

École Supérieure d'Optique

Traitement du signal avancé

Théorie logique des probabilités

Théorie de la décision

Estimation de paramètres

Filtrage de Kalman

Table des matières

1	Vecteurs aléatoires	7
1.1	Axiomes	7
1.2	Variables aléatoires	8
1.3	Vecteurs aléatoires	9
1.4	Vecteurs aléatoires gaussiens	11
2	Théorie logique des probabilités	15
2.1	Introduction	15
2.2	Les postulats de la théorie	16
2.3	Le principe d'indifférence	21
2.4	Test d'hypothèse	25
2.4.1	Notations et problématique	25
2.4.2	Méthode générale de résolution	25
2.4.3	Critère de sélection d'une hypothèse	26
2.4.4	Cas particulier : l'alternative	26
2.4.5	Test d'hypothèses multiples	29
2.4.6	Test d'un nombre infini d'hypothèses	31
2.5	Principe du maximum d'entropie	33
2.5.1	Entropie de Shannon	33
2.5.2	Méthode de Wallis	33
2.5.3	Généralisation : fonction de partition	35
2.6	Annexe : méthode des multiplicateurs de Lagrange	37
3	Théorie de la décision	41
3.1	Introduction	41
3.2	Définitions	41
3.3	Cas particulier : détection	42
3.4	Théorie de la décision : le point de vue classique	43
3.5	Stratégie de Neyman-Pearson	46
3.6	Exemple : détection de bits	49
4	Estimation de paramètres	51
4.1	Introduction – Problématique	51
4.1.1	Approche classique	51
4.1.2	Analogie avec le test d'hypothèse	52
4.2	Signal dans un bruit gaussien	52
4.3	Fonction d'ambiguïté	53
4.3.1	Moyenne et covariance de l'estimation des moindres carrés	54

4.3.2	Relation à la fonction d'ambiguïté	56
4.3.3	Fonction d'ambiguïté modifiée	57
4.4	Mesure d'un temps	57
4.4.1	Mesure avec un signal quelconque	57
4.4.2	Mesure avec un signal passe-bande: enveloppe complexe	59
4.4.3	Conception de signaux adaptés à la mesure d'un temps	61
4.5	Mesure d'une fréquence	63
4.6	Inégalité de Cramer-Rao	64
4.7	Ambiguïté temps-fréquence	68
4.7.1	Définition	68
4.7.2	Estimation conjointe temps-fréquence	70
4.8	Exemples de fonctions d'ambiguïté temps-fréquence	72
4.8.1	Impulsion gaussienne	72
4.8.2	Impulsion gaussienne "chirp"	74
4.8.3	Compression d'impulsion	75
4.8.4	Train d'impulsions	76
4.9	Inégalité de Fréchet-Darmonis-Cramer-Rao	77
5	Filtrage de Kalman	83
5.1	Problématique	83
5.1.1	Définitions	83
5.1.2	Estimation – prédiction	84
5.1.3	Modèle de mesure	85
5.1.4	Modèle d'évolution	86
5.1.5	Modèles linéaires de mesure et d'évolution	88
5.1.6	Espérance conditionnelle	89
5.1.7	Choix de l'estimateur	91
5.2	Solution générale	92
5.3	Solution gaussienne – cas linéaire	94
5.3.1	Propagation du caractère gaussien	94
5.3.2	Calcul du gain et de la covariance	96
5.3.3	Résumé et commentaires	98
5.4	Solution des moindres carrés	98
5.4.1	Prédiction de l'état	100
5.4.2	Prédiction de la mesure	101
5.4.3	Estimation de l'état	102
5.5	Solution approchée du cas général	105
5.5.1	Modèle linéaire généralisé	105
5.5.2	Approximation au premier ordre du modèle	106
5.6	Compléments	109

Avertissement

Le but de ce cours est de présenter une introduction à l'application des méthodes du traitement du signal, et en particulier de celles qui font appel aux probabilités, à l'étude des systèmes physiques. Les applications concernées sont par exemple le radar, le sonar, le lidar et de façon générale tous les systèmes de mesure, actifs ou non.

En physique, on formule des modèles permettant de décrire mathématiquement le comportement d'un système. Ces modèles incluent le plus souvent des paramètres inconnus dont il s'agit de déterminer la valeur, on parle alors d'estimation, ou dépendent d'hypothèses entre lesquelles il faut trancher.

Il est essentiel de remarquer que le résultat de l'estimation d'un paramètre ou de la sélection d'une hypothèse dépend en partie de l'observateur lui-même, c'est-à-dire de celui qui pose la question. En effet, la description que l'observateur peut faire du système physique (et non le système physique lui-même) dépend de la connaissance dont il dispose. C'est pourquoi nous décrirons la connaissance que l'observateur a d'un système sous la forme d'une probabilité qui sera toujours conditionnelle à l'information disponible. Cette description ne sera bien sûr pas complète, mais sera la meilleure que l'observateur puisse faire sachant ce qu'il sait.

La méthode scientifique nous enseigne que l'expérience est seule juge de la qualité d'un modèle. C'est pourquoi les méthodes d'estimation d'un paramètre ou de sélection d'une hypothèse que nous décrirons consisteront principalement en une comparaison probabiliste de ce que prévoit le modèle et de ce qui est effectivement mesuré.

Dans le premier chapitre, nous rappelons certaines propriétés concernant les variables et les **vecteurs aléatoires** qui seront utiles pour la suite. Le second chapitre présente une introduction à la **théorie logique des probabilités** qui permet de décrire l'information du point de vue de l'observateur d'un système. Le troisième chapitre expose les principes de la **théorie de la décision**, qui tente de répondre à une question difficile : comment agir de façon optimale dans des situations décrites en termes de probabilités. Dans le quatrième chapitre, nous aborderons la théorie de l'**estimation de paramètres**, qui fournit une technique puissante pour la conception de systèmes de mesures physiques (radar, sonar, impulsions lumineuses...). Finalement, le cinquième chapitre proposera un cas particulier important de l'estimation de paramètre, le **filtrage de Kalman**, qui est conçu pour des systèmes dynamiques (dont les paramètres évoluent au cours du temps).

Chapitre 1

Vecteurs aléatoires

1.1 Axiomes

La théorie classique des probabilités a été établie à partir des axiomes dits de Kolmogorov, et est généralement désignée sous ce nom. C'est celle qui est enseignée quasi-exclusivement en France, et qui fait l'objet des rappels de ce chapitre. Il faut cependant être conscient que ce n'est pas la seule définition possible d'une théorie des probabilités.

Définition 1.1

(Espace des épreuves)

Le triplet (Ω, \mathcal{B}, P) constitue une expérience, où :

- Ω est l'espace des épreuves. Un élément $\omega \in \Omega$ est une épreuve (en physique, le résultat d'une mesure ou d'une expérience).
- $\mathcal{B} = \mathcal{P}(\Omega)$ est l'ensemble des parties de Ω , c'est-à-dire des événements (\mathcal{B} est donc l'ensemble des questions que l'on peut se poser sur le résultat d'une expérience). \mathcal{B} est une tribu borélienne.
- P est une mesure sur (Ω, \mathcal{B}) vérifiant les axiomes suivants :

$$(A1) \quad p(\Omega) = 1$$

$$(A2) \quad \forall A, B \in \mathcal{B}, \text{ tels que } A \cap B = \emptyset, P(A \cup B) = P(A) + P(B)$$

$$(A3) \quad \forall A_i \in \mathcal{B}, i \in \mathbb{N}, \text{ tels que } A_i \cap A_j = \emptyset, P\left(\bigcup_{i=0}^{\infty} A_i\right) = \sum_{i=0}^{\infty} P(A_i)$$

Par exemple pour un jeu de dé, $\Omega = \{1, 2, 3, 4, 5, 6\}$. Les épreuves $\omega \in \Omega$ sont $\{1\}$, $\{2\}$, $\{3\}$, $\{4\}$, $\{5\}$ et $\{6\}$. Les événements sont des parties de Ω du type $\{2, 3, 5\}$ ou encore $\{1, 4, 5, 6\}$...

Axiome 1.2

(Axiome de la probabilité conditionnelle)

$$P(A|B) \doteq \frac{P(A \cap B)}{P(B)}$$

Cet axiome n'est pas présent dans toutes les versions de la théorie des probabilités dites de Kolmogorov, suivant les auteurs. La notion de probabilité conditionnelle est cependant fondamentale en physique, et sera une notion centrale de ce cours.

Axiome 1.3**(Fréquence d'apparition)**

La probabilité d'un événement s'interprète physiquement comme la fréquence d'apparition de cet événement. Si sur N épreuves l'événement A apparaît N_A fois, alors :

$$P(A) \doteq \lim_{N \rightarrow \infty} \frac{N_A}{N}$$

Cet axiome peut paraître une définition parfaitement intuitive de la mesure d'une probabilité à partir de mesures physiques. Il contient cependant une difficulté conceptuelle très forte qui est liée au passage à la limite. Pour vérifier qu'une pièce de monnaie donne bien pile ou face avec une probabilité $1/2$, il suffit suivant cet axiome de jeter la pièce un nombre très grand de fois et de mesurer les fréquences de pile et face. Mais quand doit-on s'arrêter? Est-ce que dans un cas plus complexe la probabilité que l'on cherche à mesurer ne va pas changer au cours du temps, donc au cours des expériences? En fait, il est facile de se convaincre qu'en physique il n'est pas possible de réaliser une mesure un nombre infini de fois en un temps fini. Nous reviendrons sur ces difficultés au chapitre suivant.

1.2 Variables aléatoires**Définition 1.4****(Variable aléatoire)**

Une variable aléatoire est une fonction X définie de l'espace des épreuves Ω dans \mathbb{R} ,

$$X \begin{cases} \Omega & \rightarrow \mathbb{R} \\ \omega & \mapsto X(\omega) \end{cases}$$

telle que $\forall x \in \mathbb{R}$, l'inégalité $X(\omega) \leq x$ définisse une partie de \mathcal{B} notée $(X(\omega) \leq x)$ (cette partie est donc $X^{-1}(]-\infty, x])$).

Définition 1.5**(Fonction de répartition)**

La fonction de répartition est définie à partir de la mesure P sur (Ω, \mathcal{B}) par

$$F_X(x) \doteq P(X(\omega) \leq x), \forall x \in \mathbb{R}$$

Définition 1.6**(Densité de probabilité)**

$$p_X(x) \doteq \frac{dF_X(x)}{dx}$$

$$\left\{ \begin{array}{l} p_X(x) dx = P(x \leq X(\omega) < x + dx) \\ p_X(x) \geq 0, \forall x \in \mathbb{R} \\ \int_{\mathbb{R}} p_X(x) dx = 1 \end{array} \right.$$

Définition 1.7

(Espérance)

Soit f une fonction d'une variable réelle. L'espérance de f est définie par

$$\langle f(X) \rangle \doteq \int f(X(\omega)) dP(x \leq X(\omega) < x + dx) = \int_{\mathbb{R}} f(x) p_X(x) dx$$

La première intégrale est définie au sens de la mesure, la seconde au sens d'une variable réelle.

Cette opération est linéaire, puisque :

$$\langle af(X) + bg(X) \rangle = a \langle f(X) \rangle + b \langle g(X) \rangle$$

de façon évidente.

Cas particuliers :

– Moyenne :

$$\langle X \rangle = m_X = \int_{\mathbb{R}} x p_X(x) dx$$

– Moments :

$$\langle X^n \rangle = \int_{\mathbb{R}} x^n p_X(x) dx$$

– Moments centrés :

$$\langle (X - \langle X \rangle)^n \rangle = \int_{\mathbb{R}} (x - \langle X \rangle)^n p_X(x) dx$$

– Fonction caractéristique :

$$\varphi_X(u) = \langle \exp(iuX) \rangle = \int_{\mathbb{R}} \exp(iux) p_X(x) dx$$

avec la relation

$$\langle X^n \rangle = \left(\frac{1}{i^n} \right) \frac{d^n}{du^n} \varphi_X(0)$$

1.3 Vecteurs aléatoires**Définition 1.8**

(Vecteur aléatoire)

Un vecteur aléatoire à n dimensions est une fonction vectorielle de Ω dans \mathbb{R}^n ,

$$X \left| \begin{array}{l} \Omega \rightarrow \mathbb{R}^n \\ \omega \mapsto X(\omega) = \begin{pmatrix} X_1(\omega) \\ \vdots \\ X_n(\omega) \end{pmatrix} \end{array} \right.$$

telle que $\forall x \in \mathbb{R}^n$, les n inégalités $X_k(\omega) \leq x_k$, $k \in [1, n]$ définissent une partie de \mathcal{B} notée $(X(\omega) \leq x)$.

Définition 1.9

(Fonction de répartition)

$$F_X(x) \doteq P(X(\omega) \leq x) = P(\{X_1(\omega) \leq x_1\} \cap \dots \cap \{X_n(\omega) \leq x_n\}), \forall x \in \mathbb{R}^n$$

avec $x^T = (x_1, \dots, x_n)$.

Définition 1.10

(Densité de probabilité)

$$p_X(x) \doteq \frac{\partial^n F_X(x)}{\partial x_1 \cdots \partial x_n}$$

Les variables aléatoires X_1, \dots, X_n sont mutuellement indépendantes si et seulement si :

$$p_X(x) = \prod_{k=1}^n p_{X_k}(x_k)$$

Notation :

$$\int_{\mathbb{R}^n} p_X(x) dx \doteq \underbrace{\int_{\mathbb{R}} \cdots \int_{\mathbb{R}}}_{n \text{ fois}} p_X(x_1, \dots, x_n) dx_1 \cdots dx_n$$

Densité marginale (par exemple) :

$$p_{X_k}(x_k) = \underbrace{\int_{\mathbb{R}} \cdots \int_{\mathbb{R}}}_{n-1 \text{ fois}} p_X(x_1, \dots, x_n) dx_1 \cdots dx_{k-1} dx_{k+1} \cdots dx_n$$

EXERCICE - 1.1 _____ Vecteur aléatoire uniformément distribué

Soit le vecteur aléatoire

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \in \mathbb{R}^2$$

On suppose $F_X(x) = x_1 \cdot x_2$ avec $x_1, x_2 \in [0, 1]$. Vérifier que

$$\begin{cases} F_X(0, 0) = 0 \\ F_X(1, 1) = 1 \end{cases},$$

que F_X est croissante, et que $p_X(x) = p_{X_1}(x_1) = p_{X_2}(x_2) = 1$ sur $[0, 1]^2$.

Définition 1.11

(Espérance)

Soit f une fonction vectorielle. Par définition :

$$\begin{aligned} \langle f(X) \rangle &\doteq \int f(X(\omega)) dP(x \leq X(\omega) < x + dx) \\ &= \int_{\mathbb{R}^n} f(x) p_X(x) dx \\ &= \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} f(x_1, \dots, x_n) p_X(x_1, \dots, x_n) dx_1 \cdots dx_n \end{aligned}$$

– Vecteur moyenne :

$$\langle X \rangle = m_X = \begin{pmatrix} \langle X_1 \rangle \\ \vdots \\ \langle X_n \rangle \end{pmatrix}$$

– Matrice de covariance :

$$\begin{aligned}\Gamma_{XX} &= \langle X.X^T \rangle - \langle X \rangle . \langle X \rangle^T \\ &= \begin{pmatrix} \langle X_1^2 \rangle - \langle X_1 \rangle^2 & \dots & \langle X_1.X_n \rangle - \langle X_1 \rangle . \langle X_n \rangle \\ \vdots & & \vdots \\ \langle X_1.X_n \rangle - \langle X_1 \rangle . \langle X_n \rangle & \dots & \langle X_n^2 \rangle - \langle X_n \rangle^2 \end{pmatrix}\end{aligned}$$

– Fonction caractéristique :

$$\varphi_X(u) \doteq \langle \exp(iu^T.X) \rangle = \int_{\mathbb{R}^n} \exp\left(i \sum_{k=1}^n u_k x_k\right) p_X(x) dx$$

avec $u^T = (u_1, \dots, u_n)$.

Propriété 1.1

(Factorisation de la fonction caractéristique)

Si les variables aléatoires X_1, \dots, X_n sont indépendantes, alors

$$\varphi_X(u) = \prod_{k=1}^n \varphi_{X_k}(u_k)$$

EXERCICE - 1.2 _____ Vecteur aléatoire uniformément distribué (suite)

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \in \mathbb{R}^2$$

On suppose $F_X(x) = x_1.x_2$ avec $x_1, x_2 \in [0,1]$. Montrer que :

$$\begin{aligned}\langle X \rangle^T &= (1/2, 1/2) \\ \Gamma_{XX} &= \begin{pmatrix} \frac{1}{12} & 0 \\ 0 & \frac{1}{12} \end{pmatrix} \\ \varphi_X(u) &= \exp\left(i \frac{u_1 + u_2}{2}\right) \operatorname{sinc}\left(\frac{u_1}{2\pi}\right) \operatorname{sinc}\left(\frac{u_2}{2\pi}\right)\end{aligned}$$

1.4 Vecteurs aléatoires gaussiens

Définition 1.12

(Variable aléatoire gaussienne)

Une variable aléatoire gaussienne admet pour densité de probabilité :

$$p_X(x) = \mathcal{N}(x, m, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x - m)^2}{\sigma^2}\right)$$

$$\begin{aligned}\langle X \rangle &= m_X = m \\ \langle (X - \langle X \rangle)^2 \rangle &= \sigma^2 \\ \varphi_X(u) &= \exp\left(ium - \frac{1}{2}u^2\sigma^2\right)\end{aligned}$$

Définition 1.13**(Vecteur aléatoire gaussien)**

Un vecteur aléatoire gaussien admet pour densité de probabilité :

$$p_X(x) = \mathcal{N}(x, m, \Sigma) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(x - m)^T \Sigma^{-1} (x - m)\right)$$

avec

$$\begin{aligned}\langle X \rangle &= m \quad (\text{vecteur moyenne}) \\ \Sigma &= \text{matrice de covariance} \\ |\Sigma| &= \text{déterminant de } \Sigma \\ \varphi_X(u) &= \exp\left(iu^T \cdot m - \frac{1}{2}u^T \cdot \Sigma \cdot u\right)\end{aligned}$$

Σ est diagonale si les X_k sont soit décorrélés, soit indépendants.

Les vecteurs aléatoires gaussiens jouent un rôle très important car ils se rencontrent très souvent à cause de diverses propriétés qu'ils possèdent.

Propriété 1.2**(Théorème limite central)**

Soient $(X_k)_{k=1,2,\dots,n,\dots}$ des variables aléatoires indépendantes telles que

$$\begin{aligned}\langle X_k \rangle &= 0 \\ \lim_{n \rightarrow \infty} \sum_{k=1}^n \langle X_k^2 \rangle &= \lim_{n \rightarrow \infty} \sum_{k=1}^n \sigma_k^2 = \sigma^2 < +\infty\end{aligned}$$

alors la variable aléatoire $Y = \sum_{k=1}^{\infty} X_k$ est $\mathcal{N}(y, 0, \sigma^2)$.

Propriété 1.3**(Propriété de stabilité gaussienne)**

- Le produit de deux gaussiennes est une gaussienne.
- La convolution de deux gaussiennes est une gaussienne.
- La transformée de Fourier d'une gaussienne est une gaussienne.

Propriété 1.4**(Addition de deux gaussiennes)**

Soit $Z = X + Y$ avec X et Y suivant les lois normales $\mathcal{N}(x, m_X, \Sigma_X)$ et $\mathcal{N}(y, m_Y, \Sigma_Y)$.
Si les vecteurs aléatoires X et Y sont indépendants alors Z suit une loi normale $\mathcal{N}(z = x + y, m_X + m_Y, \Sigma_X + \Sigma_Y)$

Propriété 1.5**(Transformation linéaire d'une gaussienne)**

Soit $Z = C.X + a$ avec X (à m composantes) suivant la loi normale $\mathcal{N}(x, m_X, \Sigma_X)$, C une matrice rectangulaire de taille $n \times m$, et a un vecteur de taille n , alors Z (à n composantes) suit une loi normale $\mathcal{N}(z = C.x + a, C.m_X + a, C.\Sigma_X.C^T)$.

EXEMPLE - 1.1 ————— Le problème de Herschel – 1850

Comme astronome, John Herschel se posa la question de la précision de localisation d'une étoile, c'est-à-dire de trouver la densité de probabilité conjointe $P_{XY}(x, y)$ où x et y sont les erreurs de mesure dans deux directions orthogonales (Nord-Sud et Est-Ouest par exemple). Il émit deux postulats :

- (P1) La probabilité d'erreur est indépendante de la direction.
- (P2) Les erreurs dans deux directions orthogonales sont indépendantes.

De (P1) on déduit en passant en coordonnées polaires

$$p_{XY}(x, y) \, dx dy = f(r) \, r dr d\theta = f(r) \, dx dy$$

De (P2) on déduit que

$$p_{XY}(x, y) \, dx dy = p_X(x) dx \, p_Y(y) dy$$

De plus

$$p_X(x) = \int p_{XY}(x, y) \, dy = \int f(r) dy = \int p_{XY}(y, x) \, dy = p_Y(x)$$

car $r = \sqrt{x^2 + y^2}$. En posant $p_X(x) = g(x)$, on a donc $p_Y(y) = g(y)$, et

$$g(x)g(y) = f(r) = f\left(\sqrt{x^2 + y^2}\right)$$

En faisant $y = 0$, on a $g(x)g(0) = f(x)$, donc

$$g(x)g(y) = g(r)g(0)$$

on en déduit

$$\ln\left(\frac{g(x)}{g(0)}\right) + \ln\left(\frac{g(y)}{g(0)}\right) = \ln\left(\frac{g\left(\sqrt{x^2 + y^2}\right)}{g(0)}\right)$$

d'où en dérivant par rapport à x

$$\frac{1}{x} \frac{d}{dx} \ln(g(x)) = \frac{1}{\sqrt{x^2 + y^2}} \frac{d}{dx} \ln\left(g\left(\sqrt{x^2 + y^2}\right)\right)$$

cette égalité doit être vraie quel que soit y , d'où

$$\frac{1}{x} \frac{d}{dx} \ln(g(x)) = a$$

où a est une constante. Il vient donc

$$\ln(g(x)) = \frac{a}{2} x^2 + b$$

soit encore $g(x) = \mathcal{N}(x, 0, \sigma^2)$ en posant $\sigma^2 = -1/a$ et après normalisation, et

$$P_{XY}(x, y) = \mathcal{N}(\sqrt{x^2 + y^2}, 0, \sigma^2)$$

Chapitre 2

Théorie logique des probabilités

2.1 Introduction

Les variables et vecteurs aléatoires constituent les fondements de la théorie mathématique des probabilités, construite à partir des axiomes de Kolmogorov. Cette théorie est extrêmement développée et sûre.

Un problème se pose cependant quand on tente d'appliquer cette théorie en physique. Le problème n'est pas tant de calculer avec les probabilités que de leur assigner des valeurs préalable au calcul.

En théorie classique, dite "fréquentiste", des probabilités, on attribue un caractère "aléatoire" aux phénomènes physiques à travers la notion de fréquence. Pour être définie, une "fréquence" doit être mesurée à partir d'une expérience renouvelée un nombre infini de fois, conformément à l'axiome du chapitre précédent. Mais cela correspond-il vraiment au sens commun, c'est-à-dire est-ce comme cela que nous assignons des probabilités?

EXEMPLE - 2.1 _____ Pile ou face?

Je lance une pièce de monnaie. Quelle est la probabilité d'obtenir pile (ou face)?

Si je n'ai aucune information sur l'équité de la pièce, je dirai quand même que pile et face sont équiprobables, car je n'ai aucune raison de préférer l'un ou l'autre :

$$p(\text{pile}) = p(\text{face}) = \frac{1}{2}$$

Pour décréter cela, je n'ai fait encore aucune expérience. Cette "probabilité" n'est donc pas une fréquence. Elle est la mesure de mon état d'esprit sur la pièce. On peut écrire cela de la façon suivante. Soit la proposition logique

$$A \doteq \text{"La pièce tombera sur pile"}$$

Alors $P(A) = \frac{1}{2} = P(\bar{A})$ où \bar{A} est la négation de A , c'est-à-dire $\bar{A} = \text{"La pièce tombera sur face"}$.

Il s'agit ici de la probabilité d'une proposition logique, et non de celle d'une variable aléatoire.

L'exemple précédent souligne deux points importants :

- une probabilité n'est pas toujours une fréquence ;
- la probabilité que nous attribuons à la pièce est en fait la probabilité de ce que nous pensons de la pièce.

EXEMPLE - 2.2 **Dépendance logique contre causalité physique**

Soient les propositions logiques :

$A \doteq$ “Il commencera à pleuvoir avant 10h00.”

$B \doteq$ “Le ciel se couvrira avant 10h00.”

Il est clair que B n'implique pas A . Pourtant si B est vraie, A devient plus probable. Comment mesurer ce degré de certitude?

Attention, les exemples précédents ne signifient en aucun cas que la théorie classique des probabilités soit fausse. Simplement, avant de l'appliquer en physique, il faut prendre garde d'être sûr qu'elle décrit correctement le problème considéré.

Une théorie probabiliste du raisonnement telle que nous l'avons laissée entendre par les exemples précédents doit agir à la base sur des propositions logiques plutôt que sur des variables aléatoires, c'est pourquoi nous l'appelons théorie logique des probabilités. Dans le domaine d'application des variables aléatoires, les deux théories doivent donner le même résultat, par exemple :

$X_{log} \doteq$ “La variable aléatoire X prend une valeur entre x et $x + dx$.”

doit conduire à :

$$p(X_{log}) = p_X(x) dx$$

2.2 Les postulats de la théorie

Dans la suite de ce chapitre, les lettres capitales $A, B, C \dots$ désignent des propositions logiques. On commence par définir la notion de degré de plausibilité ou simplement plausibilité.

Axiome 2.1

(Degrés de plausibilité)

Les degrés de plausibilité sont représentés par des nombres réels. Une plausibilité plus forte correspond à une valeur plus grande. On note $(A|B)$ le degré de plausibilité de A sachant B , ou “la plausibilité conditionnelle que A soit vraie, sachant que B est vraie.”

Par exemple, $(A + B|CD)$ signifie “le degré de plausibilité conditionnelle que A ou B soient vraies, sachant que C et D sont vraies.” À la différence de la théorie classique des probabilités, il n'existe en théorie logique des probabilités que des plausibilités (probabilités) conditionnelles.

Axiome 2.2

(Correspondance qualitative avec le sens commun)

Si $(A|C') > (A|C)$ et $(B|AC') = (B|AC)$, alors $(AB|C') \geq (AB|C)$ et $(\bar{A}|C') < (\bar{A}|C)$.

En langage courant cette écriture devient :

- Si sachant C' , A est plus plausible que sachant C , et si B est indifférente entre C et C' sachant A , alors AB devient plus plausible sachant C' que sachant C .
- Si sachant C' , A est plus plausible que sachant C , alors \bar{A} devient moins plausible sachant C' que sachant C .

Axiome 2.3

(Le résultat doit être conséquent)

- (a) *S'il y a plusieurs façons de raisonner, elles doivent toutes conduire à la même conclusion.*
- (b) *Il faut toujours prendre en compte toute l'information disponible.*
- (c) *Le même état de connaissance doit conduire aux mêmes valeurs du degré de plausibilité.*

Les trois axiomes précédents conduisent au résultat suivant :

Théorème 2.1

(Règles du produit et de la somme)

Quelle que soit la fonction $p(x)$ continue, strictement croissante et telle que $0 \leq p(x) \leq 1$,

$$(a) \quad p(AB|C) = p(A|C)p(B|AC) = p(B|C)p(A|BC)$$

$$(b) \quad p(A|B) + p(\bar{A}|B) = 1$$

(a) est la règle du produit, (b) la règle de la somme.

Propriété 2.2

(Règle de la somme généralisée)

$$p(A + B|C) + p(AB|C) = p(A|C) + p(B|C)$$

La règle de la somme est un cas particulier de la règle de la somme généralisée pour $B = \bar{A}$.

Démonstration

$$\begin{aligned}
p(A + B|C) &= 1 - p(\bar{A}\bar{B}|C) && \text{règle de la somme} \\
&= 1 - p(\bar{A}|C)p(\bar{B}|\bar{A}C) && \text{règle du produit} \\
&= 1 - p(\bar{A}|C)(1 - p(B|\bar{A}C)) && \text{règle de la somme} \\
&= p(A|C) + p(\bar{A}|C)p(B|\bar{A}C) && \text{règle de la somme} \\
&= p(A|C) + p(\bar{A}B|C) && \text{règle du produit} \\
&= p(A|C) + p(B|C)p(\bar{A}|BC) && \text{règle du produit} \\
&= p(A|C) + p(B|C)(1 - p(A|BC)) && \text{règle de la somme} \\
&= p(A|C) + p(B|C) - p(AB|C) && \text{règle du produit}
\end{aligned}$$

Propriété 2.3

(Découplage des informations exclusives)

Soient n propositions $\{A_1 \dots A_n\}$ qui sont mutuellement exclusives sachant I , c'est-à-dire que

$$p(A_k A_l | I) = p(A_k | I) \delta_{kl}$$

alors

$$p\left(\sum_{k=1}^n A_k \mid I\right) = \sum_{k=1}^n p(A_k | I)$$

Démonstration

$$\begin{aligned}
p\left(\sum_{k=1}^n A_k \mid I\right) &= p\left(\sum_{k=1}^{n-1} A_k \mid I\right) + p(A_n | I) - p\left(\sum_{k=1}^{n-1} A_k A_n \mid I\right) \\
p\left(\sum_{k=1}^{n-1} A_k A_n \mid I\right) &= p\left(\sum_{k=1}^{n-2} A_k A_n \mid I\right) + \underbrace{p(A_{n-1} A_n | I)}_{=0} - \underbrace{p\left(\sum_{k=1}^{n-2} A_k A_{n-1} A_n \mid I\right)}_{\propto p(A_{n-1} A_n | I) = 0} \\
&\vdots \\
&= p(A_1 A_n | I) \\
&= 0
\end{aligned}$$

Il reste donc

$$p\left(\sum_{k=1}^n A_k \mid I\right) = p\left(\sum_{k=1}^{n-1} A_k \mid I\right) + p(A_n | I)$$

D'où le résultat par récurrence.

Propriété 2.4 (Découplage des apports d'informations exclusives)

Soient n propositions $\{A_1 \dots A_n\}$ qui sont mutuellement exclusives sachant I , c'est-à-dire que

$$p(A_k A_l | I) = p(A_k | I) \delta_{kl}$$

alors

$$p\left(C \left| \left(\sum_{k=1}^n A_k \right) I \right.\right) = \frac{\sum_{k=1}^n p(C | A_k I) p(A_k | I)}{\sum_{k=1}^n p(A_k | I)}$$

Si de plus $I = \sum_{k=1}^n A_k$, ou si une des propositions A_k est vraie sachant I ($\exists k, p(A_k | I) = 1$), alors

$$p(C | I) = \sum_{k=1}^n p(C | A_k I) p(A_k | I)$$

Démonstration

– Par la règle du produit

$$p\left(C \left| \left(\sum_{k=1}^n A_k \right) I \right.\right) = \frac{p\left(C \left(\sum_{k=1}^n A_k \right) \middle| I\right)}{p\left(\sum_{k=1}^n A_k \middle| I\right)}$$

D'après la propriété (2.3), le dénominateur devient directement

$$p\left(\sum_{k=1}^n A_k \middle| I\right) = \sum_{k=1}^n p(A_k | I)$$

Il reste à transformer le numérateur. Une nouvelle application de la règle du produit donne

$$p\left(C \left(\sum_{k=1}^n A_k \right) \middle| I\right) = p(C | I) p\left(\sum_{k=1}^n A_k \middle| CI\right)$$

À partir de

$$p(A_k A_l | CI) = \frac{p(C A_k A_l | I)}{p(C | I)} = \frac{p(C | A_k A_l I)}{p(C | I)} p(A_k A_l | I) = 0$$

on voit que les n propositions $\{A_1 \dots A_n\}$ sont mutuellement exclusives sachant CI , d'où par application de la propriété (2.3)

$$p\left(\sum_{k=1}^n A_k \middle| CI\right) = \sum_{k=1}^n p(A_k | CI)$$

et le numérateur devient

$$\begin{aligned} p\left(C\left(\sum_{k=1}^n A_k\right)\middle|I\right) &= p(C|I) \sum_{k=1}^n p(A_k|CI) \\ &= \sum_{k=1}^n p(CA_k|I) \\ &= \sum_{k=1}^n p(A_k|I)p(C|A_kI) \end{aligned}$$

Ce qui démontre la première partie de la propriété.

– Si $I = \sum_{k=1}^n A_k$, alors $\left(\sum_{k=1}^n A_k\right)I = II = I$, et

$$\sum_{k=1}^n p(A_k|I) = p\left(\sum_{k=1}^n A_k\middle|I\right) = p(I|I) = 1$$

ce qui démontre la seconde partie de la propriété dans ce cas.

– Si une des propositions A_k est vraie sachant I ($\exists k, p(A_k|I) = 1$), alors $\left(\sum_{k=1}^n A_k\right) = I$ et l'on est ramené au cas précédent.

Propriété 2.5

(Syllogismes forts et faibles)

1. Si $A \Rightarrow B$ et A est vraie, alors B est vraie.
2. Si $A \Rightarrow B$ et B est fausse, alors A est fausse.
3. Si $A \Rightarrow B$ et B est vraie, alors A devient plus plausible.
4. Si $A \Rightarrow B$ et A est fausse, alors B devient moins plausible.

Démonstration

Soit $C \doteq "A \Rightarrow B"$

1.

$$p(B|AC) = \frac{p(AB|C)}{p(A|C)} = \frac{p(A|C)}{p(A|C)} = 1$$

2.

$$p(A|\bar{B}C) = \frac{p(A\bar{B}|C)}{p(\bar{B}|C)} = \frac{0}{p(\bar{B}|C)} = 0$$

3.

$$p(A|BC) = p(A|C) \frac{p(B|AC)}{p(B|C)}$$

Or $p(B|AC) = 1$ et $p(B|C) \leq 1$, d'où $p(A|BC) \geq p(A|C)$

4.

$$p(B|\bar{A}C) = p(B|C) \frac{p(\bar{A}|BC)}{p(\bar{A}|C)}$$

Or $p(\bar{A}|BC) \leq p(\bar{A}|C)$ d'après le syllogisme précédent, donc $p(B|\bar{A}C) \leq p(B|C)$

2.3 Le principe d'indifférence

Dans la démonstration de la propriété (2.4) nous avons montré au passage le résultat suivant

Propriété 2.6

(Propositions exclusives et exhaustives)

Si les n propositions $\{A_1 \dots A_n\}$ sont mutuellement exclusives sachant B

$$p(A_k A_l | B) = p(A_k | B) \delta_{kl}$$

alors

$$p(A_1 + \dots + A_m | B) = \sum_{k=1}^m p(A_k | B) \quad \text{avec } m \leq n$$

Si de plus les propositions $\{A_1 \dots A_n\}$ sont exhaustives pour B (une et une seule est vraie sachant B) alors

$$\sum_{k=1}^n p(A_k | B) = 1$$

Remarque : La “relation de normalisation” ainsi démontrée est identique au deuxième postulat de la théorie classique des probabilités. Elle découle dans le cadre de la théorie logique des probabilités des seules règles du produit et de la somme.

Jusqu'à présent, nous n'avons pas assigné de valeurs numériques aux probabilités que nous avons utilisées. Ces valeurs ne peuvent pas être obtenues “par dénombrement de tous les cas possibles”, c'est-à-dire par calcul d'une fréquence d'apparition, mais par un raisonnement logique.

Cherchons comment assigner des valeurs dans le cas de n propositions A_k mutuellement exclusives et exhaustives, et de plus indifférentes suivant B :

Définition 2.4

(Indifférence)

Deux propositions A_1 et A_2 sont dites indifférentes suivant B si B dit la même chose sur A_1 et A_2 .

Définissons le problème désigné par un ' (prime) comme celui d'assigner des valeurs aux n propositions A_k quand A_1 et A_2 sont permutées :

$$\left\{ \begin{array}{l} A'_1 \doteq A_2 \\ A'_2 \doteq A_1 \\ A'_3 \doteq A_3 \\ \vdots \\ A'_n \doteq A_n \end{array} \right.$$

On désigne les probabilités assignées dans le premier problème par un indice I , et dans le second par un indice II . La permutation de A_1 et A_2 ne changeant pas la nature du problème, mais

représentant deux façons d'aboutir au même résultat, on doit avoir en vertu du troisième axiome

$$\begin{cases} p(A_1|B)_I = p(A'_2|B)_{II} \\ p(A_2|B)_I = p(A'_1|B)_{II} \end{cases}$$

Par ailleurs, puisque B est indifférente entre A_1 et A_2 , l'état de connaissance par rapport à B est identique pour les deux problèmes, d'où par application du troisième axiome

$$p(A_k|B)_I = p(A'_k|B)_{II}, \forall k$$

Il vient donc en combinant ces deux résultats

$$p(A_1|B)_I = p(A_2|B)_I$$

En généralisant ce raisonnement à toutes les permutations possibles (une permutation cyclique suffit), on obtient

Théorème 2.7

(Principe d'indifférence)

Si B est indifférente entre les n propositions exclusives et exhaustives A_k ,

$$p(A_k|B) = \frac{1}{n}$$

Autre démonstration (par l'absurde)

Si on n'assigne pas les mêmes plausibilités, par une simple permutation on conserve le même état de connaissance en affectant des plausibilités différentes, violant ainsi le troisième axiome.

EXEMPLE - 2.3 Tirage d'une urne

Une urne contient 10 boules numérotées. Si on définit

$A_k \doteq$ "La $k^{\text{ième}}$ boule est tirée."

$B \doteq$ "L'urne contient 10 boules numérotées."

Alors (indifférence) $p(A_k|B) = \frac{1}{10}$.

EXEMPLE - 2.4 Tirage sans remise

Soit $I \doteq$ "Une urne contient n boules, dont m sont blanches et le reste noires. On tire une boule sans la replacer."

$B_k \doteq$ "Le $k^{\text{ième}}$ tirage donne une boule blanche."

$N_k \doteq$ "Le $k^{\text{ième}}$ tirage donne une boule noire."

Quelle est la probabilité de tirer une boule blanche au $k^{\text{ième}}$ tirage?

$$\begin{aligned} p(B_1|I) &= \frac{m}{n} \quad (\text{principe d'indifférence}) \\ p(N_1|I) &= \frac{n-m}{n} \end{aligned}$$

Pour obtenir ces équations, il faut supposer subdiviser les propositions B_1 et N_1 en respectivement m et $n - m$ propositions élémentaires, mutuellement exclusives et exhaustives, donc de plausibilité $1/n$. Il suffit ensuite de sommer les propositions élémentaires correspondant à B_1 et N_1 .

- Comme $B_2 = B_1B_2 + N_1B_2$ sachant I , on peut écrire

$$\begin{aligned} p(B_2|I) &= p(B_1B_2 + N_1B_2|I) \\ &= p(B_1B_2|I) + p(N_1B_2|I) \\ &= p(B_1|I)p(B_2|B_1I) + p(N_1|I)p(B_2|N_1I) \\ &= \frac{m}{n} \frac{m-1}{n-1} + \frac{n-m}{n} \frac{m}{n-1} \\ &= \frac{m}{n} \end{aligned}$$

Mais n'était-ce pas évident dès le départ? En effet, en l'absence d'information sur ce qui s'est passé lors du premier tirage (retrait soit d'une boule blanche, soit d'une boule noire), on est exactement dans la même situation que lors du premier tirage. Et cela reste vrai pour les tirages suivants

$$p(B_k|I) = \frac{m}{n}, \forall k \leq n$$

Et pourtant, si on dépasse le $n^{\text{ième}}$ tirage, il ne subsiste plus de boules dans l'urne et

$$p(B_k|I) = 0, \forall k > n$$

Il est important de remarquer ici est l'information donnée par I ne permet pas de distinguer les tirages tant que $k \leq n$. La donnée d'une information complémentaire J = "Les trois premiers tirages ont donné une boule blanche." changerait complètement ce résultat :

$$\begin{aligned} p(B_k|IJ) &= 1, 1 \leq k \leq 3 \\ p(B_k|IJ) &= \frac{m-3}{n-3}, 4 \leq k \leq n \\ p(B_k|IJ) &= 0, k > n \end{aligned}$$

Si maintenant au lieu de l'information J on donnait l'information J' = "Les trois derniers tirages ($n-2, n-1, n$) ont donné une boule blanche.", ces résultats doivent encore être changés :

$$\begin{aligned} p(B_k|IJ') &= \frac{m-3}{n-3}, 1 \leq k \leq n-3 \\ p(B_k|IJ') &= 1, n-2 \leq k \leq n \\ p(B_k|IJ') &= 0, k > n \end{aligned}$$

En effet, d'un point de vue purement *logique*, tout se passe pour l'observateur comme si l'urne ne contenait au départ que $n-3$ boules dont $m-3$ blanches, bien que les trois boules blanches manquantes aient été enlevées *postérieurement*. Cet exemple montre clairement que la conditionnalité du degré de plausibilité est d'ordre purement logique, et en aucun cas une causalité physique.

EXERCICE - 2.1 Pièces et boîtes de conserve

On considère le jeu décrit par la proposition logique I suivante :

I = "On dispose de trois boîtes de conserve vides (notées A, B et C), rangées de gauche à droite, et de deux pièces de monnaie. Les deux pièces sont placées sous deux des boîtes de conserve, la troisième restant vide (une boîte ne peut contenir qu'une pièce à la fois)."

- a) - On définit les propositions logiques A, B et C par :

A = “Une pièce se trouve sous la boîte A.”

B = “Une pièce se trouve sous la boîte B.”

C = “Une pièce se trouve sous la boîte C.”

Il est intuitif que sous la seule information I , les trois propositions logiques A , B et C sont équiprobables et que

$$p(A|I) = p(B|I) = p(C|I) = 2/3$$

Ce résultat est-il obtenu par application du principe d'indifférence sur A , B et C ? Pourquoi? Montrer que sachant I on a les trois identités:

$$\begin{aligned} AB + BC + AC &= 1 \\ ABC &= 0 \\ A &= AB + AC \end{aligned}$$

Appliquer alors le principe d'indifférence pour trouver les valeurs de $p(A|I)$, $p(B|I)$ et $p(C|I)$.

– **b)** - On donne maintenant l'information complémentaire suivante:

J = (toute l'information de I) et “Il y a deux fois plus de chance qu'une pièce soit entourée de vide (une boîte vide ou rien), que de chance qu'elle soit voisine d'une autre pièce.”

Que valent alors $p(A|J)$, $p(B|J)$ et $p(C|J)$?

Réponse

– **a)** - D'après I , il y a seulement trois possibilités: AB , BC et AC

$$\begin{array}{lll} A & B & C \\ \odot & \odot & \circ \rightarrow AB \\ \circ & \odot & \odot \rightarrow BC \\ \odot & \circ & \odot \rightarrow AC \end{array}$$

A , B et C ne sont pas mutuellement exclusives et exhaustives sachant I , on ne peut donc pas appliquer le principe d'indifférence. Par contre, AB , BC et AC sont mutuellement exclusives et exhaustives, et donc $AB + BC + AC = 1$ et $ABC = 0$ sachant I . De plus, pour avoir une pièce sous la boîte A, il faut avoir soit AB soit AC , donc $A = AB + AC$ sachant I . (On a de même $B = AB + BC$ et $C = AC + BC$.) Par application du principe d'indifférence,

$$p(AB|I) = p(BC|I) = p(AC|I) = 1/3$$

et

$$p(A|I) = p(AB + AC|I) = p(AB|I) + p(AC|I) = 2/3$$

car $p(ABC|I) = 0$. On trouve de même $p(B|I) = 2/3$ et $p(C|I) = 2/3$.

– **b)** - Dans les cas AB et BC , les pièces sont toujours voisines. Dans le cas AC , les pièces n'ont pas de voisines. On en déduit donc:

$$p(AC|J) = 2p(AB|J) = 2p(BC|J)$$

AB , BC et AC sont toujours mutuellement exclusives et exhaustives sachant J , donc

$$p(AB|J) + p(BC|J) + p(AC|J) = 1$$

On en déduit que $p(AC|J) = 1/2$ et $p(AB|J) = p(BC|J) = 1/4$, puis

$$\begin{cases} p(A|J) &= p(AB|J) + p(AC|J) = 3/4 \\ p(B|J) &= p(AB|J) + p(BC|J) = 1/2 \\ p(C|J) &= p(AC|J) + p(BC|J) = 3/4 \end{cases}$$

2.4 Test d'hypothèse

2.4.1 Notations et problématique

Supposons que nous désirions tester la véracité d'une assertion concernant un système physique. Comment s'y prendre? La méthode scientifique orthodoxe veut que l'on réalise une expérience, que l'on mesure un certain nombre de grandeurs, puis que l'on essaye de voir dans quelle mesure les données acquises soutiennent l'hypothèse formulée.

Comment écrire ce principe mathématiquement? Nous pouvons chercher la plausibilité de l'hypothèse sachant l'expérience effectuée et les résultats (données) obtenus :

Définition 2.5

(Probabilité postérieure d'une hypothèse)

$$p(H|VI)$$

$I \doteq$ "Toute l'information disponible décrivant les conditions expérimentales."

$V \doteq$ "Les données obtenues."

$H \doteq$ "L'hypothèse à tester."

$p(H|VI) =$ "Probabilité que l'hypothèse H soit vraie, sachant que pour l'expérience décrite par I on a obtenu les données V ."

2.4.2 Méthode générale de résolution

La règle du produit fournit

$$p(VH|I) = p(H|I) p(V|HI) = p(V|I) p(H|VI)$$

soit encore

Propriété 2.8

(Évaluation de la probabilité postérieure d'une hypothèse)

$$p(H|VI) = p(H|I) \frac{p(V|HI)}{p(V|I)}$$

- $p(H|I)$: probabilité antérieure de H (antérieure à la mesure de V);
- $p(V|I)$: probabilité antérieure de V (sans tenir compte de H);
- $p(V|HI)$: probabilité de mesurer V si on fait confiance à H , ou vraisemblance de H (ce terme représente toute l'information sur H que contiennent les données V).

Comment obtenir ces différents termes ?

- Les probabilités antérieures $p(H|I)$ et $p(V|I)$ doivent être assignées par l'observateur suivant l'information dont il dispose, c'est-à-dire I (en pratique à l'aide du principe d'indifférence, ou du principe du maximum d'entropie décrit plus loin).
- La vraisemblance de H , $p(V|HI)$, provient en général d'un modèle physique de l'obtention des données V (voir exemples à suivre).

2.4.3 Critère de sélection d'une hypothèse

Une fois que $p(H|VI)$ a été obtenue, comment décider si l'hypothèse est valable ou non ?

- Si $p(H|VI) \approx 1$, H est très vraisemblablement vraie.
- Si $p(H|VI) \approx 0$, H est très vraisemblablement fausse.
- Mais si $p(H|VI) \approx \frac{1}{2}$, que conclure ? Simplement que les résultats de l'expérience effectuée ne permettent pas de trancher quand à la véracité de l'hypothèse.

Toutes les situations intermédiaires sont bien sûr possibles.

2.4.4 Cas particulier : l'alternative

Définition 2.6

(Alternative)

Une alternative est un choix binaire entre une hypothèse H et sa négation \bar{H} .

H et \bar{H} sont mutuellement exclusives et exhaustives pour I

$$p(H\bar{H}|I) = 0 \quad \text{et} \quad p(H|I) + p(\bar{H}|I) = 1$$

Le problème est donc de comparer H et \bar{H} , ou de décider si H est plus vraisemblablement vraie ou fausse.

$$\begin{aligned} p(H|VI) &= p(H|I) \frac{p(V|HI)}{p(V|I)} \\ p(\bar{H}|VI) &= p(\bar{H}|I) \frac{p(V|\bar{H}I)}{p(V|I)} \end{aligned}$$

Définition 2.7

(Chances d'une hypothèse)

“Les chances” de H sachant V et I sont définies par

$$O(H|VI) \doteq \frac{p(H|VI)}{p(\bar{H}|VI)}$$

$$O(H|VI) = O(H|I) \frac{p(V|HI)}{p(V|\bar{H}I)}$$

Cette expression permet de calculer “les chances” de H postérieurement à la mesure de V , connaissant “les chances” antérieures de H , et les vraisemblances de H et \bar{H} .

Définition 2.8**(Évidence d'une hypothèse)**

L'évidence de H sachant V et I est définie par

$$e(H|VI) \doteq 10\text{Log}O(H|VI)$$

$$e(H|VI) = e(H|I) + 10\text{Log} \left[\frac{p(V|HI)}{p(V|\bar{H}I)} \right]$$

Supposons maintenant que l'on accumule des données, par exemple en répétant l'expérience, ou en fusionnant des données différentes, alors $V = V_1V_2\dots$ et

$$\begin{aligned} e(H|VI) = e(H|I) &+ 10\text{Log} \left[\frac{p(V_1|HI)}{p(V_1|\bar{H}I)} \right] \\ &+ 10\text{Log} \left[\frac{p(V_2|V_1HI)}{p(V_2|V_1\bar{H}I)} \right] \\ &+ \dots \end{aligned}$$

Un cas intéressant est celui pour lequel les données V_k sont logiquement indépendantes sachant I . Alors

$$e(H|VI) = e(H|I) + 10 \sum_k \text{Log} \left[\frac{p(V_k|HI)}{p(V_k|\bar{H}I)} \right]$$

On a ainsi découpé les apports des paquets de données indépendants V_k , pour le jugement de l'hypothèse H .

p (probabilité)	O (les chances)	e (évidence)
$\frac{1}{2}$	1:1	0 dB
$\frac{2}{3}$	2:1	3 dB
$\frac{3}{4}$	4:1	6 dB
$\frac{4}{5}$	5:1	10 dB
$\frac{10}{11}$	10:1	10 dB
$\frac{100}{101}$	100:1	20 dB
0.999	1000:1	30 dB
0.9999	10000:1	40 dB

TAB. 2.1 – Les échelles de jugement d'une hypothèse.

EXEMPLE - 2.5

Contrôle qualité

$I \doteq$ “11 machines automatiques fabriquent des transistors, qui tombent dans des caisses sans

étiquette. Les caisses sont stockées sans qu'on puisse les différencier. 10 des 11 machines fabriquent 1 transistor sur 6 défectueux, et la onzième 1 transistor sur 3 défectueux."

On prend une des caisses, qui contiennent un très grand nombre de transistors, et on veut en testant quelques uns des transistors déterminer la machine de provenance. Sachant I on considère l'alternative :

- $H \doteq$ "La caisse provient de la onzième machine."
- $\bar{H} \doteq$ "La caisse provient de l'une des 10 premières machines."

avec $p(H|I) = \frac{1}{11}$ et $p(\bar{H}|I) = \frac{10}{11}$, soit $e(H|I) = -10$ dB.

On définit $V \doteq$ "On tire un transistor de la caisse. On le teste. On note le résultat : bon (1) ou mauvais (0)." On a donc

$$\begin{cases} p(V_0|HI) &= \frac{1}{3} \\ p(V_1|HI) &= \frac{2}{3} \end{cases}$$

et

$$\begin{cases} p(V_0|\bar{H}I) &= \frac{1}{6} \\ p(V_1|\bar{H}I) &= \frac{5}{6} \end{cases}$$

Le nombre de transistors est très grand, on ne pourra donc en tester qu'un petit nombre. Les vraisemblances précédentes (le modèle de mesure) sont donc indépendantes de l'ordre des tests.

$$10\text{Log} \left[\frac{p(V_0|HI)}{p(V_0|\bar{H}I)} \right] = \text{Log} \left[\frac{\frac{1}{3}}{\frac{1}{6}} \right] = +3\text{dB}$$

$$10\text{Log} \left[\frac{p(V_1|HI)}{p(V_1|\bar{H}I)} \right] = \text{Log} \left[\frac{\frac{2}{3}}{\frac{5}{6}} \right] = -1\text{dB}$$

Donc chaque fois qu'on trouve un transistor défectueux l'évidence de H croît de 3 dB, tandis qu'elle décroît de 1 dB pour chaque transistor correct :

$$e(H|VI) = e(H|I) + 3n_0 - n_1$$

où n_0 est le nombre de transistors défectueux et n_1 le nombre de transistors corrects. Par exemple sur 12 tests, si on trouve 5 transistors défectueux

$$e(H|VI) = -10 + 3 \times 5 - (12 - 5) = -2 \text{ dB}$$

Et sur 24 tests, si on trouve 10 transistors défectueux

$$e(H|VI) = -10 + 3 \times 10 - (24 - 10) = +6 \text{ dB}$$

On voit donc que pour une même proportion d'échec l'évidence change de valeur voire même de signe suivant le nombre de tests. Dans l'exemple numérique précédent on part de -10 dB avant les tests, on passe par -2 dB à 12 tests, puis à +6 dB à 24 tests. Partant d'une hypothèse peu plausible (-10 dB), on change ainsi d'avis jusqu'à la juger relativement plausible (+6 dB), en oubliant au fur et à mesure l'*a priori* de départ quand on accumule des données.

Mais comment prendre une décision ? Par exemple, si l'évidence dépasse 10 dB on accepte que la caisse testée provient (probablement) de la onzième machine. Si au contraire l'évidence descend au dessous de -15 dB, on décrète que la caisse provient d'une des 10 premières machines. Suivant ce type de critère, on voit que le nombre de tests n'est pas fixé.

Attention! Dans le cas où on ne considère pas une alternative, mais n hypothèses ($H_k, k = 1 \dots n$) mutuellement exclusives et exhaustives suivant I , avec $n > 2$

$$p(H_k H_l | I) = p(H_k | I) \delta_{kl} \quad \text{et} \quad \sum_{k=1}^n p(H_k | I) = 1$$

et que les données V peuvent se factoriser sous la forme $V = V_1 V_2 \dots V_m$ où les V_k sont indépendantes logiquement pour toutes les hypothèses

$$p(V_1 \dots V_m | H_k I) = \prod_{l=1}^m p(V_l | H_k I)$$

alors en général

$$p(V_1 \dots V_m | \bar{H}_k I) \neq \prod_{l=1}^m p(V_l | \bar{H}_k I)$$

Ainsi on ne peut pas dans ce cas découpler les apports des V_l pour évaluer l'évidence de H_k .

2.4.5 Test d'hypothèses multiples

Dans le cas où on doit comparer un nombre fini d'hypothèses, il convient tout d'abord d'obtenir

$$p(H_k | VI), k = 1 \dots n$$

par application de

$$p(H_k | VI) = p(H_k | I) \frac{p(V | H_k I)}{p(V | I)}$$

Il reste qu'il faut faire un choix parmi ces hypothèses. Il est raisonnable (mais arbitraire) de choisir l'hypothèse la plus probable :

Définition 2.9

(Principe du maximum a posteriori (MAP))

$$H_{k_0} = \underset{H_k}{\text{Argmax}} \{p(H_k | VI)\}$$

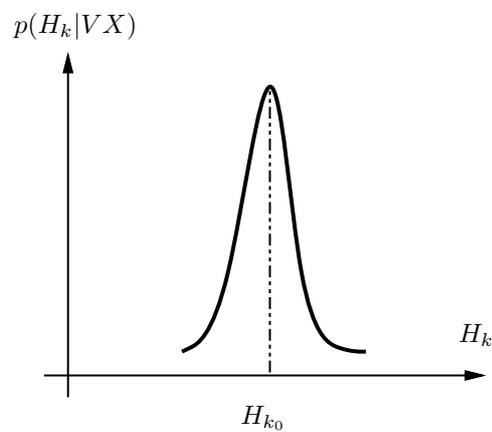
En général, $p(V | I)$ n'intervient pas dans ce choix. Si de plus toutes les hypothèses sont équiprobables sachant I ($p(H_k | I) = \frac{1}{n}$), alors le principe MAP se réduit au principe du maximum de vraisemblance (MV) :

Définition 2.10

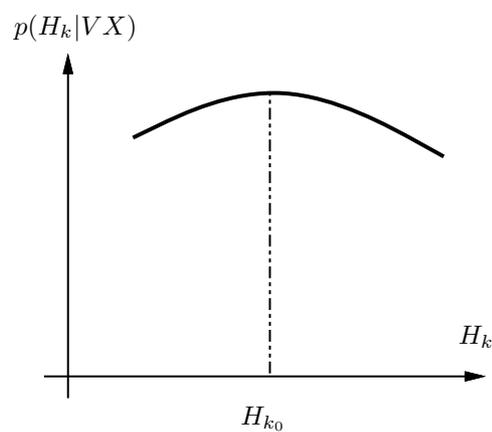
(Principe du maximum de vraisemblance (MV))

$$H_{k_0} = \underset{H_k}{\text{Argmax}} \{p(V | H_k I)\}$$

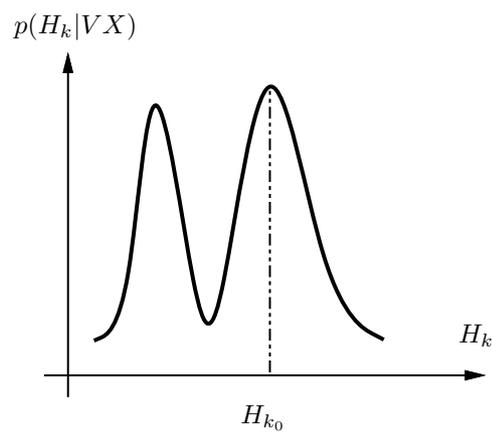
Il est en tout cas toujours possible de contrôler la validité de la solution MAP ou MV obtenue en étudiant la forme de $p(H_k | VI)$ (figure 2.1).



Courbe piquée :
solution crédible



Courbe "molle" :
solution imprécise



Que penser ?

FIG. 2.1 – Représentation graphique du critère MAP ou MV.

2.4.6 Test d'un nombre infini d'hypothèses

La théorie logique des probabilités ne travaille que sur des propositions discrètes. Il est cependant facile d'inclure des problèmes concernant un paramètre variant continûment, par la procédure suivante.

Soit f un paramètre variant continûment.

$$F \doteq (f \leq q) \quad \text{et} \quad F' \doteq (f > q)$$

F et F' sont des propositions discrètes, exclusives et exhaustives. On peut définir :

$$G(q) \doteq p(F'|I)$$

$G(q)$ est une fonction positive et croissante de q .

$$A \doteq (f \leq a), \quad B \doteq (f \leq b) \quad \text{et} \quad W \doteq (a < f \leq b)$$

$B = A + W$, et A et W sont exclusives et exhaustives pour B , donc

$$p(B|I) = p(A|I) + p(W|I)$$

soit

$$p(W|I) = p(a < f \leq b|I) = G(b) - G(a)$$

soit encore

$$p(a < f \leq b|I) = \int_a^b g(f)df$$

Il est clair que G est la fonction de répartition, et g la densité de probabilité pour f .

Soit $A_f \doteq$ "Le paramètre a une valeur entre f et $f + df$."

$$p(A_f|VI) = p(A_f|I) \frac{p(V|A_fI)}{p(V|I)}$$

$$\begin{cases} p(A_f|I) & = g(f|I)df \\ p(A_f|VI) & = g(f|VI)df \end{cases}$$

donc

$$g(f|VI) = g(f|I) \frac{p(V|A_fI)}{p(V|I)}$$

Soit $H_f \doteq$ "Le paramètre a pour valeur f ." Quand $df \rightarrow 0$, $p(V|A_fI) \rightarrow p(V|H_fI)$, et donc au final

$$g(f|VI) = g(f|I) \frac{p(V|H_fI)}{p(V|I)}$$

Par abus de notation, on écrit $p(f|I)$ et $p(f|VI)$ pour $g(f|I)$ et $g(f|VI)$.

EXEMPLE - 2.6 Contrôle qualité

On reprend l'exemple du contrôle des transistors, mais les hypothèses sont maintenant :

$$H_f \doteq \text{"La machine produit une fraction } f \text{ de transistors défectueux."}$$

On en déduit

$$p(V|H_fI) = f^n (1 - f)^{N-n}$$

avec $V \doteq$ “On teste N transistors, n sont défectueux.”

Il vient donc

$$p(f|VI) = \frac{f^n (1-f)^{N-n} p(f|I)}{\int_0^1 f^n (1-f)^{N-n} p(f|I) df}$$

– Dans l'exemple précédent, nous avons :

$$p(f|I) = \frac{10}{11} \delta\left(f - \frac{1}{6}\right) + \frac{1}{11} \delta\left(f - \frac{1}{3}\right)$$

d'où

$$p(f|VI) = \frac{\frac{10}{11} \left(\frac{1}{6}\right)^n \left(\frac{5}{6}\right)^{N-n} \delta\left(f - \frac{1}{6}\right) + \frac{1}{11} \left(\frac{1}{3}\right)^n \left(\frac{2}{3}\right)^{N-n} \delta\left(f - \frac{1}{3}\right)}{\frac{10}{11} \left(\frac{1}{6}\right)^n \left(\frac{5}{6}\right)^{N-n} + \frac{1}{11} \left(\frac{1}{3}\right)^n \left(\frac{2}{3}\right)^{N-n}}$$

– Supposons maintenant que $p(f|I) = 1$, $0 \leq f \leq 1$ (on ne dispose d'aucune information donnant une préférence à une fraction f particulière). Alors

$$p(f|VI) = \frac{(N+1)!}{n!(N-n)!} f^n (1-f)^{N-n}$$

Quelle est la solution MAP (identique à la solution MV ici) ?

$$\frac{dg}{df}(f|VI) \propto n f^{n-1} (1-f)^{N-n} - (N-n) f^n (1-f)^{N-n-1}$$

donc

$$\frac{dg}{df}(f|VI) = 0 \Leftrightarrow \hat{f} = \frac{n}{N}$$

Que devient $p(f|VI)$ quand on fait un grand nombre de tests ?

$$\ln p(f|VI) \approx n \ln f + (N-n) \ln(1-f) + \text{Const.}$$

$$\ln p(f|VI) \approx \ln g(\hat{f}|VI) - \frac{(f - \hat{f})^2}{2\sigma^2} + O((f - \hat{f})^2)$$

avec

$$\sigma^2 = \frac{\hat{f}(1-\hat{f})}{N}$$

d'où

$$p(f|VI) \approx \mathcal{N}(f, \hat{f}, \sigma^2)$$

Donc plus on fait de tests (N augmente), plus $p(f|VI)$ s'affine. En fait :

$$\lim_{N \rightarrow \infty} p(f|VI) = \delta(f - \hat{f})$$

et σ diminue comme $\frac{1}{\sqrt{N}}$.

2.5 Principe du maximum d'entropie

Comment faire pour aller au delà du principe d'indifférence?

EXEMPLE - 2.7

 Fenêtres japonaises

On a fait des statistiques au Japon lors du dernier tremblement de terre. On nous dit "La fenêtre moyenne s'est cassée en 10 morceaux." Sachant cela, quelle est la probabilité qu'une fenêtre japonaise se casse en n morceaux?

Quelle est la distribution qui présuppose le moins, tout en satisfaisant les contraintes? Ou quelle est la distribution la plus probable?

2.5.1 Entropie de Shannon

L'entropie d'une distribution discrète p_1, \dots, p_n est introduite à partir des postulats suivants :

- (1) Il existe une mesure $H_n(p_1, \dots, p_n)$ de la "quantité d'incertitude."
- (2) H_n est continue.
- (3) $h(n) \doteq H_n(\frac{1}{n}, \dots, \frac{1}{n})$ est une fonction croissante de n .
- (4) H_n est conséquente: s'il y a plusieurs façons de calculer sa valeur, toutes doivent conduire à la même valeur. En particulier, si $p_1 + q = 1$, on mesure l'incertitude par $H_2(p_1, q)$. Si on apprend de plus que $q = p_2 + p_3$, alors $p_1 + p_2 + p_3 = 1$ et l'incertitude est mesurée par $H_3(p_1, p_2, p_3)$. On doit alors postuler que

$$H_3(p_1, p_2, p_3) = H_2(p_1, q) + qH_2\left(\frac{p_2}{q}, \frac{p_3}{q}\right)$$

qui est une loi de décomposition en sous-systèmes.

On montre que la seule solution satisfaisant les conditions précédentes est

$$H_n(p_1, \dots, p_n) \doteq - \sum_{k=1}^n p_k \ln p_k$$

2.5.2 Méthode de Wallis

La méthode de Wallis est une expérience de pensée conduisant au principe du maximum d'entropie (on pourra noter qu'elle est en quelque sorte liée à l'axiome des "fréquences" de la théorie classique des probabilités).

Soit l'information I imposant des contraintes sur la distribution discrète p_1, \dots, p_n . On choisit un nombre $m \gg n$. On suppose qu'on a n boîtes, dans lesquelles on lance les m pièces. Cette distribution des pièces génère une distribution p_k en comptant les pièces dans chaque boîte

$$p_k = \frac{m_k}{m}$$

où m_k est le nombre de pièces dans la boîte numéro k ($k = 1 \dots n$). De toutes les distributions ainsi obtenues, on ne garde que celles qui sont compatibles avec l'information I .

La probabilité d'une répartition particulière est

$$W = n^{-m} \frac{m!}{m_1! \dots m_n!}$$

Quelle est la distribution la plus probable? Celle qui maximise W :

$$\max_{p_k} \frac{m!}{m_1! \cdots m_n!}$$

D'après la formule de Stirling,

$$\ln(m!) = m \ln m - m + \sqrt{2\pi m} + \frac{1}{12m} + O\left(\frac{1}{m}\right)$$

donc

$$\frac{1}{m} \ln(m!) = \ln m - 1 + O\left(\frac{1}{\sqrt{m}}\right)$$

$$\begin{aligned} \frac{1}{m} \ln\left(\frac{m!}{m_1! \cdots m_n!}\right) &= \frac{1}{m} \left(\ln(m!) - \sum_{k=1}^n \ln(m p_k)! \right) \\ &\approx \ln m - 1 - \sum_{k=1}^n p_k (\ln(m p_k) - 1) \\ &\approx \ln m - 1 - \underbrace{\left(\sum_{k=1}^n p_k \right)}_{=1} \ln m - \sum_{k=1}^n p_k \ln p_k + \underbrace{\sum_{k=1}^n p_k}_{=1} \\ &\approx - \sum_{k=1}^n p_k \ln p_k = H_n(p_1, \dots, p_n) \end{aligned}$$

D'où le résultat :

Théorème 2.9

(Principe du maximum d'entropie)

|| La distribution la plus probable est celle qui maximise l'entropie tout en étant compatible avec l'information I (les contraintes).

Nous venons de faire une "démonstration" de ce théorème dans le cas des distributions discrètes. Dans le cas continu, il se généralise en prenant pour définition de l'entropie

$$H(p) = - \int_{\mathbb{R}} p(x) \ln p(x) dx$$

EXEMPLE - 2.8 Moyenne fixée – cas discret

$n = 3$ et $\langle n \rangle = \bar{n}$ fixée. Méthode des multiplicateurs de Lagrange (on trouvera un bref rappel de ce principe en annexe de ce chapitre) :

$$\begin{aligned} \Psi(p) &= H_3(p) - \lambda \sum_{k=1}^3 k p_k - (\mu - 1) \sum_{k=1}^3 p_k \\ \frac{\partial \Psi}{\partial p_k}(p) &= -\ln p_k - 1 - \lambda k - (\mu - 1) = 0, k = 1 \dots 3 \end{aligned}$$

d'où

$$p_k = \exp(-\mu) \exp(-\lambda k)$$

Par identification de λ et μ on obtient (calcul difficile)

$$\begin{cases} p_1 &= \frac{3 - \bar{n} - p_2}{2} \\ p_2 &= \frac{1}{3} \left(\sqrt{4 - 3(\bar{n} - 2)^2} - 1 \right) \\ p_1 &= \frac{\bar{n} - 1 - p_2}{2} \end{cases}$$

EXEMPLE - 2.9 Moyenne fixée – cas continu

On suppose la moyenne de la distribution fixée, et égale à m . La méthode des multiplicateurs de Lagrange donne encore :

$$\begin{aligned} \Psi(p) &= H(p) - \lambda \int xp(x) - (\mu - 1) \int p(x) \\ \frac{\partial \Psi(p)}{\partial p(x)} &= -\ln p(x) - 1 - \lambda x - (\mu - 1) = 0, \forall x \in \mathbb{R} \end{aligned}$$

avec $\lambda > 0$ et $\mu > 0$, d'où

$$p(x) = \exp(-\mu) \exp(-\lambda x)$$

Pour éviter une divergence dans $\int p(x) dx$, il faut restreindre le domaine de définition de x , par exemple à $[0, +\infty[$. En identifiant les paramètres de Lagrange, on voit alors facilement que $\lambda = \exp(-\mu) = 1/m$, et donc

$$p(x) = \frac{1}{m} \exp\left(-\frac{x}{m}\right)$$

Cette distribution est par exemple celle de l'intensité du speckle.

2.5.3 Généralisation : fonction de partition

Le principe du maximum d'entropie est bien adapté quand l'information I consiste en un certain nombre de moyennes.

Soit x une variable qui peut prendre les valeurs $x_k, k = 1 \dots n$. m moyennes sont données, de la forme

$$\langle f_l(x) \rangle = \sum_{k=1}^n p_k f_l(x_k), l = 1 \dots m$$

Avec la contrainte additionnelle que $\sum_{k=1}^n p_k = 1$, la fonction de Lagrange associée au problème de maximisation de l'entropie est :

$$\begin{aligned} \Psi(p) &= H_n(p) - (\lambda_0 - 1) \sum_{k=1}^n p_k - \sum_{l=1}^m \lambda_l \sum_{k=1}^n p_k f_l(x_k) \\ \frac{\partial \Psi}{\partial p_k}(p) &= -\ln p_k - 1 - (\lambda_0 - 1) - \sum_{l=1}^m \lambda_l f_l(x_k) \end{aligned}$$

En annulant toutes ces dérivées partielles

$$\frac{\partial \Psi}{\partial p_k}(p) = 0, \forall k = 1 \dots n$$

on obtient

$$p_k = \exp \left(-\lambda_0 - \sum_{l=1}^m \lambda_l f_l(x_k) \right)$$

La contrainte $\sum_{k=1}^n p_k = 1$ donne

$$\exp(\lambda_0) \doteq Z(\lambda) = \sum_{k=1}^n \exp \left(- \sum_{l=1}^m \lambda_l f_l(x_k) \right)$$

$Z(\lambda)$ est la fonction de partition. On vérifie que :

$$\begin{cases} \lambda_0 & = \ln Z(\lambda) \\ \langle f_l(x) \rangle & = -\frac{\partial Z}{\partial \lambda_l}(\lambda) \end{cases}$$

On définit S comme l'entropie maximum, donc correspondant à la distribution la plus probable trouvée ci-dessus, et l'on vérifie que

$$S \doteq (H)_{\max} = \lambda_0 + \sum_{l=1}^m \lambda_l \langle f_l(x) \rangle$$

S ne dépend que des données : c'est une mesure objective de l'incertitude (il s'agit par analogie de l'entropie de la thermodynamique).

Il existe bien d'autres relations entre λ_l , $\langle f_l(x) \rangle$, $Z(\lambda)$, $S \dots$, qui sont à la base de la mécanique statistique. Par exemple :

- ensemble canonique de Gibbs (énergie moyenne fixée) ;
- ensemble grand canonique de Gibbs (énergie moyenne fixée + nombre de particules fixé)

Vérification que S est bien un maximum absolu :

$$\ln x \leq x - 1, 0 \leq x \leq \infty$$

donc

$$\sum_{k=1}^n p_k \ln \frac{u_k}{p_k} \leq \sum_{k=1}^n p_k \left(\frac{u_k}{p_k} - 1 \right) = 0$$

avec $\sum_{k=1}^n p_k = \sum_{k=1}^n u_k = 1$, donc

$$H_n(p) \leq \sum_{k=1}^n p_k \ln \frac{1}{u_k}$$

avec égalité si et seulement si les distributions p et u sont identiques.

Soit maintenant

$$u_k \doteq \frac{1}{Z(\lambda)} \exp \left(- \sum_{l=1}^m \lambda_l f_l(x_k) \right)$$

alors

$$\begin{aligned} H_n(p) &\leq \sum_{k=1}^n p_k \left(\ln Z(\lambda) + \sum_{l=1}^m \lambda_l f_l(x_k) \right) \\ H_n(p) &\leq \ln Z(\lambda) + \sum_{l=1}^m \lambda_l \langle f_l(x) \rangle \\ H_n(p) &\leq S \end{aligned}$$

S est donc bien le maximum absolu.

EXEMPLE - 2.10 _____ où l'on retrouve la gaussienne

Quelle est la distribution d'entropie maximum si on fixe la moyenne et la variance?

– **Cas discret** – Les contraintes sont :

$$\langle k \rangle = \sum_{k=1}^n k p_k \quad \text{et} \quad \sigma^2 + \langle k \rangle^2 = \sum_{k=1}^n k^2 p_k$$

On trouve donc immédiatement :

$$p_k = \exp(-\lambda_0 - \lambda_1 k - \lambda_2 k^2)$$

et il s'agit d'une gaussienne (l'identification des paramètres de Lagrange n'est pas facile).

– **Cas continu** – Les contraintes sont :

$$m = \int x p(x) dx \quad \text{et} \quad \sigma^2 + m^2 = \int x^2 p(x) dx$$

On trouve donc immédiatement :

$$p(x) = \exp(-\lambda_0 - \lambda_1 x - \lambda_2 x^2)$$

et par identification des paramètres de Lagrange :

$$p(x) = \mathcal{N}(x, m, \sigma^2)$$

2.6 Annexe : méthode des multiplicateurs de Lagrange

La méthode des multiplicateurs de Lagrange est un outil mathématique largement employé pour les problèmes d'optimisation sous contraintes. Nous allons la décrire dans le cas de l'optimisation sous contrainte d'une fonction d'une variable vectorielle, mais elle est utilisable dans un cadre bien plus général, et en particulier pour les fonctions de fonctions.

Soit $f(x)$ une fonction du vecteur x à n composantes, qui est définie sur D , sous-ensemble de \mathbb{R}^n . On cherche à obtenir le vecteur \hat{x} qui minimise la fonction $f(x)$ sur D , sous contrainte que :

$$g(x) = 0$$

où g est une fonction quelconque du vecteur x (la contrainte). On suppose que les fonctions f et g sont dérivables en tout point de D . On suppose de plus qu'au voisinage de chaque point x de D , $g(x) = 0$ définit une solution implicite :

$$x_n = h(x_1, \dots, x_{n-1})$$

Le minimum de $f(x)$ sous contrainte que $g(x) = 0$ pourra alors être obtenu en injectant la valeur $x_n = h(x_1, \dots, x_{n-1})$ dans $f(x)$, puis en minimisant le résultat

$$f(x_1, \dots, x_{n-1}, h(x_1, \dots, x_{n-1}))$$

par rapport aux $(n - 1)$ variables x_1, \dots, x_{n-1} , soit

$$\frac{\partial f(x_1, \dots, x_{n-1}, h(x_1, \dots, x_{n-1}))}{\partial x_i} + \frac{\partial h(x_1, \dots, x_{n-1})}{\partial x_i} \cdot \frac{\partial f(x)}{\partial x_n} = 0, \quad \forall i = 1, \dots, n - 1$$

De plus, puisque

$$g(x_1, \dots, x_{n-1}, h(x_1, \dots, x_{n-1})) = 0$$

on a également

$$\frac{\partial g(x_1, \dots, x_{n-1}, h(x_1, \dots, x_{n-1}))}{\partial x_i} + \frac{\partial h(x_1, \dots, x_{n-1})}{\partial x_i} \cdot \frac{\partial g(x)}{\partial x_n} = 0, \quad \forall i = 1, \dots, n - 1$$

On peut écrire ces équations sous forme matricielle :

$$\begin{pmatrix} \frac{\partial f}{\partial x_i} & \frac{\partial f}{\partial x_n} \\ \frac{\partial g}{\partial x_i} & \frac{\partial g}{\partial x_n} \end{pmatrix} \begin{pmatrix} 1 \\ \frac{\partial h}{\partial x_i} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \forall i = 1, \dots, n - 1$$

Les lignes de la matrice sont donc proportionnelles, et l'on en déduit

$$\frac{\partial f}{\partial x_i}(\hat{x}) = \lambda \frac{\partial g}{\partial x_i}(\hat{x}), \quad \forall i = 1, \dots, n$$

Le minimum de $f(x)$ sous contrainte que $g(x) = 0$ est donc obtenu en minimisant la fonction :

$$\Psi(x) = f(x) - \lambda g(x)$$

où λ est appelé le multiplicateur de Lagrange. Il reste cependant à obtenir la valeur de λ . Puisque la solution est obtenue en minimisant la fonction de Lagrange $\Psi(x) = f(x) - \lambda g(x)$, l'optimum trouvé, $\hat{x}(\lambda)$, dépend de λ . Il faut alors identifier la valeur du multiplicateur de Lagrange en résolvant l'équation

$$g(\hat{x}(\lambda)) = 0.$$

opération que l'on désigne sous le nom d'identification du paramètre.

La démonstration précédente se généralise dans le cas où plusieurs contraintes sont appliquées, et on obtient la propriété suivante.

Propriété 2.10**(Multipliateurs de Lagrange)**

Soit $f(x)$ une fonction du vecteur x à n composantes, qui est définie sur D , sous-ensemble de \mathbb{R}^n . Le vecteur \hat{x} qui minimise la fonction $f(x)$ sur D , sous les m contraintes :

$$g_k(x) = 0, \quad k = 1 \dots m$$

peut être obtenu en formant la fonction de Lagrange

$$\Psi(x, \lambda) = f(x) - \sum_{k=1}^m \lambda_k g_k(x)$$

où $\lambda_k > 0$, puis en cherchant $\hat{x}(\lambda)$ qui optimise $\Psi(x, \lambda)$ sur D par

$$\frac{\partial \Psi}{\partial x_i}(\hat{x}(\lambda), \lambda) = 0, \quad i = 1 \dots n$$

et enfin en identifiant les contraintes

$$g_k(\hat{x}(\lambda)) = 0, \quad k = 1 \dots m.$$

Chapitre 3

Théorie de la décision

3.1 Introduction

Jusqu'à présent, nous avons vu comment calculer la probabilité d'un certain nombre d'hypothèses, au sens large. Comment convertir une mesure de plausibilité (une probabilité) en une décision? c'est le but de la théorie de la décision.

La donnée des probabilités de toutes les hypothèses possibles concernant un problème réel n'est pas suffisante, car toutes les hypothèses ne sont pas équivalentes pour la décision finale. Pour pouvoir prendre une décision, il faut définir un critère mesurant l'opportunité de chaque décision, à la lumière des hypothèses plus ou moins plausibles.

EXEMPLE - 3.1 _____ **Espérance de profit**

On définit l'espérance de profit par

$$\langle M \rangle = \sum_{k=1}^n p_k M_k$$

où il y a n possibilités, chacune de probabilité p_k , et de profit financier M_k .

Pour un capitaliste, il peut sembler intéressant de toujours agir de façon à optimiser son espérance de profit.

Soit maintenant une pièce truquée légèrement (par nos soins), pour laquelle on a de bonnes raisons de croire que face sort avec une probabilité de 0.51. On propose à un quidam dans la rue de lancer la pièce, et on parie la somme M_0 sur face à 50/50:

jouer $\langle M \rangle = 0.51M_0 + 0.49(-M_0) = 0.02M_0 > 0$

ne pas jouer $\langle M \rangle = 0$

Est-ce que vous jouez tout votre argent, ce qui conduit à l'espérance de profit maximum?

3.2 Définitions

Un problème de décision est très similaire à un problème de test d'hypothèse, et généralise celui-ci.

Voici la procédure générale :

1. Rassembler toute l'information disponible sur le problème ; ceci fournit le contexte I (une proposition logique).

- Énumérer les états possibles $\{\theta_j, j = 1 \dots m\}$ du système (c'est-à-dire les hypothèses). Déterminer les probabilités antérieures

$$p(\theta_j|I), j = 1 \dots m$$

(par application du principe d'indifférence, du maximum d'entropie...)

- Collecter les données $V = V_1 V_2 \dots$. Digérer ces données pour obtenir les probabilités postérieures :

$$p(\theta_j|VI), j = 1 \dots m$$

par application de la règle du produit :

$$p(\theta_j|VI) = p(\theta_j|I) \frac{p(V|\theta_j I)}{p(V|I)}$$

où $p(V|\theta_j I)$ décrit le modèle de mesure des données.

- Énumérer les décisions possibles $\{D_i, i = 1 \dots n\}$.
- Définir une fonction de perte

$$L_{ij} = L(D_i, \theta_j)$$

“ce que coûte de prendre la décision D_i quand l'état du système est θ_j ”. C'est ici que s'introduit l'arbitraire de la méthode : il n'y a aucun principe général permettant de choisir L_{ij} .

- Pour chaque décision possible D_i , calculer l'espérance de perte :

$$K(D_i) \doteq \sum_{j=1}^m L(D_i, \theta_j) p(\theta_j|VI)$$

Dans cette expression, l'état du système θ_j apparaît comme une variable d'intégration (somme). C'est un paramètre de nuisance (on ne cherche pas à déterminer l'état dans lequel se trouve le système, mais à prendre une décision).

- Règle de discrimination entre décisions (ou règle de décision) :

choisir la décision D_{i_0} qui minimise l'espérance de perte :

$$D_{i_0} = \underset{D_i}{\operatorname{Argmin}} K(D_i)$$

3.3 Cas particulier : détection

Considérons le cas particulier simple suivant, mais très important en pratique, de la détection. Nous l'exposons sur un exemple afin que les équations soient plus “parlantes”, mais il se généralise aisément à tout problème comportant seulement deux décisions possibles.

- $I \doteq$ “Une batterie anti-aérienne est automatisée. Une caméra permet l'observation du ciel alentour. Si on détecte un avion, il faut prendre la décision de tirer ou non. On a remarqué qu'en moyenne la probabilité qu'un avion passe est p .”

2. Les états possibles du système sont pour nous

$\theta_1 \doteq$ “Il y a un avion”

$\theta_0 \doteq$ “Il y a pas d’avion”

avec $p(\theta_1|I) = p$ et $p(\theta_0|I) = 1 - p = q$.

3. $V \doteq$ “On a acquis une image du ciel.”

$$\begin{cases} p(\theta_0|VI) = p(\theta_0|I) \frac{p(V|\theta_0I)}{p(V|I)} = q \frac{p(V|\theta_0I)}{p(V|I)} \\ p(\theta_1|VI) = p(\theta_1|I) \frac{p(V|\theta_1I)}{p(V|I)} = p \frac{p(V|\theta_1I)}{p(V|I)} \end{cases}$$

$p(V|\theta_0I) \rightarrow$ vraisemblance qu’il n’y ait pas d’avion dans l’image acquise.

$p(V|\theta_1I) \rightarrow$ vraisemblance qu’il y ait un avion dans l’image acquise.

4. Décisions possibles :

$D_0 \doteq$ “On ne tire pas.”

$D_1 \doteq$ “On tire.”

5.

$$L = \begin{pmatrix} L_{00} & L_{01} \\ L_{10} & L_{11} \end{pmatrix} = \begin{pmatrix} 0 & L_r \\ L_a & 0 \end{pmatrix}$$

$L_{01} = L_r$ est le “coût d’un faux repos,” c’est-à-dire s’il y a un avion mais qu’on ne tire pas.

$L_{10} = L_a$ est le “coût d’une fausse alarme,” c’est-à-dire s’il n’y a pas d’avion mais qu’on tire.

$L_{00} = L_{11} = 0$ car dans ces deux cas on ne s’est pas trompé, et il n’y a donc aucune raison de pénaliser ces situations.

6. Espérance de perte :

$$\begin{cases} K(D_0) = L_r p(\theta_1|VI) = L_r p \frac{p(V|\theta_1I)}{p(V|I)} \\ K(D_1) = L_a p(\theta_0|VI) = L_a q \frac{p(V|\theta_0I)}{p(V|I)} \end{cases}$$

7. Tirer si $K(D_1) < K(D_0)$, c’est-à-dire si

$$\boxed{\frac{p(\theta_1|VI)}{p(\theta_0|VI)} > \frac{L_a}{L_r}} \quad \text{ou} \quad \boxed{\frac{p(V|\theta_1I)}{p(V|\theta_0I)} > \frac{q L_a}{p L_r}}$$

3.4 Théorie de la décision : le point de vue classique

La règle de décision qui a été présentée précédemment ($D_{i_0} = \text{Argmin}_{D_i} K(D_i)$) intervient tout naturellement dans le cadre de la théorie logique des probabilités.

Du point de vue de la théorie classique des probabilités (c'est-à-dire quand les probabilités sont interprétés comme des fréquences), l'expression de l'espérance de perte

$$K(D_i) \doteq \sum_{j=1}^m L(D_i, \theta_j) p(\theta_j | VI)$$

pose un problème conceptuel.

En effet, $p(\theta_j | VI)$ n'est pas une distribution que l'on peut obtenir à l'issue d'une expérience aléatoire de mesure d'une quantité physique (θ_j est une hypothèse, pas une quantité mesurable). Seules des quantités du type $p(V|I)$ ou $p(V|\theta_j I)$ peuvent être interprétées de cette façon.

Il est cependant possible de contourner cette difficulté de la façon suivante.

Définition 3.1 (Régions de décision et fonctions discriminantes)

La décision finale ne peut être prise que sur la base des données uniquement. Une règle de décision est donc constituée de la donnée de n fonctions discriminantes $p(D_i | VI)$ déterminées de la façon suivante.

Si \mathcal{R} est l'espace de définition des données V , \mathcal{R} est la réunion disjointe des régions de décisions $\mathcal{R}_1 \dots \mathcal{R}_n$ où sont vraies respectivement $D_1 \dots D_n$:

$$\mathcal{R} = \bigcup_{i=1}^n \mathcal{R}_i \quad \text{avec} \quad \mathcal{R}_i \cap \mathcal{R}_j = \emptyset \quad \text{si} \quad i \neq j$$

et $V \in \mathcal{R}_i$ implique qu'il faille choisir la décision D_i , c'est-à-dire par définition

$$p(D_i | VI) = \delta(D_i - D_{i_V}) \quad \text{si} \quad V \in \mathcal{R}_{i_V}$$

$p(D | VI) = \sum_{i=1}^n p(D_i | VI)$ est donc une fonction de partition entre les régions de décisions puisqu'un seul des symboles de Kronecker est non nul à la fois.

En conséquence, l'influence de toute proposition Y sur la décision finale D ne peut se faire qu'à travers l'influence de Y sur les données V :

$$p(D_i | YI) = \sum_V p(D_i | VI) p(V | YI)$$

On définit ensuite la perte conditionnelle.

Définition 3.2 (Perte conditionnelle)

Il s'agit de la perte encourue si le système est dans l'état θ_j , moyennée sur toutes les décisions possibles :

$$L(\theta_j) \doteq \sum_{D_i} L(D_i, \theta_j) p(D_i | \theta_j I)$$

D'où on déduit en prenant $Y = \theta_j$,

$$L(\theta_j) = \sum_{D_i} \sum_V L(D_i, \theta_j) p(D_i | VI) p(V | \theta_j I)$$

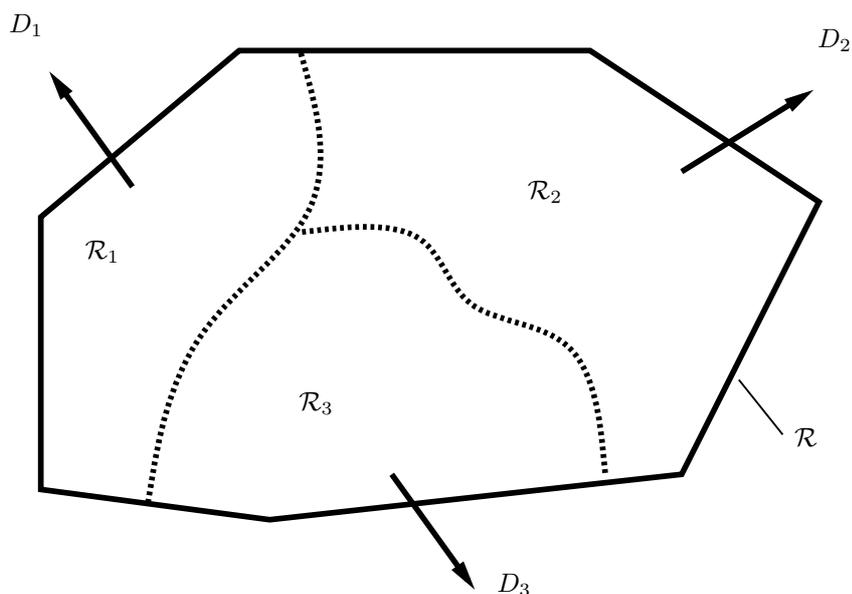


FIG. 3.1 – Représentation schématique des régions de décision.

Cette expression ne fait intervenir que des quantités bien définies au sens classique, à savoir les fonctions discriminantes et le modèle $p(V|\theta_j I)$.

Définition 3.3**(Critère Minimax)**

Pour une règle de décision fixe, on cherche le maximum de $L(\theta_j)$. On cherche alors la règle de décision qui minimise $L(\theta_j)_{max}$:

$$p_m(D|VI) \doteq \underset{p(D|VI)}{\operatorname{Argmin}} \left\{ \max_{\theta_j} \{L(\theta_j)\} \right\}$$

Ce critère revient à concentrer toute l'attention sur le cas le plus défavorable, sans tenir compte de la probabilité d'occurrence de ce cas.

La perte en moyenne est définie par

Définition 3.4**(Perte en moyenne)**

$$\langle L \rangle \doteq \sum_{\theta_j} L(\theta_j) p(\theta_j|I)$$

Attention, cette quantité est naturellement "bayésienne" (c'est-à-dire doit être comprise au sens de la théorie logique des probabilités) ; elle fait intervenir une probabilité antérieure, $p(\theta_j|I)$. Cette difficulté est contournée dans les ouvrages qui ne font pas référence à la théorie logique des probabilités en disant que θ est l'état du système défini par des nombres mesurables (la position, la vitesse...).

Définition 3.5

(Critère de Bayes)

$$p_b(D|VI) \doteq \underset{p(D|VI)}{\text{Argmin}} \{ \langle L \rangle \}$$

$$\begin{aligned} \langle L \rangle &= \sum_{D_i} \sum_V \left[\sum_{\theta_j} L(D_i, \theta_j) p(V|\theta_j I) p(\theta_j|I) \right] p(D_i|VI) \\ \langle L \rangle &= \sum_{D_i} \sum_V \left[\sum_{\theta_j} L(D_i, \theta_j) p(V\theta_j|I) \right] p(D_i|VI) \end{aligned}$$

On voit donc que $\langle L \rangle$ est minimum si pour chaque V on choisit la décision D_{i_V} pour laquelle la quantité entre crochet est minimum, puisqu'alors

$$p(D_i|VI) = \delta(D_i - D_{i_V})$$

et

$$\sum_{\theta_j} L(D_{i_V}, \theta_j) p(V\theta_j|I) \leq \sum_{\theta_j} L(D_i, \theta_j) p(V\theta_j|I)$$

impliquent que

$$\langle L \rangle = \sum_V \left[\sum_{\theta_j} L(D_{i_V}, \theta_j) p(V\theta_j|I) \right] \leq \sum_{D_i} \sum_V \left[\sum_{\theta_j} L(D_i, \theta_j) p(V\theta_j|I) \right] p(D_i|VI)$$

Mais cette solution est identique à la règle de minimisation de l'espérance de perte introduite au début de ce chapitre, puisque

$$K(D_i) \doteq \sum_{j=1}^m L(D_i, \theta_j) p(\theta_j|VI) = \frac{1}{p(V|I)} \left[\sum_{\theta_j} L(D_i, \theta_j) p(V\theta_j|I) \right]$$

soit

$$\langle L \rangle = \sum_{D_i} \sum_V K(D_i) p(V|I) p(D_i|VI)$$

Donc quel que soit V , le minimum de $\langle L \rangle$ est obtenu en choisissant la règle de décision $p(D_i|VI)$, donc la décision D_i , qui minimise $K(D_i)$.

3.5 Stratégie de Neyman-Pearson

La stratégie de Neyman-Pearson concerne plus particulièrement le problème de la détection étudiée précédemment sur un exemple. C'est un cas particulier de la section précédente.

Dans un problème à deux états (signal présent ou non) et à deux décisions (déclencher le tir ou non), il y a deux sources d'erreur :

$$D_1\theta_0 : \text{fausse alarme avec une probabilité } p(D_1\theta_0|I)$$

$$D_0\theta_1 : \text{faux repos avec une probabilité } p(D_0\theta_1|I)$$

Ces probabilités de fausse alarme et de faux repos peuvent être évaluées connaissant la règle de décision $p(D|VI)$ (voir calcul plus bas).

La stratégie de Neyman-Pearson consiste à fixer par exemple la probabilité de fausse alarme, et à minimiser la probabilité de faux repos sous cette contrainte. Il s'agit donc d'un problème de Lagrange décrit par les équations suivantes :

$$\begin{aligned} p(D_1\theta_0|I) &= \epsilon \\ \psi_\lambda(p(D|VI)) &= p(D_0\theta_1|I) + \lambda p(D_1\theta_0|I) \end{aligned}$$

ψ_λ est la fonction de Lagrange à minimiser, et λ le paramètre de Lagrange correspondant à la contrainte. Pour chaque valeur de λ il faut obtenir la règle de décision $p_\lambda(D|VI)$ qui minimise ψ_λ .

S'il l'on connaît précisément la valeur de ϵ , il faut ensuite identifier la contrainte, c'est-à-dire trouver la valeur λ_0 pour laquelle $p(D_1\theta_0|I) = \epsilon$. Dans la stratégie de Neyman-Pearson, on ne considère pas la valeur de ϵ fixée, et l'on cherche à obtenir plutôt tous les compromis possibles entre probabilité de fausse alarme et probabilité de faux repos. Ceux-ci sont obtenus en faisant varier ϵ de 0 à 1, ou de façon équivalente λ de $+\infty$ à 0 (voir figure 3.2, et calculs ci-dessous).

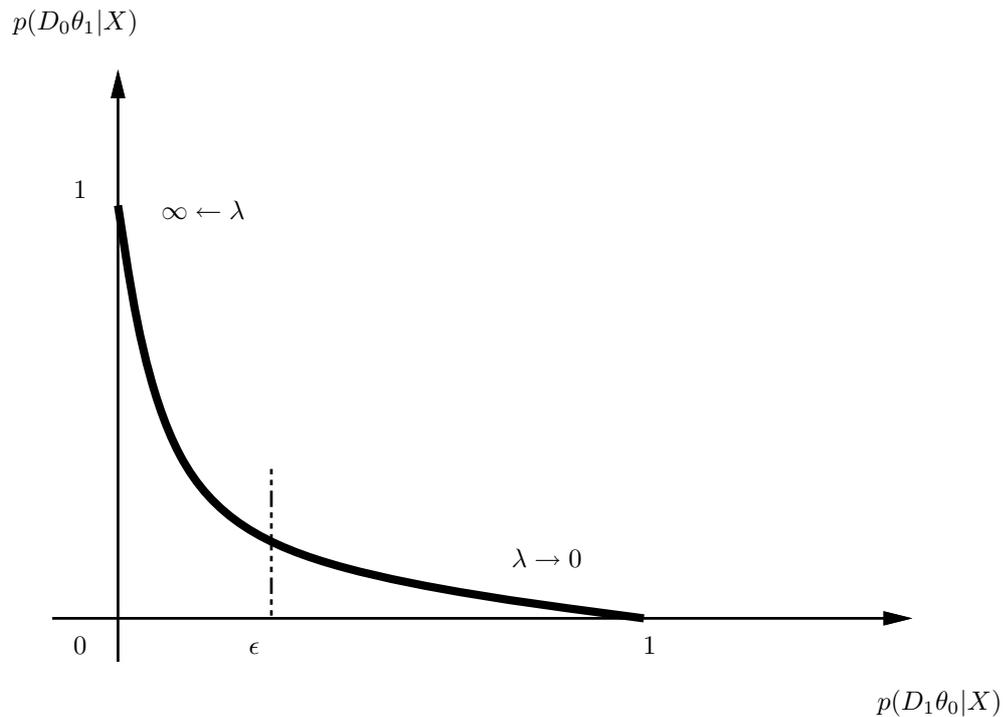


FIG. 3.2 – Le paramètre λ permet de décrire l'ensemble de la courbe.

Calculons les probabilités de fausse alarme et de faux repos :

$$\begin{aligned} p(D_0\theta_1|I) &= p(D_0|\theta_1I)p(\theta_1|I) \\ &= \sum_V p(D_0|VI)p(V|\theta_1I)p(\theta_1|I) \\ p(D_0\theta_1|I) &= \sum_V p(D_0|VI)p(V\theta_1|I) \end{aligned}$$

On montre de même facilement que

$$p(D_1\theta_0|I) = \sum_V p(D_1|VI)p(V\theta_0|I)$$

ou plus généralement (ce qui inclue les cas de vraie alarme et vrai repos)

$$p(D_i\theta_j|I) = \sum_V p(D_i|VI)p(V\theta_j|I)$$

pour i et j égaux à 0 ou 1. On en déduit les résultats suivants. D'une part en posant

$$L_{ij} = \begin{pmatrix} 0 & L_r \\ L_a & 0 \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ \lambda & 0 \end{pmatrix} \quad \text{soit} \quad \lambda = \frac{L_a}{L_r}$$

on voit que

$$\langle L \rangle = p(D_0\theta_1|I) + \lambda p(D_1\theta_0|I) = \psi_\lambda$$

ce qui prouve que la stratégie de Neyman-Pearson est un cas particulier du risque de Bayes. Il n'en reste pas moins que de calculer les probabilités de fausse alarme et de faux repos donne une bonne idée des performances du système, et permet une représentation simple des compromis possibles.

D'autre part on a de façon évidente

$$\sum_{D_i} \sum_{\theta_j} p(D_i\theta_j|I) = p(D_0\theta_0|I) + p(D_0\theta_1|I) + p(D_1\theta_0|I) + p(D_1\theta_1|I) = 1$$

ce qui implique que si $p(D_0\theta_1|I) = 1$ alors $p(D_1\theta_0|I) = 0$ et réciproquement, ce qui permet d'obtenir les points extrêmes de la courbe de la figure 3.2.

Précisons en outre la relation entre ϵ et λ . Si $\lambda = 0$, alors $\psi_0 = p(D_0\theta_1|I)$, et la minimisation aboutira à $p(D_0\theta_1|I) = 0$ soit $\epsilon = 1$. Inversement, si $\lambda \rightarrow \infty$, $\psi_\lambda \rightarrow \lambda p(D_1\theta_0|I)$, et la minimisation aboutira à $p(D_1\theta_0|I) \rightarrow 0$, soit $\epsilon \rightarrow 0$.

3.6 Exemple : détection de bits

I \doteq “On considère un système linéaire, pour lequel on mesure une tension v à un instant isolé. On doit décider si un signal d’amplitude fixe s est présent ou non, afin d’inscrire un 1 ou un 0 dans une mémoire (échantillonnage). On sait que le bruit d’acquisition est centré et de variance σ^2 . p , q , L_a et L_r sont supposés donnés.”

θ_1 \doteq “Le signal d’amplitude s est présent.”

θ_0 \doteq “Le signal d’amplitude s est absent.”

D_1 \doteq “On inscrit un 1.”

D_0 \doteq “On inscrit un 0.”

V \doteq On mesure la tension v .”

Quel modèle de mesure devons nous employer? Nous savons que le bruit d’acquisition est centré et de variance σ^2 , et que le signal attendu est s si θ_1 est vraie et 0 sinon (ou θ_0 est vraie). La tension mesurée a donc pour moyenne s si θ_1 est vraie et 0 sinon, et pour variance σ^2 . En appliquant le principe du maximum d’entropie on obtient donc :

$$\begin{cases} p(V|\theta_1 I) = \mathcal{N}(v, s, \sigma^2) \\ p(V|\theta_0 I) = \mathcal{N}(v, 0, \sigma^2) \end{cases}$$

D’où l’on déduit

$$\frac{p(V|\theta_1 I)}{p(V|\theta_0 I)} = \frac{\mathcal{N}(v, s, \sigma^2)}{\mathcal{N}(v, 0, \sigma^2)} = \exp\left(\frac{2vs - s^2}{2\sigma^2}\right)$$

La règle de décision est donc de choisir D_1 si

$$\frac{p(V|\theta_1 I)}{p(V|\theta_0 I)} = \exp\left(\frac{2vs - s^2}{2\sigma^2}\right) > \frac{q L_a}{p L_r}$$

et D_0 sinon. En prenant le logarithme de cette dernière expression, on voit que la décision est prise en comparant la tension mesurée v à un seuil v_b prédéterminé (donné seulement par l’information I), soit choisir D_1 si

$$v > v_b = \frac{s}{2} + \frac{\sigma^2}{s} \ln\left(\frac{q L_a}{p L_r}\right)$$

et D_0 sinon.

On peut maintenant calculer les probabilités de fausse alarme et de faux repos.

$$\begin{aligned} p(D_0\theta_1|I) &= p \sum_V p(D_0|VI)p(V|\theta_1 I) \\ &= p \int_{-\infty}^{v_b} \mathcal{N}(v, s, \sigma^2) dv \\ p(D_0\theta_1|I) &= p \Phi\left(\frac{v_b - s}{\sigma}\right) \end{aligned}$$

où $\Phi(x)$ est la fonction erreur définie par :

$$\Phi(x) \doteq \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp(-t^2/2) dt$$

De même la probabilité de fausse alarme est

$$\begin{aligned} p(D_1\theta_0|I) &= q \sum_V p(D_1|VI)p(V|\theta_0I) \\ &= q \int_{v_b}^{\infty} \mathcal{N}(v,0,\sigma^2)dv \\ p(D_1\theta_0|I) &= q \left(1 - \Phi\left(\frac{v_b}{\sigma}\right)\right) \end{aligned}$$

Par exemple pour $q = 10p$ et $L_r = 10L_a$, $v_b = s/2$, $p(D_0\theta_1|I) = 2.7\%$, et $p(D_1\theta_0|I) = 0.27\%$.

Chapitre 4

Estimation de paramètres

4.1 Introduction – Problématique

4.1.1 Approche classique

L'objectif général de l'estimation de paramètres est d'extraire d'un signal les informations qu'il contient. Cette information à obtenir est constituée du vecteur paramètre θ (à n composantes). Le signal attendu $|x_\theta\rangle$ est supposé dépendre du vecteur paramètre θ :

$$\left\{ \begin{array}{l} \mathbb{R}^n \rightarrow \text{espace des signaux (discret ou continu)} \\ \theta \mapsto |x_\theta\rangle \end{array} \right.$$

avec $|x_\theta\rangle = x_\theta(t), t \in \mathbb{R}$ pour un signal continu et $|x_\theta\rangle = x_\theta(i), i = 1 \dots m$ pour un signal discret.

EXEMPLE - 4.1 --- Radar ou sonar

Pour la mesure d'un temps, un signal connu $x(t)$ (une impulsion électromagnétique, optique ou sonore) est émis, et réfléchi par la cible dont on veut connaître la distance. Le signal reçu est proportionnel à $x(t - \Delta\tau)$, où $\Delta\tau$ est le paramètre à déterminer. On mesure ainsi la distance $d = \frac{1}{2}c\Delta\tau$, où c est la célérité des ondes utilisées.

Si la cible est mobile, on peut de plus mesurer sa vitesse radiale par effet Doppler : si le signal $x(t)$ émis est à bande étroite autour de la fréquence ν_0 , la fréquence des signaux reçus sera $\nu_0 + \Delta\nu$, avec $\frac{\Delta\nu}{\nu_0} = 2\frac{v_r}{c}$ (v_r : vitesse radiale de la cible).

Il est clair que les signaux adaptés à la mesure d'un temps (signaux large-bande, donc de courte durée) ne le sont pas nécessairement pour la mesure d'une fréquence (signaux à bande étroite, donc de grande durée). C'est l'ambiguïté temps-fréquence.

De façon générale, le signal observé $|r\rangle$ (ou $r(t)$) est formé du signal attendu $|x_{\theta_0}\rangle$ et de bruit $|b\rangle$, ce qui constitue le modèle de mesure

$$|r\rangle = |x_{\theta_0}\rangle + |b\rangle$$

- Le bruit $|b\rangle$ est un signal aléatoire, dont il convient de fixer les caractéristiques au mieux pour améliorer l'estimation.
- θ_0 est la valeur inconnue du paramètre, et $|x_{\theta_0}\rangle$ est un signal certain mais inconnu.

- $|r \rangle$ est le signal observé, c'est donc un signal certain.

Définition 4.1**(Estimateur)**

Un estimateur $\hat{\theta}(r)$ est une fonction certaine de l'observation $|r \rangle$, dont on voudrait que la valeur soit proche de θ_0 quand $|r \rangle = |x_{\theta_0} \rangle + |b \rangle$, où $|b \rangle$ est le bruit de mesure.

Nous précisons plus loin les propriétés désirables d'un estimateur.

4.1.2 Analogie avec le test d'hypothèse

Le problème de l'estimation de paramètres est formellement identique au test d'hypothèses. En définissant

- $I \doteq$ "Le signal mesuré est en moyenne $|r \rangle = |x_{\theta_0} \rangle$, où θ_0 est une valeur particulière du paramètre θ ."
- $H_\theta \doteq$ "Le paramètre a pour valeur θ ."
- $V \doteq$ "On a mesuré $|r \rangle$."

Il est clair que la spécification des propriétés statistiques de $|b \rangle$ dans l'approche classique équivaut à la donnée de $p(V|H_\theta I)$.

Le problème est alors de choisir une des hypothèses H_θ . Nous avons vu que pour cela nous pouvons utiliser le principe du Maximum a Posteriori :

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\text{Argmax}} p(H_\theta|VI)$$

ou si toutes les hypothèses sont équiprobables ($p(H_\theta|I) = \text{constante}$) le principe du **Maximum de Vraisemblance** :

$$\hat{\theta}_{\text{MV}} = \underset{\theta}{\text{Argmax}} p(V|H_\theta I)$$

Dans la suite, nous utiliserons l'estimateur du maximum de vraisemblance (nous considérerons donc toutes les hypothèses équiprobables).

Notation abusive :

$$p(r|\theta) \doteq p(V|H_\theta I)$$

4.2 Signal dans un bruit gaussien

Comme souvent, l'hypothèse que le bruit de mesure suit une statistique gaussienne (loi normale) permet de simplifier les calculs.

Les calculs qui suivent sont faits en notation vectorielle (signaux discrets). La généralisation aux signaux continus est immédiate.

$$|b \rangle = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}$$

On suppose que le bruit est centré, et que l'on connaît sa matrice de covariance¹ $\Gamma_{ij} = \langle b_i^* b_j \rangle$. Alors le principe du maximum d'entropie nous conduit à attribuer à ce bruit une distribution gaussienne :

$$p(b) = (\sqrt{2\pi})^{-m} |\Gamma|^{-1/2} \exp\left(-\frac{1}{2} b^\dagger \Gamma^{-1} b\right)$$

On a alors

$$p(r|\theta) = (\sqrt{2\pi})^{-m} |\Gamma|^{-1/2} \exp\left(-\frac{1}{2} (r - x_\theta)^\dagger \Gamma^{-1} (r - x_\theta)\right)$$

puisque $|b\rangle = |r\rangle - |x_\theta\rangle$ sous l'hypothèse H_θ .

Définissons la log-vraisemblance :

$$V(r|\theta) \doteq \ln p(r|\theta) = C^{ste} - \frac{1}{2} (r - x_\theta)^\dagger \Gamma^{-1} (r - x_\theta)$$

Chercher le maximum de $p(r|\theta)$ revient à chercher le maximum de $V(r|\theta)$, ou encore le minimum de $(r - x_\theta)^\dagger \Gamma^{-1} (r - x_\theta)$, qui est une forme quadratique.

Si de plus le bruit est blanc, c'est-à-dire que $\Gamma_{ij} = N_0 \delta_{ij}$, où N_0 est la puissance de bruit, alors

$$\begin{aligned} (r - x_\theta)^\dagger \Gamma^{-1} (r - x_\theta) &= \frac{1}{N_0} (r - x_\theta)^\dagger (r - x_\theta) \\ &= \frac{1}{N_0} \left(\sum_{i=1}^m |r_i - (x_\theta)_i|^2 \right) \\ &= \frac{1}{N_0} \|r - x_\theta\|^2 \end{aligned}$$

On en conclue donc

Propriété 4.1

(Moindres carrés)

|| Dans le cas de la mesure dans un bruit blanc gaussien, chercher le maximum de vraisemblance revient à minimiser l'écart quadratique entre le modèle $|x_\theta\rangle$ et l'observée $|r\rangle$

$$\hat{\theta}_{MC} \doteq \underset{\theta}{\text{Argmin}} \|r - x_\theta\|^2 = \underset{\theta}{\text{Argmin}} \langle r - x_\theta | r - x_\theta \rangle$$

Il s'agit dans ce cas de minimiser la distance entre $|r\rangle$ et $|x_\theta\rangle$.

4.3 Fonction d'ambiguïté

Nous nous plaçons toujours dans le cas du bruit blanc gaussien. Quelle est la qualité de l'estimateur des moindres carrés $\hat{\theta}_{MC}$?

1. Les calculs qui suivent sont effectués pour des signaux à valeurs complexes car ils serviraient également pour le formalisme de l'enveloppe complexe décrit plus loin. En particulier, le symbole \dagger représente le complexe conjugué transposé.

4.3.1 Moyenne et covariance de l'estimation des moindres carrés

Si on suppose $|r \rangle = |x_{\theta_0} \rangle + |b \rangle$,

$$\begin{aligned} \langle r - x_{\theta} | r - x_{\theta} \rangle &= \langle x_{\theta} - x_{\theta_0} - b | x_{\theta} - x_{\theta_0} - b \rangle \\ &= \langle x_{\theta} - x_{\theta_0} | x_{\theta} - x_{\theta_0} \rangle + \langle b | b \rangle - 2\Re \{ \langle x_{\theta} - x_{\theta_0} | b \rangle \} \end{aligned}$$

On suppose de plus que les perturbations apportées par le bruit restent faibles. Écrivons les termes de l'équation précédente en notations continue :

$$\begin{cases} \langle x_{\theta} - x_{\theta_0} | x_{\theta} - x_{\theta_0} \rangle = \int_{\mathbb{R}} |x(t, \theta) - x(t, \theta_0)|^2 dt \\ \langle b | b \rangle = \int_{\mathbb{R}} |b(t)|^2 dt \\ \langle x_{\theta} - x_{\theta_0} | b \rangle = \int_{\mathbb{R}} (x(t, \theta) - x(t, \theta_0))^* b(t) dt \end{cases}$$

Posons $\theta = \theta_0 + d\theta$, soit en notation développée

$$\begin{pmatrix} \theta_1 \\ \vdots \\ \theta_n \end{pmatrix} = \begin{pmatrix} (\theta_0)_1 \\ \vdots \\ (\theta_0)_n \end{pmatrix} + \begin{pmatrix} d\theta_1 \\ \vdots \\ d\theta_n \end{pmatrix}$$

Alors au premier ordre :

$$\begin{aligned} x(t, \theta) - x(t, \theta_0) &\approx \sum_{i=1}^n \frac{\partial x}{\partial \theta_i}(t, \theta_0) d\theta_i \\ |x(t, \theta) - x(t, \theta_0)|^2 &\approx \sum_{i=1}^n \sum_{j=1}^n \frac{\partial x^*}{\partial \theta_i}(t, \theta_0) \frac{\partial x}{\partial \theta_j}(t, \theta_0) d\theta_i d\theta_j \end{aligned}$$

On pose

$$P_i \doteq \Re \left\{ \int_{\mathbb{R}} \frac{\partial x^*}{\partial \theta_i}(t, \theta_0) b(t) dt \right\}$$

d'où

$$2\Re \{ \langle x_{\theta} - x_{\theta_0} | b \rangle \} \approx 2 \sum_{i=1}^n P_i d\theta_i = 2P^T \cdot d\theta$$

On pose de plus

$$A_{ij} \doteq - \int_{\mathbb{R}} \Re \left\{ \frac{\partial x^*}{\partial \theta_i}(t, \theta_0) \frac{\partial x}{\partial \theta_j}(t, \theta_0) dt \right\}$$

alors

$$\langle x_{\theta} - x_{\theta_0} | x_{\theta} - x_{\theta_0} \rangle \approx - \sum_{i=1}^n \sum_{j=1}^n A_{ij} d\theta_i d\theta_j = -d\theta^T \cdot A \cdot d\theta$$

Au total, au voisinage de θ_0 nous avons :

$$\begin{aligned} \langle r - x_{\theta} | r - x_{\theta} \rangle &= \langle b | b \rangle - d\theta^T \cdot A \cdot d\theta - 2P^T \cdot d\theta \\ &= \langle b | b \rangle + P^T \cdot A^{-1} \cdot P - (d\theta + A^{-1} \cdot P)^T \cdot A \cdot (d\theta + A^{-1} \cdot P) \end{aligned}$$

En effet, compte tenu du fait que A est symétrique, $A^T = A$ et $(A^{-1})^T = A^{-1}$. D'autre part $(A^{-1}.P)^T = P^T.(A^{-1})^T = P^T.A^{-1}$, et $P^T.d\theta = d\theta.P^T$, d'où

$$\begin{aligned} (d\theta + A^{-1}.P)^T.A.(d\theta + A^{-1}.P) &= d\theta^T.A.d\theta + (A^{-1}.P)^T.A.d\theta \\ &\quad + d\theta^T.A.(A^{-1}.P) + (A^{-1}.P)^T.A.(A^{-1}.P) \\ &= d\theta^T.A.d\theta + P^T.A^{-1}.A.d\theta \\ &\quad + d\theta^T.A.A^{-1}.P + P^T.A^{-1}.A.A^{-1}.P \\ (d\theta + A^{-1}.P)^T.A.(d\theta + A^{-1}.P) &= d\theta^T.A.d\theta + 2P^T.d\theta + P^T.A^{-1}.P \end{aligned}$$

La matrice $-A$ est définie positive², donc $\langle r - x_\theta | r - x_\theta \rangle$ est minimum si $d\theta + A^{-1}.P = 0$, soit

$$\boxed{d\theta = -A^{-1}.P}$$

C'est-à-dire que sous l'effet de la perturbation P due au bruit, le minimum se décale de $d\theta$:

$$\hat{\theta}_{\text{MC}} \approx \theta_0 - A^{-1}.P$$

On peut maintenant chercher ce qui se passe en moyenne quand on considère tous les bruits possibles, ou tous les signaux attendus possibles. L'espérance $E[\hat{\theta}_{\text{MC}}]$ est la valeur prédite en moyenne par l'estimateur des moindres carrés.

$$E[\hat{\theta}_{\text{MC}}] = \int \hat{\theta}_{\text{MC}} p(r|\theta) dr = \theta_0 - A^{-1}.E[P] = \theta_0$$

car l'opérateur valeur moyenne est linéaire et le bruit centré. On dit que l'estimateur est non-biaisé.

On peut de même définir la covariance de l'estimation par :

$$\Gamma_{\hat{\theta}_{\text{MC}}} = E[(\hat{\theta}_{\text{MC}} - \theta_0).(\hat{\theta}_{\text{MC}} - \theta_0)^T]$$

d'où

$$\Gamma_{\hat{\theta}_{\text{MC}}} = A^{-1}.E[P.P^T].A^{-1}$$

avec $E[P.P^T]$ la covariance des perturbations apportées par le bruit.

$$E[P_i.P_j^*] = \int_{\mathbb{R}} \int_{\mathbb{R}} \frac{\partial x^*}{\partial \theta_i}(t, \theta_0) \frac{\partial x}{\partial \theta_j}(u, \theta_0) E[b(t)b^*(u)] dt du$$

Or $E[b(t)b^*(u)] = N_0\delta(t - u)$ puisque le bruit est blanc, donc

$$E[P_i.P_j^*] = N_0 \int_{\mathbb{R}} \frac{\partial x^*}{\partial \theta_i}(t, \theta_0) \frac{\partial x}{\partial \theta_j}(t, \theta_0) dt = -N_0 A_{ij}$$

et donc

$$\boxed{\Gamma_{\hat{\theta}_{\text{MC}}} = -N_0.A^{-1}}$$

2. Si une matrice M est définie positive alors $u^\dagger.M.u \geq 0$ quel que soit le vecteur u .

4.3.2 Relation à la fonction d'ambiguïté

La fonction d'ambiguïté pour un paramètre θ est généralement définie dans le cas où l'énergie du modèle du signal ne dépend pas de θ :

$$\langle x_\theta | x_\theta \rangle = C^{ste}$$

On peut écrire

$$\langle r - x_\theta | r - x_\theta \rangle = \langle r | r \rangle + \langle x_\theta | x_\theta \rangle - 2\Re \{ \langle x_\theta | r \rangle \}$$

Le minimum de cette forme quadratique correspond au maximum de $\Re \{ \langle x_\theta | r \rangle \}$, puisque l'énergie du modèle du signal est constante et que $\langle r | r \rangle$ ne dépend pas de θ . On écrit encore

$$\Re \{ \langle x_\theta | r \rangle \} = \chi(\theta, \theta_0) + \Re \{ \langle x_\theta | b \rangle \}$$

en tenant compte de ce que $|r\rangle = |x_{\theta_0}\rangle + |b\rangle$.

Définition 4.2

(Fonction d'ambiguïté)

Dans le cas où l'énergie du modèle du signal ne dépend pas de θ , la fonction d'ambiguïté est définie par :

$$\chi(\theta, \theta_0) \doteq \Re \{ \langle x_\theta | x_{\theta_0} \rangle \} = \Re \left\{ \int_{\mathbb{R}} x^*(t, \theta) x(t, \theta_0) dt \right\}$$

La fonction d'ambiguïté est maximum pour $\theta = \theta_0$. En effet, en utilisant l'inégalité de Schwartz (produit scalaire) :

$$|\chi(\theta, \theta_0)|^2 = | \langle x_\theta | x_{\theta_0} \rangle |^2 \leq \langle x_\theta | x_\theta \rangle \langle x_{\theta_0} | x_{\theta_0} \rangle = |\chi(\theta_0, \theta_0)|^2$$

D'autre part, on a :

$$A_{ij} = \frac{\partial^2 \chi}{\partial \theta_i \partial \theta_j}(\theta_0, \theta_0)$$

ce qui fournit un autre moyen de calculer la covariance de l'estimation des moindres carrés.

En effet,

$$\begin{aligned} \int_{\mathbb{R}} |x(t, \theta)|^2 dt = C^{ste} &\Rightarrow 2\Re \left\{ \int_{\mathbb{R}} \left(\frac{\partial x^*}{\partial \theta_i}(t, \theta) \frac{\partial x}{\partial \theta_j}(t, \theta) + \frac{\partial^2 x^*}{\partial \theta_i \partial \theta_j}(t, \theta) x(t, \theta) \right) dt \right\} = 0 \\ &\Rightarrow -A_{ij} + \frac{\partial^2 \chi}{\partial \theta_i \partial \theta_j}(\theta_0, \theta_0) = 0 \end{aligned}$$

pour $\theta = \theta_0$.

En particulier, dans le cas de l'estimation d'un paramètre scalaire :

$$\sigma_{\hat{\theta}_{MC}}^2 = - \frac{N_0}{\frac{d^2 \chi}{d\theta^2}(\theta_0, \theta_0)}$$

Puisque la dérivée de la fonction d'ambiguïté est nulle en θ_0 (au maximum), la dérivée seconde définit un rayon de courbure de la parabole osculatrice en θ_0 . On en conclue donc qu'une bonne estimation correspond à une fonction d'ambiguïté très piquée, c'est-à-dire à une forte valeur de la dérivée seconde.

4.3.3 Fonction d'ambiguïté modifiée

Dans le cas où l'énergie du modèle du signal n'est pas constante, donc dépend de θ , il est toujours possible de définir une fonction d'ambiguïté par :

Définition 4.3

(Fonction d'ambiguïté modifiée)

Dans le cas où l'énergie du modèle du signal dépend de θ , la fonction d'ambiguïté est définie par :

$$\chi(\theta, \theta_0) \doteq \Re \{ \langle x_\theta | x_{\theta_0} \rangle \} - \frac{1}{2} \langle x_\theta | x_\theta \rangle$$

Il est facile de montrer que les propriétés de la fonction d'ambiguïté démontrées précédemment restent vraies. En effet, on a toujours $\chi(\theta, \theta_0) \leq \chi(\theta_0, \theta_0)$, puisque

$$\begin{aligned} \chi(\theta, \theta_0) - \chi(\theta_0, \theta_0) &= \Re \{ \langle x_\theta | x_{\theta_0} \rangle \} - \frac{1}{2} \langle x_\theta | x_\theta \rangle - \frac{1}{2} \langle x_{\theta_0} | x_{\theta_0} \rangle \\ &= -\frac{1}{2} \langle x_\theta - x_{\theta_0} | x_\theta - x_{\theta_0} \rangle \\ \chi(\theta, \theta_0) - \chi(\theta_0, \theta_0) &\leq 0 \end{aligned}$$

D'autre part

$$\begin{aligned} \frac{\partial^2 \chi}{\partial \theta_i \partial \theta_j}(\theta, \theta) &= \Re \left\{ \int_{\mathbb{R}} \frac{\partial^2 x^*}{\partial \theta_i \partial \theta_j}(t, \theta) x(t, \theta) dt \right\} \\ &\quad - \frac{1}{2} \times 2 \Re \left\{ \int_{\mathbb{R}} \left(\frac{\partial x^*}{\partial \theta_i}(t, \theta) \frac{\partial x}{\partial \theta_j}(t, \theta) + \frac{\partial^2 x^*}{\partial \theta_i \partial \theta_j}(t, \theta) x(t, \theta) \right) dt \right\} \end{aligned}$$

d'où pour $\theta = \theta_0$:

$$\frac{\partial^2 \chi}{\partial \theta_i \partial \theta_j}(\theta_0, \theta_0) = A_{ij}$$

4.4 Mesure d'un temps

4.4.1 Mesure avec un signal quelconque

Le signal émis $x(t)$ est reçu avec un retard $\Delta\tau$, qui est le paramètre à estimer, de sorte que (on ne prend pas en compte une atténuation éventuelle)

$$|x_{\Delta\tau}\rangle = x(t - \Delta\tau)$$

L'énergie de ce signal est

$$\langle x_{\Delta\tau} | x_{\Delta\tau} \rangle = \int_{\mathbb{R}} |x(t - \Delta\tau)|^2 dt = \int_{\mathbb{R}} |x(t)|^2 dt = E_x$$

L'énergie est donc indépendante du paramètre, et la fonction d'ambiguïté est :

$$\chi(\Delta\tau, \Delta\tau_0) = \Re \left\{ \int_{\mathbb{R}} x^*(t - \Delta\tau) x(t - \Delta\tau_0) dt \right\}$$

Soit la fonction de corrélation de $x(t)$:

$$C_{XX}(\tau) \doteq \int_{\mathbb{R}} x^*(t - \tau) \cdot x(t) dt$$

alors

$$\chi(\Delta\tau, \Delta\tau_0) = \Re \{ C_{XX}(\Delta\tau - \Delta\tau_0) \}$$

Le principal intérêt de la fonction de corrélation est qu'elle est la transformée de Fourier inverse de la densité spectrale :

$$C_{XX}(\tau) = \text{TF}^{-1} [|\tilde{x}(\nu)|^2]$$

où $\tilde{x}(\nu)$ est la TF de $x(t)$ définie par

$$\begin{cases} \tilde{x}(\nu) = \int_{\mathbb{R}} x(t) \exp(-2i\pi\nu t) dt \\ x(t) = \int_{\mathbb{R}} \tilde{x}(\nu) \exp(+2i\pi\nu t) d\nu \end{cases}$$

La variance de l'estimation du temps de retour est donnée par la dérivée seconde de la fonction d'ambiguïté :

$$\frac{d^2\chi}{d\Delta\tau^2}(\Delta\tau_0, \Delta\tau_0) = \Re \left\{ \frac{d^2 C_{XX}}{d\tau^2}(0) \right\}$$

$$C_{XX}(\tau) = \int_{\mathbb{R}} |\tilde{x}(\nu)|^2 \exp(2i\pi\nu\tau) d\nu$$

$$\frac{dC_{XX}}{d\tau}(\tau) = (2i\pi) \int_{\mathbb{R}} \nu |\tilde{x}(\nu)|^2 \exp(2i\pi\nu\tau) d\nu$$

$$\frac{d^2 C_{XX}}{d\tau^2}(\tau) = (-4\pi^2) \int_{\mathbb{R}} \nu^2 |\tilde{x}(\nu)|^2 \exp(2i\pi\nu\tau) d\nu$$

d'où l'on tire

$$\frac{d^2\chi}{d\Delta\tau^2}(\Delta\tau_0, \Delta\tau_0) = \Re \left\{ \frac{d^2 C_{XX}}{d\tau^2}(0) \right\} = (-4\pi^2) \int_{\mathbb{R}} \nu^2 |\tilde{x}(\nu)|^2 d\nu$$

On a donc

$$\sigma_{\Delta\tau}^2 = - \frac{N_0}{\frac{d^2\chi}{d\Delta\tau^2}(\Delta\tau_0, \Delta\tau_0)} = \frac{N_0}{4\pi^2 \int_{\mathbb{R}} \nu^2 |\tilde{x}(\nu)|^2 d\nu}$$

Propriété 4.2

(Estimation d'un temps de retour)

La variance de l'estimation d'un temps de retour est

$$\sigma_{\Delta\tau}^2 = \frac{1}{4\pi^2 \overline{\nu^2}} \frac{N_0}{E_x}$$

où $\overline{\nu^2}$ est l'épanouissement fréquentiel défini par

$$\overline{\nu^2} \doteq \frac{\int_{\mathbb{R}} \nu^2 |\tilde{x}(\nu)|^2 d\nu}{\int_{\mathbb{R}} |\tilde{x}(\nu)|^2 d\nu}$$

E_x est l'énergie du signal, et E_x/N_0 le rapport signal à bruit.

4.4.2 Mesure avec un signal passe-bande : enveloppe complexe

En pratique, le signal $x(t)$ utilisé est souvent constitué à partir d'une porteuse modulée. Le signal utilisé pour l'estimation ne sera pas en général le signal émis $x(t)$ décalé en temps, mais l'enveloppe de ce signal. On doit dans ce cas obtenir la qualité de l'estimation à partir de l'enveloppe complexe du signal plutôt qu'à partir du signal lui-même. Ce procédé intuitif peut être introduit rigoureusement à l'aide de la notion de signal analytique complexe.

Définition 4.4

(Signal analytique)

Le signal analytique $z(t)$ (fonction à valeurs complexes) est associé au signal réel $x(t)$ par :

$$\tilde{z}(\nu) = 2H(\nu)\tilde{x}(\nu)$$

où $H(\nu)$ est la fonction de Heaviside (nulle pour $\nu < 0$, égale à 1 pour $\nu \geq 0$).

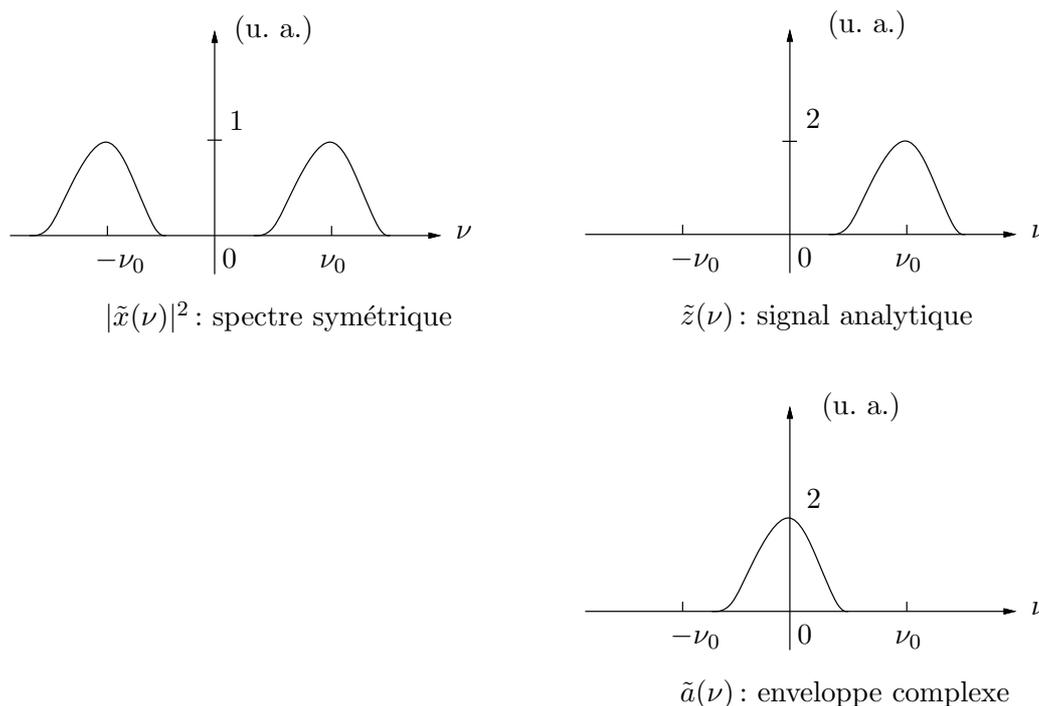


FIG. 4.1 – Extraction de l'enveloppe complexe d'un signal réel.

Pour calculer le signal analytique, il suffit donc de ne conserver que les fréquences positives du spectre initial, et de multiplier par deux le résultat. Cette définition qui peut apparaître arbitraire est éclairée par la propriété suivante :

Propriété 4.3

(Théorème de Bedrosian)

Soit le signal réel $x(t) = \alpha(t) \cos(2\pi\nu_0 t + \varphi(t))$, où α et φ sont des fonctions réelles. Si la fonction complexe $a(t) = \alpha(t) \exp(i\varphi(t))$ est telle que $|\tilde{a}(\nu)| = 0$ si $\nu < -\nu_0$, alors

$$z(t) = \alpha(t) \exp(2i\pi\nu_0 t + i\varphi(t)) = a(t) \exp(2i\pi\nu_0 t)$$

$a(t)$ est l'enveloppe complexe du signal réel $x(t)$, et $\exp(2i\pi\nu_0 t)$ est la porteuse.

En principe, la valeur de ν_0 n'est pas imposée. Il est d'usage de choisir ν_0 telle que

$$\int \nu |\tilde{a}(\nu)|^2 d\nu = 0$$

c'est-à-dire que l'on centre le spectre de l'enveloppe complexe (signal passe-bas).

Remarque : Le signal analytique fournit une définition rigoureuse du procédé usuel en physique qui consiste à remplacer une fonction harmonique réelle du type $\cos(\omega t)$ par $\exp(i\omega t)$ (ondes planes en optique, en électromagnétisme ou en acoustique, signal sinusoïdal en électricité,...).

Calculons la fonction de corrélation du signal $x(t)$ en fonction de celle de son enveloppe complexe $a(t)$:

$$\begin{aligned} C_{XX}(\tau) &= \int |\tilde{x}(\nu)|^2 \exp(2i\pi\nu\tau) d\nu \\ &= \int_{\mathbb{R}^+} |\tilde{x}(\nu)|^2 \exp(2i\pi\nu\tau) d\nu + \int_{\mathbb{R}^+} |\tilde{x}(\nu)|^2 \exp(-2i\pi\nu\tau) d\nu \\ &= \frac{1}{4} \int_{\mathbb{R}} |\tilde{z}(\nu)|^2 \exp(2i\pi\nu\tau) d\nu + \frac{1}{4} \int_{\mathbb{R}} |\tilde{z}(\nu)|^2 \exp(-2i\pi\nu\tau) d\nu \\ &= \frac{1}{4} C_{ZZ}(\tau) + \frac{1}{4} C_{ZZ}(-\tau) \\ &= \frac{1}{4} \int z^*(t-\tau).z(t) dt + \frac{1}{4} \int z^*(t+\tau).z(t) dt \\ &= \frac{1}{4} \int z^*(t-\tau).z(t) dt + \frac{1}{4} \int z^*(t).z(t-\tau) dt \\ &= \frac{1}{4} \int a^*(t-\tau).a(t) dt \exp(2i\pi\nu_0\tau) + \frac{1}{4} \int a^*(t).a(t-\tau) dt \exp(-2i\pi\nu_0\tau) \\ C_{XX}(\tau) &= \frac{1}{2} \Re \{ C_{AA}(\tau) \exp(2i\pi\nu_0\tau) \} \end{aligned}$$

Cette expression signifie que le terme lentement variable $C_{AA}(\tau)$ est rapidement modulé le terme exponentiel. En pratique, au voisinage d'un temps τ donné, la valeur de $C_{XX}(\tau)$ varie très rapidement entre $\frac{1}{2}|C_{AA}(\tau)|$ et $-\frac{1}{2}|C_{AA}(\tau)|$, de sorte que l'on peut écrire (en négligeant l'influence de la phase de $C_{XX}(\tau)$, qui décale le système de franges d'une quantité négligeable devant le nombre de périodes décrites) :

$$C_{XX}(\tau) \approx \frac{1}{2} |C_{AA}(\tau)| \cos(2\pi\nu_0\tau)$$

$\frac{1}{2}|C_{AA}(\tau)|$ est donc l'enveloppe de la fonction de corrélation du signal. En pratique, les systèmes d'analyse ne possèdent pas une résolution suffisante pour distinguer les franges modulatrices, et

ne permettent de détecter que l'enveloppe de la fonction de corrélation. On en déduit immédiatement que dans ce cas :

Propriété 4.4 (Estimation d'un temps de retour par l'enveloppe complexe)

La variance de l'estimation d'un temps de retour par l'enveloppe complexe est

$$\sigma_{\Delta\tau}^2 = \frac{1}{4\pi^2 \overline{\nu^2}} \frac{N_0}{E_x}$$

où $\overline{\nu^2}$ est l'épanouissement fréquentiel

$$\overline{\nu^2} \doteq \frac{\int_B \nu^2 |\tilde{a}(\nu)|^2 d\nu}{\int_B |\tilde{a}(\nu)|^2 d\nu}$$

$E_x = E_a/2$ est l'énergie du signal, et E_x/N_0 le rapport signal à bruit (N_0 est la puissance de bruit dans la bande B autour de la fréquence ν_0 de la porteuse).

On peut remarquer que dans le cas où l'on est capable de résoudre les franges de modulation, l'épanouissement fréquentiel est bien plus grand :

$$\overline{\nu^2} = \frac{\int_B \nu^2 |\tilde{a}(\nu)|^2 d\nu}{\int_B |\tilde{a}(\nu)|^2 d\nu} + \nu_0^2$$

Cependant cette situation théorique se heurte à de graves difficultés techniques !

4.4.3 Conception de signaux adaptés à la mesure d'un temps

D'après l'étude précédente, il est clair que pour une énergie du signal fixée et pour une bande donnée (donc pour un rapport signal à bruit donné) il faut utiliser un signal qui maximise l'épanouissement fréquentiel

$$\frac{\int_B \nu^2 |\tilde{a}(\nu)|^2 d\nu}{\int_B |\tilde{a}(\nu)|^2 d\nu}$$

sous contrainte que

$$\int_B \nu |\tilde{a}(\nu)|^2 d\nu = 0$$

où B est la bande disponible. Ce problème n'a pas de solution simple, et il faut ajouter en pratique des contraintes de réalisation.

Puissance constante sur la bande – Si $|\tilde{a}(\nu)|$ est constante, alors

$$\tilde{a}(\nu) = \sqrt{\frac{E_a}{B}} \exp(i\varphi(\nu))$$

où la phase $\varphi(\nu)$ n'est pas imposée. Il est donc possible de jouer sur cette phase pour améliorer les caractéristiques de la détection.

Puissance et phase constantes – Si $\varphi(\nu) = 0$, alors $\tilde{a}(\nu) = \sqrt{\frac{E_a}{B}}$ et

$$a(t) = \int_{-B/2}^{B/2} \tilde{a}(\nu) \exp(2i\pi\nu t) d\nu = \sqrt{\frac{E_a}{B}} \text{TF} [\text{rect}(\nu/B)] = \sqrt{BE_a} \text{sinc}(Bt)$$

De plus

$$\int_{-B/2}^{B/2} \nu^2 |\tilde{a}(\nu)|^2 d\nu = \frac{E_a}{B} \left[\frac{1}{3} \left(2 \left(\frac{B}{2} \right)^3 \right) \right] = \frac{E_a B^2}{12}$$

d'où

$$\overline{\nu^2} = \frac{B^2}{12}$$

et

$$\sigma_{\Delta\tau}^2 = \frac{3}{\pi^2 B^2} \frac{N_0}{E_x}$$

Le défaut de ce signal est que sa puissance temporelle $|a(t)|^2 = BE_a \text{sinc}^2(Bt)$ est très concentrée à l'origine (cette fonction tend vers un Dirac quand B tend vers l'infini), ce qui rend sa production et son émission plus difficile. Pour éviter cette concentration de l'énergie on peut jouer sur la phase $\varphi(\nu)$.

Signal “chirp” – Un signal chirp est une gaussienne de phase.

$$\varphi(\nu) = \pi\alpha\nu^2$$

soit

$$\tilde{a}(\nu) = \sqrt{\frac{E_a}{B}} \exp(i\pi\alpha\nu^2)$$

On en déduit

$$a(t) = \sqrt{\frac{E_a}{B}} \text{TF} [\text{rect}(\nu/B) \exp(i\pi(\sqrt{\alpha}\nu)^2)] \approx \sqrt{\frac{E_a}{\alpha B}} \exp\left(i\pi \frac{t^2}{\alpha}\right)$$

pourvu que $\alpha \gg 1/B^2$, c'est-à-dire si la gaussienne oscille beaucoup dans la bande³. Le domaine de définition de t où cette approximation est valable est limité à l'intervalle $[-\alpha B/2, \alpha B/2]$, en dehors duquel $a(t) \approx 0$. La variance de l'estimation du temps $\sigma_{\Delta\tau}^2$ est inchangée par rapport au signal à phase constante (puisque on a toujours $|\tilde{a}(\nu)| = 1$), mais l'énergie temporelle du signal est répartie sur l'ensemble de sa durée :

$$|a(t)|^2 \approx \frac{E_a}{\alpha B}$$

3. Cette situation possède une analogie optique immédiate. L'ombre portée à la distance d derrière une ouverture rectangulaire est donnée par la diffraction de Fresnel, soit exactement l'expression précédente avec B la largeur de l'ouverture, $\alpha = 1/\lambda d$, et t la fréquence spatiale. Si l'ouverture est très grande devant la longueur d'onde ou la distance d'observation suffisamment proche, l'ombre conserve la forme rectangulaire.

Une application commerciale – Si l'enveloppe $a(t)$ est modulée par un code $\beta(t)$ (par exemple des sauts "aléatoires" entre 1 et -1)

$$z(t) = \beta(t)a(t) \exp(2i\pi\nu_0 t)$$

pour réaliser la corrélation de l'enveloppe, il faudra connaître le code.

Un satellite émettant ce signal codé permettra par exemple de se positionner sur Terre, mais le prix de vente dépendra de la longueur de code achetée, puisque la précision de localisation est inversement proportionnelle au carré de la durée acquise, ou inversement proportionnelle au carré de la bande. C'est par exemple le cas du GPS (*Global Positioning System*), où la longueur de code vendue pour les applications civiles (navigation, cartographie...) est plus petite que celle que se réserve les militaires.

4.5 Mesure d'une fréquence

Le problème de la mesure d'une fréquence est similaire à celui de la mesure d'un temps. On émet le signal $x(t)$. Après réflexion sur une cible en mouvement, le signal reçu est translaté en fréquence de :

$$\frac{\Delta f}{f} = 2 \frac{v_r}{c}$$

où v_r est la vitesse radiale de la cible et c la célérité des ondes utilisées (effet Doppler classique). Il est clair que la translation de fréquence n'est proportionnelle à v_r que lorsque le signal $x(t)$ est constitué d'une enveloppe à variations lentes modulée sur une porteuse, car alors

$$\frac{\Delta f}{f} \approx \frac{\Delta f}{\nu_0} \approx 2 \frac{v_r}{c}$$

De plus, dans ce cas les fréquences positives sont translatées de $+\Delta f$, et les fréquences négatives de $-\Delta f$, ce qui rend le formalisme lourd sans utilisation du signal analytique $z(t) = a(t) \exp(2i\pi\nu_0 t)$, qui n'utilise que les fréquences positives du signal de départ. On écrit le signal reçu sous la forme $|\tilde{Z}_{\Delta f} \rangle$ ou $\tilde{z}(\nu - \Delta f)$. L'énergie de ce signal est indépendante du paramètre à estimer (Δf) :

$$\langle \tilde{Z}_{\Delta f} | \tilde{Z}_{\Delta f} \rangle = \int |\tilde{z}(\nu - \Delta f)|^2 d\nu = \int |\tilde{z}(\nu)|^2 d\nu = \langle Z_0 | Z_0 \rangle$$

La fonction d'ambiguïté est donc

$$\chi(\Delta f, \Delta f_0) \doteq \Re \left\{ \int \tilde{z}^*(\nu - \Delta f) \tilde{z}(\nu - \Delta f_0) d\nu \right\}$$

soit en utilisant la fonction de corrélation en fréquence

$$\begin{aligned} C_{\tilde{Z}\tilde{Z}}(f) &\doteq \int \tilde{z}^*(\nu - f) \tilde{z}(\nu) d\nu = \langle \tilde{Z}_f | \tilde{Z}_0 \rangle \\ C_{\tilde{Z}\tilde{Z}}(f) &= \int |z(t)|^2 \exp(2i\pi f t) dt \\ \chi(\Delta f, \Delta f_0) &= \Re \{ C_{\tilde{Z}\tilde{Z}}(\Delta f - \Delta f_0) \} \end{aligned}$$

On obtient donc de façon similaire à l'estimation d'un temps de retour :

Propriété 4.5 (Estimation d'une fréquence par l'enveloppe complexe)

La variance de l'estimation d'une fréquence par l'enveloppe complexe est

$$\sigma_{\widehat{\Delta f}}^2 = \frac{1}{4\pi^2 \overline{\Delta t^2}} \frac{N_0}{E_x}$$

où $\overline{\Delta t^2}$ est l'épanouissement temporel défini par

$$\overline{\Delta t^2} \doteq \frac{\int (t - t_0)^2 |a(t)|^2 dt}{\int |a(t)|^2 dt}$$

et

$$t_0 \doteq \frac{\int t |a(t)|^2 dt}{\int |a(t)|^2 dt}$$

Il est donc évident qu'à l'inverse des signaux adaptés à la mesure d'un temps, les signaux adaptés à la mesure d'une fréquence doivent avoir une grande durée, donc un spectre étroit.

Remarque – On peut toujours imposer $t_0 = 0$ en décalant l'origine des temps. Ainsi les équations sont complètement symétriques de celles de la mesure d'un temps.

4.6 Inégalité de Cramer-Rao

Jusqu'à présent, nous avons considéré le cas particulier du bruit blanc gaussien, et nous avons obtenu la variance de l'estimation du maximum de vraisemblance dans ce cas :

$$\Gamma_{\hat{\theta}_{\text{MC}}} = -N_0 \cdot A^{-1}$$

avec

$$A_{ij} = \frac{\partial^2 \chi}{\partial \theta_i \partial \theta_j}(\theta_0, \theta_0) .$$

Existe-t'il une relation semblable dans le cas où la vraisemblance $p(r|\theta)$ n'est pas nécessairement une gaussienne? C'est l'inégalité de Cramer-Rao.

Propriété 4.6**(Inégalité de Cramer-Rao (paramètre scalaire))**

$$E \left[|\hat{\theta}(r) - \theta|^2 \right] \geq \frac{\left| \frac{d}{d\theta} E \left[\hat{\theta}(r) \right] \right|^2}{E \left[\left| \frac{d}{d\theta} \ln p(r|\theta) \right|^2 \right]} = \frac{\left| \frac{d}{d\theta} E \left[\hat{\theta}(r) \right] \right|^2}{E \left[\frac{d^2}{d\theta^2} \ln p(r|\theta) \right]}$$

où $\hat{\theta}(r)$ est un estimateur quelconque du paramètre scalaire θ ($\hat{\theta}(r)$ ne dépend que des données r), et $E[.]$ est l'espérance relative à la vraisemblance $p(r|\theta)$:

$$E[f(r)] \doteq \int f(r)p(r|\theta)dr$$

où $f(r)$ est une fonction quelconque de r . $E \left[\frac{d^2}{d\theta^2} \ln p(r|\theta) \right]$ est appelée courbure de la log-vraisemblance.

Preuve – Par définition :

$$E[\hat{\theta}(r)] = \int \hat{\theta}(r)p(r|\theta)dr$$

donc

$$\frac{d}{d\theta} E[\hat{\theta}(r)] = \int \hat{\theta}(r) \frac{d}{d\theta} p(r|\theta)dr$$

De plus

$$\int \frac{d}{d\theta} p(r|\theta)dr = 0$$

donc

$$\int \theta \frac{d}{d\theta} p(r|\theta)dr = 0$$

d'où

$$\underline{\underline{\frac{d}{d\theta} E[\hat{\theta}(r)] = \int (\hat{\theta}(r) - \theta) \frac{d}{d\theta} p(r|\theta)dr}}$$

Maintenant

$$\frac{d}{d\theta} \ln p(r|\theta) = \frac{1}{p(r|\theta)} \frac{d}{d\theta} p(r|\theta)$$

donc

$$\underline{\underline{\frac{d}{d\theta} E[\hat{\theta}(r)] = \int (\hat{\theta}(r) - \theta)p(r|\theta) \frac{d}{d\theta} \ln p(r|\theta)dr}}$$

Si on pose :

$$\begin{aligned} f(r) &\doteq (\hat{\theta}(r) - \theta) \sqrt{p(r|\theta)} \\ g(r) &\doteq \sqrt{p(r|\theta)} \frac{d}{d\theta} \ln p(r|\theta) \end{aligned}$$

on voit alors que

$$\frac{d}{d\theta} E[\hat{\theta}(r)] = \langle f|g \rangle .$$

L'inégalité de Cauchy-Schwartz permet d'écrire

$$|\langle f|g \rangle|^2 \leq \langle f|f \rangle \cdot \langle g|g \rangle ,$$

donc

$$\left| \frac{d}{d\theta} E [\hat{\theta}(r)] \right|^2 \leq \left(\int |\hat{\theta}(r) - \theta|^2 p(r|\theta) dr \right) \cdot \left(\int \left| \frac{d}{d\theta} \ln p(r|\theta) \right|^2 p(r|\theta) dr \right)$$

soit

$$\left| \frac{d}{d\theta} E [\hat{\theta}(r)] \right|^2 \leq E [|\hat{\theta}(r) - \theta|^2] \cdot E \left[\left| \frac{d}{d\theta} \ln p(r|\theta) \right|^2 \right]$$

ce qui montre la première partie de l'inégalité.

Puisque

$$\int \frac{d}{d\theta} p(r|\theta) dr = 0$$

on a aussi

$$\int p(r|\theta) \frac{d}{d\theta} \ln p(r|\theta) dr = 0$$

et en dérivant encore une fois

$$\underbrace{\int \frac{d}{d\theta} p(r|\theta) \frac{d}{d\theta} \ln p(r|\theta) dr}_{\int p(r|\theta) \left| \frac{d}{d\theta} \ln p(r|\theta) \right|^2 dr} + \int p(r|\theta) \frac{d^2}{d\theta^2} \ln p(r|\theta) dr = 0$$

et donc

$$E \left[\left| \frac{d}{d\theta} \ln p(r|\theta) \right|^2 \right] = -E \left[\frac{d^2}{d\theta^2} \ln p(r|\theta) \right]$$

ce qui montre la seconde partie de l'inégalité. Par définition, $E \left[\frac{d^2}{d\theta^2} \ln p(r|\theta) \right]$ est la courbure de la log-vraisemblance, et joue un rôle similaire à celui de la dérivée seconde de la fonction d'ambiguïté dans le cas de l'estimation des moindres carrés.

L'inégalité de Cramer-Rao est valable quel que soit l'estimateur $\hat{\theta}(r)$, c'est-à-dire pour une fonction quelconque des données ! Elle ne devient intéressante que si l'on impose des contraintes supplémentaires sur la qualité de cette estimation, ce qui conduit aux définitions suivantes.

Définition 4.5

(Estimateur non biaisé)

Un estimateur non biaisé est tel que $E[\hat{\theta}(r)] = \theta$ si θ est la vraie valeur du paramètre. C'est-à-dire que pour θ fixée, la moyenne de $\hat{\theta}(r)$ prise sur toutes les réalisations possibles de la mesure est égale à θ .

L'estimateur des moindres carrés $\hat{\theta}_{MC}$ dans le cas du bruit blanc gaussien est non biaisé, ainsi que nous l'avons montré plus haut. Dans le cas général, rien ne prouve que l'estimateur du maximum de vraisemblance soit non biaisé : cela dépend de la forme de la vraisemblance $p(r|\theta)$.

Définition 4.6

(Estimateur efficace)

Un estimateur efficace est un estimateur non-biaisé qui atteint les bornes de Cramer-Rao.

Dans le cas d'un estimateur non biaisé, $\frac{d}{d\theta} E[\hat{\theta}(r)] = 1$, et donc

$$E[|\hat{\theta}(r) - \theta|^2] \geq -\frac{1}{E\left[\frac{d^2}{d\theta^2} \ln p(r|\theta)\right]}$$

L'égalité (estimateur efficace) ne peut être vérifiée que si $f(r)$ est colinéaire à $g(r)$ d'après la démonstration précédente, soit $g(r) = k(\theta)f(r)$ où $k(\theta)$ est une fonction quelconque de θ ne dépendant pas de r , soit

$$\frac{d}{d\theta} \ln p(r|\theta) = k(\theta)(\hat{\theta}(r) - \theta)$$

Considérons maintenant $\theta = \hat{\theta}_{MV}$. Pour cette valeur, on a par définition

$$\frac{d}{d\theta} \ln p(r|\theta = \hat{\theta}_{MV}) = 0$$

puisque $\theta = \hat{\theta}_{MV}$ est le maximum de la vraisemblance $p(r|\theta)$, et donc

$$\hat{\theta}(r) = \hat{\theta}_{MV}$$

si $\hat{\theta}(r)$ est un estimateur efficace, ce qui montre la propriété suivante.

Propriété 4.7

(Estimateur efficace et maximum de vraisemblance)

S'il existe un estimateur efficace pour un problème d'estimation de paramètre, alors il s'identifie avec l'estimateur du maximum de vraisemblance :

$$\hat{\theta}(r) = \hat{\theta}_{MV}$$

Dans ce cas, l'estimateur du maximum de vraisemblance est donc également efficace.

Il faut remarquer que cette propriété est seulement une implication et n'admet pas de réciproque en général. Par ailleurs, dans tout ce qui précède, on peut remplacer la vraisemblance $p(r|\theta)$ par la probabilité postérieure $p(\theta|r)$, et l'on obtiendra qu'un estimateur efficace au sens de $p(\theta|r)$ s'identifie avec l'estimateur MAP.

Correspondance avec le cas du bruit blanc gaussien – Vérifions que l'inégalité de Cramer-Rao permet bien de retrouver le résultat obtenu directement en supposant le bruit blanc gaussien.

$$\begin{aligned} \ln p(r|\theta) &= C^{ste} - \frac{1}{2N_0} \|r - x_\theta\|^2 \\ &= C^{ste} - \frac{1}{2N_0} \langle r - x_\theta | r - x_\theta \rangle \\ &= C^{ste} - \frac{1}{2N_0} \langle x_\theta - x_{\theta_0} - b | x_\theta - x_{\theta_0} - b \rangle \\ \ln p(r|\theta) &= C^{ste} - \frac{1}{2N_0} (\langle x_\theta - x_{\theta_0} | x_\theta - x_{\theta_0} \rangle + \langle b | b \rangle - 2\Re\{\langle x_\theta - x_{\theta_0} | b \rangle\}) \end{aligned}$$

Donc si le bruit est blanc et centré

$$E[\ln p(r|\theta)] = C^{ste} - \frac{1}{2N_0} (\langle x_\theta - x_{\theta_0} | x_\theta - x_{\theta_0} \rangle + N_0)$$

ou encore

$$E[\ln p(r|\theta)] = C^{ste} - \frac{1}{2N_0} \int |x(t,\theta) - x(t,\theta_0)|^2 dt$$

D'autre part il faut calculer

$$E \left[\frac{d^2}{d\theta^2} \ln p(r|\theta) \right] = \frac{d^2}{d\theta^2} E [\ln p(r|\theta)]$$

$$\begin{aligned} \frac{d}{d\theta} E[\ln p(r|\theta)] &= -\frac{1}{N_0} \Re \left\{ \int (x(t,\theta) - x(t,\theta_0)) \frac{\partial x^*}{\partial \theta}(t,\theta) dt \right\} \\ \frac{d^2}{d\theta^2} E[\ln p(r|\theta)] &= -\frac{1}{N_0} \Re \left\{ \int (x(t,\theta) - x(t,\theta_0)) \frac{\partial^2 x^*}{\partial \theta^2}(t,\theta) dt + \int \left| \frac{\partial x}{\partial \theta}(t,\theta) \right|^2 dt \right\} \end{aligned}$$

Soit pour $\theta = \theta_0$:

$$E \left[\frac{d^2}{d\theta^2} \ln p(r|\theta) \right] = -\frac{1}{N_0} \int \left| \frac{\partial x}{\partial \theta}(t,\theta_0) \right|^2 dt = \frac{A_{11}}{N_0}$$

Donc l'inégalité de Cramer-Rao (égalité ici) donne bien

$$E \left[|\hat{\theta}_{MC}(r) - \theta|^2 \right] = -\frac{A_{11}}{N_0}$$

4.7 Ambiguïté temps-fréquence

4.7.1 Définition

Un signal $x(t)$ se réfléchissant sur une cible en mouvement subit à la fois un retard $\Delta\tau$ et une translation de fréquence Δf . On utilise toujours le signal analytique, et on suppose que la détection se fait sur l'enveloppe de la fonction d'ambiguïté. Au lieu de $\frac{1}{2} \Re \{ \langle A_{\Delta\tau\Delta f} | A_{\Delta\tau_0\Delta f_0} \rangle \}$, on utilise parfois pour des raisons historiques la fonction d'ambiguïté temps-fréquence de Woodward :

$$\chi(\Delta\tau - \Delta\tau_0, \Delta f - \Delta f_0) \doteq \Re \{ \langle a_{\Delta\tau\Delta f_0} | a_{\Delta\tau_0\Delta f} \rangle \}$$

c'est-à-dire le produit scalaire de l'enveloppe complexe décalée en temps par l'enveloppe complexe décalée en fréquence. On "oublie" de plus le facteur 1/2 et l'opérateur partie réelle. Cela ne porte pas à conséquence (voir plus loin).

Nous utiliserons la définition suivante de la fonction d'ambiguïté temps-fréquence, dite de Wigner-Ville

Définition 4.7**(Fonction d'ambiguïté temps-fréquence)**

$$\begin{aligned}\chi(\tau, f) &\doteq \int a^*(t - \tau/2)a(t + \tau/2) \exp(2i\pi ft) dt \\ &= \int \tilde{a}^*(\nu + f/2)\tilde{a}(\nu - f/2) \exp(2i\pi\nu\tau) d\nu\end{aligned}$$

en posant $\tau = \Delta\tau - \Delta\tau_0$ et $f = \Delta f - \Delta f_0$. On suppose de plus que l'enveloppe est centrée à la fois en temps et en fréquence (cela revient simplement à décaler les origines en temps et fréquence) :

$$\int t|a(t)|^2 dt = 0 \quad \text{et} \quad \int \nu|\tilde{a}(\nu)|^2 d\nu = 0$$

Remarque – Les différentes fonctions d'ambiguïté possibles ne diffèrent que par un facteur du type $\exp(i\pi f\tau)$, dont l'influence disparaît quand on prend la partie réelle de la fonction d'ambiguïté (Il vaut mieux croire cette affirmation sur parole : les calculs pour la vérifier sont un peu fastidieux...). Notre définition (la plus usuelle) correspond par exemple à multiplier la fonction d'ambiguïté temps-fréquence de Woodward (définition historique) par $\exp(-i\pi f\tau)$; elle conduit à des calculs plus simples en général, à cause de la propriété (ii) suivante.

Propriété 4.8**(Propriétés de la fonction d'ambiguïté temps-fréquence)**

(i) Valeur centrale

$$\chi(0,0) = \int |a(t)|^2 dt = \int |\tilde{a}(\nu)|^2 d\nu = 2E_x$$

(ii) Symétrie hermitienne

$$\chi(-\tau, -f) = \chi^*(\tau, f)$$

(iii) χ est maximum à l'origine

$$|\chi(\tau, f)| \leq \chi(0,0)$$

(iv) Le module carré de la fonction d'ambiguïté temps-fréquence est sa propre double TF :

$$|\chi(\eta, \xi)|^2 = \iint |\chi(\tau, f)|^2 \exp(2i\pi(\eta f - \xi\tau)) d\tau df$$

En particulier, "l'amplitude est égale au volume" :

$$|\chi(0,0)|^2 = \iint |\chi(\tau, f)|^2 d\tau df$$

Cette propriété implique que la fonction d'ambiguïté doit avoir une certaine extension en temps et en fréquence, et ne peut être concentrée complètement à l'origine.

4.7.2 Estimation conjointe temps-fréquence

Propriété 4.9

(Estimation conjointe temps-fréquence)

L'estimation conjointe des moindres carrés conduit à chercher le maximum de la fonction d'ambiguïté temps-fréquence. La matrice de covariance de l'erreur est

$$\Gamma_{\widehat{\Delta\tau}, \widehat{\Delta f}} = \frac{1}{4\pi} \frac{N_0}{E_X} \begin{pmatrix} \overline{\Delta f^2} & \overline{\Delta\tau \cdot \Delta f} \\ \overline{\Delta\tau \cdot \Delta f} & \overline{\Delta\tau^2} \end{pmatrix}^{-1} = \frac{1}{4\pi} \frac{N_0}{E_X} \frac{1}{\Delta} \begin{pmatrix} \overline{\Delta\tau^2} & -\overline{\Delta\tau \cdot \Delta f} \\ -\overline{\Delta\tau \cdot \Delta f} & \overline{\Delta f^2} \end{pmatrix}$$

$$\text{avec } \Delta = \overline{\Delta\tau^2} \cdot \overline{\Delta f^2} - (\overline{\Delta\tau \cdot \Delta f})^2.$$

Il est aisé de vérifier cette propriété à partir de

$$\Gamma_{\widehat{\Delta\tau}, \widehat{\Delta f}} = -N_0 \cdot A^{-1}$$

et

$$A = \begin{pmatrix} \frac{\partial^2 \chi}{\partial \tau^2}(0,0) & \frac{\partial^2 \chi}{\partial \tau \partial f}(0,0) \\ \frac{\partial^2 \chi}{\partial \tau \partial f}(0,0) & \frac{\partial^2 \chi}{\partial f^2}(0,0) \end{pmatrix}$$

avec les définitions suivantes :

– Épanouissement temporel

$$\overline{\Delta\tau^2} = \frac{\int t^2 |a(t)|^2 dt}{\int |a(t)|^2 dt} = \frac{-1}{4\pi^2} \frac{1}{\chi(0,0)} \frac{\partial^2 \chi}{\partial f^2}(0,0)$$

– Épanouissement fréquentiel

$$\overline{\Delta f^2} = \frac{\int \nu^2 |\tilde{a}(\nu)|^2 d\nu}{\int |\tilde{a}(\nu)|^2 d\nu} = \frac{-1}{4\pi^2} \frac{1}{\chi(0,0)} \frac{\partial^2 \chi}{\partial \tau^2}(0,0)$$

– Produit temps-fréquence moyen

$$\overline{\Delta\tau \cdot \Delta f} = \frac{-1}{4\pi^2} \frac{1}{\chi(0,0)} \frac{\partial^2 \chi}{\partial \tau \partial f}(0,0)$$

Fonction d'ambiguïté temps-fréquence au voisinage de l'origine – On peut vérifier que du fait de la condition de centrage en temps et en fréquence de l'enveloppe complexe les dérivées premières de χ sont nulles. Un développement au second ordre donne alors :

$$\begin{aligned} \chi(\tau, f) &= \chi(0,0) + \frac{1}{2} \frac{\partial^2 \chi}{\partial \tau^2}(0,0) \tau^2 + \frac{1}{2} \frac{\partial^2 \chi}{\partial f^2}(0,0) f^2 + \frac{\partial^2 \chi}{\partial \tau \partial f}(0,0) \tau f \\ &= \chi(0,0) \left(1 - 2\pi^2 (\overline{\Delta\tau^2} f^2 + \overline{\Delta f^2} \tau^2 + 2\overline{\Delta\tau \cdot \Delta f} \tau f) \right) \end{aligned}$$

Les courbes de niveau autour de l'origine sont donc des ellipses d'équation :

$$\overline{\Delta\tau^2}f^2 + \overline{\Delta f^2}\tau^2 + 2\overline{\Delta\tau.\Delta f}\tau f = C^{ste}$$

ce que représente la figure 4.2. Dans le cas d'une impulsion gaussienne ce résultat est exact pour toutes valeurs de τ et f . Les axes de l'ellipse ne sont pas orientés suivant les axes des abscisses et des ordonnées si $\overline{\Delta\tau.\Delta f} \neq 0$.

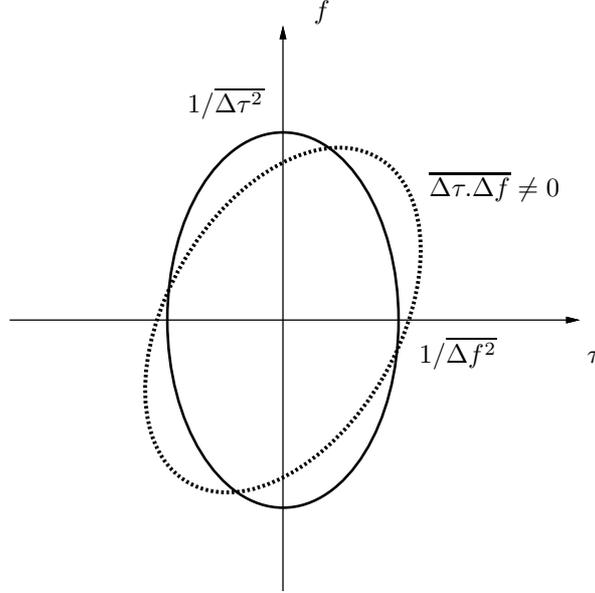


FIG. 4.2 – Fonction d'ambiguïté temps-fréquence au voisinage de l'origine.

Supposons que l'on mesure le temps de retour en faisant l'hypothèse que la cible est immobile, alors qu'en réalité elle est en mouvement, ce qui crée un décalage en fréquence δf . Au lieu de chercher le maximum de la fonction d'ambiguïté en temps seul $\chi(\tau, 0)$ comme on croit le faire, on va faire en fait la recherche du temps de retour sur la fonction $\chi(\tau, -\delta f)$ pour un certain δf fixé mais inconnu (le signe moins vient de ce que l'on se trompe sur l'origine des fréquences en faisant l'hypothèse que la cible est immobile). En supposant δf et τ petits, on a

$$\chi(\tau, -\delta f) = \chi(0, 0) \left(1 - 2\pi^2(\overline{\Delta\tau^2}\delta f^2 + \overline{\Delta f^2}\tau^2 - 2\overline{\Delta\tau.\Delta f}\tau\delta f) \right)$$

En conséquence le maximum sera trouvé lorsque

$$\frac{\partial}{\partial \tau} \chi(\tau, -\delta f) = 0$$

soit

$$\tau = \frac{\overline{\Delta\tau.\Delta f}}{\overline{\Delta f^2}} \delta f$$

Ainsi, le coefficient de couplage $\overline{\Delta\tau.\Delta f}$ introduit un biais systématique sur l'estimation du temps de retour si l'estimation de la fréquence (ou l'hypothèse sur la fréquence) est erronée.

De même, si l'on mesure une fréquence avec une hypothèse erronée sur le temps de retour ($\delta\tau$), un biais est introduit par le couplage :

$$f = \frac{\overline{\Delta\tau.\Delta f}}{\overline{\Delta\tau^2}} \delta\tau$$

Le terme $\overline{\Delta\tau \cdot \Delta f}$ provoque un couplage des erreurs de mesure sur le temps et la fréquence. Pour éviter ce couplage, il faut utiliser des signaux tels que $\overline{\Delta\tau \cdot \Delta f} = 0$. On peut montrer que cette condition revient à imposer :

$$\int t|a(t)|^2 \frac{d\varphi(t)}{dt} dt = 0$$

où $\varphi(t)$ est la phase de l'enveloppe complexe⁴. En particulier si cette phase est constante l'intégrale est nulle. Dans ce cas la matrice de covariance des erreurs est diagonale, c'est-à-dire que les erreurs sont découplées. On peut alors pratiquer sans risque une estimation séparée du temps et de la fréquence :

$$\begin{cases} \sigma_{\widehat{\Delta\tau}}^2 &= \frac{1}{4\pi^2} \frac{N_0}{E_X} \frac{1}{\overline{\Delta f^2}} \\ \sigma_{\widehat{\Delta f}}^2 &= \frac{1}{4\pi^2} \frac{N_0}{E_X} \frac{1}{\overline{\Delta\tau^2}} \end{cases}$$

En conclusion, il faut un signal le plus large possible en fréquence, le plus long possible en temps, et qui découple les mesures de temps et de fréquence.

4.8 Exemples de fonctions d'ambiguïté temps-fréquence

4.8.1 Impulsion gaussienne

Soit l'impulsion gaussienne suivante :

$$a(t) = \sqrt{\frac{A}{\sqrt{2\pi}T}} \exp\left(-\frac{t^2}{4T^2}\right)$$

$$|a(t)|^2 = \frac{A}{\sqrt{2\pi}T} \exp\left(-\frac{t^2}{2T^2}\right)$$

qui est telle que (par analogie avec une loi normale centrée de variance T^2)

$$\int |a(t)|^2 dt = A$$

Calcul direct – On voit que

$$\int t|a(t)|^2 dt = 0$$

et (toujours par analogie avec une loi normale centrée de variance T^2)

$$\int t^2|a(t)|^2 dt = A \cdot T^2$$

donc

$$\overline{\Delta\tau^2} = T^2$$

4. Exercice : montrer que

$$\frac{\partial^2 \chi}{\partial \tau \partial f}(0,0) = -2\pi \int t|a(t)|^2 \frac{d\varphi(t)}{dt} dt$$

Cherchons la transformée de Fourier de l'enveloppe :

$$\begin{aligned}
\tilde{a}(\nu) &= \sqrt{\frac{A}{\sqrt{2\pi}T}} \text{TF} \left[\exp \left(-\pi \left(\frac{t}{\sqrt{4\pi T^2}} \right)^2 \right) \right] (\nu) \\
&= \sqrt{\frac{A(4\pi T^2)}{\sqrt{2\pi}T}} \exp \left(-\pi (\sqrt{4\pi T^2} \nu)^2 \right) \\
\tilde{a}(\nu) &= \sqrt{\frac{A(4\pi T)}{\sqrt{2\pi}}} \exp \left(-4\pi^2 T^2 \nu^2 \right) \\
|\tilde{a}(\nu)|^2 &= \frac{A(4\pi T)}{\sqrt{2\pi}} \exp \left(-8\pi^2 T^2 \nu^2 \right) \\
&= \frac{A(4\pi T)}{\sqrt{2\pi}} \exp \left(-\frac{\nu^2}{2\sigma^2} \right) \quad \text{avec } \sigma^2 = \frac{1}{16\pi^2 T^2} \\
|\tilde{a}(\nu)|^2 &= \frac{A}{\sqrt{2\pi}\sigma} \exp \left(-\frac{\nu^2}{2\sigma^2} \right)
\end{aligned}$$

donc (toujours par analogie avec une loi normale centrée de variance σ^2)

$$\begin{aligned}
\int |\tilde{a}(\nu)|^2 d\nu &= A \\
\int \nu |\tilde{a}(\nu)|^2 d\nu &= 0 \\
\int \nu^2 |\tilde{a}(\nu)|^2 d\nu &= A \cdot \sigma^2 = \frac{A}{16\pi^2 T^2}
\end{aligned}$$

d'où l'on déduit

$$\overline{\Delta f^2} = \frac{1}{16\pi^2 T^2}$$

Calcul par la fonction d'ambiguïté

$$\begin{aligned}
\chi(\tau, f) &= \frac{A}{\sqrt{2\pi}T} \int \exp \left(-\frac{(t - \tau/2)^2 + (t + \tau/2)^2}{4T^2} \right) \exp(2i\pi ft) dt \\
&= \frac{A}{\sqrt{2\pi}T} \exp \left(-\frac{\tau^2}{8T^2} \right) \int \exp \left(-\frac{t^2}{2T^2} \right) \exp(2i\pi ft) dt \\
&= \frac{A}{\sqrt{2\pi}T} \exp \left(-\frac{\tau^2}{8T^2} \right) \text{TF} \left[\exp \left(-\pi \left(\frac{t}{\sqrt{2\pi}T} \right)^2 \right) \right] (-f) \\
\chi(\tau, f) &= A \cdot \exp \left(-\frac{\tau^2}{8T^2} \right) \cdot \exp \left(-2\pi^2 T^2 f^2 \right)
\end{aligned}$$

On peut maintenant évaluer les dérivées partielles :

$$\begin{aligned}
\frac{\partial \chi}{\partial \tau}(\tau, f) &= \left(-\frac{\tau}{4T^2} \right) \chi(\tau, f) \\
\frac{\partial \chi}{\partial \tau}(0, 0) &= 0 \\
\frac{\partial \chi}{\partial f}(\tau, f) &= \left(-4\pi^2 T^2 f \right) \chi(\tau, f)
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \chi}{\partial f}(0,0) &= 0 \\
\frac{\partial^2 \chi}{\partial \tau^2}(\tau, f) &= \left(-\frac{1}{4T^2}\right) \chi(\tau, f) + \left(\frac{\tau}{4T^2}\right)^2 \chi(\tau, f) \\
\frac{\partial^2 \chi}{\partial \tau^2}(0,0) &= -\frac{1}{4T^2} \chi(0,0) \\
\frac{\partial^2 \chi}{\partial f^2}(\tau, f) &= \left(-4\pi^2 T^2\right) \chi(\tau, f) + \left(-4\pi^2 T^2 f\right)^2 \chi(\tau, f) \\
\frac{\partial^2 \chi}{\partial f^2}(0,0) &= -4\pi^2 T^2 \chi(0,0) \\
\frac{\partial^2 \chi}{\partial \tau \partial f}(\tau, f) &= \left(-\frac{\tau}{4T^2}\right) \left(-4\pi^2 T^2 f\right) \chi(\tau, f) \\
\frac{\partial^2 \chi}{\partial \tau \partial f}(0,0) &= 0
\end{aligned}$$

D'où l'on voit que l'on retrouve les expressions des épanouissements temporel et fréquentiel, et que le produit temps-fréquence moyen est nul (ce qui était évident au départ puisque la phase de l'impulsion est nulle).

4.8.2 Impulsion gaussienne “chirp”

Nous reprenons l'exemple précédent mais avec un terme de phase quadratique additionnel (modulation linéaire de fréquence) :

$$a(t) = \sqrt{\frac{A}{\sqrt{2\pi T}}} \exp\left(-\frac{t^2}{4T^2} + ibt^2\right)$$

On a toujours

$$|a(t)|^2 = \frac{A}{\sqrt{2\pi T}} \exp\left(-\frac{t^2}{2T^2}\right)$$

et donc

$$\begin{aligned}
\int |a(t)|^2 dt &= A \\
\int t |a(t)|^2 dt &= 0 \\
\int t^2 |\tilde{a}(\nu)|^2 d\nu &= A.T^2 \\
\overline{\Delta\tau^2} &= T^2
\end{aligned}$$

La précision de la mesure de temps n'est donc pas modifiée.

$$\begin{aligned}
\chi(\tau, f) &= \frac{A}{\sqrt{2\pi T}} \int \exp\left(-\left(t - \tau/2\right)^2 \left(\frac{1}{4T^2} + ib\right) - \left(t + \tau/2\right)^2 \left(\frac{1}{4T^2} - ib\right)\right) \exp(2i\pi ft) dt \\
&= \frac{A}{\sqrt{2\pi T}} \exp\left(-\frac{\tau^2}{8T^2}\right) \int \exp\left(-\frac{t^2}{2T^2} + 2ib\tau t\right) \exp(2i\pi ft) dt \\
&= \frac{A}{\sqrt{2\pi T}} \exp\left(-\frac{\tau^2}{8T^2}\right) \text{TF} \left[\exp\left(-\pi \left(\frac{t}{\sqrt{2\pi T}}\right)^2\right) \right] \left(-f - \frac{b\tau}{\pi}\right)
\end{aligned}$$

$$\begin{aligned}\chi(\tau, f) &= A \cdot \exp\left(-\frac{\tau^2}{8T^2}\right) \cdot \exp\left(-2\pi^2 T^2 \left(f + \frac{b\tau}{\pi}\right)^2\right) \\ \chi(\tau, f) &= A \cdot \exp\left(-\left(\frac{1}{8T^2} + 2b^2 T^2\right)\tau^2 - 2\pi^2 T^2 f^2 - 4\pi b T^2 f \tau\right)\end{aligned}$$

Après évaluation des dérivées partielles à l'origine, on obtient

$$\overline{\Delta f^2} = \frac{1}{16\pi^2 T^2} + \frac{b^2 T^2}{\pi^2}$$

et

$$\overline{\Delta \tau \cdot \Delta f} = \frac{b T^2}{\pi}$$

et donc un terme de phase quadratique conduit à l'apparition d'un couplage des erreurs de mesure en temps et fréquence, même s'il permet d'améliorer la précision de la mesure en fréquence seule ($\overline{\Delta f^2}$ augmente avec b^2). Dans la pratique une impulsion "chirp" a donc un intérêt certain.

4.8.3 Compression d'impulsion

La fonction d'ambiguïté pour la mesure d'un temps est :

$$\chi(\tau, 0) = \int a^*(t - \tau) a(t) dt$$

Si on pose $h(t) = a^*(-t)$ alors

$$\chi(\tau, 0) = [h \star a](\tau)$$

c'est-à-dire le produit de convolution de a par h . $h(t)$ est appelé filtre adapté à l'impulsion d'enveloppe complexe $a(t)$. Le système de traitement des signaux pour la mesure d'un temps filtre le signal de retour avec le filtre adapté à l'impulsion. Si on écrit maintenant

$$\tilde{a}(\nu) = |\tilde{a}(\nu)| \exp(i\phi(\nu))$$

alors

$$a(t) = \text{TF}^{-1}[|\tilde{a}(\nu)|](t) \star \text{TF}^{-1}[\exp(i\phi(\nu))](t)$$

La phase $\phi(\nu)$ ne joue pas sur l'énergie du signal ($\int |\tilde{a}(\nu)|^2 d\nu$), mais "élargit" l'impulsion $\text{TF}^{-1}[|\tilde{a}(\nu)|](t)$ par sa convolution avec $\text{TF}^{-1}[\exp(i\phi(\nu))](t)$. En effet, si $\phi(\nu) = 0$, alors

$$\text{TF}^{-1}[\exp(i\phi(\nu))](t) = \delta(t)$$

et l'impulsion n'est pas élargie par convolution. Par contre si $\phi(\nu) \neq 0$, l'extension de la fonction $\text{TF}^{-1}[\exp(i\phi(\nu))](t)$ devient non négligeable.

On dit que le filtre adapté, de transformée de Fourier $\tilde{h}(\nu) = |\tilde{a}(\nu)| \exp(-i\phi(\nu))$, réalise une compression de l'impulsion $a(t)$ en remettant en phase les fréquences, puisque

$$\tilde{h}(\nu) \cdot \tilde{a}(\nu) = |\tilde{a}(\nu)|^2$$

Pour la mesure, tout se passe donc comme si on avait utilisé l'impulsion comprimée $\text{TF}^{-1}[|\tilde{a}(\nu)|](t)$ à la place de $a(t)$.

4.8.4 Train d'impulsions

Supposons qu'à partir d'une impulsion $a(t)$ de durée T , on produise le train de $(2n + 1)$ impulsions défini par

$$a_n(t) = \frac{1}{\sqrt{2n+1}} \sum_{k=-n}^{+n} a(t - kT_r)$$

où T_r est le temps de répétition ($T \ll T_r$). Que devient la fonction d'ambiguïté?

$$\chi_n(\tau, f) = \frac{1}{2n+1} \sum_{k=-n}^{+n} \sum_{l=-n}^{+n} \int a^*(t - kT_r - \tau/2) a(t - lT_r + \tau/2) \exp(2i\pi ft) dt$$

soit en faisant le changement de variable $t = t' + 1/2(k+l)T_r$:

$$\chi_n(\tau, f) = \frac{1}{2n+1} \sum_{k=-n}^{+n} \sum_{l=-n}^{+n} \chi(\tau + (k-l)T_r, f) \exp(i\pi(k+l)fT_r)$$

où $\chi(\tau, f)$ est la fonction d'ambiguïté associée à l'impulsion $a(t)$. $\chi(\tau, f)$ a une largeur temporelle de l'ordre de $2T$ (elle est nulle en dehors de $\tau \in [-T, T]$). Par conséquent, $\chi(\tau + (k-l)T_r, f)$ n'est différente de 0 que dans l'intervalle

$$-T \leq \tau + (k-l)T_r \leq T$$

ou

$$-T + (l-k)T_r \leq \tau \leq T + (l-k)T_r$$

$\chi(\tau + (k-l)T_r, f)$ est non nulle pour τ dans des intervalles de la forme $[-T + mT_r, T + mT_r], m \in \mathbb{Z}$. Donc si $\tau \in [-T + mT_r, T + mT_r]$, $l - k = m$, et

$$\chi_n(\tau + mT_r, f) = \frac{1}{2n+1} \left(\sum_{k=p(m)}^{q(m)} \exp(2i\pi k f T_r) \right) \chi(\tau, f) \exp(i\pi m f T_r), \quad \tau \in [-T, T]$$

avec

$$\begin{cases} p(m) &= \max(-n - m, -n) \\ q(m) &= \min(n - m, n) \end{cases}$$

puisque'on doit avoir à la fois $k \in [-n, n]$ et $k + m \in [-n, n]$. On vérifie facilement que $p(m) + q(m) + m = 0$, et que $q(m) - p(m) = 2n - |m|$. La somme est celle d'une suite géométrique, on montre donc que

$$\sum_{k=p(m)}^{q(m)} \exp(2i\pi k f T_r) = \exp(i\pi(q(m) + p(m))fT_r) \frac{\sin(\pi(q(m) - p(m) + 1)fT_r)}{\sin(\pi f T_r)}$$

et donc

$$\chi_n(\tau + mT_r, f) = \frac{\sin(\pi(2n + 1 - |m|)fT_r)}{(2n + 1) \sin(\pi f T_r)} \chi(\tau, f), \quad \tau \in [-T, T]$$

En clair, la fonction d'ambiguïté de l'impulsion originale est découpée par une fonction de type Fabry-Perrot qui l'affine en fréquence d'un facteur $1/(2n + 1 - |m|)$. Par contre la fonction d'ambiguïté de l'impulsion originale est également périodisée en temps, avec une période T_r . Par conséquent la précision de la mesure de fréquence est augmentée, mais en contrepartie l'ambiguïté en temps est plus grande car on peut confondre les maxima successifs $\chi_n(mT_r, 0)$ qui sont difficilement distinguables si n est grand (en présence de bruit bien sûr) :

$$\chi_n(mT_r, 0) = \frac{2n + 1 - |m|}{2n + 1} \chi(0, 0)$$

4.9 Inégalité de Fréchet-Darmois-Cramer-Rao

L'inégalité de Fréchet-Darmois-Cramer-Rao généralise l'inégalité de Cramer-Rao (valable pour un paramètre scalaire) pour un paramètre vectoriel à n composantes (ou à l'estimation conjointe de n paramètres scalaires).

Propriété 4.10

(Inégalité de Fréchet-Darmois-Cramer-Rao)

$$E \left[(\hat{\theta}(r) - \theta)(\hat{\theta}(r) - \theta)^T \right] \geq M.I_F^{-1}.M^T$$

où $\hat{\theta}(r)$ est un vecteur estimateur du vecteur paramètre θ . I_F est la matrice d'information de Fisher, définie par :

$$(I_F)_{ij} \doteq E \left[\frac{\partial}{\partial \theta_i} \ln p(r|\theta) \cdot \frac{\partial}{\partial \theta_j} \ln p(r|\theta) \right]$$

soit encore

$$I_F = E \left[\nabla \ln p(r|\theta) \cdot (\nabla \ln p(r|\theta))^T \right]$$

où ∇ est l'opérateur gradient défini par

$$\nabla \ln p(r|\theta) \doteq \begin{pmatrix} \frac{\partial}{\partial \theta_1} \ln p(r|\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_n} \ln p(r|\theta) \end{pmatrix}$$

et M est la matrice des dérivées partielles du vecteur $E[\hat{\theta}(r)]$

$$M_{ij} \doteq \frac{\partial}{\partial \theta_j} E[\hat{\theta}_i(r)]$$

et enfin l'espérance est prise sur $p(r|\theta)$

$$E[f(r)] = \int f(r)p(r|\theta)dr$$

où $f(r)$ est une fonction quelconque de r (L'application de l'opérateur espérance fait disparaître la dépendance en r , mais pas en θ).

Preuve – La démonstration n'est pas compliquée, mais nécessite quelques manipulations d'algèbre linéaire.

$$\begin{aligned} E[\hat{\theta}_i(r)] &= \theta_i + (E[\hat{\theta}_i(r)] - \theta_i) \\ &= \theta_i + \int (\hat{\theta}_i(r) - \theta_i) p(r|\theta)dr \end{aligned}$$

donc

$$\begin{aligned} M_{ij} &= \frac{\partial}{\partial \theta_j} E[\hat{\theta}_i(r)] \\ &= \frac{\partial}{\partial \theta_j} \theta_i + \int (\hat{\theta}_i(r) - \theta_i) \frac{\partial}{\partial \theta_j} p(r|\theta)dr - \frac{\partial}{\partial \theta_j} \theta_i \int p(r|\theta)dr \end{aligned}$$

$$\begin{aligned}
&= \int (\hat{\theta}_i(r) - \theta_i) \frac{\partial}{\partial \theta_j} p(r|\theta) dr \\
&= \int (\hat{\theta}_i(r) - \theta_i) p(r|\theta) \frac{\partial}{\partial \theta_j} \ln p(r|\theta) dr \\
&= E \left[(\hat{\theta}_i(r) - \theta_i) \frac{\partial}{\partial \theta_j} \ln p(r|\theta) \right]
\end{aligned}$$

soit encore

$$M = E \left[(\hat{\theta}(r) - \theta) \cdot (\nabla \ln p(r|\theta))^T \right]$$

Posons

$$\delta \doteq \hat{\theta}(r) - \theta$$

et

$$V \doteq \delta - M.I_F^{-1} \cdot (\nabla \ln p(r|\theta))$$

Alors

$$E[V.V^T] = E \left[(\delta - M.I_F^{-1} \cdot (\nabla \ln p(r|\theta))) \cdot (\delta^T - (\nabla \ln p(r|\theta))^T \cdot I_F^{-1} \cdot M^T) \right]$$

car

$$(\delta - M.I_F^{-1} \cdot (\nabla \ln p(r|\theta)))^T = \delta^T - (\nabla \ln p(r|\theta))^T \cdot I_F^{-1} \cdot M^T$$

et I_F est symétrique (donc I_F^{-1} aussi).

$$\begin{aligned}
E[V.V^T] &= E[\delta.\delta^T] - 2E[\delta.(\nabla \ln p(r|\theta))^T \cdot I_F^{-1} \cdot M^T] \\
&\quad + E[M.I_F^{-1} \cdot (\nabla \ln p(r|\theta)) \cdot (\nabla \ln p(r|\theta))^T \cdot I_F^{-1} \cdot M^T] \\
E[V.V^T] &= E[\delta.\delta^T] - 2E[\delta.(\nabla \ln p(r|\theta))^T] \cdot I_F^{-1} \cdot M^T \\
&\quad + M.I_F^{-1} \cdot E[(\nabla \ln p(r|\theta)) \cdot (\nabla \ln p(r|\theta))^T] \cdot I_F^{-1} \cdot M^T
\end{aligned}$$

car M et I_F ne dépendent pas de r (ce sont déjà des espérances!), et donc

$$\begin{aligned}
E[V.V^T] &= E[\delta.\delta^T] - 2M.I_F^{-1} \cdot M^T \\
&\quad + M.I_F^{-1} \cdot I_F \cdot I_F^{-1} \cdot M^T \\
E[V.V^T] &= E[\delta.\delta^T] - M.I_F^{-1} \cdot M^T
\end{aligned}$$

Pour conclure il suffit de remarquer que $E[V.V^T] \geq 0$, et donc

$$E[\delta.\delta^T] \geq M.I_F^{-1} \cdot M^T$$

ce qui constitue l'inégalité F.D.C.R. (Fréchet-Darmonis-Cramer-Rao).

Comme dans le cas scalaire on peut définir les notions d'estimateurs non biaisés et efficaces.

Définition 4.8

(Estimateur vectoriel non biaisé)

Un estimateur vectoriel non biaisé est tel que $E[\hat{\theta}(r)] = \theta$ si θ est la vraie valeur du paramètre vectoriel. C'est-à-dire que pour θ fixée, la moyenne de $\hat{\theta}(r)$ prise sur toutes les réalisations possibles de la mesure est égale à θ . Dans ce cas l'inégalité F.D.C.R. devient

$$E[(\hat{\theta}(r) - \theta)(\hat{\theta}(r) - \theta)^T] \geq I_F^{-1}$$

Définition 4.9**(Estimateur vectoriel efficace)**

Un estimateur vectoriel efficace est un estimateur vectoriel non-biaisé qui atteint les bornes de Fréchet-Darmois-Cramer-Rao.

EXERCICE - 4.1

Réfectomètre optique

On veut étudier la précision d'un réfectomètre pour fibres optiques. Le principe est d'injecter une impulsion lumineuse connue de puissance instantanée $y(t)$ quasi-monochromatique dans une fibre optique défectueuse, qui est rompue à une distance d de l'extrémité connectée au réfectomètre. On cherche alors à déterminer le temps de retour de l'impulsion réfléchi sur l'interface due à la cassure. On donne l'information suivante, contenue dans la proposition logique I suivante :

$I =$ "On suppose en première approximation que l'impulsion se propage sans se déformer. On connaît l'indice effectif de propagation dans la fibre, que l'on note n , et on pose $\beta = n/c$, où c est la célérité de la lumière dans le vide. Par ailleurs, on connaît les pertes par unité de longueur de la fibre qui sont mesurées par le coefficient α en m^{-1} ; c'est-à-dire qu'après propagation sur une distance d , l'impulsion est atténuée d'un facteur $\exp(-\alpha d)$. Le détecteur utilisé pour enregistrer le signal de retour est affecté d'un bruit blanc centré de puissance connue σ^2 ."

- a) - On définit le vecteur paramètre θ par

$$\theta = \begin{pmatrix} d \\ \rho \end{pmatrix}$$

où ρ est le coefficient de réflexion en intensité à l'interface (cassure), et d est la distance à laquelle se trouve cette interface. On désigne par V la proposition logique :

$V =$ "Le signal de retour mesuré est $v(t)$."

Justifier que sous l'information I , on peut écrire le modèle de mesure sous la forme :

$$\begin{aligned} p(V|\theta I) &= \mathcal{N}(v(t), x(t, \theta), \sigma^2) = A \exp\left(-\frac{1}{2\sigma^2} \int (v(t) - x(t, \theta))^2 dt\right) \\ x(t, \theta) &= \rho \exp(-2\alpha d) y(t - 2\beta d) \end{aligned}$$

où A est une constante que l'on ne cherchera pas à déterminer, et que $p(\theta|VI)$ est directement proportionnelle à $p(V|\theta I)$.

- b) - On rappelle que dans le cas gaussien, les principes d'estimation MAP et MC s'identifient. On suppose dans un premier temps que la valeur de ρ est connue. On pose

$$J[f] = \int f(t) dt$$

et

$$[f \otimes g](z) = \int f(t) g(t - z) dt$$

où f et g sont des fonctions quelconques de t , et $a = J[y^2]$; $b = J\left[y \frac{dy}{dt}\right]$ et $c = J\left[\left(\frac{dy}{dt}\right)^2\right]$.

Calculer $\int (v(t) - x(t, \theta))^2 dt$ en utilisant les notations précédentes. Montrer alors que la valeur de l'estimation \hat{d} au sens MAP ou MC est solution de l'équation implicite :

$$\beta \frac{d}{dt}[v \otimes y](2\beta \hat{d}) = \alpha \cdot ([v \otimes y](2\beta \hat{d}) - a\rho \exp(-2\alpha \hat{d}))$$

Calculer $\frac{\partial x(t, \theta)}{\partial d}$, et en déduire la matrice de covariance $\Gamma_{\hat{d}}$ (ici réduite à un scalaire) de l'estimation de la distance d .

- c) - On suppose maintenant que la valeur de ρ est inconnue. Montrer que l'on a (toujours au sens MAP ou MC) :

$$\hat{\rho} = \frac{\exp(2\alpha\hat{d})}{a} [v \otimes y](2\beta\hat{d})$$

$$\frac{d}{dt} [v \otimes y](2\beta\hat{d}) = 0$$

Interpréter simplement ces deux équations. Montrer que la matrice A intervenant dans le calcul de la matrice de covariance de l'estimation MC a pour expression :

$$A = -\exp(-4\alpha\hat{d}) \begin{pmatrix} a & -2\hat{\rho}(\alpha a + \beta b) \\ -2\hat{\rho}(\alpha a + \beta b) & 4\hat{\rho}^2(\alpha^2 a + 2\alpha\beta b + \beta^2 c) \end{pmatrix}$$

(on ne cherchera pas à inverser cette matrice pour obtenir la matrice de covariance de l'estimation).

Réponse

- a) - Déterminons tout d'abord la moyenne du modèle de mesure, c'est-à-dire $x(t, \theta)$. L'impulsion parcourt une longueur $2d$ dans un milieu d'indice n . Elle revient donc au bout d'un temps $n \times 2d/c = 2\beta d$. Les pertes sont dues à la réflexion (ρ) et à la propagation ($\exp(-2\alpha d)$), donc au total

$$x(t, \theta) = \rho \exp(-2\alpha d) y(t - 2\beta d)$$

Le bruit de détection est centré, blanc, et de puissance σ^2 . Sous ces contraintes, le principe du maximum d'entropie conduit à attribuer une distribution gaussienne à $p(V|\theta I)$, de moyenne $x(t, \theta)$ et de variance σ^2 .

Rien dans I ne permet de privilégier une hypothèse sur θ , donc $p(\theta|I)$ est constante, et $p(\theta|VI)$ est directement proportionnelle à $p(V|\theta I)$.

- b) -

$$\int (v(t) - x(t, \theta))^2 dt = J[v^2] + a\rho^2 \exp(-4\alpha d) - 2\rho \exp(-2\alpha d) [v \otimes y](2\beta d)$$

puisque

$$\int y^2(t - 2\beta d) dt = \int y^2(t) dt = a$$

\hat{d} est obtenue en cherchant le minimum de $\int (v(t) - x(t, \theta))^2 dt$, donc pour

$$\frac{\partial}{\partial d} \int (v(t) - x(t, \theta))^2 dt = 0$$

soit

$$-4\alpha a \rho^2 \exp(-4\alpha d) + 4\alpha \rho \exp(-2\alpha d) [v \otimes y](2\beta d) - 4\beta \rho \exp(-2\alpha d) \frac{d}{dt} [v \otimes y](2\beta d) = 0$$

donc \hat{d} est solution de l'équation implicite

$$\beta \frac{d}{dt} [v \otimes y](2\beta\hat{d}) = \alpha \cdot ([v \otimes y](2\beta\hat{d}) - a\rho \exp(-2\alpha\hat{d}))$$

$$\frac{\partial x(t, \theta)}{\partial d} = -2\alpha \rho \exp(-2\alpha d) y(t - 2\beta d) - 2\beta \rho \exp(-2\alpha d) \frac{d}{dt} y(t - 2\beta d)$$

$$\Gamma_{\hat{d}} = -\sigma^2 \cdot A^{-1} = -\frac{\sigma^2}{A_{11}}$$

avec

$$A_{11} = -\int \left(\frac{\partial x(t, \theta)}{\partial d} \right)^2 dt$$

$$A_{11} = -4\rho^2 \exp(-4\alpha d) (\alpha^2 a + 2\alpha\beta b + \beta^2 c)$$

donc

$$\Gamma_{\hat{d}} = \frac{\sigma^2 \exp(4\alpha d)}{4\rho^2(\alpha^2 a + 2\alpha\beta b + \beta^2 c)}$$

– c) - $\hat{\rho}$ et \hat{d} sont solutions de

$$\begin{cases} \frac{\partial}{\partial \rho} \int (v(t) - x(t, \theta))^2 dt = 0 = 2\rho \exp(-4\alpha d)a - 2 \exp(-2\alpha d)[v \otimes y](2\beta d) \\ \frac{\partial}{\partial d} \int (v(t) - x(t, \theta))^2 dt = 0 \end{cases}$$

La seconde équation est identique à celle de la question précédente. La première équation donne

$$\hat{\rho} = \frac{\exp(2\alpha \hat{d})}{a} [v \otimes y](2\beta \hat{d})$$

et en injectant cette valeur dans l'expression de \hat{d} trouvée à la question précédente

$$\frac{d}{dt} [v \otimes y](2\beta \hat{d}) = 0$$

Pour trouver \hat{d} , il faut donc chercher le maximum de la corrélation de v avec y . On obtient ensuite $\hat{\rho}$ en normalisant la valeur au maximum comme si l'acquisition avait été faite sans bruit :

$$[v \otimes y]_{max} = \hat{\rho} \exp(-2\alpha \hat{d}) \underbrace{[y \otimes y](0)}_{=a}$$

$$\begin{cases} \frac{\partial x(t, \theta)}{\partial d} = -2\alpha \rho \exp(-2\alpha d)y(t - 2\beta d) - 2\beta \rho \exp(-2\alpha d) \frac{d}{dt} y(t - 2\beta d) \\ \frac{\partial x(t, \theta)}{\partial \rho} = \exp(-2\alpha d)y(t - 2\beta d) \end{cases}$$

donc

$$A_{22} = -\exp(-4\alpha d)a$$

et

$$A_{12} = A_{21} = 2\rho \exp(-4\alpha d)(\alpha a + \beta b)$$

d'où l'expression de la matrice A .

Chapitre 5

Filtrage de Kalman

5.1 Problématique

5.1.1 Définitions

On suppose que l'on observe l'évolution d'un système au cours du temps, en effectuant des mesures sur ce système.

- L'état du système est décrit par le **vecteur paramètre** θ_t , qui évolue au cours du temps. La valeur exacte de θ_t n'est pas connue : il faut l'estimer.
- Afin de rendre cette estimation possible, on effectue des mesures, dont le résultat est décrit par le **vecteur mesure** v_t . Ces mesures portent sur des grandeurs physiques liées à la valeur du vecteur paramètre, c'est-à-dire que l'on formule un modèle reliant les observations à l'état du système.
- On suppose également **un modèle de l'évolution temporelle de l'état du système**, soit du vecteur paramètre θ_t .

Le but du filtrage de Kalman est l'estimation temporelle de l'état du système θ_t , au fur et à mesure que l'on accumule des données.

Notations : On considérera par la suite une évolution à temps discret, c'est-à-dire que t est un entier (un indice).

$$\theta_t \in \mathbb{R}^n$$

$$v_t \in \mathbb{R}^m$$

$$V_t = v_1 \cdots v_t \quad (\text{produit de propositions logiques : } V_t \text{ est constituée de toutes les données acquises jusqu'à l'instant } t.)$$

$$I : \quad \text{Proposition logique décrivant l'ensemble de notre état de connaissance sur le système considéré (tout ce qui n'est pas contenu dans les } v_t \text{ et } \theta_t).$$

EXEMPLE - 5.1 Poursuite d'un avion

Supposons un problème de poursuite d'un point brillant (un avion) dans une séquence d'images (ce qui est un problème difficile en pratique!). L'avion apparaît comme un point car il est en limite de portée de la caméra. On acquière une séquence d'images.

- Comment décrire l'état du système (l'avion)?
 - position (3 composantes)
 - direction (2 composantes)
 - module de vitesse (1 composante)
 - ...

- position du manche à balai
- psychologie du pilote (?)
- ...

On peut rendre le problème d'estimation arbitrairement compliqué, mais il est clair que certains paramètres sont plus importants que d'autres.

- Que mesure t'on? les images (donc $m =$ nombre de pixels des images). Il est clair que l'on peut estimer à partir de ces mesures la position, la direction et le module de la vitesse de l'avion, le tout avec une précision finie.

Le filtrage de Kalman est particulièrement utile dans toutes les applications où le signal attendu peut disparaître pendant plusieurs incréments de temps. Cette situation peut se produire par exemple lors de la poursuite d'un hélicoptère par un équipement optronique. Si l'hélicoptère disparaît derrière une rangée d'arbres, il sera nécessaire de prédire approximativement l'endroit où il va réapparaître afin de pouvoir continuer la poursuite. Cette prédiction de la trajectoire devra être faite sur la base des mouvements précédents. Un autre exemple est fourni par la restauration d'enregistrements détériorés, par exemple de vieux disques vinyle. Il faudra détecter les endroits abîmés et les remplacer par des enregistrements prédits à l'aide des parties intactes précédentes.

5.1.2 Estimation – prédiction

Au bout d'un temps t , on a acquis les données $V_t = v_1 \cdots v_t$. Connaissant de plus toute l'information disponible sur le système, contenue dans la proposition logique I , que peut-on dire de l'état du système à un instant t' quelconque?

La réponse à cette question est contenue dans la densité de probabilité

$$p(\theta_{t'} | V_t I)$$

Il est d'usage de différencier les trois cas suivants :

$$\begin{cases} t' < t & \rightarrow \text{“lissage”} \\ t' = t & \rightarrow \text{“filtrage”} \\ t' > t & \rightarrow \text{“prédiction”} \end{cases}$$

Il est clair que plus on dispose d'information (plus t est grand) plus l'estimation sera bonne : l'estimation de $\theta_{t'}$ que l'on peut faire à partir de $p(\theta_{t'} | V_{t_2} I)$ doit être de qualité supérieure ou égale à celle que l'on peut faire à partir de $p(\theta_{t'} | V_{t_1} I)$ si $t_2 > t_1$.

Plus précisément, l'apport d'une donnée complémentaire est le suivant. De $V_{t+1} = v_{t+1} \cdot V_t$, la règle du produit donne :

$$p(\theta_{t'} | V_{t+1} I) = \frac{p(\theta_{t'} | V_t I)}{p(v_{t+1} | V_t I)} p(v_{t+1} | \theta_{t'} V_t I)$$

soit encore

$$p(\theta_{t'} | V_{t+1} I) = \underbrace{\frac{p(v_{t+1} | \theta_{t'} V_t I)}{p(v_{t+1} | V_t I)}}_{\text{apport complémentaire}} p(\theta_{t'} | V_t I)$$

On peut chercher de même à déterminer des quantités du type :

$$p(v_{t'} | V_t I)$$

Dans ce cas, seule la prédiction ($t' > t$) a un sens puisque :

$$t' \leq t \quad p(v_{t'}|V_t I) = \delta(v_{t'} - \text{“mesure à l’instant } t’\text{”})$$

Cette équation signifie simplement qu’il n’y a aucun intérêt à estimer ce qui a déjà été mesuré !

5.1.3 Modèle de mesure

De quoi peut dépendre (logiquement) le résultat de la mesure à l’instant t ?

- du temps t
- de la succession des états $\Theta_t = \theta_0 \dots \theta_t$
- des mesures précédentes $V_{t-1} = v_1 \dots v_{t-1}$

On peut donc définir le modèle de mesure comme étant la donnée de :

$$\boxed{p(v_t|\Theta_t V_{t-1} I)} \tag{5.1}$$

Il est clair qu’un modèle de mesure de ce type est d’une complexité excessive en pratique.

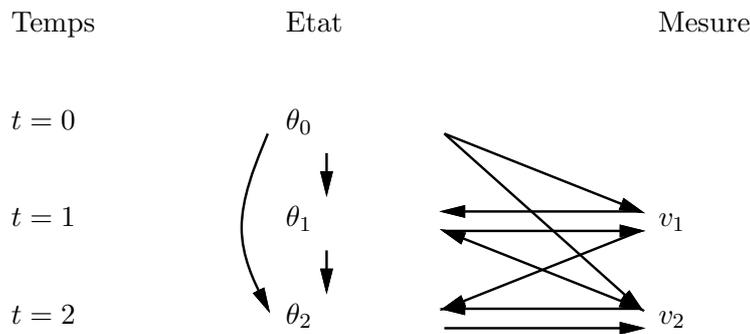


FIG. 5.1 – *Dépendance logique pour le modèle de mesure.*

En général on fait l’hypothèse que la mesure ne dépend que de l’état présent du système et du temps présent, soit :

$$\boxed{p(v_t|\Theta_t V_{t-1} I) = p(v_t|\theta_t I)} \tag{5.2}$$

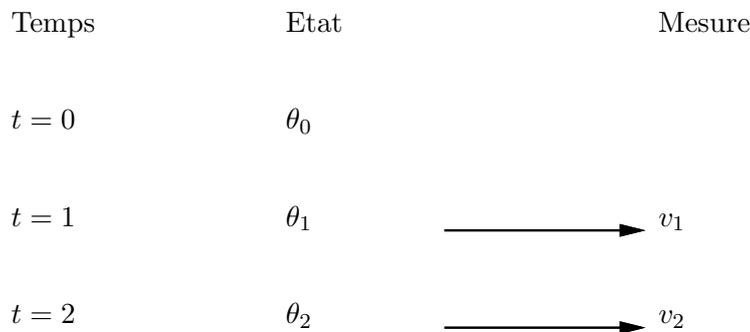


FIG. 5.2 – *Modèle de mesure, dans l’hypothèse où la mesure ne dépend que de l’état présent.*

Pour simplifier plus encore, on se contente en général de spécifier ce modèle à l'ordre deux seulement :

$$\boxed{p(v_t|\theta_t I) = f_t(v_t, x(t, \theta_t), r(t, \theta_t))} \quad (5.3)$$

Cette notation signifie que f_t est une fonction de la variable v_t , de moyenne $x(t, \theta_t)$ (vecteur de \mathbb{R}^m) et de covariance $r(t, \theta_t)$ (matrice carrée de $\mathbb{R}^{m \times m}$) :

$$\begin{aligned} \int p(v_t|\theta_t I) dv_t &= \int f_t(\omega, -, -) d\omega = 1 \\ \int v_t p(v_t|\theta_t I) dv_t &= \int \omega f_t(\omega, -, -) d\omega = x(t, \theta_t) \\ \int (v_t - x(t, \theta_t)) \cdot (v_t - x(t, \theta_t))^T p(v_t|\theta_t I) dv_t &= \int (\omega - \langle \omega \rangle) \cdot (\omega - \langle \omega \rangle)^T f_t(\omega, -, -) d\omega \\ &= r(t, \theta_t) \end{aligned}$$

Par exemple, ce peut être une gaussienne :

$$\begin{aligned} p(v_t|\theta_t I) &= \mathcal{N}(v_t, x(t, \theta_t), r(t, \theta_t)) \\ &= \frac{1}{(2\pi)^{m/2} \sqrt{|r(t, \theta_t)|}} \exp\left(-\frac{1}{2} (v_t - x(t, \theta_t))^T r(t, \theta_t)^{-1} (v_t - x(t, \theta_t))\right) \end{aligned}$$

Une autre écriture est encore souvent employée dans la littérature :

$$\boxed{v_t = x(t, \theta_t) + b_t, \text{ où } b_t \text{ est un bruit centré de covariance } r(t, \theta_t)} \quad (5.4)$$

Remarque : L'hypothèse simplificatrice classique que la mesure ne dépend que de l'état présent nous prive de la possibilité de rendre le filtrage (de Kalman) adaptatif.

En effet, s'il est assez naturel que la mesure ne dépende *physiquement* que de l'état présent du système, l'accumulation des données devrait permettre *logiquement* de préciser la qualité de l'approximation du modèle, c'est-à-dire la forme de la distribution autour de la valeur moyenne $x(t, \theta_t)$.

On pourrait par exemple faire dépendre la covariance du modèle de mesure des mesures précédentes ($r(t, \theta_t, V_{t-1})$), afin de profiter de l'information apportée par ces mesures à notre confiance dans le modèle de mesure. Cela rendrait en ce sens le filtrage adaptatif. Une telle théorie est cependant complexe, et dépasse le cadre de ce cours.

5.1.4 Modèle d'évolution

Comme pour la mesure à l'instant t , l'état du système dépend du temps, des états et des mesures précédentes. Le modèle d'évolution est constitué par :

$$\boxed{p(\theta_{t+1}|\Theta_t V_t I)} \quad (5.5)$$

Attention : La dépendance de θ_{t+1} avec V_t , soit avec les mesures précédentes, ne résulte en aucun cas d'une causalité physique (une mesure optronique, lidar, radar, sonar..., ne modifie pas ou très peu l'état du système observé). La dépendance est d'ordre purement logique : les données peuvent modifier notre prédiction pour θ_{t+1} , en excluant certains cas par exemple, mais en général elles ne modifient pas la "valeur vraie" de cet état.

Comme dans le cas du modèle de mesure, ce modèle d'évolution est trop complexe en pratique. On fait alors de façon similaire l'hypothèse que l'état présent ne dépend que de l'état précédent :

$$\boxed{p(\theta_{t+1}|\Theta_t V_t I) = p(\theta_{t+1}|\theta_t I)} \quad (5.6)$$

Pour simplifier plus encore, on se contente en général de spécifier ce modèle à l'ordre deux seulement :

$$\boxed{p(\theta_{t+1}|\theta_t I) = g_t(\theta_{t+1}, \varphi(t, \theta_t), q(t, \theta_t))} \quad (5.7)$$

Cette notation signifie que g_t est une fonction de la variable θ_{t+1} , de moyenne $\varphi(t, \theta_t)$ (vecteur de \mathbb{R}^n) et de covariance $q(t, \theta_t)$ (matrice carrée de $\mathbb{R}^{n \times n}$) :

$$\begin{aligned} \int p(\theta_{t+1}|\theta_t I) d\theta_t &= \int g_t(\omega, -, -) d\omega = 1 \\ \int \theta_{t+1} p(\theta_{t+1}|\theta_t I) d\theta_{t+1} &= \int \omega g_t(\omega, -, -) d\omega = \varphi(t, \theta_t) \\ \int (\theta_{t+1} - \varphi(t, \theta_t)) \cdot (\theta_{t+1} - \varphi(t, \theta_t))^T p(\theta_{t+1}|\theta_t I) d\theta_{t+1} &= \int (\omega - \langle \omega \rangle) \cdot (\omega - \langle \omega \rangle)^T g_t(\omega, -, -) d\omega \\ &= q(t, \theta_t) \end{aligned}$$

Par exemple, ce peut être une gaussienne :

$$\begin{aligned} p(\theta_{t+1}|\theta_t I) &= \mathcal{N}(\theta_{t+1}, \varphi(t, \theta_t), q(t, \theta_t)) \\ &= \frac{1}{(2\pi)^{n/2} \sqrt{|q(t, \theta_t)|}} \exp\left(-\frac{1}{2} (\theta_{t+1} - \varphi(t, \theta_t))^T q(t, \theta_t)^{-1} (\theta_{t+1} - \varphi(t, \theta_t))\right) \end{aligned}$$

Temps	Etat	Mesure
$t = 0$	θ_0	
	↓	
$t = 1$	θ_1	v_1
	↓	
$t = 2$	θ_2	v_2

FIG. 5.3 – Modèle d'évolution, dans l'hypothèse où l'état présent ne dépend que de l'état précédent.

Une autre écriture est souvent employée dans la littérature :

$$\boxed{\theta_{t+1} = \varphi(t, \theta_t) + w_t, \text{ où } w_t \text{ est un bruit centré de covariance } q(t, \theta_t)} \quad (5.8)$$

Remarque : Une autre écriture est encore souvent employée dans la littérature :

$$\boxed{\theta_{t+1} = \varphi(t, \theta_t) + \gamma(t, \theta_t) \cdot w_t} \quad (5.9)$$

où $w_t (\in \mathbb{R}^p)$ est un bruit centré de covariance $q(t, \theta_t) (\in \mathbb{R}^{p \times p})$ et $\gamma(t, \theta_t)$ est une matrice rectangle de $\mathbb{R}^{n \times p}$. Il suffit de poser $w'_t = \gamma(t, \theta_t) \cdot w_t$ pour se ramener au cas précédent, puisqu'alors w'_t est centré, et de covariance $\gamma(t, \theta_t) \cdot q(t, \theta_t) \cdot \gamma(t, \theta_t)^T$.

En effet, si on pose à t et θ_t donnés, $q(t, \theta_t) = q$, $\gamma(t, \theta_t) = \gamma$, $w_t = w$ et $w'_t = w'$, alors par définition :

$$q = \int w.w^T p(w)dw$$

donc

$$\begin{aligned} \gamma.q.\gamma^T &= \int \gamma.w.w^T.\gamma^T p(w)dw \\ &= \int (\gamma.w).(\gamma.w)^T p(w)dw \end{aligned}$$

comme de plus $p(w')dw' = p(w)dw$, il vient $\gamma.q.\gamma^T = \int w'.w'^T p(w')dw'$.

5.1.5 Modèles linéaires de mesure et d'évolution

Un cas particulier fondamental est obtenu en supposant le modèle linéaire suivant :

modèle de mesure			
moyenne	$x(t, \theta_t) = X_t.\theta_t$	X_t	: matrice de $\mathbb{R}^{m \times n}$
covariance	$r(t, \theta_t) = R_t$	R_t	: matrice de $\mathbb{R}^{m \times m}$
modèle d'évolution			
moyenne	$\varphi(t, \theta_t) = \Phi_t.\theta_t$	Φ_t	: matrice de $\mathbb{R}^{n \times n}$
covariance	$q(t, \theta_t) = Q_t$	Q_t	: matrice de $\mathbb{R}^{n \times n}$

Il est essentiel de remarquer que dans ce cas les covariances des modèles ne dépendent plus de l'état θ_t , ce qui simplifie beaucoup les calculs.

EXEMPLE - 5.2 ————— Déplacement d'un mobile à vitesse constante dans le plan

On étudie le déplacement d'un mobile (solide indéformable) dans un plan. On repère la position du centre de gravité du mobile par le vecteur $(z_1(t), z_2(t))^T$. La vitesse est donnée par $\left(\frac{dz_1}{dt}(t), \frac{dz_2}{dt}(t)\right)^T$. On suppose pour le modèle d'évolution que cette vitesse est constante, soit en discrétisant (l'incrément de temps est supposé égal à 1) :

$$\begin{cases} z_1(t+1) &= z_1(t) + \frac{dz_1}{dt}(t) \\ z_2(t+1) &= z_2(t) + \frac{dz_2}{dt}(t) \\ \frac{dz_1}{dt}(t+1) &= \frac{dz_1}{dt}(t) \\ \frac{dz_2}{dt}(t+1) &= \frac{dz_2}{dt}(t) \end{cases}$$

Si l'on prend :

$$\theta_t \doteq \begin{pmatrix} z_1(t) \\ z_2(t) \\ \frac{dz_1}{dt}(t) \\ \frac{dz_2}{dt}(t) \end{pmatrix}$$

on aura alors :

$$\Phi_t = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

L'incertitude dans ce modèle sera mesurée par une matrice de covariance Q_t (matrice 4×4).

On suppose par ailleurs qu'on ne mesure que la position, c'est-à-dire que :

$$v_t \doteq \begin{pmatrix} z_1(t) \\ z_2(t) \end{pmatrix}$$

on aura alors :

$$X_t = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

L'incertitude dans ce modèle sera mesurée par une matrice de covariance R_t (matrice 2×2).

5.1.6 Espérance conditionnelle

Définition 5.1

(Espérance conditionnelle)

L'opérateur espérance conditionnelle est défini par :

$$\langle \cdot \rangle_y \doteq \int (\cdot) p(x|y) dx$$

En particulier la moyenne conditionnelle de x sachant y est

$$\langle x \rangle_y = \int xp(x|y) dx$$

et la fonction caractéristique conditionnelle est

$$\varphi_{x|y}(u) = \langle \exp(iu^T \cdot x) \rangle_y = \int \exp(iu^T \cdot x) p(x|y) dx$$

Propriété 5.1**(La trace de la covariance est un produit scalaire)**

La trace de la covariance (conditionnelle ou non) de deux vecteurs aléatoires centrés est un produit scalaire de ces deux vecteurs aléatoires. L'autocovariance définit donc une norme, qui permet de définir les classes d'équivalence de vecteurs aléatoires de même norme (v.a. égaux en moyenne quadratique).

$$\begin{aligned} \text{Tr} \langle a.b^T \rangle_y &= \langle a^T.b \rangle_y = \iint a^T.b p(ab|y) da db \\ \text{Tr} \langle a.a^T \rangle_y &= \langle a^T.a \rangle_y = \int a^T.a p(a|y) da \doteq \|a\|_y^2 \end{aligned}$$

Démonstration : Il faut montrer que la trace de la covariance possède les propriétés de symétrie et de linéarité (pour le produit scalaire), et qu'elle est de plus définie positive (pour la norme).

- Symétrie : $\langle a^T.b \rangle_y = \langle b^T.a \rangle_y$
ceci est évident puisque $a^T.b = b^T.a$ et que $p(ab|y) = p(ba|y)$.
- Linéarité : Il faut montrer que $\forall a, b, c$ vecteurs aléatoires, $\forall \lambda \in \mathbb{R}$,

$$\begin{aligned} (a) \quad \langle (\lambda a)^T.b \rangle_y &= \lambda \langle a^T.b \rangle_y \\ (b) \quad \langle (a+b)^T.c \rangle_y &= \langle a^T.c \rangle_y + \langle b^T.c \rangle_y \end{aligned}$$

- (a) Soit $a' = \lambda a$, alors $p(a'b|y) da' db = p(ab|y) da db$, d'où :

$$\begin{aligned} \langle (\lambda a)^T.b \rangle_y &= \int a'^T.b p(a'b|y) da' db \\ &= \lambda \int a^T.b p(ab|y) da db \\ \langle (\lambda a)^T.b \rangle_y &= \lambda \langle a^T.b \rangle_y \end{aligned}$$

- (b) Soit $s = a + b$, alors $p(s|cy) = \int p(A = s - x|cy)p(B = x|cy) dx$ (théorème classique des probabilités : la distribution d'une somme de v.a. est la convolution des distributions de ces v.a.). Donc :

$$\begin{aligned} \langle s^T.c \rangle_y &= \iint s^T.c p(sc|y) ds dc \\ &= \iint s^T.c p(s|cy) p(c|y) ds dc \\ \langle s^T.c \rangle_y &= \iiint s^T.c p(A = s - x|cy)p(B = x|cy) p(c|y) ds dc dx \end{aligned}$$

Soit le changement de variable suivant :

$$\begin{cases} a &= s - x \\ b &= x \\ c &= c \end{cases}$$

1. Attention, le produit scalaire $a^T.b$ est celui d'un espace vectoriel classique ; le produit scalaire que nous définissons est défini pour un espace de vecteurs *aléatoires*.

dont le jacobien² est unitaire :

$$J = \begin{vmatrix} \begin{pmatrix} 1 & & \\ & \searrow & \\ & & 1 \end{pmatrix} & 0 & 0 \\ \begin{pmatrix} -1 & & \\ & \searrow & \\ & & -1 \end{pmatrix} & \begin{pmatrix} 1 & & \\ & \searrow & \\ & & 1 \end{pmatrix} & 0 \\ 0 & 0 & \begin{pmatrix} 1 & & \\ & \searrow & \\ & & 1 \end{pmatrix} \end{vmatrix} = 1$$

Il vient donc :

$$\begin{aligned} \langle s^T . c \rangle_y &= \iiint (a+b)^T . c p(a|cy) p(b|cy) p(c|y) da db dc \\ &= \iint a^T . c p(a|cy) p(c|y) da dc + \iint b^T . c p(b|cy) p(c|y) db dc \\ &= \iint a^T . c p(ac|y) da dc + \iint b^T . c p(bc|y) db dc \\ \langle s^T . c \rangle_y &= \langle a^T . c \rangle_y + \langle b^T . c \rangle_y \end{aligned}$$

- La trace de la covariance est définie positive si $\langle a^T . a \rangle_y = 0 \Rightarrow a = 0$. Cette condition est réalisée puisque $a^T . a \geq 0$ et $p(a|y) \geq 0$ sans que cette distribution puisse être uniformément nulle (puisque $\int p(a|y) da = 1$).

5.1.7 Choix de l'estimateur

Ainsi que nous l'avons dit en introduction, le but du filtrage de Kalman est d'obtenir la densité de probabilité :

$$p(\theta_t | V_t I)$$

À partir de cette densité de probabilité, on peut alors chercher un estimateur $\hat{\theta}_t$ de l'état du système. Mais il faut garder à l'esprit que c'est dans le choix de l'estimateur que s'introduit l'arbitraire de cette méthode. En aucun cas, la donnée de l'estimateur $\hat{\theta}_t$ ne peut remplacer celle de $p(\theta_t | V_t I)$.

Quels sont les estimateurs possibles? Les plus utilisés sont ceux définis au sens du MAP (maximum *a posteriori*), au sens du MV (maximum de vraisemblance) ou au sens des MC (moindres carrés).

- MAP

$$\hat{\theta}_t \doteq \underset{\theta_t}{\text{Argmax}} \{p(\theta_t | V_t I)\}$$

- MV

$$\hat{\theta}_t \doteq \underset{\theta_t}{\text{Argmax}} \{p(V_t | \theta_t I)\}$$

- MC

$$\hat{\theta}_t = \langle \theta_t \rangle_{V_t I} = \int \theta_t p(\theta_t | V_t I) d\theta_t$$

2. Jacobien : déterminant de la matrice des dérivées partielles des nouvelles variables par rapport aux anciennes, ici une matrice carrée $3n \times 3n$.

Détaillons plus précisément le cas MC. Par définition, l'estimateur des moindres carrés minimise la trace de la covariance conditionnelle de θ_t (pour les calculs qui suivent, se souvenir que $\hat{\theta}_t$ est une fonction *a priori* quelconque des données, donc de $V_t I$, mais pas de θ_t):

$$\begin{aligned} \langle (\theta_t - \hat{\theta}_t)^T \cdot (\theta_t - \hat{\theta}_t) \rangle_{V_t I} &\doteq \int (\theta_t - \hat{\theta}_t)^T \cdot (\theta_t - \hat{\theta}_t) p(\theta_t | V_t I) d\theta_t \\ &= \langle \theta_t^T \cdot \theta_t \rangle_{V_t I} + \hat{\theta}_t^T \cdot \hat{\theta}_t - 2 \cdot \langle \theta_t \rangle_{V_t I}^T \cdot \hat{\theta}_t \\ &= \|\theta_t\|_{V_t I}^2 + \|\hat{\theta}_t - \langle \theta_t \rangle_{V_t I}\|_{V_t I}^2 - \|\langle \theta_t \rangle_{V_t I}\|_{V_t I}^2 \end{aligned}$$

Il est donc clair que cette quantité est minimale pour $\hat{\theta}_t = \langle \theta_t \rangle_{V_t I}$, c'est-à-dire que l'estimateur des moindres carrés est la moyenne de la distribution. Dans ces conditions, la trace de la covariance conditionnelle de l'estimateur est égale à celle de la distribution :

$$\langle (\theta_t - \hat{\theta}_t)^T \cdot (\theta_t - \hat{\theta}_t) \rangle_{V_t I} = \|\theta_t - \langle \theta_t \rangle_{V_t I}\|_{V_t I}^2$$

5.2 Solution générale

On suppose donnés les modèles de mesure et d'évolution (simplifiés) :

$$p(v_t | \Theta_t V_{t-1} I) = p(v_t | \theta_t I)$$

$$p(\theta_{t+1} | \Theta_t V_t I) = p(\theta_{t+1} | \theta_t I)$$

Par ailleurs, on suppose donnée $p(\theta_0 | I)$.

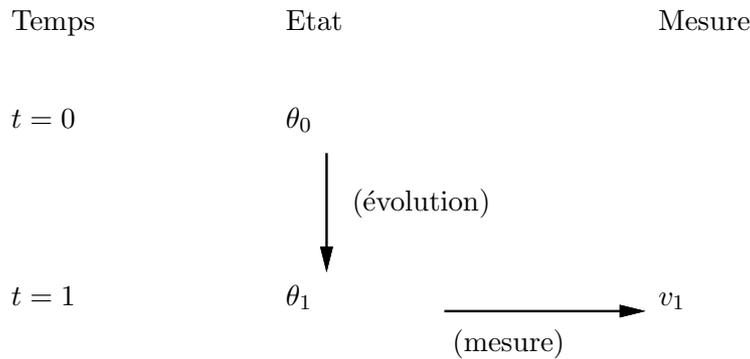


FIG. 5.4 – Départ de la prédiction (conditions initiales).

Avant l'acquisition des premières mesures v_1 , on peut prédire les distributions de θ_1 et de v_1 .

– prédiction de l'état suivant du système

$$p(\theta_1 | I) = \int p(\theta_1 | \theta_0 I) p(\theta_0 | I) d\theta_0$$

– prédiction de la mesure suivante

$$\begin{aligned} p(v_1 | I) &= \int p(v_1 | \theta_1 I) p(\theta_1 | I) d\theta_1 \\ &= \iint p(v_1 | \theta_1 I) p(\theta_1 | \theta_0 I) p(\theta_0 | I) d\theta_0 d\theta_1 \end{aligned}$$

Après l'acquisition des premières mesures v_1 , on peut calculer la densité de probabilité pour l'estimation :

$$p(\theta_1|v_1I) = \frac{p(\theta_1|I)}{p(v_1|I)} p(v_1|\theta_1I)$$

Cette expression combine donc par la règle du produit les prédictions avec la vraisemblance (ou modèle de mesure).

Cette structure de prédiction/estimation se répète récursivement pour les instants suivants. Avant l'acquisition de la mesure v_{t+1} , on peut prédire les distributions de θ_{t+1} et de v_{t+1} .

– prédiction de l'état suivant du système

$$p(\theta_{t+1}|V_tI) = \int \underbrace{p(\theta_{t+1}|\theta_tI)}_{\text{évolution}} \underbrace{p(\theta_t|V_tI)}_{\text{estimation précédente}} d\theta_t$$

– prédiction de la mesure suivante

$$\begin{aligned} p(v_{t+1}|V_tI) &= \int p(v_{t+1}|\theta_{t+1}I) p(\theta_{t+1}|V_tI) d\theta_{t+1} \\ &= \iiint \underbrace{p(v_{t+1}|\theta_{t+1}I)}_{\text{modèle de mesure}} \underbrace{p(\theta_{t+1}|\theta_tI)}_{\text{évolution}} \underbrace{p(\theta_t|V_tI)}_{\text{estimation précédente}} d\theta_t d\theta_{t+1} \end{aligned}$$

Pour écrire les relations précédentes, nous avons utilisé le fait que :

$$\begin{cases} p(v_{t+1}|\theta_{t+1}V_tI) = p(v_{t+1}|\theta_{t+1}I) \\ p(\theta_{t+1}|\theta_tV_tI) = p(\theta_{t+1}|\theta_tI) \end{cases}$$

Ces relations sont contenues dans les hypothèses qui permettent de simplifier les modèles de mesure et d'évolution.

Après l'acquisition de la mesure v_{t+1} , on peut calculer la densité de probabilité pour l'estimation :

$$\begin{aligned} p(\theta_{t+1}|V_{t+1}I) &= p(\theta_{t+1}|v_{t+1}V_tI) \\ &= \frac{p(\theta_{t+1}|V_tI)}{p(v_{t+1}|V_tI)} p(v_{t+1}|\theta_{t+1}V_tI) \end{aligned}$$

soit finalement

$$\boxed{p(\theta_{t+1}|V_{t+1}I) = \frac{p(\theta_{t+1}|V_tI)}{p(v_{t+1}|V_tI)} p(v_{t+1}|\theta_{t+1}I)}$$

Remarque : À partir de cette dernière expression, on retrouve aisément la formule de prédiction de la mesure, puisque

$$\int p(\theta_{t+1}|V_{t+1}I) d\theta_{t+1} = 1$$

et donc

$$p(v_{t+1}|V_tI) = \int p(v_{t+1}|\theta_{t+1}I) p(\theta_{t+1}|V_tI) d\theta_{t+1}$$

Cette expression est une relation de normalisation.

5.3 Solution gaussienne – cas linéaire

On se place dans le cas linéaire et on suppose des densités de probabilité gaussiennes pour :

- les conditions initiales $p(\theta_0|I) = \mathcal{N}(\theta_0, \hat{\theta}_0, P_0)$
- le modèle de mesure $p(v_t|\theta_t I) = \mathcal{N}(v_t, X_t.\theta_t, R_t)$
- le modèle d'évolution $p(\theta_{t+1}|\theta_t I) = \mathcal{N}(\theta_{t+1}, \Phi_t.\theta_t, Q_t)$

5.3.1 Propagation du caractère gaussien

Propriété 5.2

(Propagation du caractère gaussien)

Soient α et β deux vecteurs aléatoires. α est un vecteur aléatoire gaussien à n composantes, de moyenne q (vecteur de dimension n) et de covariance C (matrice $n \times n$), soit $p(\alpha) = \mathcal{N}(\alpha, q, C)$. Le vecteur aléatoire β est de dimension m , et tel que $p(\beta|\alpha) = \mathcal{N}(\beta, A.\alpha, B)$, avec A une matrice $m \times n$ et B une matrice $m \times m$. Alors β est un vecteur aléatoire gaussien :

$$p(\beta) = \int \mathcal{N}(\beta, A.\alpha, B) \mathcal{N}(\alpha, q, C) d\alpha = \mathcal{N}(\beta, A.q, A.C.A^T + B)$$

Preuve : La fonction caractéristique de β est

$$\begin{aligned} \varphi_\beta(u) &\doteq \langle \exp(iu^T.\beta) \rangle \\ \varphi_\beta(u) &= \int \exp(iu^T.\beta) p(\beta) d\beta \end{aligned}$$

Puisque

$$p(\beta) = \int p(\beta|\alpha) p(\alpha) d\alpha$$

il vient

$$\begin{aligned} \varphi_\beta(u) &= \int p(\alpha) d\alpha \int \exp(iu^T.\beta) p(\beta|\alpha) d\beta \\ &= \int p(\alpha) d\alpha \langle \exp(iu^T.\beta) \rangle_\alpha \\ \varphi_\beta(u) &= \int p(\alpha) d\alpha \varphi_{\beta|\alpha}(u) \end{aligned}$$

Or

$$\varphi_{\beta|\alpha}(u) = \exp\left(iu^T.(A.\alpha) - \frac{1}{2}u^T.B.u\right) = \exp\left(i(A^T.u)^T.\alpha - \frac{1}{2}u^T.B.u\right)$$

donc

$$\begin{aligned} \varphi_\beta(u) &= \exp\left(-\frac{1}{2}u^T.B.u\right) \int p(\alpha) d\alpha \exp\left(i(A^T.u)^T.\alpha\right) \\ &= \exp\left(-\frac{1}{2}u^T.B.u\right) \varphi_\alpha(A^T.u) \\ &= \exp\left(-\frac{1}{2}u^T.B.u\right) \exp\left(i(A^T.u)^T.q - \frac{1}{2}(A^T.u)^T.C.(A^T.u)\right) \\ \varphi_\beta(u) &= \exp\left(iu^T.(A.q)\right) \exp\left(-\frac{1}{2}u^T.(B + A.C.A^T).u\right) \end{aligned}$$

Cette propriété de stabilité des gaussiennes (une de plus!) permet de montrer par récurrence que toutes les densités de probabilité intervenant sont gaussiennes. Nous pouvons reprendre les calculs de la section précédente en explicitant le caractère gaussien des densités de probabilité.

– Prédiction initiale de l'état :

$$\begin{aligned} p(\theta_1|I) &= \int p(\theta_1|\theta_0 I) p(\theta_0|I) d\theta_0 \\ &= \int \mathcal{N}(\theta_1, \Phi_0, \theta_0, Q_0) \mathcal{N}(\theta_0, \hat{\theta}_0) d\theta_0 \\ p(\theta_1|I) &= \mathcal{N}(\theta_1, \Phi_0, \hat{\theta}_0, Q_0 + \Phi_0 \cdot P_0 \cdot \Phi_0^T) \end{aligned}$$

On pose donc

$$\boxed{\begin{aligned} p(\theta_1|I) &= \mathcal{N}(\theta_1, \hat{\theta}_1^-, P_1^-) \\ \text{avec} \quad &\begin{cases} \hat{\theta}_1^- &= \Phi_0 \cdot \hat{\theta}_0 \\ P_1^- &= Q_0 + \Phi_0 \cdot P_0 \cdot \Phi_0^T \end{cases} \end{aligned}}$$

– Prédiction initiale de la mesure :

$$\begin{aligned} p(v_1|I) &= \int p(v_1|\theta_1 I) p(\theta_1|I) d\theta_1 \\ &= \int \mathcal{N}(v_1, X_1, \theta_1, R_1) \mathcal{N}(\theta_1, \hat{\theta}_1^-, P_1^-) d\theta_1 \\ p(v_1|I) &= \mathcal{N}(v_1, X_1, \hat{\theta}_1^-, R_1 + X_1 \cdot P_1^- \cdot X_1^T) \end{aligned}$$

On pose donc

$$\boxed{\begin{aligned} p(v_1|I) &= \mathcal{N}(v_1, \hat{v}_1^-, S_1) \\ \text{avec} \quad &\begin{cases} \hat{v}_1^- &= X_1 \cdot \hat{\theta}_1^- \\ S_1 &= R_1 + X_1 \cdot P_1^- \cdot X_1^T \end{cases} \end{aligned}}$$

– Estimation de l'état :

$$p(\theta_1|v_1 I) = \frac{p(\theta_1|I)}{p(v_1|I)} p(v_1|\theta_1 I) = \frac{\mathcal{N}(\theta_1, \hat{\theta}_1^-, P_1^-)}{\mathcal{N}(v_1, \hat{v}_1^-, S_1)} \mathcal{N}(v_1, X_1, \theta_1, R_1)$$

Un produit de gaussiennes étant une gaussienne, on notera

$$\boxed{p(\theta_1|v_1 I) = \mathcal{N}(\theta_1, \hat{\theta}_1, P_1)}$$

L'identification de $\hat{\theta}_1$ et P_1 n'est pas triviale et sera faite plus loin.

On voit par récurrence que toutes les densités de probabilité nécessaires sont gaussiennes, et l'on note :

– Prédiction de l'état :

$$\boxed{\begin{array}{l} p(\theta_{t+1}|V_t I) = \mathcal{N}(\theta_{t+1}, \hat{\theta}_{t+1}^-, P_{t+1}^-) \\ \text{avec} \quad \left\{ \begin{array}{l} \hat{\theta}_{t+1}^- = \Phi_t \cdot \hat{\theta}_t \\ P_{t+1}^- = Q_t + \Phi_t \cdot P_t \cdot \Phi_t^T \end{array} \right. \end{array}}$$

– Prédiction de la mesure :

$$\boxed{\begin{array}{l} p(v_{t+1}|V_t I) = \mathcal{N}(v_{t+1}, \hat{v}_{t+1}^-, S_{t+1}) \\ \text{avec} \quad \left\{ \begin{array}{l} \hat{v}_{t+1}^- = X_{t+1} \cdot \hat{\theta}_{t+1}^- \\ S_{t+1} = R_{t+1} + X_{t+1} \cdot P_{t+1}^- \cdot X_{t+1}^T \end{array} \right. \end{array}}$$

– Estimation de l'état :

$$\boxed{p(\theta_{t+1}|V_{t+1} I) = \mathcal{N}(\theta_{t+1}, \hat{\theta}_{t+1}, P_{t+1})}$$

5.3.2 Calcul du gain et de la covariance

Il reste à calculer la moyenne $\hat{\theta}_{t+1}$ et la covariance P_{t+1} à partir de la relation

$$p(\theta_{t+1}|V_{t+1} I) = \frac{p(\theta_{t+1}|V_t I)}{p(v_{t+1}|V_t I)} p(v_{t+1}|\theta_{t+1} I)$$

ou

$$\mathcal{N}(\theta_{t+1}, \hat{\theta}_{t+1}, P_{t+1}) = \frac{\mathcal{N}(\theta_{t+1}, \hat{\theta}_{t+1}^-, P_{t+1}^-)}{\mathcal{N}(v_{t+1}, \hat{v}_{t+1}^-, S_{t+1})} \mathcal{N}(v_{t+1}, X_{t+1} \cdot \hat{\theta}_{t+1}, R_{t+1})$$

Nous oublierons l'indice $(t+1)$ dans les calculs à suivre.

$$\begin{aligned} (\theta - \hat{\theta})^T \cdot P^{-1} \cdot (\theta - \hat{\theta}) &= (\theta - \hat{\theta}^-)^T \cdot P^{-1} \cdot (\theta - \hat{\theta}^-) + (v - X \cdot \theta)^T \cdot R^{-1} \cdot (v - X \cdot \theta) \\ &\quad - (v - v^-)^T \cdot S^{-1} \cdot (v - v^-) \end{aligned}$$

$$\text{On pose} \quad \left\{ \begin{array}{l} \Delta\theta = \theta - \hat{\theta}^- \\ \Delta v = v - v^- \end{array} \right.$$

Alors $v - X \cdot \theta = \Delta v - X \cdot \Delta\theta$, et donc

$$\begin{aligned} (\theta - \hat{\theta})^T \cdot P^{-1} \cdot (\theta - \hat{\theta}) &= \Delta\theta^T \cdot P^{-1} \cdot \Delta\theta + (\Delta v - X \cdot \Delta\theta)^T \cdot R^{-1} \cdot (\Delta v - X \cdot \Delta\theta) \\ &\quad - \Delta v^T \cdot S^{-1} \cdot \Delta v \end{aligned}$$

Cette dernière expression semble indiquer que $(\theta - \hat{\theta})$ doit être une combinaison linéaire de $\Delta\theta$ et de Δv (ce que nous devons vérifier par la suite). On pose donc :

$$\boxed{\theta - \hat{\theta} = \Delta\theta - K \cdot \Delta v}$$

où K est une matrice (“le gain du filtre de Kalman”) à déterminer.

$$\begin{aligned} (\theta - \hat{\theta})^T . P^{-1} . (\theta - \hat{\theta}) &= \Delta\theta^T . P^{-1} . \Delta\theta - \Delta\theta^T . P^{-1} . K . \Delta v \\ &\quad - \Delta v^T . K^T . P^{-1} . \Delta\theta + \Delta v^T . K^T . P^{-1} . K . \Delta v \\ (\theta - \hat{\theta})^T . P^{-1} . (\theta - \hat{\theta}) &= \Delta\theta^T . [(P^-)^{-1} + X^T . R^{-1} . X] . \Delta\theta - \Delta\theta^T . X^T . R^{-1} . \Delta v \\ &\quad - \Delta v^T . R^{-1} . X . \Delta\theta + \Delta v^T . [R^{-1} - S^{-1}] . \Delta v \end{aligned}$$

On en tire le système d'équations :

$$\left\{ \begin{array}{ll} P^{-1} = (P^-)^{-1} + X^T . R^{-1} . X & (a) \\ P^{-1} . K = X^T . R^{-1} & (b) \\ K^T . P^{-1} = R^{-1} . X & (c) \\ K^T . P^{-1} . K = R^{-1} - S^{-1} & (d) \end{array} \right.$$

$$\begin{aligned} X^T . (d) &\Rightarrow X^T . K^T . P^{-1} . K = X^T . R^{-1} - X^T . S^{-1} \\ (b) &\Rightarrow X^T . K^T . P^{-1} . K = P^{-1} . K - X^T . S^{-1} \\ (c) &\Rightarrow X^T . R^{-1} . X . K = P^{-1} . K - X^T . S^{-1} \\ (a) &\Rightarrow X^T . R^{-1} . X . K = (P^-)^{-1} . K + X^T . R^{-1} . X . K - X^T . S^{-1} \end{aligned}$$

donc $(P^-)^{-1} . K = X^T . S^{-1}$, ou

$$\boxed{K = P^- . X^T . S^{-1}}$$

De plus

$$(a) \text{ et } (b) \Rightarrow P^{-1} = (P^-)^{-1} + P^{-1} . K . X$$

donc $P^{-1} . (I_d - K . X) = (P^-)^{-1}$, soit

$$\boxed{P = (I_d - K . X) . P^-}$$

Il reste à vérifier que K et P ainsi obtenus sont bien solutions de l'identité de départ. Il suffit de vérifier que les expressions obtenues pour K et P impliquent (a), (b), (c) et (d).

– (b) :

$$\begin{aligned} P . X^T - K . R &= (I_d - K . X) . P^- . X^T - K . R && \text{définition de } P \\ &= P^- . X^T - K . (X . P^- . X^T + R) \\ &= P^- . X^T - K . S && \text{définition de } S \end{aligned}$$

$$P . X^T - K . R = 0 \quad \text{définition de } K$$

donc $P^{-1} . K = X^T . R^{-1}$.

– (c) : Il suffit de transposer (b), compte tenu du fait que P et R sont des matrices symétriques par définition.

– (a) :

$$P^{-1} \cdot (I_d - K \cdot X) = (P^-)^{-1} \quad \text{définition de } P$$

$$P^{-1} - (P^{-1} \cdot K) \cdot X = (P^-)^{-1}$$

$$P^{-1} = X^T \cdot R^{-1} \cdot X + (P^-)^{-1} \quad \text{d'après (b)}$$

– (d) :

$$K^T \cdot P^{-1} \cdot K = R^{-1} \cdot X \cdot K \quad \text{d'après (b)}$$

$$= R^{-1} \cdot (X \cdot P^- \cdot X^T) \cdot S^{-1} \quad \text{définition de } K$$

$$= R^{-1} \cdot (S - R) \cdot S^{-1} \quad \text{définition de } S$$

$$K^T \cdot P^{-1} \cdot K = R^{-1} - S^{-1}$$

5.3.3 Résumé et commentaires

Les matrices de covariance dont on suit l'évolution au cours du temps, P_t^- , S_t et P_t (qu'il faut évaluer dans cet ordre), ne dépendent pas des données V_t que l'on accumule au cours de l'évolution du système; elles ne dépendent que des modèles (X_t , R_t , Φ_t , Q_t) et des données initiales (P_0).

Cette remarque rejoint celle sur le filtrage adaptatif (page 5.1.3) : c'est l'absence d'une dépendance des covariances sur l'état du système ou sur les mesures qui empêche que le filtrage soit adaptatif, mais qui rend la résolution analytique précédente possible. Si par exemple X_t , R_t , Φ_t ou Q_t dépendait de θ_t , la propriété de stabilité des gaussiennes ne pourrait être employée.

P_t^- , S_t et P_t ne dépendant pas des données acquises au cours de l'évolution, ces matrices peuvent être calculées par avance (hors ligne), ce qui représente un allègement important des calculs nécessaires en temps réel. La question se pose alors de savoir si ces matrices de covariance tendent vers une limite (stabilité), si elles oscillent, ou si elles divergent (instabilité : notre confiance dans l'estimation va en s'amenuisant). Il est en général très difficile de répondre à cette question même dans des cas simples, car les équations de Kalman sont fortement couplées.

5.4 Solution des moindres carrés

Nous avons discuté précédemment le choix de l'estimateur. Dans le cas gaussien, les critères MAP (maximum *a posteriori*) et MC (moindres carrés) conduisent au même estimateur : la moyenne de la distribution. La covariance des moindres carrés est également dans ce cas la covariance de la distribution gaussienne.

Nous nous plaçons ici hors de l'hypothèse gaussienne, et nous cherchons à déterminer le filtrage de Kalman pour des distributions quelconques, caractérisées seulement par leur moyenne et leur covariance, pour l'estimateur des moindres carrés. Cette étude, nous le verrons, redonnera exactement les résultats du cas gaussien traité précédemment.

Les modèles de mesure sont spécifiés au second ordre :

$$p(v_t | \theta_t I) = f_t(v_t, x(t, \theta_t), r(t, \theta_t))$$

Modèle linéaire de mesure	
$p(v_t \theta_t I) = \mathcal{N}(v_t, X_t \cdot \theta_t, R_t)$	
X_t	matrice $m \times n$
R_t	matrice $m \times m$
Modèle linéaire d'évolution	
$p(\theta_{t+1} \theta_t I) = \mathcal{N}(\theta_{t+1}, \Phi_t \cdot \theta_t, Q_t)$	
Φ_t	matrice $n \times n$
Q_t	matrice $n \times n$
Conditions initiales	
$p(\theta_0 I) = \mathcal{N}(\theta_0, \hat{\theta}_0, P_0)$	
P_0	matrice $n \times n$
Prédiction de l'état	
$p(\theta_{t+1} V_t I) = \mathcal{N}(\theta_{t+1}, \hat{\theta}_{t+1}^-, P_{t+1}^-)$	
$\hat{\theta}_{t+1}^- = \Phi_t \cdot \hat{\theta}_t$	vecteur n
$P_{t+1}^- = Q_t + \Phi_t \cdot P_t \cdot \Phi_t^T$	matrice $n \times n$
Prédiction de la mesure	
$p(v_{t+1} V_t I) = \mathcal{N}(v_{t+1}, \hat{v}_{t+1}^-, S_{t+1})$	
$\hat{v}_{t+1}^- = X_{t+1} \cdot \hat{\theta}_{t+1}^-$	vecteur m
$S_{t+1} = R_{t+1} + X_{t+1} \cdot P_{t+1}^- \cdot X_{t+1}^T$	matrice $m \times m$
Estimation de l'état	
$p(\theta_{t+1} V_{t+1} I) = \mathcal{N}(\theta_{t+1}, \hat{\theta}_{t+1}, P_{t+1})$	
$\hat{\theta}_{t+1} = \hat{\theta}_{t+1}^- + K_{t+1} \cdot (v_{t+1} - \hat{v}_{t+1}^-) = (I_d - K_{t+1} \cdot X_{t+1}) \cdot \hat{\theta}_{t+1}^- + K_{t+1} \cdot v_{t+1}$	vecteur n
$K_{t+1} = P_{t+1}^- \cdot X_{t+1}^T \cdot S_{t+1}^{-1}$	matrice $n \times m$
$P_{t+1} = (I_d - K_{t+1} \cdot X_{t+1}) \cdot P_{t+1}^-$	matrice $n \times n$

TAB. 5.1 – Filtrage de Kalman dans le cas linéaire gaussien.

où $x(t, \theta_t)$ (vecteur de \mathbb{R}^m) est la moyenne et $r(t, \theta_t)$ (matrice carrée de $\mathbb{R}^{m \times m}$) la covariance, et

$$p(\theta_{t+1} | \theta_t I) = g_t(\theta_{t+1}, \varphi(t, \theta_t), q(t, \theta_t))$$

où $\varphi(t, \theta_t)$ (vecteur de \mathbb{R}^n) est la moyenne et $q(t, \theta_t)$ (matrice carrée de $\mathbb{R}^{n \times n}$) la covariance. De même, la condition initiale est que $p(\theta_0 | I)$ a pour moyenne $\hat{\theta}_0$ et covariance P_0 .

5.4.1 Prédiction de l'état

On a par définition

$$\begin{cases} \hat{\theta}_t &= \langle \theta_t \rangle_{V_t I} = \int \theta_t p(\theta_t | V_t I) d\theta_t \\ P_t &= \langle (\theta_t - \hat{\theta}_t) \cdot (\theta_t - \hat{\theta}_t)^T \rangle_{V_t I} = \int (\theta_t - \hat{\theta}_t) \cdot (\theta_t - \hat{\theta}_t)^T p(\theta_t | V_t I) d\theta_t \end{cases}$$

$\hat{\theta}_t$ et P_t sont *a priori* des fonctions des données jusqu'à l'instant t et des données initiales (I).

Cherchons tout d'abord :

$$\begin{cases} \hat{\theta}_{t+1}^- &= \langle \theta_{t+1} \rangle_{V_t I} = \int \theta_{t+1} p(\theta_{t+1} | V_t I) d\theta_{t+1} \\ P_{t+1}^- &= \langle (\theta_{t+1} - \hat{\theta}_{t+1}^-) \cdot (\theta_{t+1} - \hat{\theta}_{t+1}^-)^T \rangle_{V_t I} \\ &= \int (\theta_{t+1} - \hat{\theta}_{t+1}^-) \cdot (\theta_{t+1} - \hat{\theta}_{t+1}^-)^T p(\theta_{t+1} | V_t I) d\theta_{t+1} \end{cases}$$

Puisque $p(\theta_{t+1} | V_t I) = \int p(\theta_{t+1} | \theta_t I) p(\theta_t | V_t I) d\theta_t$,

$$\int (\cdot) p(\theta_{t+1} | V_t I) d\theta_{t+1} = \int p(\theta_t | V_t I) d\theta_t \int (\cdot) p(\theta_{t+1} | \theta_t I) d\theta_{t+1}$$

d'où l'on tire :

$$\boxed{\hat{\theta}_{t+1}^- = \int \varphi(t, \theta_t) p(\theta_t | V_t I) d\theta_t}$$

soit l'espérance (conditionnelle) de $\varphi(t, \theta_t)$. En fait, puisqu'on ne connaît $p(\theta_t | V_t I)$ qu'au second ordre, on est bien incapable d'évaluer cette intégrale, sauf dans le cas du modèle linéaire, c'est-à-dire si $\varphi(t, \theta_t) = \Phi_t \cdot \theta_t$, auquel cas on trouve comme dans le cas gaussien :

$$\boxed{\hat{\theta}_{t+1}^- = \Phi_t \cdot \hat{\theta}_t}$$

Écrivons $\theta_{t+1} - \hat{\theta}_{t+1}^- = (\theta_{t+1} - \varphi(t, \theta_t)) - (\hat{\theta}_{t+1}^- - \varphi(t, \theta_t))$, alors

$$\int (\theta_{t+1} - \hat{\theta}_{t+1}^-) \cdot (\theta_{t+1} - \hat{\theta}_{t+1}^-)^T p(\theta_{t+1} | \theta_t I) d\theta_{t+1} = q(t, \theta_t) + (\hat{\theta}_{t+1}^- - \varphi(t, \theta_t)) \cdot (\hat{\theta}_{t+1}^- - \varphi(t, \theta_t))^T$$

et donc

$$\boxed{P_{t+1}^- = \int q(t, \theta_t) p(\theta_t | V_t I) d\theta_t + \int (\hat{\theta}_{t+1}^- - \varphi(t, \theta_t)) \cdot (\hat{\theta}_{t+1}^- - \varphi(t, \theta_t))^T p(\theta_t | V_t I) d\theta_t}$$

Ici encore, on ne sait pas en général évaluer ces intégrales, sauf dans le cas linéaire, c'est-à-dire si $\varphi(t, \theta_t) = \Phi_t \cdot \theta_t$ et $q(t, \theta_t) = Q_t$, et on retrouve alors la solution linéaire gaussienne :

$$\boxed{P_{t+1}^- = Q_t + \Phi_t \cdot P_t \cdot \Phi_t^T}$$

En effet,

$$\int q(t, \theta_t) p(\theta_t | V_t I) d\theta_t = Q_t$$

et

$$\begin{aligned} & \int (\hat{\theta}_{t+1}^- - \varphi(t, \theta_t)) \cdot (\hat{\theta}_{t+1}^- - \varphi(t, \theta_t))^T p(\theta_t | V_t I) d\theta_t \\ &= \Phi_t \cdot \left(\int (\hat{\theta}_t - \theta_t) \cdot (\hat{\theta}_t - \theta_t)^T p(\theta_t | V_t I) d\theta_t \right) \cdot \Phi_t^T \\ &= \Phi_t \cdot P_t \cdot \Phi_t^T \end{aligned}$$

5.4.2 Prédiction de la mesure

Cherchons maintenant :

$$\left\{ \begin{aligned} \hat{v}_{t+1}^- &= \langle v_{t+1} \rangle_{V_t I} = \int v_{t+1} p(v_{t+1} | V_t I) dv_{t+1} \\ S_{t+1} &= \langle (v_{t+1} - \hat{v}_{t+1}^-) \cdot (v_{t+1} - \hat{v}_{t+1}^-)^T \rangle_{V_t I} \\ &= \int (v_{t+1} - \hat{v}_{t+1}^-) \cdot (v_{t+1} - \hat{v}_{t+1}^-)^T p(v_{t+1} | V_t I) dv_{t+1} \end{aligned} \right.$$

Puisque $p(v_{t+1} | V_t I) = \int p(v_{t+1} | \theta_{t+1} I) p(\theta_{t+1} | V_t I) d\theta_{t+1}$,

$$\int (\cdot) p(v_{t+1} | V_t I) dv_{t+1} = \int p(\theta_{t+1} | V_t I) d\theta_{t+1} \int (\cdot) p(v_{t+1} | \theta_{t+1} I) dv_{t+1}$$

d'où l'on tire :

$$\boxed{\hat{v}_{t+1}^- = \int x(t+1, \theta_{t+1}) p(\theta_{t+1} | V_t I) d\theta_{t+1}}$$

que l'on sait calculer dans le cas du modèle linéaire, c'est-à-dire si $x(t+1, \theta_{t+1}) = X_{t+1} \cdot \theta_{t+1}$, auquel cas on trouve comme dans le cas gaussien :

$$\boxed{\hat{v}_{t+1}^- = X_{t+1} \cdot \hat{\theta}_{t+1}^-}$$

À partir de $v_{t+1} - \hat{v}_{t+1}^- = (v_{t+1} - x(t+1, \theta_{t+1})) - (\hat{v}_{t+1}^- - x(t+1, \theta_{t+1}))$, on montre comme précédemment :

$$\boxed{\begin{aligned} S_{t+1} &= \int r(t+1, \theta_{t+1}) p(\theta_{t+1} | V_t I) d\theta_{t+1} \\ &\quad + \int (\hat{v}_{t+1}^- - x(t+1, \theta_{t+1})) \cdot (\hat{v}_{t+1}^- - x(t+1, \theta_{t+1}))^T p(\theta_{t+1} | V_t I) d\theta_{t+1} \end{aligned}}$$

Et dans le cas linéaire, on retrouve alors la solution linéaire gaussienne :

$$\boxed{S_{t+1} = R_{t+1} + X_{t+1} \cdot P_{t+1}^- \cdot X_{t+1}^T}$$

5.4.3 Estimation de l'état

Il faut maintenant calculer :

$$\begin{cases} \hat{\theta}_t &= \langle \theta_t \rangle_{V_t I} = \int \theta_t p(\theta_t | V_t I) d\theta_t \\ P_t &= \langle (\theta_t - \hat{\theta}_t) \cdot (\theta_t - \hat{\theta}_t)^T \rangle_{V_t I} = \int (\theta_t - \hat{\theta}_t) \cdot (\theta_t - \hat{\theta}_t)^T p(\theta_t | V_t I) d\theta_t \end{cases}$$

sachant que

$$p(\theta_t | V_t I) = \frac{p(\theta_t | V_{t-1} I)}{p(v_t | V_{t-1} I)} p(v_t | \theta_t I)$$

Remarque : Pour utiliser un calcul du type de celui employé pour la prédiction de l'état et de la mesure, il faudrait connaître la dépendance de $\hat{\theta}_t$ et P_t avec v_t :

$$\int (\cdot) p(\theta_t | V_t I) d\theta_t = \frac{1}{p(v_t | V_{t-1} I)} \int (\cdot) p(\theta_t | V_{t-1} I) p(v_t | \theta_t I) d\theta_t$$

or cette identité fait apparaître une intégration couplée dans le membre de droite. Pour éliminer ce couplage, il faut intégrer de plus sur v_t :

$$\int p(v_t | V_{t-1} I) dv_t \int (\cdot) p(\theta_t | V_t I) d\theta_t = \int p(\theta_t | V_{t-1} I) d\theta_t \int (\cdot) p(v_t | \theta_t I) dv_t$$

Hypothèse de linéarité de l'estimation : Nous avons vu lors de la résolution du cas linéaire gaussien que l'estimation de l'état dépendait linéairement des données :

$$\hat{\theta}_t = \hat{\theta}_t^- + K_t \cdot (v_t - \hat{v}_t^-)$$

$\hat{\theta}_t^-$ et \hat{v}_t^- ne dépendaient alors que de V_{t-1} (données jusqu'à l'instant $t-1$), et K_t ne dépendait pas des données :

$$\hat{\theta}_t = \underbrace{(I_d - K_t \cdot X_t) \cdot \hat{\theta}_t^-}_{\text{dépend de } V_{t-1}} + \underbrace{K_t \cdot v_t}_{\text{dépend de } v_t}$$

avec $\hat{\theta}_t^- = \Phi_{t-1} \cdot \hat{\theta}_{t-1}$, donc

$$\hat{\theta}_t = (I_d - K_t \cdot X_t) \cdot \Phi_{t-1} \cdot \hat{\theta}_{t-1} + K_t \cdot v_t$$

Par récurrence, on voit que l'estimée de l'état dépend linéairement des données dans l'hypothèse linéaire gaussienne.

Il est logique (mais arbitraire) de supposer cette propriété de linéarité de l'estimée avec les données vraie également dans le cas des moindres carrés. C'est ce que nous supposerons à partir de maintenant afin de pouvoir poursuivre le calcul.

Orthogonalité des innovations : Soit l'innovation à l'instant t ,

$$\Delta v_t \doteq v_t - v_t^- ,$$

avec

$$v_t^- = \int v_t p(v_t | V_{t-1} I) dv_t$$

On a donc :

$$\int (v_t - v_t^-) p(v_t | V_{t-1} I) dv_t = 0$$

Soit $t' < t$, on a encore (puisque $v_{t'}$ est contenu dans V_{t-1}) :

$$\int (v_{t'} - v_{t'}^-) \cdot (v_t - v_t^-)^T p(v_t | V_{t-1} I) dv_t = 0$$

et donc

$$\langle \Delta v_{t'}^T \cdot \Delta v_t \rangle_{V_{t-1} I} = 0$$

En utilisant la notion de produit scalaire de vecteurs aléatoires centrés, on voit donc que $\Delta v_{t'} \perp \Delta v_t$ au sens de la probabilité conditionnelle par rapport à $V_{t-1} I$.

Soient $\text{ev}[V_{t-1}] \doteq$ "sous-espace vectoriel engendré par la famille $\Delta v_1 \dots \Delta v_{t-1}$ ", et $\text{ev}[v_t] \doteq$ "sous-espace vectoriel engendré par Δv_t ". Ces deux sous-espaces vectoriels sont orthogonaux, de dimension respective $t-1$ et 1, et forment une somme directe :

$$\text{ev}[V_t] = \text{ev}[V_{t-1}] \oplus \text{ev}[v_t]$$

Cette propriété conduit à supposer une fois de plus la forme :

$$\hat{\theta}_t = \hat{\theta}_t^- + K_t \cdot \Delta v_t$$

Calcul du gain :

$$\int (\theta_t - \hat{\theta}_t) \cdot v_t^T p(\theta_t | V_t I) d\theta_t = 0$$

car V_t contient v_t . Puisque

$$p(\theta_t | V_t I) = \frac{p(\theta_t | V_{t-1} I)}{p(v_t | V_{t-1} I)} p(v_t | \theta_t I)$$

il vient donc

$$\int (\theta_t - \hat{\theta}_t) \cdot v_t^T p(\theta_t | V_{t-1} I) p(v_t | \theta_t I) d\theta_t = 0$$

En intégrant de plus sur v_t :

$$\int p(\theta_t | V_{t-1} I) d\theta_t \int (\theta_t - \hat{\theta}_t) \cdot v_t^T p(v_t | \theta_t I) dv_t = 0$$

Or

$$(\theta_t - \hat{\theta}_t) \cdot v_t^T = (\theta_t - \hat{\theta}_t^-) \cdot v_t^T - K_t \cdot v_t \cdot v_t^T + K_t \cdot X_t \cdot \hat{\theta}_t^- \cdot v_t^T$$

L'intégration sur v_t donne :

$$(\theta_t - \hat{\theta}_t^-) \cdot (X_t \cdot \theta_t)^T - K_t \cdot (R_t + X_t \cdot \theta_t \cdot \theta_t^T \cdot X_t^T) + K_t \cdot X_t \cdot \hat{\theta}_t^- \cdot \theta_t^T \cdot X_t^T$$

ou

$$\begin{aligned} & (\theta_t - \hat{\theta}_t^-) \cdot (\theta_t - \hat{\theta}_t^-)^T \cdot X_t^T + (\theta_t - \hat{\theta}_t^-) \cdot (X_t \cdot \hat{\theta}_t^-)^T \\ & - K_t \cdot R_t - K_t \cdot X_t \cdot (\theta_t - \hat{\theta}_t^-) \cdot (\theta_t - \hat{\theta}_t^-)^T \cdot X_t^T - K_t \cdot X_t \cdot (\theta_t - \hat{\theta}_t^-) \cdot (X_t \cdot \hat{\theta}_t^-)^T \end{aligned}$$

L'intégration sur θ_t donne :

$$P_t^- \cdot X_t^T - K_t \cdot R_t - K_t \cdot X_t \cdot P_t^- \cdot X_t^T = P_t^- \cdot X_t^T - K_t \cdot S_t = 0$$

Et on retrouve bien le résultat du calcul fait dans le cas gaussien :

$$K_t = P_t^- . X_t^T . S_t^{-1}$$

Calcul de la covariance de l'estimation :

$$P_t = \int (\theta_t - \hat{\theta}_t) . (\theta_t - \hat{\theta}_t)^T p(\theta_t | V_t I) d\theta_t$$

avec

$$p(\theta_t | V_t I) = \frac{p(\theta_t | V_{t-1} I)}{p(v_t | V_{t-1} I)} p(v_t | \theta_t I)$$

d'où l'on déduit

$$\int P_t p(v_t | V_{t-1} I) dv_t = \int p(\theta_t | V_{t-1} I) d\theta_t \int (\theta_t - \hat{\theta}_t) . (\theta_t - \hat{\theta}_t)^T p(v_t | V_{t-1} I) dv_t$$

car P_t peut *a priori* dépendre de v_t . Il faut donc supposer que P_t ne dépend pas de v_t (ce qui s'introduisait naturellement dans le cas gaussien) pour pouvoir poursuivre le calcul. Écrivons :

$$(\theta_t - \hat{\theta}_t) . (\theta_t - \hat{\theta}_t)^T = ((I_d - K_t . X_t) . \Delta\theta_t - K_t (v_t - X_t . \theta_t)) . ((I_d - K_t . X_t) . \Delta\theta_t - K_t (v_t - X_t . \theta_t))^T$$

avec

$$\Delta\theta_t = \theta_t - \theta_t^-$$

L'intégration sur v_t donne :

$$P_t = \int p(\theta_t | V_{t-1} I) d\theta_t \left[(I_d - K_t . X_t) . \Delta\theta_t . \Delta\theta_t^T . (I_d - K_t . X_t)^T + K_t . R_t . K_t^T \right]$$

L'intégration sur θ_t donne :

$$P_t = (I_d - K_t . X_t) . P_t^- . (I_d - K_t . X_t)^T + K_t . R_t . K_t^T$$

Or

$$K_t = P_t^- . X_t^T . (R_t + X_t . P_t^- . X_t^T)^{-1}$$

donc

$$K_t . R_t = (I_d - K_t . X_t) . P_t^- . X_t^T$$

et

$$K_t . R_t . K_t^T = (I_d - K_t . X_t) . P_t^- . (K_t . X_t)^T$$

donc au final on retrouve la forme obtenue dans le cas linéaire gaussien :

$$P_t = (I_d - K_t . X_t) . P_t^-$$

On peut noter que l'expression de P_t trouvée plus haut :

$$P_t = (I_d - K_t . X_t) . P_t^- . (I_d - K_t . X_t)^T + K_t . R_t . K_t^T$$

bien qu'elle paraisse plus compliquée permet en pratique un calcul numérique plus précis, car la forme même de l'expression définit une matrice symétrique. Si la matrice P_t doit être évaluée pour de nombreuses valeurs de t , les erreurs qui s'accumulent au cours du calcul récursif risquent de poser des problèmes de précision. Les erreurs commises sur le calcul de K_t , qui suppose une

inversion de la matrice S_t , n'interviendront qu'au second ordre dans l'expression symétrique alors qu'elles interviennent au premier ordre dans l'expression compacte. On peut noter que les expressions servant à calculer S_t et P_t^- présentent également une forme symétrique.

Justification de la forme de $\hat{\theta}_t$: Nous avons supposé la forme $\hat{\theta}_t = \hat{\theta}_t^- + K_t \cdot \Delta v_t$, où $\hat{\theta}_t^- \in \text{ev}[V_{t-1}]$ et $K_t \cdot \Delta v_t \in \text{ev}[\Delta v_t]$. En fait, nous n'avons pas montré qu'il fallait prendre $\hat{\theta}_t^-$ comme projection de θ_t sur $\text{ev}[V_{t-1}]$. On peut justifier ce choix en précisant le calcul de P_t précédent. Posons $\hat{\theta}_t = \alpha + K_t \cdot \Delta v_t$, avec $\alpha \in \text{ev}[V_{t-1}]$. Alors

$$\theta_t - \hat{\theta}_t = (I_d - K_t \cdot X_t) \cdot \Delta \theta_t - K_t (v_t - X_t \cdot \theta_t) - (\alpha - \hat{\theta}_t^-)$$

En reprenant le calcul de P_t on voit que

$$P_t = (I_d - K_t \cdot X_t) \cdot P_t^- + \int p(\theta_t | V_{t-1} I) d\theta_t \int p(v_t | V_{t-1} I) dv_t (\alpha - \hat{\theta}_t^-) \cdot (\alpha - \hat{\theta}_t^-)^T$$

ou puisque ni α ni $\hat{\theta}_t^-$ ne dépendent de θ_t ou dv_t

$$P_t = (I_d - K_t \cdot X_t) \cdot P_t^- + (\alpha - \hat{\theta}_t^-) \cdot (\alpha - \hat{\theta}_t^-)^T$$

La valeur minimale du critère des moindres carrés ($\doteq \text{Tr} P_t$) est donc atteinte pour $\alpha = \hat{\theta}_t^-$.

5.5 Solution approchée du cas général

5.5.1 Modèle linéaire généralisé

Nous modifions légèrement le modèle linéaire de la page 5.1.5 de la façon suivante :

modèle de mesure	
$x(t, \theta_t)$	$= X_t \cdot \theta_t + x_t$
$r(t, \theta_t)$	$= R_t$
modèle d'évolution	
$\varphi(t, \theta_t)$	$= \Phi_t \cdot \theta_t + \varphi_t$
$q(t, \theta_t)$	$= Q_t$

où φ_t et x_t sont des vecteurs de dimensions respectives n et m ne dépendant que du temps t (mais plus de θ_t). Comment est modifiée la solution linéaire gaussienne (ou linéaire des moindres carrés) ?

Nous allons modifier légèrement une propriété vue précédemment :

Propriété 5.3**(Propagation du caractère gaussien (2))**

Soient α et β deux vecteurs aléatoires. α est un vecteur aléatoire gaussien à n composantes, de moyenne q (vecteur de dimension n) et de covariance C (matrice $n \times n$), soit $p(\alpha) = \mathcal{N}(\alpha, q, C)$. Le vecteur aléatoire β est de dimension m , et tel que $p(\beta|\alpha) = \mathcal{N}(\beta, A.\alpha + a, B)$, avec A une matrice $m \times n$, B une matrice $m \times m$, et a un vecteur de dimension m .

Alors β est un vecteur aléatoire gaussien :

$$p(\beta) = \int \mathcal{N}(\beta, A.\alpha + a, B) \mathcal{N}(\alpha, q, C) d\alpha = \mathcal{N}(\beta, A.q + a, A.C.A^T + B)$$

On obtient alors de façon presque immédiate le filtrage de Kalman dans le cas linéaire gaussien généralisé résumé par la table 5.2.

5.5.2 Approximation au premier ordre du modèle

Nous nous plaçons toujours dans le cadre des modèles du second ordre pour la mesure et l'évolution du système. Ceux-ci sont donc spécifiés par $x(t, \theta_t)$, $r(t, \theta_t)$, $\varphi(t, \theta_t)$ et $q(t, \theta_t)$. On peut obtenir une solution approchée du filtrage de Kalman dans ce cas en linéarisant moyennes et covariances autour des prédictions et estimations obtenues à chaque étape, et se ramener ainsi au cas linéaire généralisé. Reprenons les différentes étapes.

Prédiction de l'état du système au temps $t+1$: De l'étape précédente du filtrage (estimation au temps t), on connaît la valeur de $\hat{\theta}_t$. On pose en développant au premier ordre autour de $\hat{\theta}_t$:

$$\varphi(t, \theta_t) = \varphi(t, \hat{\theta}_t) + \nabla \varphi(t, \hat{\theta}_t) \cdot (\theta_t - \hat{\theta}_t) \quad (+ \text{ termes du second ordre})$$

où $\nabla \varphi(t, \hat{\theta}_t)$ est la matrice des dérivées partielles définie par :

$$\nabla \varphi(t, \hat{\theta}_t) \doteq \begin{pmatrix} \frac{\partial \varphi_1}{\partial (\theta_t)_1}(t, \hat{\theta}_t) & \cdots & \frac{\partial \varphi_1}{\partial (\theta_t)_n}(t, \hat{\theta}_t) \\ \vdots & & \vdots \\ \frac{\partial \varphi_n}{\partial (\theta_t)_1}(t, \hat{\theta}_t) & \cdots & \frac{\partial \varphi_n}{\partial (\theta_t)_n}(t, \hat{\theta}_t) \end{pmatrix}$$

on a donc

$$\varphi(t, \theta_t) \approx \Phi_t \cdot \theta_t + \varphi_t$$

avec

$$\begin{cases} \Phi_t &= \nabla \varphi(t, \hat{\theta}_t) \\ \varphi_t &= \varphi(t, \hat{\theta}_t) - \nabla \varphi(t, \hat{\theta}_t) \cdot \hat{\theta}_t \end{cases}$$

On prend de plus en développant à l'ordre 0 la covariance :

$$Q_t \approx q(t, \hat{\theta}_t)$$

Modèle linéaire de mesure	
$p(v_t \theta_t I) = \mathcal{N}(v_t, X_t \cdot \theta_t + x_t, R_t)$	
x_t	vecteur m
X_t	matrice $m \times n$
R_t	matrice $m \times m$
Modèle linéaire d'évolution	
$p(\theta_{t+1} \theta_t I) = \mathcal{N}(\theta_{t+1}, \Phi_t \cdot \theta_t + \varphi_t, Q_t)$	
φ_t	vecteur n
Φ_t	matrice $n \times n$
Q_t	matrice $n \times n$
Conditions initiales	
$p(\theta_0 I) = \mathcal{N}(\theta_0, \hat{\theta}_0, P_0)$	
P_0	matrice $n \times n$
Prédiction de l'état	
$p(\theta_{t+1} V_t I) = \mathcal{N}(\theta_{t+1}, \hat{\theta}_{t+1}^-, P_{t+1}^-)$	
$\hat{\theta}_{t+1}^- = \Phi_t \cdot \hat{\theta}_t + \varphi_t$	vecteur n
$P_{t+1}^- = Q_t + \Phi_t \cdot P_t \cdot \Phi_t^T$	matrice $n \times n$
Prédiction de la mesure	
$p(v_{t+1} V_t I) = \mathcal{N}(v_{t+1}, \hat{v}_{t+1}^-, S_{t+1})$	
$\hat{v}_{t+1}^- = X_{t+1} \cdot \hat{\theta}_{t+1}^- + x_t$	vecteur m
$S_{t+1} = R_{t+1} + X_{t+1} \cdot P_{t+1}^- \cdot X_{t+1}^T$	matrice $m \times m$
Estimation de l'état	
$p(\theta_{t+1} V_{t+1} I) = \mathcal{N}(\theta_{t+1}, \hat{\theta}_{t+1}, P_{t+1})$	
$\hat{\theta}_{t+1} = \hat{\theta}_{t+1}^- + K_{t+1} \cdot (v_{t+1} - \hat{v}_{t+1}^-) = (I_d - K_{t+1} \cdot X_{t+1}) \cdot \hat{\theta}_{t+1}^- + K_{t+1} \cdot v_{t+1}$	vecteur n
$K_{t+1} = P_{t+1}^- \cdot X_{t+1}^T \cdot S_{t+1}^{-1}$	matrice $n \times m$
$P_{t+1} = (I_d - K_{t+1} \cdot X_{t+1}) \cdot P_{t+1}^-$	matrice $n \times n$

TAB. 5.2 – Filtrage de Kalman dans le cas linéaire gaussien généralisé.

À partir de cette approximation, on obtient alors :

$$\begin{cases} \hat{\theta}_{t+1}^- &= \Phi_t \cdot \hat{\theta}_t + \varphi_t &= \varphi(t, \hat{\theta}_t) \\ P_{t+1}^- &= Q_t + \Phi_t \cdot P_t \cdot \Phi_t^T &= q(t, \hat{\theta}_t) + \nabla \varphi(t, \hat{\theta}_t) \cdot P_t \cdot \nabla \varphi(t, \hat{\theta}_t)^T \end{cases}$$

Prédiction de la mesure au temps $t + 1$: De l'étape précédente, on connaît la valeur de $\hat{\theta}_{t+1}^-$. On pose en développant au premier ordre autour de $\hat{\theta}_{t+1}^-$:

$$x(t + 1, \theta_{t+1}) = x(t + 1, \hat{\theta}_{t+1}^-) + \nabla x(t + 1, \hat{\theta}_{t+1}^-) \cdot (\theta_{t+1} - \hat{\theta}_{t+1}^-) \quad (+ \text{ termes du second ordre})$$

où $\nabla x(t + 1, \hat{\theta}_{t+1}^-)$ est la matrice des dérivées partielles définie par :

$$\nabla x(t + 1, \hat{\theta}_{t+1}^-) \doteq \begin{pmatrix} \frac{\partial x_1}{\partial (\theta_{t+1})_1}(t + 1, \hat{\theta}_{t+1}^-) & \cdots & \frac{\partial x_1}{\partial (\theta_{t+1})_n}(t + 1, \hat{\theta}_{t+1}^-) \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial (\theta_{t+1})_1}(t + 1, \hat{\theta}_{t+1}^-) & \cdots & \frac{\partial x_n}{\partial (\theta_{t+1})_n}(t + 1, \hat{\theta}_{t+1}^-) \end{pmatrix}$$

on a donc

$$x(t + 1, \theta_{t+1}) \approx X_{t+1} \cdot \theta_{t+1} + x_{t+1}$$

avec

$$\begin{cases} X_{t+1} &= \nabla x(t + 1, \hat{\theta}_{t+1}^-) \\ x_{t+1} &= x(t + 1, \hat{\theta}_{t+1}^-) - \nabla x(t + 1, \hat{\theta}_{t+1}^-) \cdot \hat{\theta}_{t+1}^- \end{cases}$$

On prend de plus en développant à l'ordre 0 la covariance :

$$R_{t+1} \approx r(t + 1, \hat{\theta}_{t+1}^-)$$

À partir de cette approximation, on obtient alors :

$$\begin{cases} \hat{v}_{t+1}^- &\approx X_{t+1} \cdot \hat{\theta}_{t+1}^- + x_{t+1} &= x(t + 1, \hat{\theta}_{t+1}^-) \\ S_{t+1} &\approx R_t + X_{t+1} \cdot P_{t+1}^- \cdot X_{t+1}^T &= r(t + 1, \hat{\theta}_{t+1}^-) + \nabla x(t + 1, \hat{\theta}_{t+1}^-) \cdot P_t \cdot \nabla x(t + 1, \hat{\theta}_{t+1}^-)^T \end{cases}$$

Estimation de l'état au temps $t + 1$: À partir des résultats des approximations précédentes, on peut écrire simplement :

$$\hat{\theta}_{t+1} = \hat{\theta}_{t+1}^- + K_{t+1} \cdot (v_{t+1} - \hat{v}_{t+1}^-) = \varphi(t, \hat{\theta}_t) + K_{t+1} \cdot (v_{t+1} - x(t + 1, \hat{\theta}_{t+1}^-))$$

avec

$$K_{t+1} = P_{t+1}^- \cdot X_{t+1}^T \cdot S_{t+1}^{-1}$$

et

$$P_{t+1} = (I_d - K_{t+1} \cdot X_{t+1}^T) \cdot P_{t+1}^-$$

où P_{t+1}^- , X_{t+1} et S_{t+1} ont été obtenues précédemment.

On voit donc que l'évaluation des matrices des dérivées partielles est nécessaire pour l'obtention des matrices de covariance et du gain du filtre de Kalman. Ces matrices dépendent maintenant des données acquises au cours de l'évolution (à travers $\hat{\theta}_t$ et $\hat{\theta}_{t+1}^-$), ce qui interdit leur calcul par avance.

5.6 Compléments

EXEMPLE - 5.3 Cas scalaire : étude de la stabilité

On suppose que l'état et la mesure sont des scalaires, et que les modèles de mesure et d'évolution sont linéaires et indépendants du temps, c'est-à-dire que $X_t = X \in \mathbb{R}$, $\Phi_t = \Phi \in \mathbb{R}$, $R_t = R \in \mathbb{R}$ et $Q_t = Q \in \mathbb{R}$. On pose

$$\begin{cases} a &= \Phi^2/Q \\ b &= X^2/R \end{cases}$$

de sorte que

$$\begin{cases} P_{t+1}^- &= Q + \Phi \cdot P_t \cdot \Phi = Q(aP_t + 1) \\ S_{t+1} &= R + X \cdot P_{t+1}^- \cdot X = R(bP_{t+1}^- + 1) \end{cases}$$

- Gain :

$$\begin{aligned} K_{t+1} &= \frac{P_{t+1}^- \cdot X}{S_{t+1}} \\ &= \frac{X}{R} \frac{(aP_t + 1)}{(bP_{t+1}^- + 1)} \\ &= \frac{XQ}{R} \frac{(aP_t + 1)}{(b(aP_t + 1) + 1)} \\ &= \frac{bQ}{X} \frac{(aP_t + 1)}{Qb(aP_t + 1) + 1} \\ K_{t+1} &= \frac{1}{X} \frac{(aP_t + 1)}{(aP_t + 1) + \frac{1}{Qb}} \end{aligned}$$

- Covariance de l'estimation :

$$\begin{aligned} P_{t+1} &= (1 - K_{t+1}X)P_{t+1}^- \\ &= Q \left(1 - \frac{(aP_t + 1)}{(aP_t + 1) + \frac{1}{Qb}} \right) (aP_t + 1) \\ P_{t+1} &= \frac{1}{b} \frac{(aP_t + 1)}{(aP_t + 1) + \frac{1}{Qb}} \end{aligned}$$

et

$$P_{t+1}^{-1} = b + \frac{1}{Q} \frac{1}{aP_t + 1}$$

et donc

$$\frac{P_t}{P_{t+1}} = bP_t + \frac{1}{Q} \frac{P_t}{aP_t + 1}$$

Si P_t tend vers une limite, notée P_∞ , celle-ci doit vérifier :

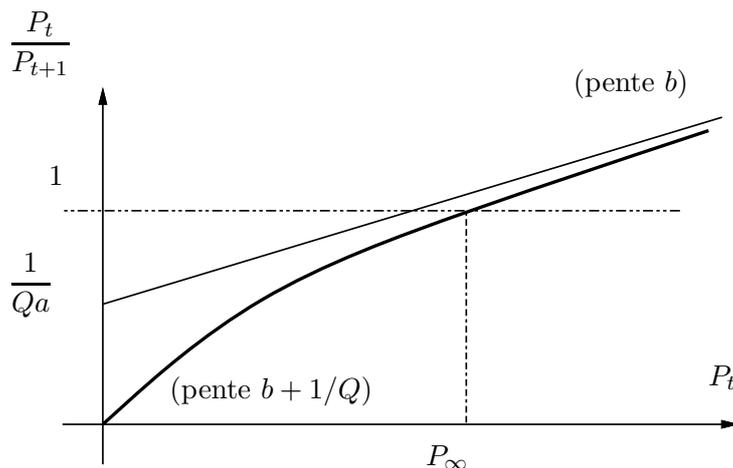
$$bP_\infty + \frac{1}{Q} \frac{P_\infty}{aP_\infty + 1} = 1$$

P_∞ est donc racine du polynôme du second degré

$$abQP_\infty^2 + (1 + bQ - aQ)P_\infty - Q = 0$$

dont le discriminant est positif ($\Delta = (1 + bQ - aQ)^2 + 4abQ^2$), et donc

$$P_\infty = \frac{\sqrt{\Delta} - (1 + bQ - aQ)}{2abQ}$$

FIG. 5.5 – *Attracteur.*

(l'autre racine est négative). $\left(\frac{P_{t+1}}{P_t}\right)$ est une fonction strictement décroissante de P_t , donc si $P_t > P_\infty$, alors $P_{t+1} < P_t$, et si $P_t < P_\infty$, alors $P_{t+1} > P_t$. P_∞ est attracteur.

EXERCICE - 5.1 Filtrage de Kalman étendu

On rappelle les équations de linéarisation du modèle de mesure pour la solution approchée du filtrage de Kalman :

$$\begin{cases} X_t &= \nabla x(t, \hat{\theta}_t^-) \\ \hat{v}_t^- &= x(t, \hat{\theta}_t^-) \end{cases}$$

ainsi que l'équation pour la moyenne de l'estimation de l'état du système :

$$\hat{\theta}_t = \varphi(t, \hat{\theta}_{t-1}) + K_t \cdot (v_t - x(t, \hat{\theta}_t^-))$$

On se propose de remplacer les étapes de prédiction de la mesure et d'estimation de l'état par l'algorithme itératif suivant :

$$y_{k+1} = \varphi(t, \hat{\theta}_{t-1}) + (K_t)_k \cdot (v_t - x(t, y_k))$$

$$\text{avec } \begin{cases} (X_t)_k &= \nabla x(t, y_k) \\ (S_t)_k &= R_t + (X_t)_k \cdot P_t^- \cdot (X_t)_k^T \\ (K_t)_k &= P_t^- \cdot (X_t)_k^T \cdot (S_t)_k^{-1} \end{cases}$$

condition initiale: $y_0 = \varphi(t, \hat{\theta}_{t-1}) = \hat{\theta}_t^-$
condition d'arrêt: $y_{k+1} = y_k$

Que permet d'obtenir cet algorithme (quel sens donner à la limite de la suite (y_k) quand k tend vers l'infini) ?

Quel est l'avantage de cet algorithme par rapport aux étapes classiques de prédiction de la mesure et de l'estimation de l'état ? Son inconvénient ?

Est-il intéressant d'employer le même type d'algorithme itératif pour remplacer également l'étape de prédiction de l'état ?

Réponse : La suite (y_k) tend vers une estimation de l'état $\hat{\theta}_t$ telle que

$$\hat{\theta}_t = \varphi(t, \hat{\theta}_{t-1}) + (K_t)_k \cdot (v_t - x(t, \hat{\theta}_t))$$

Il s'agit d'une équation implicite dans laquelle tout se passe comme si le modèle de mesure avait été développé autour de l'estimation $\hat{\theta}_t$ au lieu de la prédiction $\hat{\theta}_t^-$ (qui est la condition initiale). On améliore ainsi à la fois la prédiction de la mesure (qui devient $x(t, \hat{\theta}_t)$) et l'estimation de l'état. L'inconvénient est le temps de calcul qui peut être considérablement rallongé.

Il n'est pas intéressant de remplacer également l'étape de prédiction de l'état par un algorithme itératif de ce type, puisque le développement du modèle d'évolution est déjà fait autour de l'estimation à l'instant précédent $\hat{\theta}_{t-1}$, et non autour d'une valeur prédite.

EXERCICE - 5.2 Mesure de la polarisation

On mesure l'état de polarisation d'une onde plane. On sait que la polarisation est rectiligne, mais on ne connaît ni sa direction, ni la puissance lumineuse. On écrit

$$\theta = \begin{pmatrix} a \\ \psi \end{pmatrix}$$

où a est la puissance lumineuse et ψ l'angle que fait la direction de polarisation avec l'horizontale. Deux détecteurs sont placés derrière un cube séparateur de polarisation et donneraient en l'absence de bruit d'acquisition $v_1 = a \cos^2 \psi$ et $v_2 = a \sin^2 \psi$. v_1 et v_2 sont les composantes du vecteur mesure v , $\psi = 0$ en polarisation \odot et $\psi = \pi/2$ en polarisation \rightarrow .

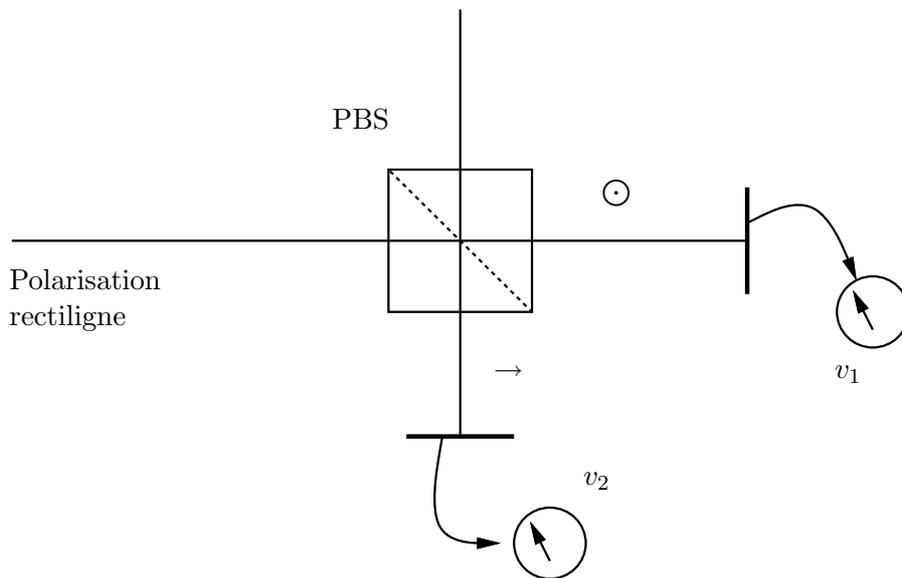


FIG. 5.6 – *Expérience de mesure de la polarisation.*

1. On sait que les bruits de détection sont centrés, et de variances respectives σ_1^2 et σ_2^2 . Justifier que les mesures v_1 et v_2 sont indépendantes. Que vaut $p(v|\theta I)$?
2. On suppose a connu. La polarisation peut être seulement \odot ou \rightarrow (hypothèses H_1 et H_2), sans qu'on sache rien *a priori* sur leur répartition. Quelle est la règle de choix optimale (au sens du MAP) entre H_1 et H_2 ?

3. On suppose que $\sigma_1^2 = \sigma_2^2 = \sigma^2$ et que a est connu. On ne sait rien *a priori* sur la répartition de ψ (paramètre continu).

Montrer que

$$p(\theta|vI) \propto \exp\left(-\frac{1}{2\sigma^2}f(\theta)\right)$$

avec

$$f(\theta) = v_1^2 + v_2^2 + a^2(\cos^4 \psi + \sin^4 \psi) - 2a(v_1 \cos^2 \psi + v_2 \sin^2 \psi)$$

puis que

$$\hat{\psi} = \frac{1}{2} \arccos\left(\frac{v_1 - v_2}{a}\right)$$

4. On reprend la question précédente, mais en supposant qu'on ne sait rien *a priori* sur a . Montrer que

$$\hat{a} = \frac{v_1 \cos^2 \hat{\psi} + v_2 \sin^2 \hat{\psi}}{\cos^4 \hat{\psi} + \sin^4 \hat{\psi}}$$

où $\hat{\psi}$ est solution implicite de l'équation

$$v_1 - v_2 = \frac{v_1 \cos^2 \psi + v_2 \sin^2 \psi}{\cos^4 \psi + \sin^4 \psi} \cos(2\psi)$$

Montrer alors que

$$\begin{cases} \tan^2 \hat{\psi} &= \frac{v_1}{v_2} \\ \hat{a} &= v_1 + v_2 \end{cases}$$

Commenter.

5. On suppose que la polarisation tourne à la vitesse angulaire constante ω , et que l'intensité a reste constante. On réalise une acquisition de v_1 et v_2 à intervalles de temps réguliers τ . Écrire les modèles de mesure et d'évolution (au second ordre statistique) pour le filtrage de Kalman. Donner la linéarisation au premier ordre du filtrage de Kalman étendu.

Réponse

1. Connaissant seulement les moyennes et les variances de v_1 et v_2 , on assigne par le principe du maximum d'entropie :

$$\begin{cases} p(v_1|\theta I) &= \mathcal{N}(v_1, a \cos^2 \psi, \sigma_1^2) \\ p(v_2|\theta I) &= \mathcal{N}(v_2, a \sin^2 \psi, \sigma_2^2) \end{cases}$$

Il est clair que les moyennes sont liées puisque $a \cos^2 \psi + a \sin^2 \psi = a$. Cependant puisque les détecteurs sont complètement déconnectés, et que rien ne laisse supposer que les bruit d'acquisition dépendent du signal, les mesures v_1 et v_2 sont logiquement indépendantes³, et donc

$$p(v|\theta I) = p(v_1|\theta I) p(v_2|\theta I)$$

2. $p(H_1|I) = p(H_2|I) = 1/2$ par application du principe d'indifférence, donc $p(H|vI) \propto p(v|HI)$. De plus d'après la question précédente

$$\begin{cases} p(v|H_1 I) &= \mathcal{N}(v_1, a, \sigma_1^2) \mathcal{N}(v_2, 0, \sigma_2^2) \\ p(v|H_2 I) &= \mathcal{N}(v_1, 0, \sigma_1^2) \mathcal{N}(v_2, a, \sigma_2^2) \end{cases}$$

3. Attention : cette indépendance est logique, pas physique. Elle s'entend du point de vue de l'observateur. Il ne s'agit pas d'une hypothèse sur la nature physique des capteurs, mais du mieux que l'on puisse en dire avec l'information dont on dispose.

La règle de choix au sens du MAP est de choisir H_1 si $p(H_1|vI) > p(H_2|vI)$, et H_2 sinon ; soit encore choisir H_1 si $p(v|H_1I) > p(v|H_2I)$, et H_2 sinon :

$$\begin{aligned} \mathcal{N}(v_1, a, \sigma_1^2) \mathcal{N}(v_2, 0, \sigma_2^2) &> \mathcal{N}(v_1, 0, \sigma_1^2) \mathcal{N}(v_2, a, \sigma_2^2) \\ \exp\left(-\frac{1}{2\sigma_1^2}(v_1 - a)^2\right) \exp\left(\frac{1}{2\sigma_1^2}v_1^2\right) &> \exp\left(-\frac{1}{2\sigma_2^2}(v_2 - a)^2\right) \exp\left(\frac{1}{2\sigma_2^2}v_2^2\right) \\ \frac{a^2 - 2av_1}{\sigma_1^2} &< \frac{a^2 - 2av_2}{\sigma_2^2} \\ \frac{v_1}{\sigma_1^2} - \frac{v_2}{\sigma_2^2} &> \frac{a}{2} \left(\frac{1}{\sigma_1^2} - \frac{1}{\sigma_2^2} \right) \end{aligned}$$

En particulier si $\sigma_1 = \sigma_2$, la règle de choix au sens du MAP se simplifie en $v_1 > v_2$.

3. De $p(v|\theta I) = p(v_1|\theta I) p(v_2|\theta I)$ avec $\sigma_1 = \sigma_2 = \sigma$ on déduit puisque $p(\theta|I)$ est constante :

$$\begin{aligned} p(\theta|vI) &\propto p(v|\theta I) \\ &\propto \mathcal{N}(v_1, a \cos^2 \psi, \sigma^2) \mathcal{N}(v_2, a \sin^2 \psi, \sigma^2) \\ &\propto \exp\left(-\frac{1}{2\sigma^2}((v_1 - a \cos^2 \psi)^2 + (v_2 - a \sin^2 \psi)^2)\right) \\ p(\theta|vI) &\propto \exp\left(-\frac{1}{2\sigma^2}f(\theta)\right) \end{aligned}$$

Puisque

$$\hat{\theta} = \underset{\theta}{\text{Argmax}} p(\theta|vI) = \underset{\theta}{\text{Argmin}} f(\theta)$$

$$\begin{aligned} \frac{\partial f(\theta)}{\partial \psi} &= a^2(-4 \sin \psi \cos^3 \psi + 4 \cos \psi \sin^3 \psi) - 2a(-2v_1 \sin \psi \cos \psi + 2v_2 \sin \psi \cos \psi) \\ \frac{\partial f(\theta)}{\partial \psi} &= -a \sin 2\psi(a \cos 2\psi - v_1 + v_2) \end{aligned}$$

et cette dérivée s'annule pour

$$\hat{\psi} = \frac{1}{2} \arccos\left(\frac{v_1 - v_2}{a}\right)$$

4. Le problème est le même que celui de la question précédente, mais l'optimisation doit porter sur a également :

$$\begin{cases} \frac{\partial f(\theta)}{\partial a} = 2a(\cos^4 \psi + \sin^4 \psi) - 2(v_1 \sin^2 \psi + v_2 \sin^2 \psi) = 0 \\ \frac{\partial f(\theta)}{\partial \psi} = -a \sin 2\psi(a \cos 2\psi - v_1 + v_2) = 0 \end{cases}$$

La première équation donne

$$\hat{a} = \frac{v_1 \cos^2 \hat{\psi} + v_2 \sin^2 \hat{\psi}}{\cos^4 \hat{\psi} + \sin^4 \hat{\psi}}$$

et en injectant cette valeur dans la seconde équation

$$v_1 - v_2 = \frac{v_1 \cos^2 \hat{\psi} + v_2 \sin^2 \hat{\psi}}{\cos^4 \hat{\psi} + \sin^4 \hat{\psi}} \cos(2\hat{\psi})$$

Cette dernière expression donne

$$\begin{aligned} (v_1 - v_2)(\cos^4 \hat{\psi} + \sin^4 \hat{\psi}) &= (v_1 \cos^2 \hat{\psi} + v_2 \sin^2 \hat{\psi})(\cos^2 \hat{\psi} - \sin^2 \hat{\psi}) \\ v_1(\sin^4 \hat{\psi} + \sin^2 \hat{\psi} \cos^2 \hat{\psi}) &= v_2(\cos^4 \hat{\psi} + \sin^2 \hat{\psi} \cos^2 \hat{\psi}) \\ v_1 \sin^2 \hat{\psi} &= v_2 \cos^2 \hat{\psi} \end{aligned}$$

donc

$$\tan^2 \hat{\psi} = \frac{v_1}{v_2}$$

De

$$\hat{a} = \frac{v_1 - v_2}{\cos(2\hat{\psi})}$$

et de

$$\cos(2\hat{\psi}) = \frac{1 - \tan^2 \hat{\psi}}{1 + \tan^2 \hat{\psi}} = \frac{v_1 - v_2}{v_1 + v_2}$$

on tire

$$\hat{a} = v_1 + v_2$$

5. Moyenne du modèle d'évolution :

$$\varphi(t, \theta_t) = \theta_t + \begin{pmatrix} 0 \\ \omega\tau \end{pmatrix} = \begin{pmatrix} a_t \\ \psi_t + \omega\tau \end{pmatrix}$$

donc

$$\Phi_t = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

et

$$\varphi_t = \begin{pmatrix} 0 \\ \omega\tau \end{pmatrix}$$

Moyenne du modèle de mesure :

$$x(t, \theta_t) = \begin{pmatrix} a_t \cos^2 \psi_t \\ a_t \sin^2 \psi_t \end{pmatrix}$$

donc

$$X_t = \nabla x(t, \hat{\theta}_t^-) = \begin{pmatrix} \cos^2 \hat{\psi}_t^- & -a_t \sin 2\hat{\psi}_t^- \\ \sin^2 \hat{\psi}_t^- & a_t \sin 2\hat{\psi}_t^- \end{pmatrix}$$

Toutes les autres quantités se déduisent des précédentes par une application immédiate des formules du filtrage de Kalman étendu.
