbf{f}$ .

En effet, en notant

$$\mathbf{s} = \mathbf{H}^t \boldsymbol{\lambda}, \quad G^*(\mathbf{s}) = \log Z(\mathbf{s}) = \log \int_{\mathcal{C}} \exp[\mathbf{s}^t \mathbf{f}] \, d\mu(\mathbf{f}),$$

et

$$\Omega(\mathbf{f}) = \inf_{\mathbf{s}} \{\mathbf{s}^t \mathbf{f} - G^*(\mathbf{s})\}, \quad D(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^t \mathbf{g} - G^*(\mathbf{H}^t \boldsymbol{\lambda}),$$

on peut alors montrer que :

- Le vecteur de  $\boldsymbol{\lambda}$  optimal peut être calculé en cherchant l'optimum de  $D(\boldsymbol{\lambda})$  :

$$\hat{\boldsymbol{\lambda}} = \arg \max_{\boldsymbol{\lambda}} \{D(\boldsymbol{\lambda})\}$$

Le critère  $D(\boldsymbol{\lambda})$  est appelé *Critère dual* ;

- La solution du problème  $\hat{\mathbf{f}}$  peut être calculée en cherchant l'optimum d'un critère  $\Omega(\mathbf{f}, \mathbf{m})$  appelé *critère primal* :

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f} \in \mathcal{C}} \{H(\mathbf{f}, \mathbf{m})\} \quad \text{sous les contraintes} \quad \mathbf{g} = \mathbf{H}\mathbf{f};$$

- Il existe une relation explicite pour calculer la solution en fonction de  $G^*(\mathbf{s})$  :

$$\hat{\mathbf{f}}(\mathbf{s}) = \frac{dG^*(\mathbf{s})}{d\mathbf{s}},$$

et où :

- les fonctions  $G$  et  $\Omega$  dépendent de la mesure de référence  $\mu(\mathbf{f})$  ;
- $D(\boldsymbol{\lambda})$  est appelé le *critère dual* et dépend des données  $\mathbf{g}$  et de la fonction  $G$  ;
- $\Omega(\mathbf{f}, \mathbf{m})$  est appelé le *critère primal* et peut être considéré comme une mesure de distance entre  $\mathbf{f}$  et  $\mathbf{m}$ , ce qui signifie qu'elle satisfait les conditions suivantes :
  - $\Omega(\mathbf{f}, \mathbf{m}) \geq 0$ , et  $\Omega(\mathbf{f}, \mathbf{m}) = 0$  ssi  $\mathbf{f} = \mathbf{m}$  ;
  - $\Omega(\mathbf{f}, \mathbf{m})$  est continue et convexe sur  $\mathcal{C}$  ;
  - $\Omega(\mathbf{f}, \mathbf{m}) = \infty$  si  $\mathbf{f} \notin \mathcal{C}$ .

Lorsque la mesure de référence est séparable

$$\mu(\mathbf{f}) = \prod_{j=1}^N \mu_j(f_j),$$

on a aussi :

$$dP(\mathbf{f}, \boldsymbol{\lambda}) = \prod_{j=1}^N dP_j(f_j, \lambda_j),$$

et

$$G(\mathbf{s}) = \sum_j g_j(s_j), \quad \Omega(\mathbf{f}, \mathbf{m}) = \sum_j h_j(f_j, m_j), \quad \hat{x}_j = g'_j(s_j).$$



En remplaçant  $\mathbf{s} = \mathbf{H}^t \boldsymbol{\lambda}$  on obtient :

$$G(\boldsymbol{\lambda}) = \sum_j g_j \left( [\mathbf{H}^t \boldsymbol{\lambda}]_j \right), \quad \Omega(\mathbf{f}, \mathbf{m}) = \sum_j h_j(f_j, m_j), \quad \hat{x}_j = g'_j \left( [\mathbf{H}^t \hat{\boldsymbol{\lambda}}]_j \right).$$

On peut voir que  $\Omega(\mathbf{f}, \mathbf{m})$  est aussi séparable et que  $h_j$  et  $g_j$  dépendent de la mesure de référence  $\mu_j(f_j)$ :

–  $g_j$  est la transformée de Cramer (log de la transformée de LAPLACE) de  $\mu_j$ :

$$g_j(s) = \log \int \exp [s f_j] d\mu_j(f_j);$$

–  $h_j$  est la convexe conjuguée de  $g_j$ :  $h_j(f) = \inf_s \{s f - g_j(s)\}$ .

Prenons quelques exemples :

	$\mu_j(f)$	$g_j(s)$	$h_j(f, m)$
Gaussien :	$\exp \left[ -\frac{1}{2} (f - m)^2 \right]$	$\frac{1}{2} (s - m)^2$	$\frac{1}{2} (f - m)^2$
Poisson :	$\frac{m^f}{f!} \exp [-m]$	$\exp [m - s]$	$-f \log \frac{f}{m} + m - f$
Gamma :	$f^{\alpha-1} \exp \left[ -\frac{f}{m} \right]$	$\log (s - m)$	$-\log \frac{f}{m} + \frac{f}{m} - 1$

Ainsi, si on considère le critère primal :

$$\text{maximiser } \Omega(\mathbf{f}, \mathbf{m}) = \sum_{j=1}^N h_j(f_j, m_j) \text{ sous contraintes } \mathbf{H} \mathbf{f} = \mathbf{g},$$

on peut par exemple remarquer que lorsque la mesure de référence est une mesure poissonnienne on a  $h_j(f_j, m_j) = -f_j \ln f_j/m_j + (m_j - f_j)$ . On peut ainsi faire le parallèle avec la méthode du ME classique.

Cependant, lorsque la mesure de référence n'est pas séparable ou, lorsqu'il y a de l'incertitude sur les données  $\mathbf{g}$ , il n'est plus aussi facile d'utiliser cette approche. Mentionnons quand même que cette approche est récente, et des travaux tout à fait récents ont essayé de prendre en compte le bruit. Ici nous en présentons une approche qui consiste à remplacer  $\mathbf{g} = \mathbf{H} \mathbf{f} + \mathbf{n}$  par

$$\mathbf{g} = [\mathbf{H} | \mathbf{I}] \begin{bmatrix} \mathbf{f} \\ \mathbf{n} \end{bmatrix} \longrightarrow \mathbf{g} = \tilde{\mathbf{H}} \tilde{\mathbf{f}}$$

et faire l'hypothèse d'indépendance entre  $\mathbf{b}$  et  $\mathbf{f}$  :

$$\mu(\tilde{\mathbf{f}}) = \mu_x(\mathbf{f}) \mu_n(\mathbf{n})$$

Ensuite, suivant l'approche, on montre que la solution de l'optimisation du critère primal devient :

$$\hat{\tilde{\mathbf{f}}} = \arg \min_{\tilde{\mathbf{f}} \in \mathcal{C}} \{ \mathcal{Q}(\mathbf{g} - \mathbf{H} \tilde{\mathbf{f}}) + \alpha \Omega(\tilde{\mathbf{f}}, \mathbf{m}) \}$$

avec

$$\Omega(\mathbf{f}, \mathbf{m}) = \sum_{j=1}^N h_j(f_j, m_j), \quad \text{et} \quad \mathcal{Q}(\mathbf{z}) = \sum_{i=1}^M q_i(z_i).$$

Là encore, les fonctions  $h_j(f_j)$  et  $q_i(z_i)$  dépendent des mesures de références  $\mu_f(\mathbf{f})$  et  $\mu_n(\mathbf{n})$ .

