

Analyse et correction des Systèmes linéaires continus ou échantillonnés à l'aide des variables d'état

Gonzalo Cabodevila

gonzalo.cabodevila@femto-st.fr

2^{ème} année
Semestre vert
Automatique avancée
filière EAOI



**École Nationale Supérieure de
Mécanique et des Microtechniques**
26, chemin de l'Épitaphe
25030 Besançon cedex – FRANCE
<http://intranet-tice.ens2m.fr>

Table des matières

1 Exemple introductif : fil rouge	7
1.1 Différentes représentations d'un système physique	7
1.1.1 Equations différentielles	7
1.1.2 Fonction de transfert	8
1.1.3 Réponse impulsionnelle	8
1.1.4 Représentation d'état	8
1.2 Propriétés de la représentation d'état	10
1.2.1 Non unicité de la représentation d'état	10
1.2.2 Matrice de transfert	10
1.3 Intérêt de cette représentation	12
1.4 Résolution des équations d'état	12
1.4.1 Cas simple	12
1.4.2 Cas général	14
1.4.3 Généralisation aux systèmes variants dans le temps	14
1.4.4 Simulation sur calculateur	15
2 Obtention des équations d'état	17
2.1 Méthode directe	17
2.2 A partir de la fonction de transfert	17
2.2.1 Forme 1 : forme canonique de commandabilité	17
2.2.2 Forme 2 : forme canonique d'observabilité	18
2.2.3 Représentation modale	19
2.2.4 Forme canonique de Jordan (forme diagonale)	19
3 Commandabilité et observabilité des systèmes	25
3.1 Définitions	25
3.2 Que faire si un système n'est pas observable et/ou commandable	27
3.2.1 Retour sur conception	27
3.2.2 Réduction de modèles	27
4 Transformation en l'une des formes canoniques	29
4.1 Diagonalisation de la matrice \mathbf{A}	29
4.2 Conséquences pour la commandabilité et l'observabilité	30
4.3 Cas des valeurs propres complexes	31
4.3.1 Diagonalisation classique	31
4.3.2 Transformation modifiée : \mathbf{T}_m	31
4.4 Transformation en la forme canonique d'asservissement	32
4.5 Transformation en la forme canonique d'observabilité	32

5	Stabilité des systèmes dynamiques linéaires	35
5.1	Définition	35
5.2	Etude de la stabilité	35
5.3	Stabilité au sens de Lyapounov	37
5.3.1	Théorème	37
5.3.2	Interprétation physique	37
5.3.3	Applications aux systèmes linéaires	38
5.3.4	Fil rouge	38
6	Commande des systèmes	39
6.1	Placement de pôles	39
6.1.1	Calcul du régulateur \mathbf{L}	39
6.1.2	Calcul de la matrice de préfiltre \mathbf{S}	40
6.2	Cas d'une représentation quelconque du système à asservir	41
6.2.1	Transformation en la forme canonique de commandabilité	41
6.2.2	Théorème de Cayley-Hamilton	41
6.3	Commande Modale	42
6.3.1	Définition	42
6.3.2	Méthode de synthèse	42
6.4	Choix des pôles	44
6.4.1	Pôles complexes conjugués dominants	45
6.4.2	Maximalement plat	46
6.4.3	Pôles à partie réelle identique	46
6.4.4	Polynômes de Naslin	46
6.5	Commande optimale	48
6.5.1	Définition	48
6.5.2	Stabilité de la commande optimale	48
6.5.3	Choix des matrice \mathbf{R} et \mathbf{Q}	48
6.5.4	Exemple : fil rouge	49
7	Synthèse d'observateurs d'état	51
7.1	Introduction au problème de la reconstruction d'état	51
7.1.1	Par calcul direct	51
7.1.2	Par simulation du processus	51
7.1.3	Par simulation du processus et asservissement sur les parties connues du vecteur d'état	52
7.2	Observateurs de Luenberger	52
7.3	Observateurs d'ordre réduit	54
7.3.1	Hypothèses	54
7.4	Observateur généralisé	55
7.5	Equation d'état d'un système asservi avec observateur	57
7.5.1	Théorème de séparation	57
7.6	Filtrage de Kalman	58
8	Représentation d'état des systèmes linéaires échantillonnés	59
8.1	Système discret	59
8.2	Résolution des équations dans le domaine du temps	60
8.3	Application de la transformée en z	60
8.4	Matrice de transfert	61
8.5	Obtention d'un modèle d'état à partir de la fonction de transfert en z	61
8.6	Résolution de l'équation d'état dans le domaine de z	61
8.7	Commandabilité et observabilité	61

<i>TABLE DES MATIÈRES</i>	5
8.7.1 Commandabilité d'un système échantillonné	61
8.7.2 Observabilité d'un système échantillonné	62
8.8 Stabilité des systèmes échantillonnés	63
8.9 Commandes des systèmes échantillonnés	63
8.9.1 Calcul de la matrice de préfiltre	64
8.9.2 Commande optimale dans le cas discret	64
9 Annales d'examens	65
Devoir personnel Juin 1999	66
Examen final Juin 1999	69
Examen final Juin 2010	71
10 Travaux dirigés	75
I Annexes	89
A Quelques publications originales	91

Chapitre 1

Exemple introductif : fil rouge

1.1 Différentes représentations d'un système physique

Soit un moteur à courant continu commandé par l'inducteur.

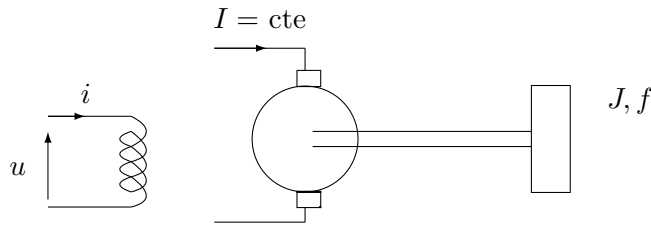


FIGURE 1.1 – Moteur à courant continu commandé par l'inducteur

- commande : u
- sortie : ω

Le système est monovariante, linéaire invariant dans le temps, il peut donc être représenté par une équation différentielle à coefficients constants.

1.1.1 Equations différentielles

Le système représenté en figure 1.1 est décrit par les équations suivantes :

$$u = Ri + L \frac{di}{dt} \quad (1.1)$$

$$J \frac{d\omega}{dt} + f\omega = \gamma \quad (1.2)$$

$$\gamma = ki \quad (1.3)$$

$$J \frac{d^2\omega}{dt^2} + f \frac{d\omega}{dt} = k \frac{di}{dt} = \frac{k}{L} (u - Ri) = \frac{k}{L} \left(u - \frac{R}{k} \left(\frac{d\omega}{dt} + f\omega \right) \right) \quad (1.4)$$

$$J \frac{d^2\omega}{dt^2} + \left(f + \frac{RJ}{L} \right) \frac{d\omega}{dt} + \frac{Rf}{L} \omega = \frac{k}{L} u \quad (1.5)$$

$$\frac{d^2\omega}{dt^2} + \left(\frac{f}{J} + \frac{R}{L} \right) \frac{d\omega}{dt} + \frac{Rf}{LJ} \omega = \frac{k}{LJ} u \quad (1.6)$$

$$\frac{d^2\omega}{dt^2} + a_1 \frac{d\omega}{dt} + a_0 \omega = b_0 u \quad (1.7)$$

Dès lors $\omega(t)$ est connu si $u(t)$ et les deux conditions initiales ($\omega(0)$ et $\frac{d\omega(0)}{dt}$) sont connues.

1.1.2 Fonction de transfert

En utilisant la transformation de Laplace :

$$\mathcal{L}[e(t)] = \int_0^{\infty} e^{-pt} e(t) dt, \quad (1.8)$$

$$\mathcal{L}\left[\frac{de(t)}{dt}\right] = pE(p) - e(0), \quad (1.9)$$

nous obtenons la transformée de (1.7) :

$$p^2\Omega(p) - p\Omega(0) - \dot{\Omega}(0) + a_1(p\Omega(p) - \Omega(0) + a_0\Omega(p)) = b_0U(p) \quad (1.10)$$

d'où :

$$\Omega(p) = \frac{b_0}{p^2 + a_1p + a_0}U(p) + \frac{\Omega(0)(a_1 + p) + \dot{\Omega}(0)}{p^2 + a_1p + a_0} \quad (1.11)$$

Si les conditions initiales sont nulles :

$$\Omega(p) = \frac{b_0}{p^2 + a_1p + a_0}U(p) \quad (1.12)$$

$$\frac{\Omega(p)}{U(p)} = \frac{b_0}{p^2 + a_1p + a_0} = H(p) \quad (1.13)$$

$H(p)$ est la fonction de transfert du système.

1.1.3 Réponse impulsionnelle

Si les conditions initiales sont nulles :

$$\Omega(p) = H(p) \times U(p) \quad (1.14)$$

En repassant en temporel, la multiplication est transformée en une convolution

$$\omega(t) = h(t) \star u(t) = \int_0^{\infty} h(\tau)u(t - \tau)d\tau \quad (1.15)$$

donc

$$h(t) = \frac{k}{RJ - Lf} \left(e^{-\frac{f}{J}t} - e^{-\frac{R}{L}t} \right) \quad (1.16)$$

1.1.4 Représentation d'état

Si l'on désire réaliser une simulation analogique du système à partir d'intégrateurs, l'équation (1.7) peut se mettre sous la forme :

$$\frac{d^2\omega}{dt^2} = b_0u - a_1\frac{d\omega}{dt} - a_0\omega \quad (1.17)$$

d'où le schéma suivant,

Les variables d'état sont les sorties des intégrateurs.

$$x_1 = \omega, \quad x_2 = \frac{d\omega}{dt} = \dot{\omega} \quad (1.18)$$

Le système peut être représenté par les deux équations suivantes,

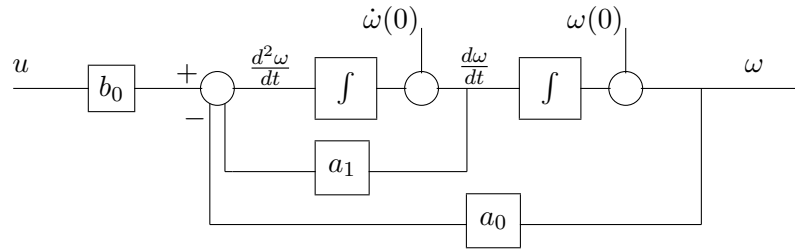


FIGURE 1.2 – Schéma d'un simulateur analogique

$$\dot{x}_1 = x_2 \quad (1.19)$$

$$\dot{x}_2 = b_0 u - a_0 x_1 - a_1 x_2 \quad (1.20)$$

La représentation utilisée classiquement est la représentation matricielle, on définit alors un vecteur d'état,

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (1.21)$$

Equation d'état :

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -a_0 & -a_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ b_0 \end{bmatrix} u \quad (1.22)$$

Equation de sortie :

$$\omega = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (1.23)$$

équation d'état + équation de sortie = représentation d'état

Dans le cas général, l'écriture de la représentation d'état est la suivante,

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}u \quad (1.24)$$

$$\underline{y} = \mathbf{C}\underline{x} + \mathbf{D}u \quad (1.25)$$

Dans les cas qui nous intéressent c'est-à-dire les systèmes à une seule entrée et une seule sortie, l'écriture générale de la représentation d'état est la suivante,

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}u \quad (1.26)$$

$$y = \mathbf{C}\underline{x} + \mathbf{D}u \quad (1.27)$$

Dimensions : si le vecteur d'état est de dimension n , alors

- \mathbf{A} est de dimension n lignes $\times n$ colonnes,
- \mathbf{B} est de dimension n lignes $\times 1$ colonne,
- \mathbf{C} est de dimension 1 ligne $\times n$ colonnes,
- \mathbf{D} est une constante (très souvent nulle).

1.2 Propriétés de la représentation d'état

1.2.1 Non unicité de la représentation d'état

La représentation d'état n'est pas unique, dans l'exemple traité jusqu'à présent nous avons choisi le vecteur d'état suivant

$$x = \begin{bmatrix} \omega \\ \dot{\omega} \end{bmatrix} \quad (1.28)$$

nous aurions pu choisir un autre vecteur d'état, par exemple

$$x = \begin{bmatrix} i \\ \omega \end{bmatrix} \quad (1.29)$$

En effet, en reprenant les équations fondamentales du système

$$u = Ri + L \frac{di}{dt} \quad (1.30)$$

$$J \frac{d\omega}{dt} + f\omega = ki \quad (1.31)$$

avec ce nouveau vecteur d'état elle peuvent être écrites sous la forme

$$U = Rx_1 + L\dot{x}_1 \quad (1.32)$$

$$J\dot{x}_2 + fx_2 = kx_1 \quad (1.33)$$

d'où la représentation d'état

Equation d'état :

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} -\frac{R}{L} & 0 \\ \frac{k}{J} & -\frac{f}{J} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} \frac{1}{L} \\ 0 \end{bmatrix} u \quad (1.34)$$

Equation de sortie :

$$\omega = \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (1.35)$$

Notez que nous n'aurions pas pu prendre i et $\frac{di}{dt}$ comme vecteur d'état car $\frac{di}{dt}$ n'est pas une sortie d'intégrateur.

En fait, on peut prendre comme vecteur d'état n'importe quelle combinaison linéaire d'un vecteur d'état valable.

Soit $x' = \mathbf{M}x$

$$\dot{x}' = \mathbf{M}^{-1}\mathbf{A}\mathbf{M}x' + \mathbf{M}^{-1}\mathbf{B}u \quad (1.36)$$

$$y = \mathbf{C}\mathbf{M}x' + \mathbf{D}u \quad (1.37)$$

1.2.2 Matrice de transfert

En prenant la transformée de Laplace de la représentation d'état, on obtient,

$$p\underline{X}(p) - \underline{x}(0) = \mathbf{A}\underline{X}(p) + \mathbf{B}U(p) \quad (1.38)$$

$$Y(p) = \mathbf{C}\underline{X}(p) + \mathbf{D}U(p) \quad (1.39)$$

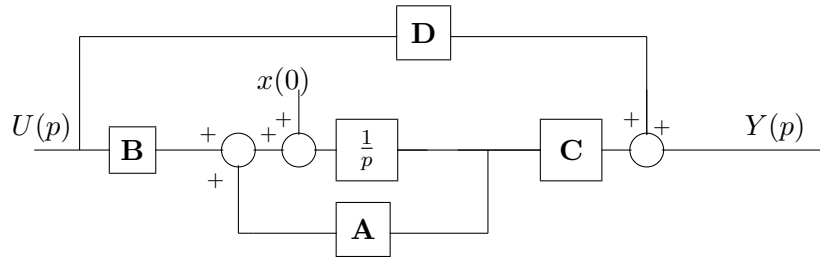


FIGURE 1.3 – Schéma général après transformation de Laplace

les équations (1.38) et (1.39) se réécrivent sous la forme

$$[p\mathbf{I} - \mathbf{A}] \underline{X}(p) = \underline{x}(0) + \mathbf{B}U(p) \quad (1.40)$$

$$Y(p) = \mathbf{C}\underline{X}(p) + \mathbf{D}U(p) \quad (1.41)$$

$$\underline{X}(p) = [p\mathbf{I} - \mathbf{A}]^{-1} \underline{x}(0) + [p\mathbf{I} - \mathbf{A}]^{-1} \mathbf{B}U(p) \quad (1.42)$$

$$Y(p) = \underbrace{\mathbf{C}[p\mathbf{I} - \mathbf{A}]^{-1} \underline{x}(0)}_{\mathbf{C} \mathbf{I}} + \underbrace{[\mathbf{C}[p\mathbf{I} - \mathbf{A}]^{-1} \mathbf{B} + \mathbf{D}]}_{\text{matrice de transfert}} U(p) \quad (1.43)$$

en posant

$$\mathbf{T} = [\mathbf{C}[p\mathbf{I} - \mathbf{A}]^{-1} \mathbf{B} + \mathbf{D}] \quad (1.44)$$

on obtient

$$Y(p) = \mathbf{T}U(p) + \mathbf{C}[p\mathbf{I} - \mathbf{A}]^{-1} \underline{x}(0) \quad (1.45)$$

Note : dans le cas multivariable (plusieurs entrées, plusieurs sorties)

$$\mathbf{T}_{(i,j)}(p) = \frac{Y_i(p)}{U_j(p)} \quad (1.46)$$

Dans le cas SISO la matrice \mathbf{T} représente la fonction de transfert du système.

La matrice \mathbf{T} est unique, elle ne dépend pas de la représentation d'état choisie, elle ne dépend que des entrées et des sorties.

Démonstration :

soit

$$p\underline{X}(p) = \mathbf{A}\underline{X}(p) + \mathbf{B}U(p) \quad (1.47)$$

$$Y(p) = \mathbf{C}\underline{X}(p) + \mathbf{D}U(p) \quad (1.48)$$

posons

$$\underline{X} = \mathbf{K}\underline{X}' \quad (1.49)$$

donc

$$p\underline{X}'(p) = \mathbf{K}^{-1} \mathbf{A} \mathbf{K} \underline{X}'(p) + \mathbf{K}^{-1} \mathbf{B}U(p) \quad (1.50)$$

$$Y(p) = \mathbf{C} \mathbf{K} \underline{X}'(p) + \mathbf{D}U(p) \quad (1.51)$$

$$\begin{aligned}
\mathbf{T} &= \mathbf{CK}[p\mathbf{I} - \mathbf{K}^{-1}\mathbf{AK}]^{-1}\mathbf{K}^{-1}\mathbf{B} + \mathbf{D} && \text{en utilisant } [\mathbf{AB}]^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1} \\
&= \mathbf{CK}[\mathbf{K}p\mathbf{I} - \mathbf{AK}]^{-1}\mathbf{B} + \mathbf{D} && \text{en utilisant } [\mathbf{AB}]^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1} \\
&= \mathbf{C}[\mathbf{K}p\mathbf{I}\mathbf{K}^{-1} - \mathbf{A}]^{-1}\mathbf{B} + \mathbf{D} && p\mathbf{I} \text{ est une matrice diagonale} \\
&= \mathbf{C}[p\mathbf{I} - \mathbf{A}]^{-1}\mathbf{B} + \mathbf{D} && \text{CQFD}
\end{aligned} \tag{1.52}$$

\mathbf{T} est bien indépendante de la représentation d'état choisie.

1.3 Intérêt de cette représentation

1. Systèmes linéaires invariants
 - → simulation du système, analogique ou numérique
 - → extension aux systèmes multivariables
 - → notions nouvelles de commandabilité et d'observabilité
2. Systèmes linéaires variants (\mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{D} dépendent de t)
 - → simulation
3. Systèmes non linéaires
 - → simulation, étude de la dynamique
4. Qu'est-ce que l'état ?
 - l'état d'un système à un instant donné représente la mémoire minimale du passé nécessaire à la détermination du futur.
 - les variables d'état doivent apporter une description interne du système et on choisit celles pour lesquelles on peut définir l'état initial, c'est-à-dire celles qui décrivent des "réservoirs" d'énergie

1.4 Résolution des équations d'état

soit le système :

$$\begin{aligned}
\dot{\underline{x}} &= \mathbf{A}\underline{x} + \mathbf{B}u \\
y &= \mathbf{C}\underline{x} + \mathbf{D}u
\end{aligned} \tag{1.53}$$

1.4.1 Cas simple

Dans le cas simple où x ne contient qu'une seule composante, ce système devient

$$\frac{dx}{dt} = ax + bu \tag{1.54}$$

En régime libre, la solution est

$$x(t) = ke^{at} \tag{1.55}$$

En régime forcé, la méthode de la variation de la constante donne

$$x(t) = k(t)e^{at} \tag{1.56}$$

$$ak(t)e^{at} + k'(t)e^{at} = ak(t)e^{at} + bu(t) \tag{1.57}$$

$$k'(t) = bu(t)e^{-at} \tag{1.58}$$

$$k(t) = \int_0^t bu(\tau)e^{-a\tau} d\tau + cte \tag{1.59}$$

d'où

$$x(t) = e^{at} \left[\int_0^t bu(\tau)e^{-a\tau} d\tau + cte \right] \quad (1.60)$$

en identifiant $x(0) = x_0$

$$x(t) = e^{at} \left[\int_0^t bu(\tau)e^{-a\tau} d\tau \right] + e^{at}x_0 \quad (1.61)$$

en généralisant

$$x(t) = \int_{t_0}^t bu(\tau)e^{-a(t-\tau)} d\tau + e^{a(t-t_0)}x(t_0). \quad (1.62)$$

1.4.2 Cas général

$$\underline{x}(t) = e^{\mathbf{A}t} \left[\int_{t_0}^t e^{-\mathbf{A}\tau} \mathbf{B}u(\tau) d\tau \right] + e^{\mathbf{A}(t-t_0)} \underline{x}(t_0) \quad (1.63)$$

Définition :

$$e^{\mathbf{A}t} = 1 + \mathbf{A}t + \frac{1}{2}(\mathbf{A}t)^2 + \frac{1}{3!}(\mathbf{A}t)^3 + \dots \quad (1.64)$$

propriétés de $e^{\mathbf{A}t}$

$$e^{\mathbf{A}t} e^{\mathbf{B}t} = e^{(\mathbf{A}+\mathbf{B})t} \quad \text{si} \quad \mathbf{A}\mathbf{B} = \mathbf{B}\mathbf{A} \quad (1.65)$$

$$\frac{d}{dt} e^{\mathbf{A}t} = \mathbf{A}e^{\mathbf{A}t} \quad (1.66)$$

$$\int_{t_1}^{t_2} e^{\mathbf{A}\tau} d\tau = \mathbf{A}^{-1} [e^{\mathbf{A}t_2} - e^{\mathbf{A}t_1}] = [e^{\mathbf{A}t_2} - e^{\mathbf{A}t_1}] \mathbf{A}^{-1} \quad (1.67)$$

1.4.3 Généralisation aux systèmes variants dans le temps

En posant $\Phi(t, t_0)$ solution de l'équation

$$\dot{\Phi}(t, t_0) = \mathbf{A}(t)\Phi(t, t_0) \quad (1.68)$$

alors

$$\underline{x}(t) = \int_{t_0}^t \Phi(t, \tau) \mathbf{B}(\tau) u(\tau) d\tau + \Phi(t, t_0) x(t_0) \quad (1.69)$$

Sauf dans quelques cas particuliers, cette équation est irrésolvable.

dans le cas particulier où la matrice \mathbf{A} est diagonale la solution est de la forme

$$\Phi(t, t_0) = \exp \left(\int_{t_0}^t \mathbf{A}(\tau) d\tau \right) \quad (1.70)$$

1.4.4 Simulation sur ordinateur

Equation générale, en posant $t = kT_e$

$$\underline{x}((k+1)T_e) = e^{\mathbf{A}(k+1)T_e} \left[\int_{kT_e}^{(k+1)T_e} e^{-\mathbf{A}\tau} \mathbf{B}u(\tau) d\tau \right] + e^{\mathbf{A}T_e} \underline{x}(kT_e) \quad (1.71)$$

Hypothèse $u(t) = cte$ entre kT_e et $(k+1)T_e$,

$$\underline{x}((k+1)T_e) = e^{\mathbf{A}(k+1)T_e} \left[e^{-\mathbf{A}(k+1)T_e} - e^{-\mathbf{A}kT_e} \right] \mathbf{A}^{-1} \mathbf{B}u(kT_e) + e^{\mathbf{A}T_e} \underline{x}(kT_e) \quad (1.72)$$

$$\underline{x}((k+1)T_e) = - \left[1 - e^{-\mathbf{A}T_e} \right] \mathbf{A}^{-1} \mathbf{B}u(kT_e) + e^{\mathbf{A}T_e} \underline{x}(kT_e) \quad (1.73)$$

$$\underline{x}((k+1)T_e) = \left(T_e + \frac{T_e^2}{2} \mathbf{A} + \frac{T_e^3}{3!} (\mathbf{A})^2 + \dots \right) \mathbf{B}u(kT_e) + e^{\mathbf{A}T_e} \underline{x}(kT_e) \quad (1.74)$$

simplification : Méthode d'Euler

$$\underline{x}((k+1)T_e) = - \left[1 - e^{-\mathbf{A}T_e} \right] \mathbf{A}^{-1} \mathbf{B}u(kT_e) + e^{\mathbf{A}T_e} \underline{x}(kT_e) \quad (1.75)$$

avec un développement au premier degré,

$$\underline{x}((k+1)T_e) = T_e \mathbf{B}u(kT_e) + (1 + \mathbf{A}T_e) \underline{x}(kT_e) \quad (1.76)$$

Vérification dans le cas simple

$$\frac{dx}{dt} = \frac{x((k+1)T_e) - x(kT_e)}{T_e} \quad (1.77)$$

et

$$\frac{dx}{dt} = ax + bu \quad (1.78)$$

$$x((k+1)T_e) - x(kT_e) = T_e ax(kT_e) + T_e bu(kT_e) \quad (1.79)$$

$$x((k+1)T_e) = (T_e a + 1)x(kT_e) + T_e bu(kT_e) \quad (1.80)$$

Vous êtes maintenant capables simuler sur un ordinateur le comportement d'un système donné sous la forme d'une représentation d'état. Notez que, si le système est invariant dans le temps, le choix de T_e n'est pas crucial. Les termes $- \left[1 - e^{-\mathbf{A}T_e} \right] \mathbf{A}^{-1} \mathbf{B}$ et $e^{\mathbf{A}T_e}$ de l'équation (1.75) peuvent être pré-calculés une fois pour toutes. Dans le cas contraire il vaut mieux choisir T_e suffisamment petit pour que l'équation (1.76) soit valable.

Chapitre 2

Obtention des équations d'état

2.1 Méthode directe

La méthode qui donne le maximum d'informations est celle qui consiste à déterminer la représentation d'état directement à partir des équations de la physique appliquées au système considéré.

2.2 A partir de la fonction de transfert

Souvent, les systèmes physiques sont donnés sous la forme d'une fonction de transfert, les paragraphes suivants traitent du passage de la fonction de transfert vers la représentation d'état.

Tout système monovarié peut être représenté sous la forme

$$\frac{Y(p)}{U(p)} = \frac{b_0 + b_1p + b_2p^2 + \dots + b_m p^m}{a_0 + a_1p + a_2p^2 + \dots + a_n p^n} \quad (2.1)$$

En général $m < n$, car les systèmes physiques sont filtrants, quelquefois, $m = n$.

2.2.1 Forme 1 : forme canonique de commandabilité

L'équation (2.1) peut se décomposer sous la forme suivante

$$\frac{Y(p)}{W(p)} \times \frac{W(p)}{U(p)} \quad \text{avec} \quad \frac{W(p)}{U(p)} = \frac{1}{a_0 + a_1p + a_2p^2 + \dots + a_n p^n} \quad (2.2)$$

sous cette forme, nous pouvons en déduire que

$$U(p) = a_0W(p) + a_1pW(p) + a_2p^2W(p) + \dots + a_n p^n W(p) \quad (2.3)$$

donc,

$$U(p) = a_0W(p) + a_1\dot{W}(p) + a_2\ddot{W}(p) + \dots + a_n \overset{n}{W}(p) \quad (2.4)$$

d'où le choix du vecteur d'état

$$\underline{x} = \begin{bmatrix} W \\ \dot{W} \\ \vdots \\ \overset{(n-1)}{W} \\ W \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad (2.5)$$

on en déduit

$$\begin{aligned} \dot{x}_1 &= x_2 \\ \dot{x}_2 &= x_3 \\ &\vdots \\ \dot{x}_n &= \frac{1}{a_n}U(p) - \frac{a_0}{a_n}x_1 - \frac{a_1}{a_n}x_2 + \dots - \frac{a_{n-1}}{a_n}x_n \end{aligned} \quad (2.6)$$

$$y = b_0x_1 + b_1x_2 + b_2x_3 + \cdots + b_{n-1}x_n + b_n \left(\frac{1}{a_n}U(p) - \frac{a_0}{a_n}x_1 - \frac{a_1}{a_n}x_2 + \cdots - \frac{a_{n-1}}{a_n}x_n \right) \quad (2.7)$$

la représentation d'état est alors, dans le cas $n = m$

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \vdots \\ \dot{x}_{n-1} \\ \dot{x}_n \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & 0 & \cdots & 0 & 1 \\ -\frac{a_0}{a_n} & -\frac{a_1}{a_n} & -\frac{a_2}{a_n} & \cdots & -\frac{a_{n-1}}{a_n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \frac{1}{a_n} \end{bmatrix} u \quad (2.8)$$

$$y = \left[b_0 - \frac{a_0}{a_n}b_n, \cdots, b_{n-1} - \frac{a_{n-1}}{a_n}b_n \right] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} + \frac{b_n}{a_n}u \quad (2.9)$$

Dans le cas nettement plus fréquent où $m < n$ et après simplification par a_n , le système est plus simple

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \vdots \\ \dot{x}_{n-1} \\ \dot{x}_n \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & 0 & \cdots & 0 & 1 \\ -a_0 & -a_1 & -a_2 & \cdots & -a_{n-1} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} u \quad (2.10)$$

$$y = [b_0, b_1, \cdots, b_m, 0, \cdots, 0] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} \quad (2.11)$$

2.2.2 Forme 2 : forme canonique d'observabilité

Dans le cas où $m < n$ et après simplification par a_n ,

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \vdots \\ \dot{x}_{n-1} \\ \dot{x}_n \end{bmatrix} = \begin{bmatrix} 0 & 0 & \cdots & 0 & -a_0 \\ 1 & 0 & \cdots & 0 & -a_1 \\ 0 & 1 & \ddots & \vdots & \vdots \\ \vdots & & \ddots & 0 & -a_{n-2} \\ 0 & 0 & \cdots & 1 & -a_{n-1} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_m \\ 0 \\ \vdots \\ 0 \end{bmatrix} u \quad (2.12)$$

$$y = [0, \cdots, 0, 1] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} \quad (2.13)$$

Noter la symétrie des deux formes canoniques, transposition par rapport à la diagonale principale.

2.2.3 Représentation modale

La représentation modale consiste à choisir le vecteur d'état tel que la matrice \mathbf{A} soit diagonale et fasse apparaître les valeurs propres du système.

Valeurs propres d'une matrice,

$$\det [\lambda \mathbf{I} - \mathbf{A}] = 0 \quad (2.14)$$

S'il y a n valeurs propres distinctes, il y a aussi n vecteurs propres \underline{v}_i tels que

$$[\lambda_i \mathbf{I} - \mathbf{A}] \underline{v}_i = 0 \quad (2.15)$$

d'où l'on en déduit une matrice de passage

$$\mathbf{P}_{(n,n)} = [\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n] \quad (2.16)$$

le nouveau vecteur d'état est x' tel que

$$x = \mathbf{P}x' \quad (2.17)$$

La nouvelle matrice d'état est donc

$$\mathbf{A}' = \mathbf{P}^{-1} \mathbf{A} \mathbf{P} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & \lambda_n \end{bmatrix} \quad (2.18)$$

2.2.4 Forme canonique de Jordan (forme diagonale)

Méthode : décomposition en éléments simples de la fraction rationnelle $G(p)$

Cas où tous les pôles sont simples

$$Y(p) = \left(\sum_{i=1}^n \frac{r_i}{p - \lambda_i} + r_0 \right) U(p)$$

où $r_0 \neq 0$ si $m = n$

Choix des variables d'état :

$$X_i(p) = \frac{1}{p - \lambda_i} U(p), \quad i = 1, 2, \dots, n$$

donc

$$Y(p) = \sum_{i=1}^n r_i X_i(p) + r_0 U(p)$$

En prenant la transformée de Laplace inverse

$$\begin{aligned} \dot{x}_i &= \lambda_i x_i + u, & i = 1, 2, \dots, n \\ y &= \sum_{i=1}^n r_i x_i + r_0 u \end{aligned}$$

Ecriture matricielle

$$\begin{pmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \vdots \\ \dot{x}_{n-1} \\ \dot{x}_n \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & \lambda_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix} u \quad (2.19)$$

$$y = (r_1, r_2, \dots, r_n) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{pmatrix} + r_0 u \quad (2.20)$$

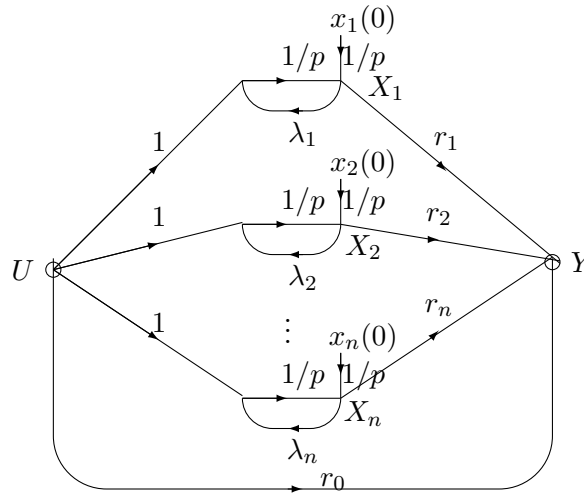


FIGURE 2.1 – Graphe de fluence de la représentation d'état sous la forme canonique de Jordan

Cas d'un pôle multiple

soit λ_1 d'ordre q donc

$$Y(p) = \left[\frac{r_1}{p - \lambda_1} + \cdots + \frac{r_q}{(p - \lambda_1)^q} + \sum_{i=q+1}^n \frac{r_i}{p - \lambda_i} + r_0 \right] U(p)$$

Choix des variables d'état

$$\begin{aligned} X_1 &= \frac{1}{p - \lambda_1} U, & X_2 &= \frac{1}{(p - \lambda_1)^2} U, \dots, X_q = \frac{1}{(p - \lambda_1)^q} U \\ X_i &= \frac{1}{p - \lambda_i} U, & i &= q + 1, \dots, n \end{aligned}$$

d'où

$$X_2 = \frac{1}{p - \lambda_1} X_1, \quad X_3 = \frac{1}{p - \lambda_1} X_2, \dots, X_q = \frac{1}{p - \lambda_1} X_{q-1}$$

En prenant la transformée de Laplace inverse

$$\begin{aligned} \dot{x}_1 &= \lambda_1 x_1 + u \\ \dot{x}_2 &= \lambda_1 x_2 + x_1 \\ &\vdots \\ \dot{x}_q &= \lambda_1 x_q + x_{q-1} \\ \dot{x}_i &= \lambda_i x_i + u, & i &= q + 1, \dots, n \\ y &= \sum_{i=1}^n r_i x_i + r_0 u \end{aligned}$$

Ecriture matricielle

$$\begin{pmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \vdots \\ \dot{x}_q \\ x_{q+1} \\ \vdots \\ \dot{x}_n \end{pmatrix} = \begin{pmatrix} \lambda_1 & & & & & & \\ 1 & \lambda_1 & & & & & \\ & \ddots & \ddots & & & & \\ & & 1 & \lambda_1 & & & \\ & & & \lambda_{q+1} & & & \\ & & & & \ddots & & \\ & & & & & \lambda_n & \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_q \\ x_{q+1} \\ \vdots \\ x_n \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix} u \quad (2.21)$$

$$y = (r_1, r_2, \dots, r_q, r_{q+1}, \dots, r_n) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_q \\ x_{q+1} \\ \vdots \\ x_n \end{pmatrix} + r_0 u \quad (2.22)$$

Matrice du système obtenue : forme réduite de Jordan

Grphe correspondant :

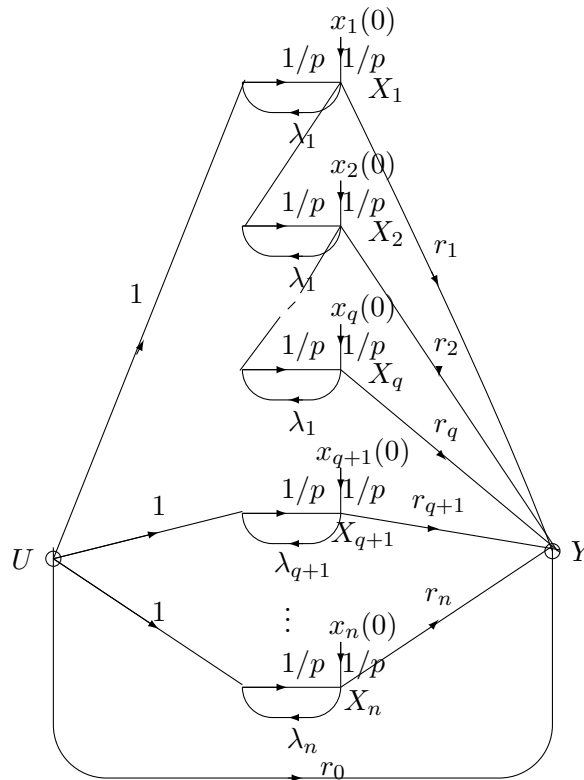


FIGURE 2.2 – Grphe de fluence de la représentation d'état sous la forme canonique de Jordan

Remarque 1 : Dans le cas de plusieurs pôles multiples, à chaque pôle correspond un sous bloc de Jordan carré $q \times q$, avec la valeur de ce pôle sur la diagonale et 1 sur le sous diagonale.

Remarque 2 : Les variables d'état ne sont plus entièrement découplées.

Cas d'une paire de pôles complexes conjugués

2 solutions :

1. Etablir une forme réduite de Jordan comme précédemment avec des variables d'état complexes.
2. Rechercher une représentation autre que diagonale, où tous les coefficients sont réels. (voir §4).

Exemple :

Equation d'état :

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} a + jb & 0 \\ 0 & a - jb \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} u \quad (2.23)$$

Equation de sortie :

$$\omega = [c \quad \bar{c}] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (2.24)$$

Dans cet exemple, les variables d'état sont des fonctions complexes conjuguées.

Le changement de variable $\underline{z} = \mathbf{M}\underline{x}$ avec $\mathbf{M} = \begin{bmatrix} 1 & 1 \\ j & -j \end{bmatrix}$ permet de transformer le système précédent en :

Equation d'état :

$$\begin{bmatrix} \dot{z}_1 \\ \dot{z}_2 \end{bmatrix} = \begin{bmatrix} a & b \\ -b & a \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} + \begin{bmatrix} 2 \\ 0 \end{bmatrix} u \quad (2.25)$$

Equation de sortie :

$$\omega = [\operatorname{Re}(c) \quad \operatorname{Im}(c)] \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad (2.26)$$

En résumé

$$\mathbf{A} = \begin{bmatrix} a & b & 0 & 0 & 0 \\ -b & a & 0 & 0 & 0 \\ 0 & 0 & c & 0 & 0 \\ 0 & 0 & 0 & d & 1 \\ 0 & 0 & 0 & 0 & d \end{bmatrix} \left. \begin{array}{l} \left. \begin{array}{l} \text{pôles complexes conjugués} \\ \text{pôle réel} \end{array} \right\} a \pm jb \\ \left. \begin{array}{l} \text{pôle réel} \\ \text{pôles réels d'ordre 2} \end{array} \right\} c \\ \left. \begin{array}{l} \text{pôles réels d'ordre 2} \end{array} \right\} d \end{array} \right. \quad (2.27)$$

;

$$\mathbf{B} = \begin{bmatrix} 2 \\ 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} \left. \begin{array}{l} \left. \begin{array}{l} \text{pôles complexes conjugués} \\ \text{pôle réel} \end{array} \right\} a \pm jb \\ \left. \begin{array}{l} \text{pôle réel} \\ \text{pôles réels d'ordre 2} \end{array} \right\} c \\ \left. \begin{array}{l} \text{pôles réels d'ordre 2} \end{array} \right\} d \end{array} \right. \quad (2.28)$$

Chapitre 3

Commandabilité et observabilité des systèmes

Problème :

- Peut-on à partir des commandes $u(t)$ contrôler l'état du système c'est-à-dire $\underline{x}(t)$?
— → problème de commandabilité.
- Peut-on à partir de l'observation de $y(t)$ remonter à l'état du système $\underline{x}(t)$?
— → problème d'observabilité.

3.1 Définitions

- Un système est totalement commandable s'il existe $u(t)$ qui conduise en un temps fini d'un état $x(t_0)$ à un état $x(t_1)$ quels que soient $x(t_0)$ et $x(t_1)$.
- Un système est observable, si l'observation de $y(t)$ entre t_0 et t_1 permet de retrouver $x(t_0)$

Exemple

$$\mathbf{A} = \begin{bmatrix} a_{1,1} & 0 & \cdots & 0 \\ 0 & a_{2,2} & & \vdots \\ \vdots & & a_{3,3} & 0 \\ 0 & \cdots & 0 & a_{4,4} \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} b_1 \\ b_2 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} 0 & c_{1,2} & 0 & 0 \\ 0 & 0 & 0 & c_{2,4} \end{bmatrix} \quad (3.1)$$

Le schéma représentatif de ce système est donné en figure 3.1.

Condition de commandabilité : \mathbf{Q}_S est de rang n avec

$$\mathbf{Q}_S = [\mathbf{B}, \mathbf{AB}, \mathbf{A}^2\mathbf{B}, \dots, \mathbf{A}^{n-1}\mathbf{B}] \quad (3.2)$$

Application à l'exemple du moteur commandé par l'inducteur,

Equation d'état :

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -a_0 & -a_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ b_0 \end{bmatrix} u \quad (3.3)$$

Calcul de \mathbf{Q}_S

$$\mathbf{B} = \begin{bmatrix} 0 \\ b_0 \end{bmatrix}, \quad \mathbf{AB} = \begin{bmatrix} b_0 \\ -b_0 a_1 \end{bmatrix} \quad (3.4)$$

d'où

$$\mathbf{Q}_S = \begin{bmatrix} 0 & b_0 \\ b_0 & -b_0 a_1 \end{bmatrix} \quad (3.5)$$

et

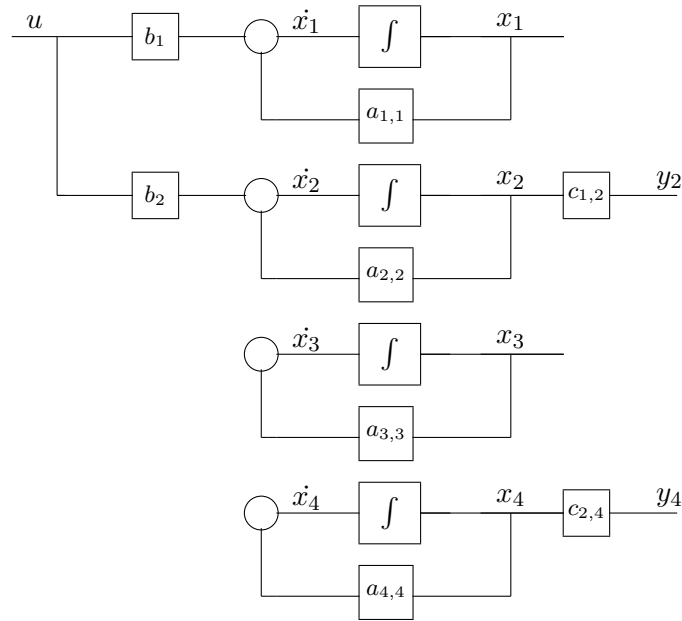


FIGURE 3.1 – Schéma d'un simulateur analogique

$$\det \mathbf{Q}_S = -b_0^2 \quad (3.6)$$

donc le système est commandable (si $b_0 \neq 0$).

Condition d'observabilité : \mathbf{Q}_B est de rang n avec

$$\mathbf{Q}_B = \begin{bmatrix} \mathbf{C} \\ \mathbf{CA} \\ \mathbf{CA}^2 \\ \vdots \\ \mathbf{CA}^{n-1} \end{bmatrix} \quad (3.7)$$

Application à l'exemple du moteur commandé par l'inducteur,

Equation de sortie :

$$\omega = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (3.8)$$

donc

$$\mathbf{C} = \begin{bmatrix} 1 & 0 \end{bmatrix}, \quad \mathbf{CA} = \begin{bmatrix} 0 & 1 \end{bmatrix} \quad (3.9)$$

donc

$$\mathbf{Q}_B = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (3.10)$$

et comme $\det \mathbf{Q}_B = 1$, \mathbf{Q}_B est de rang n et donc le système est observable.

3.2 Que faire si un système n'est pas observable et/ou commandable

Il arrive souvent qu'après une première analyse du système celui-ci s'avère être non commandable ou non observable. Deux solutions s'offrent à vous.

3.2.1 Retour sur conception

La première solution consiste à ajouter ou changer des organes de commande (non commandable) ou des capteurs (non observable).

3.2.2 Réduction de modèles

La deuxième solution n'est valable que si la maîtrise complète de l'état n'est pas importante. Dans ce cas, la réduction du modèle est envisageable. Une façon simple de procéder est de déterminer la représentation d'état à partir de l'équation différentielle (qui "élimine" les inobservabilités) ou de la fonction de transfert du système (qui "élimine" les inobservabilités et les non commandabilités).

Chapitre 4

Transformation en l'une des formes canoniques

4.1 Diagonalisation de la matrice \mathbf{A}

Hypothèses :

— le système est décrit par les équations suivantes :

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}u \quad (4.1)$$

$$\underline{y} = \mathbf{C}\underline{x} + \mathbf{D}u \quad (4.2)$$

— Les valeurs propres de \mathbf{A} ($\lambda_1, \dots, \lambda_n$) sont toutes différentes.

Il existe alors une transformation \mathbf{T} définie par

$$\underline{x} = \mathbf{T}\underline{x}^*$$

avec : \underline{x} = ancien vecteur d'état,

\underline{x}^* = nouveau vecteur d'état qui diagonalise \mathbf{A}

Dans ce cas \mathbf{T} est la matrice modale.

On en déduit une nouvelle représentation d'état.

$$\dot{\underline{x}}^* = \mathbf{\Lambda}\underline{x}^* + \widehat{\mathbf{B}}u \quad (4.3)$$

$$\underline{y} = \widehat{\mathbf{C}}\underline{x}^* + \widehat{\mathbf{D}}u \quad (4.4)$$

$$\mathbf{\Lambda} = \mathbf{T}^{-1}\mathbf{A}\mathbf{T} = \begin{pmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{pmatrix}$$

avec :

$$\begin{aligned} \widehat{\mathbf{B}} &= \mathbf{T}^{-1}\mathbf{B} \\ \widehat{\mathbf{C}} &= \mathbf{C}\mathbf{T} \\ \widehat{\mathbf{D}} &= \mathbf{D} \end{aligned}$$

• Cas où \mathbf{A} est une matrice compagne (une forme canonique) avec des valeurs propres toutes différentes, alors \mathbf{T} est une matrice de Vandermonde :

$$\mathbf{T} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ \lambda_1 & \lambda_2 & \cdots & \lambda_n \\ \lambda_1^2 & \lambda_2^2 & \cdots & \lambda_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{n-1} & \lambda_2^{n-1} & \cdots & \lambda_n^{n-1} \end{pmatrix}$$

• Cas général.

Vecteur propre \underline{v}_i associé à la valeur propre λ_i , obtenu par

$$\mathbf{A}\underline{v}_i = \lambda_i \underline{v}_i$$

c'est-à-dire

$$(\lambda_i \mathbf{I} - \mathbf{A})\underline{v}_i = 0.$$

Rappel : les λ_i sont obtenus par la résolution de l'équation caractéristique de la matrice \mathbf{A}

$$\det(\lambda_i \mathbf{I} - \mathbf{A}) = 0$$

la matrice de transformation \mathbf{T} est alors donnée par

$$\mathbf{T} = \begin{pmatrix} v_{11} & v_{12} & \cdots & v_{1n} \\ v_{21} & v_{22} & \cdots & v_{2n} \\ \vdots & & & \vdots \\ v_{n1} & v_{n2} & \cdots & v_{nn} \end{pmatrix} = (\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n)$$

4.2 Conséquences pour la commandabilité et l'observabilité

Propriété 1 : Si le système $\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}\underline{u}$ possède des valeurs propres toutes différentes, il est commandable si et seulement si aucun élément du vecteur colonne $\widehat{\mathbf{B}} = \mathbf{T}^{-1}\mathbf{B}$ n'est nul.

démonstration : voir graphe du système découplé (fig. (2.2) page 22).

note : commandabilité $\longleftrightarrow u(t)$ doit pouvoir influencer tout x_i .

Propriété 2 : Si le système

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}\underline{u} \tag{4.5}$$

$$\underline{y} = \mathbf{C}\underline{x} + \mathbf{D}\underline{u} \tag{4.6}$$

possède des valeurs propres toutes différentes, il est observable si et seulement si aucun élément du vecteur ligne $\widehat{\mathbf{C}} = \mathbf{C}\mathbf{T}$ n'est nul.

démonstration : voir graphe du système découplé (fig. (2.2) page 22).

note : observabilité \longleftrightarrow tout x_i doit pouvoir influencer $y(t)$.

Relation avec la fonction de transfert

$$\dot{x}_i^* = \lambda_i x_i^* + \widehat{b}_i u \tag{4.7}$$

$$X_i^* = \frac{1}{p - \lambda_i} \widehat{b}_i u \tag{4.8}$$

$$Y = \sum_{i=1}^n \widehat{c}_i X_i^* + dU = \underbrace{\left(\sum_{i=1}^n \frac{\widehat{b}_i \widehat{c}_i}{p - \lambda_i} + d \right)}_{G(p)} U \tag{4.9}$$

Si le système n'est pas commandable ou pas observable, au moins l'un des \widehat{b}_i ou des \widehat{c}_i sera nul, donc le terme correspondant du développement de $G(p)$ en éléments simples disparaîtra, donc la valeur propre λ_i correspondante ne sera pas un pôle de $G(p)$.

Propriété 3 : un système à 1 entrée et 1 sortie et à valeurs propres toutes différentes est à la fois commandable et observable

- si et seulement si toutes les valeurs propres sont des pôles de la fonction de transfert.
- si et seulement si la fonction de transfert $G(p)$, dans laquelle ne doit plus apparaître au numérateur et au dénominateur un terme identique, a un dénominateur de degré égal au nombre n des équations différentielles d'état.

4.3 Cas des valeurs propres complexes

S'il existe des valeurs propres complexes

$$\lambda_i = \sigma + j\omega \quad \text{et} \quad (4.10)$$

$$\lambda_{i+1} = \sigma - j\omega = \overline{\lambda_i} \quad (4.11)$$

parmi les valeurs propres de \mathbf{A} , il existe deux solutions possibles.

4.3.1 Diagonalisation classique

Matrice diagonale, à éléments diagonaux λ_i et λ_{i+1} complexes.

exemple : système du deuxième ordre.

$$\mathbf{\Lambda} = \begin{pmatrix} \sigma + j\omega & 0 \\ 0 & \sigma - j\omega \end{pmatrix}$$

On montre aisément que les vecteur propres associés sont eux aussi complexes conjugués.

$$\underline{v}_i = \underline{\alpha} + j\underline{\beta} \quad (4.12)$$

$$\underline{v}_{i+1} = \underline{\alpha} - j\underline{\beta} \quad (4.13)$$

$$(4.14)$$

4.3.2 Transformation modifiée : \mathbf{T}_m

$$\underline{x} = \mathbf{T}_m \underline{w}$$

avec \underline{w} vecteur d'état qui "diagonalise" \mathbf{A} et $\mathbf{T}_m = (\underline{\alpha} \quad \underline{\beta})$

Comme

$$\begin{aligned} \mathbf{A}\underline{v}_1 &= \lambda_1 \underline{v}_1 \\ &= (\sigma + j\omega)\underline{v}_1 \\ &= (\sigma + j\omega)(\underline{\alpha} + j\underline{\beta}) \\ &= (\sigma\underline{\alpha} - \omega\underline{\beta}) + j(\omega\underline{\alpha} + \sigma\underline{\beta}) \end{aligned}$$

et que : $\mathbf{A}\underline{v}_i = \mathbf{A}\underline{\alpha} + j\mathbf{A}\underline{\beta}$

on a :

$$\begin{aligned} \mathbf{A}\underline{\alpha} &= \sigma\underline{\alpha} - \omega\underline{\beta} \\ \mathbf{A}\underline{\beta} &= \omega\underline{\alpha} + \sigma\underline{\beta} \end{aligned}$$

donc

$$\mathbf{A} (\underline{\alpha} \quad \underline{\beta}) = (\underline{\alpha} \quad \underline{\beta}) \begin{pmatrix} \sigma & \omega \\ -\omega & \sigma \end{pmatrix} \quad (4.15)$$

$$\mathbf{A}\mathbf{T}_m = \mathbf{T}_m\mathbf{\Lambda}_m \quad (4.16)$$

d'où

$$\mathbf{\Lambda}_m = \mathbf{T}_m^{-1}\mathbf{A}\mathbf{T}_m = \begin{pmatrix} \sigma & \omega \\ -\omega & \sigma \end{pmatrix}$$

Généralisation : système d'ordre n , ayant (à titre d'exemple) deux paires de valeurs propres complexes conjuguées $\lambda_1, \lambda_2, \overline{\lambda_1}$ et $\overline{\lambda_2}$

la transformation modifiée est :

$$\underline{x} = \mathbf{T}_m \underline{w}$$

avec

$$\mathbf{T}_m = [\text{Re}(\underline{v}_1), \text{Im}(\underline{v}_1), \text{Re}(\underline{v}_2), \text{Im}(\underline{v}_2), \underline{v}_3, \dots, \underline{v}_n]$$

la nouvelle représentation est :

$$\dot{\underline{x}} = \mathbf{A}_m \underline{x} + \mathbf{B}_m \underline{u} \quad (4.17)$$

$$\underline{y} = \mathbf{C}_m \underline{x} + \mathbf{D}_m \underline{u} \quad (4.18)$$

avec :

$$\mathbf{A}_m = \mathbf{T}_m^{-1} \mathbf{A} \mathbf{T}_m = \begin{pmatrix} \sigma_1 & \omega_1 & & & & \\ -\omega_1 & \sigma_1 & & & & \\ & & \sigma_2 & \omega_2 & & \\ & & -\omega_2 & \sigma_2 & & \\ & & & & \lambda_5 & \\ & & & & & \ddots \\ & & & & & & \lambda_n \end{pmatrix}$$

et

$$\mathbf{B}_m = \mathbf{T}^{-1} \mathbf{B}, \quad \mathbf{C}_m = \mathbf{C} \mathbf{T}, \quad \mathbf{D}_m = \mathbf{D}$$

4.4 Transformation en la forme canonique d'asservissement

Hypothèses : $a_n = 1$ et $b_n = 0$.

On montre aisément que l'équation caractéristique

$$\det(p\mathbf{I} - \mathbf{A}) = 0$$

se réduit à :

$$a_0 + a_1 p + a_2 p^2 + \dots + a_{n-1} p^{n-1} + p^n = 0$$

Dans la forme canonique d'asservissement, la dernière ligne de la matrice du système est composée des coefficients de l'équation caractéristique (excepté celui du terme de plus haut degré) changés de signe.

Théorème : Si le système est commandable on peut le mettre sous la forme canonique d'asservissement au moyen de la transformation

$$\underline{z} = \mathbf{T}^{-1} \underline{x}$$

avec

$$\mathbf{T}^{-1} = \begin{pmatrix} \underline{q}_S^T & & & & \\ \underline{q}_S^T & \mathbf{A} & & & \\ \underline{q}_S^T & \mathbf{A}^2 & & & \\ \vdots & & & & \\ \underline{q}_S^T & \mathbf{A}^{n-1} & & & \end{pmatrix}$$

où \underline{q}_S^T est la dernière ligne de l'inverse de la matrice de commandabilité \mathbf{Q}_S^{-1}

L'utilisation de ce théorème peut sembler lourde, mais dans la pratique il se révèle puissant, puisqu'une grande partie des termes, notamment de \mathbf{Q}_S^{-1} , est inutile.

4.5 Transformation en la forme canonique d'observabilité

Dans la forme canonique d'observation, la dernière colonne de la matrice du système est composée des coefficients de l'équation caractéristique (excepté celui du terme de plus haut degré) changés de signe.

Théorème : Si le système est observable on peut le mettre sous la forme canonique d'asservissement au moyen de la transformation

$$\underline{z} = \mathbf{T} \underline{x}$$

avec

$$\mathbf{T} = \left(\underline{q}_B \quad \mathbf{A}\underline{q}_B \quad \mathbf{A}^2\underline{q}_B \quad \cdots \quad \mathbf{A}^{n-1}\underline{q}_B \right)$$

où \underline{q}_b est la dernière colonne de l'inverse de la matrice d'observabilité \mathbf{Q}_B^{-1}

Chapitre 5

Stabilité des systèmes dynamiques linéaires

La stabilité des systèmes linéaires, dans le cas SISO, consiste le plus souvent à revenir à des méthodes classiques des systèmes représentés sous la forme d'une fonction de transfert. Attention, la fonction de transfert d'un système présentant moins d'information que sa représentation d'état, il est possible mais rare que la fonction de transfert soit stable alors que la représentation d'état présente un mode instable !

5.1 Définition

Un système est dit asymptotiquement stable, si et seulement si, écarté de sa position d'origine et soumis à une entrée nulle, celui-ci revient à sa position d'origine.

5.2 Etude de la stabilité

En reprenant l'équation 1.63 et en supposant que $u(t) = 0$ et $t_0 = 0$, nous obtenons :

$$\underline{x}(t) = e^{\mathbf{A}t} \underline{x}(0)$$

D'après la définition de la stabilité, il faut que $\underline{x}(t) \rightarrow 0$ donc que $e^{\mathbf{A}t} \rightarrow 0$.

Si \mathbf{A} est diagonalisable ($\mathbf{\Lambda}$) alors il existe une matrice \mathbf{T} telle que :

$$\mathbf{\Lambda} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & \lambda_n \end{bmatrix} \quad (5.1)$$

donc

$$\underline{x}(t) = e^{\mathbf{T} \mathbf{\Lambda} \mathbf{T}^{-1} t} \underline{x}(0) \rightarrow 0$$

donc

$$\exp \begin{bmatrix} \lambda_1 t & 0 & \cdots & 0 \\ 0 & \lambda_2 t & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & \lambda_n t \end{bmatrix} \rightarrow 0 \quad (5.2)$$

donc il faut et il suffit que les λ_i soient à partie réelle négative, donc que les valeurs propres de la matrice \mathbf{A} soient à partie réelle négative.

En d'autres termes, pour qu'un système soit asymptotiquement stable en boucle ouverte, il faut et il suffit que tous ses modes soient stables. L'extension aux systèmes non diagonalisables (blocs de Jordan) est immédiate puisque les termes de l'exponentielle d'une matrice triangulaire sont une combinaison linéaire des $e^{\lambda_i t}$.

Forme modale

La lecture de la stabilité est immédiate, puisque les valeurs propres de la matrice sont sur la diagonale.

Forme canonique de commandabilité

La matrice \mathbf{A} se présente alors sous la forme :

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & 0 & \dots & 0 & 1 \\ -\frac{a_0}{a_n} & -\frac{a_1}{a_n} & -\frac{a_2}{a_n} & \dots & -\frac{a_{n-1}}{a_n} \end{bmatrix}$$

le calcul des valeurs propres de la matrice revient à calculer les zéros du polynôme

$$P(\lambda) = -\frac{a_0}{a_n} - \frac{a_1}{a_n}\lambda - \frac{a_2}{a_n}\lambda^2 - \dots - \frac{a_{n-1}}{a_n}\lambda^{n-1}$$

En fait, seul le signe des parties réelles des valeurs propres nous intéresse pour l'étude de la stabilité, donc il suffit d'appliquer le critère de Routh sur le polynôme P .

Forme canonique d'observabilité

Par transposition de la matrice \mathbf{A} , on revient au cas précédent.

5.3 Stabilité au sens de Lyapounov

L'origine est un état stable si pour tout $\varepsilon > 0$, il existe un nombre $\delta(\varepsilon, t_0) > 0$ tel que $\|\underline{x}(t_0)\| < \delta$ entraîne $\|\underline{x}(t)\| < \varepsilon, \forall t > t_0$. Cet état est asymptotiquement stable si, de plus, il existe un nombre $\delta_a(t_0)$ tel que $\|\underline{x}(t_0)\| < \delta_a(t_0)$ entraîne $\lim_{t \rightarrow \infty} \|\underline{x}(t)\| = 0$.

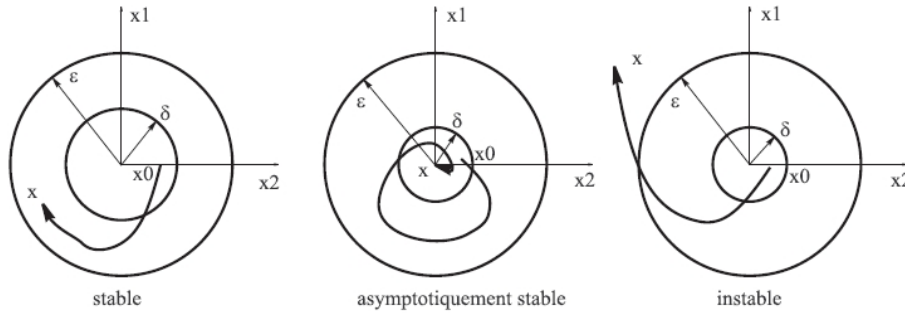


FIGURE 5.1 – Illustration de la stabilité dans le plan de phase

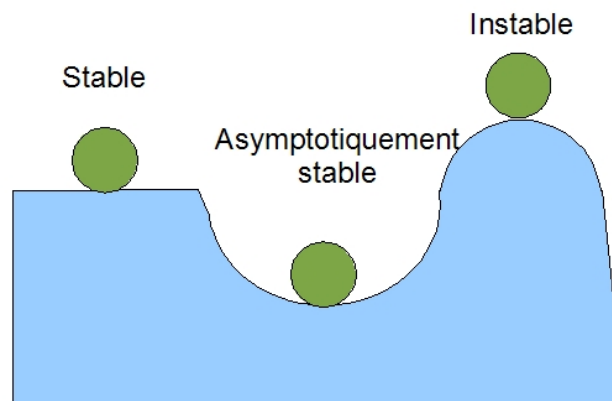


FIGURE 5.2 – Illustration de la stabilité d'une bille sur un profil

5.3.1 Théorème

S'il existe une fonction $V(\underline{x})$ telle que :

1. $V(\underline{x}) > 0, \forall \underline{x} \neq 0$,
2. $V(0) = 0$,
3. $V(\underline{x}) \rightarrow \infty$ pour $\|\underline{x}\| \rightarrow \infty$,

alors l'équilibre en $\underline{x} = 0$ est localement asymptotiquement stable.

Dans le cas des systèmes linéaires invariants dans le temps, la stabilité est globale.

5.3.2 Interprétation physique

Comme nous l'avons dit précédemment, les variables d'état représentent les réservoirs d'énergie du système ou du moins une combinaison linéaire de ceux-ci. Si, à l'origine des temps le système contient de l'énergie alors $\underline{x}(0) \neq 0$. La stabilité est alors synonyme de décroissance de la quantité d'énergie présente dans le système, celle-ci peut converger vers une valeur non nulle. La stabilité asymptotique par contre implique la convergence vers 0.

Supposons que la quantité d'énergie présente dans le système soit $V(t)$, il suffit alors de démontrer que $\dot{V}(t) < 0$ pour que le système soit asymptotiquement stable ($\dot{V}(t) \leq 0$ pour que le système soit stable).

5.3.3 Applications aux systèmes linéaires

Posons

$$V(\underline{x}) = \underline{x}^T \mathbf{P}(t) \underline{x}$$

où \mathbf{P} est une matrice symétrique définie positive, sa dérivée par rapport au temps s'écrit :

$$\dot{V}(\underline{x}) = \dot{\underline{x}}^T \mathbf{P} \underline{x} + \underline{x}^T \dot{\mathbf{P}} \underline{x} + \underline{x}^T \mathbf{P} \dot{\underline{x}} = \underline{x}^T (\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} + \dot{\mathbf{P}}) \underline{x}$$

Sa dérivée est toujours négative s'il existe une matrice définie positive \mathbf{Q} telle que :

$$-\mathbf{Q} = \mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} + \dot{\mathbf{P}}$$

Dans le cas des systèmes linéaires invariants l'équation précédente se simplifie en

$$-\mathbf{Q} = \mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} \quad (5.3)$$

Dans la pratique la méthode est la suivante :

1. Choisir une matrice \mathbf{Q} arbitraire, symétrique, définie positive (la matrice \mathbf{I} conviendra le plus souvent),
2. déterminer ensuite \mathbf{P} à partir de (5.3), cela revient à résoudre $n(n+1)/2$ équations,
3. le système est asymptotiquement stable si \mathbf{P} est définie positive.

Cette méthode peut sembler lourde au regard de celle présentées précédemment, mais c'est la seule qui reste valable dans le cas des systèmes non linéaires et/ou variants dans le temps. La recherche de la fonction de Lyapounov $V(\underline{x}, t)$ devient alors la clef du problème.

5.3.4 Fil rouge

Le système est défini par sa matrice

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -a_0 & -a_1 \end{bmatrix}$$

posons

$$\mathbf{P} = \begin{bmatrix} p_1 & p_2 \\ p_2 & p_3 \end{bmatrix} \quad \text{et} \quad \mathbf{Q} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

l'application de la relation (5.3) donne

$$\begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} = \begin{bmatrix} -2 a_0 p_2 & -a_0 p_3 + p_1 - a_1 p_2 \\ -a_0 p_3 + p_1 - a_1 p_2 & 2 p_2 - 2 a_1 p_3 \end{bmatrix}$$

la résolution de 3 équations à 3 inconnues donne la matrice suivante

$$\mathbf{P} = \begin{bmatrix} \frac{a_0^2 + a_0 + a_1^2}{2 a_0 a_1} & \frac{1}{2 a_0} \\ \frac{1}{2 a_0} & \frac{a_0 + 1}{2 a_0 a_1} \end{bmatrix}$$

sachant que a_0 et a_1 sont positifs, la matrice \mathbf{P} est bien définie positive, le système est stable.

Chapitre 6

Commande des systèmes

6.1 Placement de pôles

Hypothèses : le système est décrit par les équations suivantes

$$\begin{aligned}\dot{\underline{x}} &= \mathbf{A}\underline{x} + \mathbf{B}u \\ \underline{y} &= \mathbf{C}\underline{x}\end{aligned}$$

L'objectif de la synthèse est une régulation (donc pas un asservissement), donc il s'agit de maintenir y proche d'une valeur de consigne y_c . Le schéma de la régulation est donné figure 6.1.

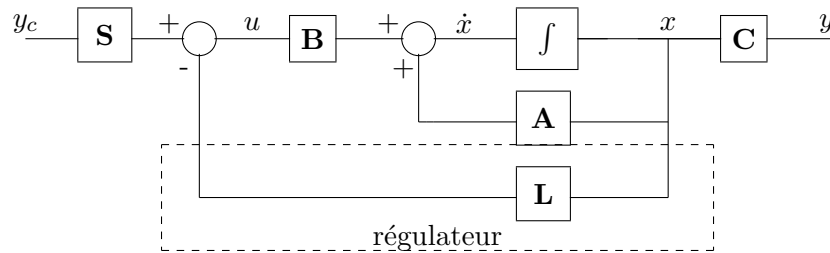


FIGURE 6.1 – Régulation continue par retour d'état

6.1.1 Calcul du régulateur L

Supposons que la fonction de transfert du système en boucle ouverte soit :

$$F_{BO}(p) = \frac{Y(p)}{U(p)} = \frac{b_0 + b_1p + b_2p^2 + \dots + b_m p^m}{a_0 + a_1p + a_2p^2 + \dots + a_n p^n},$$

avec $m < n$

alors la représentation d'état sous la forme canonique de commandabilité est (après simplification par a_n) :

$$\dot{\underline{x}} = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & 0 & \dots & 0 & 1 \\ -a_0 & -a_1 & -a_2 & \dots & -a_{n-1} \end{bmatrix} \underline{x} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} u \quad (6.1)$$

$$\underline{y} = [b_0, b_1, \dots, b_m, 0, \dots, 0] \underline{x} \quad (6.2)$$

On remarque que la dernière ligne de la matrice \mathbf{A} contient les coefficients du dénominateur de la fonction de transfert.

Le schéma 6.1 donne comme commande :

$$u = -\mathbf{L}\underline{x} + \mathbf{S}\underline{y}_c = -(l_0 \ l_1 \ l_2 \ \cdots \ l_{n-1})\underline{x} + \mathbf{S}\underline{y}_c.$$

Avec cette commande, le système en boucle fermée admet pour représentation d'état :

$$\dot{\underline{x}} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & 0 & \cdots & 0 & 1 \\ -l_0 - a_0 & -l_1 - a_1 & -l_2 - a_2 & \cdots & -l_{n-1} - a_{n-1} \end{bmatrix} \underline{x} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \mathbf{S}\underline{y}_c \quad (6.3)$$

$$y = [b_0, b_1, \cdots, b_m, 0, \cdots, 0] \underline{x}. \quad (6.4)$$

Ainsi la fonction de transfert du système en boucle fermée s'écrit :

$$F_{BF}(p) = \frac{Y(p)}{\mathbf{S}y_c(p)} = \frac{b_0 + b_1p + b_2p^2 + \cdots + b_m p^m}{(l_0 + a_0) + (l_1 + a_1)p + (l_2 + a_2)p^2 + \cdots + p^n},$$

donc, le régulateur permet d'imposer arbitrairement le polynôme caractéristique (donc la dynamique) de la fonction de transfert du système en boucle fermée.

Supposons que le polynôme choisi soit :

$$P(p) = \alpha_0 + \alpha_1 p + \alpha_2 p^2 + \cdots + p^n,$$

Le calcul du régulateur est immédiat puisque :

$$(l_i + a_i) = \alpha_i,$$

donc

$$l_i = \alpha_i - a_i.$$

Si le système n'est pas directement sous la forme de commandabilité, il suffit de s'y ramener (voir § 4.4). Le correcteur obtenu précédemment devra subir la transformation inverse à celle utilisée pour obtenir la forme de commandabilité soit :

$$\mathbf{L}^* = \mathbf{L}\mathbf{T}^{-1}$$

6.1.2 Calcul de la matrice de préfiltre \mathbf{S}

Le modèle du système en boucle fermée est :

$$\begin{aligned} \dot{\underline{x}} &= (\mathbf{A} - \mathbf{B}\mathbf{L})\underline{x} + \mathbf{B}\mathbf{S}\underline{y}_c \\ \underline{y} &= \mathbf{C}\underline{x} \end{aligned}$$

Si le système est stable, alors $\lim_{t \rightarrow \infty} \dot{\underline{x}} = 0$ donc,

$$\lim_{t \rightarrow \infty} \underline{y}(t) = -\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{L})^{-1} \mathbf{B}\mathbf{S}\underline{y}_c$$

or, on désire que $\lim_{t \rightarrow \infty} \underline{y}(t) = \underline{y}_c$, donc,

$$\begin{aligned} \mathbf{S}^{-1} &= -\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{L})^{-1} \mathbf{B} \\ \mathbf{S} &= -\left(\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{L})^{-1} \mathbf{B}\right)^{-1} \end{aligned}$$

6.2 Cas d'une représentation quelconque du système à asservir

6.2.1 Transformation en la forme canonique de commandabilité

$$\underline{z} = \mathbf{T}^{-1}\underline{x}$$

avec

$$\mathbf{T}^{-1} = \begin{pmatrix} \underline{q}_S^T \\ \underline{q}_S^T \mathbf{A} \\ \underline{q}_S^T \mathbf{A}^2 \\ \vdots \\ \underline{q}_S^T \mathbf{A}^{n-1} \end{pmatrix}$$

où \underline{q}_S^T est la dernière ligne de l'inverse de la matrice de commandabilité \mathbf{Q}_S^{-1} .

Dans cette représentation

$$\underline{u} = -\mathbf{L}_{FC}\underline{z} = -\mathbf{L}_{FC}\mathbf{T}^{-1}\underline{x} = -\mathbf{L}\underline{x}$$

avec

$$\mathbf{L} = \mathbf{L}_{FC}\mathbf{T}^{-1} = (\alpha_0 - a_0)\underline{q}_S^T + (\alpha_1 - a_1)\underline{q}_S^T \mathbf{A} + \dots + (\alpha_{n-1} - a_{n-1})\underline{q}_S^T \mathbf{A}^{n-1}$$

6.2.2 Théorème de Cayley-Hamilton

Toute matrice carrée \mathbf{A} satisfait sa propre équation caractéristique.

Si

$$P(\lambda) = \det(\lambda\mathbf{I} - \mathbf{A}) = \lambda^n + \alpha_{n-1}\lambda^{n-1} + \dots + \alpha_1\lambda + \alpha_0 = 0$$

alors

$$P(\mathbf{A}) = \det(\mathbf{A}\mathbf{I} - \mathbf{A}) = \mathbf{A}^n + \alpha_{n-1}\mathbf{A}^{n-1} + \dots + \alpha_1\mathbf{A} + \alpha_0\mathbf{I} = 0$$

En appliquant ce théorème aux résultats précédents, \mathbf{A} et \mathbf{A}_{FC} sont semblables et ont la même équation caractéristique, donc \mathbf{A} doit satisfaire l'équation caractéristique de \mathbf{A}_{FC} .

Donc

$$\mathbf{A}^n + a_{n-1}\mathbf{A}^{n-1} + \dots + a_1\mathbf{A} + a_0\mathbf{I} = 0$$

$$\underline{q}_S^T \mathbf{A}^n = -a_{n-1}\underline{q}_S^T \mathbf{A}^{n-1} - \dots - a_1\underline{q}_S^T \mathbf{A} - a_0\underline{q}_S^T$$

d'où, par substitution : Théorème d'Ackermann

$$\mathbf{L} = \underline{q}_S^T \mathbf{A}^n + \alpha_{n-1}\underline{q}_S^T \mathbf{A}^{n-1} + \dots + \alpha_1\underline{q}_S^T \mathbf{A} + \alpha_0\underline{q}_S^T$$

$$\mathbf{L} = \underline{q}_S^T P(\mathbf{A})$$

Théorème : Pour un système à une entrée et une sortie les zéros du système restent inchangés par le placement des pôles

6.3 Commande Modale

6.3.1 Définition

Soit la transformation $\underline{x} = \mathbf{T}\underline{x}^* \mathbf{T}$ est la transformation modale qui diagonalise \mathbf{A} si toutes les valeurs propres sont différentes ou du moins la met sous forme de blocs de Jordan si toutes les valeurs propres ne sont pas différentes.

Dans la nouvelle base, les équations sont

$$\dot{\underline{x}}_i^* = \lambda_i \underline{x}_i^* + u_i^*$$

si le système est diagonalisable.

Si les u_{d_i} sont indépendantes les une des autres, on peut asservir chaque variable d'état \underline{x}_{d_i} , et donc chaque mode du système.

6.3.2 Méthode de synthèse

Si l'on dispose de p grandeurs de commande indépendantes, on peut asservir p coordonnées modales.

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}\underline{u} \Rightarrow \dot{\underline{x}}^* = \mathbf{\Lambda}\underline{x}^* + \widehat{\mathbf{B}}\underline{u}$$

avec :

$$\mathbf{\Lambda} = \mathbf{T}^{-1}\mathbf{A}\mathbf{T} = \begin{pmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{pmatrix}$$

$$\widehat{\mathbf{B}} = \mathbf{T}^{-1}\mathbf{B}$$

Choix des p coordonnées ramenées en contre réaction :

- les plus proches de l'axe imaginaire (système plus rapide),
- celles qui sont le plus influencées par une perturbation (système moins sensible aux perturbations)
- ...

le système peut être mis sous la forme :

$$\begin{pmatrix} \dot{\underline{x}}_p^* \\ \dot{\underline{x}}_{n-p}^* \end{pmatrix} = \begin{pmatrix} \mathbf{\Lambda}_p & 0 \\ 0 & \mathbf{\Lambda}_{n-p} \end{pmatrix} \begin{pmatrix} \underline{x}_p^* \\ \underline{x}_{n-p}^* \end{pmatrix} + \begin{pmatrix} \widehat{\mathbf{B}}_p \\ \widehat{\mathbf{B}}_{n-p} \end{pmatrix} \underline{u}$$

avec :

$\mathbf{\Lambda}_p$: sous matrice ($p \times p$) de $\mathbf{\Lambda}$

$\widehat{\mathbf{B}}_p$: sous matrice ($p \times p$) de $\widehat{\mathbf{B}}$

Principe : déterminer la contre réaction de ces p coordonnées de telle manière que :

1. les valeurs propres du système λ_1 à λ_p soient remplacées par les valeurs propres souhaitées λ_1^* à λ_p^*
2. le système bouclé reste diagonal (en ces p variables, donc chaque mode est régulé indépendamment)

donc le système bouclé doit avoir, en mode de régulation pour les p coordonnées concernées une équation de la forme :

$$\dot{\underline{x}}_p^* = \mathbf{\Lambda}_p^* \underline{x}_p^*$$

avec $\mathbf{\Lambda}_p^* = \begin{pmatrix} \lambda_1^* & & 0 \\ & \ddots & \\ 0 & & \lambda_p^* \end{pmatrix}$

donc

$$\mathbf{\Lambda}_p \underline{x}_p^* + \underline{u}^* = \mathbf{\Lambda}_p^* \underline{x}_p^*$$

donc

$$\underline{u}^* = -(\mathbf{\Lambda}_p - \mathbf{\Lambda}_p^*) \underline{x}_p^* = - \begin{pmatrix} \lambda_1 - \lambda_1^* & & 0 \\ & \ddots & \\ 0 & & \lambda_p - \lambda_p^* \end{pmatrix} \begin{pmatrix} \underline{x}_1^* \\ \vdots \\ \underline{x}_p^* \end{pmatrix}$$

Problème : revenir à \underline{u} notre véritable vecteur de commande

\mathbf{B} ($n \times p$) est de rang p ; \mathbf{T}^{-1} est régulière donc : $\widehat{\mathbf{B}}(n \times p)$ est de rang p donc $\widehat{\mathbf{B}}_p(p \times p)$ est régulière

$$\underline{u} = \widehat{\mathbf{B}}_p^{-1} \underline{u}^* = -\widehat{\mathbf{B}}_p^{-1} (\mathbf{\Lambda}_p - \mathbf{\Lambda}_p^*) \underline{u}^*$$

d'où

$$\dot{\underline{x}}_p^* = \mathbf{\Lambda}_p^* \underline{x}_p^* \quad (6.5)$$

$$\dot{\underline{x}}_{n-p}^* = -\widehat{\mathbf{B}}_{n-p} \widehat{\mathbf{B}}_p^{-1} (\mathbf{\Lambda}_p - \mathbf{\Lambda}_p^*) \underline{x}_p^* + \mathbf{\Lambda}_{n-p} \underline{x}_{n-p}^* \quad (6.6)$$

- l'équation (6.5) montre que les p premières coordonnées modales sont bien asservies individuellement ; les valeurs propres sont bien celles choisies (λ_i^*),
- l'équation (6.6) montre que les $n - p$ autres coordonnées modales sont influencées par les p premières.

Vérification de cette deuxième constatation. Calculons les valeurs propres du système bouclé :

$$\det \begin{pmatrix} p\mathbf{I}_p - \mathbf{\Lambda}_p^* & 0 \\ \widehat{\mathbf{B}}_{n-p} \widehat{\mathbf{B}}_p^{-1} (\mathbf{\Lambda}_p - \mathbf{\Lambda}_p^*) & p\mathbf{I}_p - \mathbf{\Lambda}_{n-p} \end{pmatrix} = 0$$

$$\det(p\mathbf{I}_p - \mathbf{\Lambda}_p^*) \times \det(p\mathbf{I}_p - \mathbf{\Lambda}_{n-p}) = 0$$

Les valeurs propres du système sont donc :

$\lambda_1^*, \dots, \lambda_p^*$: les p premières qui ont été déplacées par la commande modale,

$\lambda_{p+1}, \dots, \lambda_n$: les $n - p$ autres qui restent inchangées.

La commande modale ne présente pas d'effet indésirable du fait du couplage entre les $n - p$ coordonnées modales restantes et les p premières.

Synthèse du régulateur modal

Retour à la représentation de départ :

$$\underline{x}^* = \mathbf{T}^{-1} \underline{x}$$

$$\underline{u} = -\widehat{\mathbf{B}}_p^{-1} (\mathbf{\Lambda}_p - \mathbf{\Lambda}_p^*) \underline{x}^* = -\widehat{\mathbf{B}}_p^{-1} (\mathbf{\Lambda}_p - \mathbf{\Lambda}_p^*) \mathbf{T}^{-1} \underline{x} = -\mathbf{L} \underline{x}$$

$$\mathbf{L} = \widehat{\mathbf{B}}_p^{-1} \begin{pmatrix} \lambda_1 - \lambda_1^* & & 0 \\ & \ddots & \\ 0 & & \lambda_p - \lambda_p^* \end{pmatrix} \mathbf{T}^{-1} \underline{x}$$

6.4 Choix des pôles

Le choix de pôles en boucle fermée tient plus de l'art que de la science. En effet, les phénomènes à prendre en compte sont nombreux, très dépendants du système considéré et des performances désirées.

Voici quelques règles fondamentales à respecter :

1. Les pôles choisis doivent être stables,
2. pas trop proches de l'axe des imaginaires, sinon la moindre variation de modèle peut provoquer une instabilité,
3. Les pôles complexes conjugués seront choisis pour présenter un dépassement convenable (typiquement : inférieur à 20 %) sinon le régime transitoire sera long
4. pas trop rapides (typiquement : 4 à 10 fois plus rapides que les pôles en B0), il est peu probable que le modèle que vous avez soit encore valable au-delà de ce domaine et/ou que la commande ne sature pas.

La figure 6.2 résume ces quelques règles.

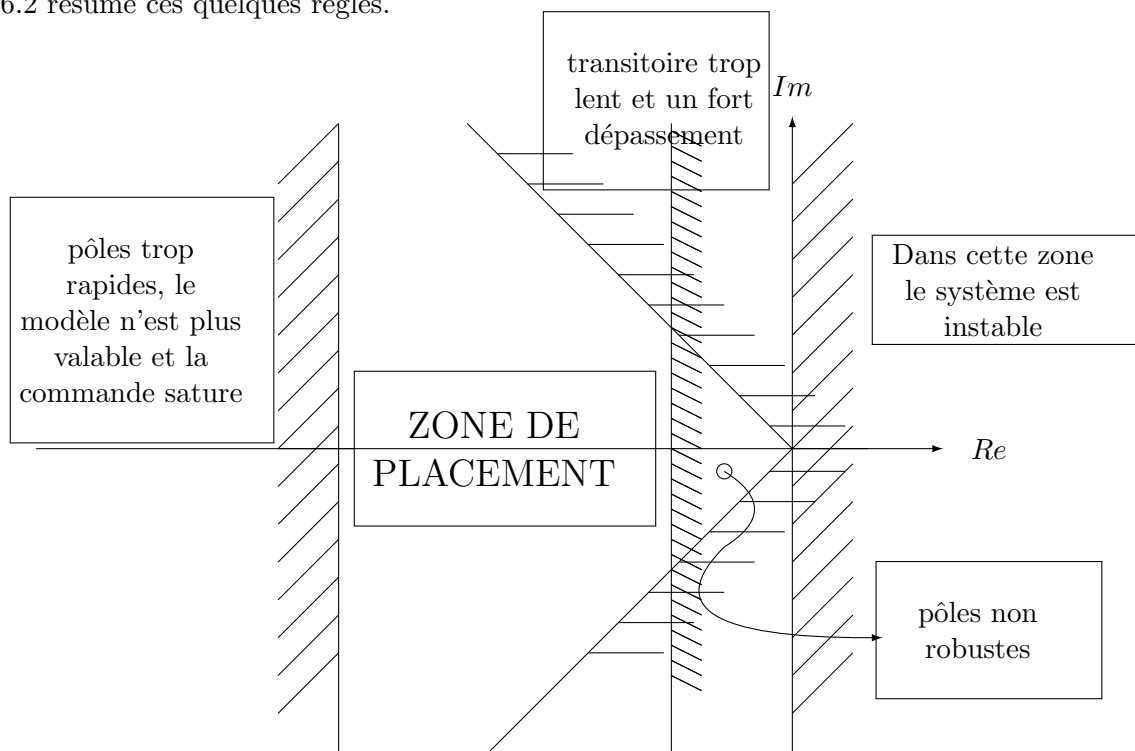
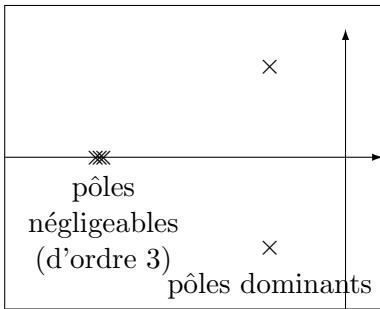


FIGURE 6.2 – Règles de placement de pôles

Le choix des pôles au sein de la zone de placement est tout autant un art et les méthodes sont nombreuses. Afin de réduire vos maux de tête lors du choix je vous propose 3 méthodes ayant fait leurs preuves, mais sachez qu'il en existe d'autres.

Les exemples sont donnés pour un système d'ordre 5.

6.4.1 Pôles complexes conjugués dominants



Le grand avantage de cette méthode est une simplification des calculs et une bonne maîtrise du comportement en boucle fermée du système (comportement comme les pôles dominants). Les résultats en boucle fermée sont quelques fois surprenants si les pôles négligés ne l'étaient pas vraiment. Bien sûr si le comportement souhaité est de type premier ordre, il suffit de choisir un pôle simple comme pôle dominant.

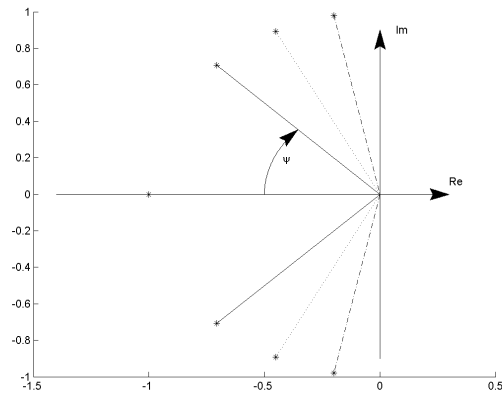
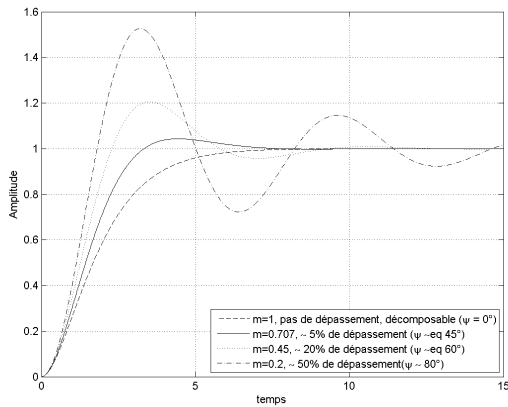


FIGURE 6.3 – Relation entre dépassement et position des pôles dans le plan complexe

Quelques relations pour la détermination des pôles :

Pour $H(p) = \frac{1}{1 + \frac{2m}{\omega_0}p + \frac{p^2}{\omega_0^2}}$

Temps de réponse à 5% :

$$T_{r5\%} = \frac{3}{m\omega_0}$$

Temps de montée (10 à 90% de la valeur finale) :

$$T_m = \frac{\pi}{2\omega_0\sqrt{1-m^2}}$$

Premier dépassement

$$D1 = 100e^{-\frac{m\pi}{\sqrt{1-m^2}}}$$

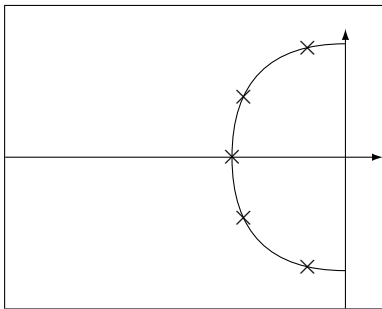
la partie réelle des pôles négligeables sera choisie égale entre 3 et 10 plus grande que la partie réelle des pôles dominants soit :

$$P_{neg} = 3 \text{ à } 10 \times m\omega_0$$

TABLE 6.1 – Polynômes normalisés de Butterworth à l'ordre n (remplacer p par p/ω_0)

n	Polynôme	$D\%$	$\omega_0 T(D\%)$	$\omega_0 T_{r5\%}$
1	$(p + 1)$	0		3
2	$(p^2 + 1.4142p + 1)$	4.3	4.5	2.8
3	$(p + 1)(p^2 + p + 1)$	8.1	4.9	5.9
4	$(p^2 + 0.7654p + 1)(p^2 + 1.8478p + 1)$	10.8	5.5	6.7
5	$(p + 1)(p^2 + 0.6180p + 1)(p^2 + 1.6180p + 1)$	12.7	6.3	7.5
6	$(p^2 + 0.5176p + 1)(p^2 + 1.4142p + 1)(p^2 + 1.9319p + 1)$	14.1	6.9	10.7

6.4.2 Maximalement plat



Premier dépassement

$$D1 = -0.20n^2 + 4n - 2.35$$

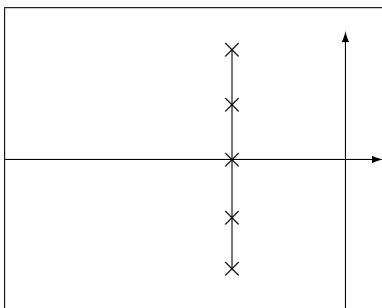
Temps du premier maximum :

$$T_{D1} = \frac{0.016n^2 + 0.52n + 3.3}{\omega_0}$$

les pôles sont donnés par :

$$\text{pour } 0 \leq k \leq n-1, \quad p_k = \omega_0 e^{\frac{j(2k+n-1)\pi}{2n}}$$

6.4.3 Pôles à partie réelle identique



Cette méthode, proche de la méthode des pôles dominants, consiste à équirépartir l'ensemble des pôles sur la même verticale. En l'absence de zéros influents dans le système, la réponse à l'échelon est toujours apériodique. Les pôles "négligeables" atténuent le dépassement des pôles "dominants". Utile en cas de bruit, car les pôles "négligeables" (qui ne le sont plus donc !) filtrent fortement ces bruits.

6.4.4 Polynômes de Naslin

C'est encore dans les vieux pots que l'on fait les meilleures soupes. Cette méthode perdure depuis les années 60 et donne de bons résultats quel que soit le domaine d'application.

Rappels : On choisit α et ω_0 en fonction des caractéristiques voulues pour la boucle fermée.

Temps de montée (10 à 90% de la valeur finale) :

$$T_m = \frac{2.2}{\omega_0}$$

Premier dépassement

$$D1\% = 10^{4.8-2\alpha}$$

les coefficients du polynôme sont alors :

$$\text{pour } 0 \leq k \leq n \quad a_k = \alpha^{\frac{-k(k-1)}{2}} \omega_0^{-k}$$

En cas de zéro dans la fonctions de transfert, cette méthode permet de "précompenser" l'influence de ce zéro. Si le numérateur de la fonction de transfert est de la forme $b_0 + b_1 p$, alors $\omega'_0 = \frac{b_0}{b_1}$. Il suffit alors de reprendre les formules précédentes en remplaçant α par α_e et ω_0 par ω_{0c} .

$$\alpha_e = 1.5 + \frac{\omega'_0}{4\omega_0}(\alpha - 1.5)$$

$$\omega_{0c} = \left(\frac{1}{\omega_0} - \frac{1}{2\omega'_0} \right)^{-1}$$

Si ω_{0c} est négatif, seule la méthode "try and error" pour le choix de α et ω_0 reste valable!

6.5 Commande optimale

Cette partie n'est donnée qu'à titre de culture générale. La commande optimale est une base fondamentale de l'automatique, pour continuer dans cette voie, une autoformation est nécessaire !

6.5.1 Définition

L'objectif de cette commande est de minimiser un critère quadratique de la forme :

$$J = \int_0^{\infty} (\underline{x}^T \mathbf{Q} \underline{x} + \underline{u}^T \mathbf{R} \underline{u}) dt \quad (6.7)$$

avec : \mathbf{Q} = matrice symétrique définie positive

\mathbf{R} = matrice symétrique non négative

la commande u est alors définie par l'équation suivante :

$$u = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{K} \underline{x}$$

où \mathbf{K} est une matrice symétrique, solution définie négative de l'équation de Riccati :

$$\mathbf{K} \mathbf{A} + \mathbf{A}^T \mathbf{K} - \mathbf{K} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{K} + \mathbf{Q} = 0$$

6.5.2 Stabilité de la commande optimale

Les commandes optimales à critère quadratique conduisent toujours à un système stable en boucle fermée. La marge de phase est toujours supérieure à 60 degrés.

6.5.3 Choix des matrices \mathbf{R} et \mathbf{Q}

La forme du critère (eq. (6.7)) représente typiquement un critère énergétique. Aussi la commande optimale vise-t'elle à minimiser l'énergie requise pour faire passer le système d'un état à un autre. Dans ce cas les matrices \mathbf{R} et \mathbf{Q} seront le plus souvent choisies diagonales, les coefficients de la diagonale ayant une signification physique (inductance, masse ...).

Plus généralement, cette méthode de synthèse s'applique lorsque l'on désire minimiser une ou plusieurs grandeurs de commande et/ou un ou plusieurs états pour des problèmes de saturation, de sécurité ou de durée de vie du système (limitation de la vitesse d'un moteur, limitation de la pression dans un réservoir ...).

6.5.4 Exemple : fil rouge.

Calcul d'une commande qui minimise les pertes énergétiques dans le système. Les pertes sont :

- $\frac{1}{2}Ri^2$: pertes Joule dans la résistance d'inducteur
- $\frac{1}{2}f_v\omega^2$: pertes mécaniques par frottements visqueux

Le critère est donc

$$J = \int_0^{\infty} (Ri^2 + f_v\omega^2)dt$$

sous forme matricielle

$$J = \int_0^{\infty} (\underline{x}^T \begin{pmatrix} 0 & 0 \\ 0 & f_v \end{pmatrix} \underline{x} + \underline{u}^T (R)\underline{u})dt \quad (6.8)$$

En posant (\mathbf{K} symétrique)

$$\mathbf{K} = \begin{pmatrix} k_{11} & k_{12} \\ k_{12} & k_{22} \end{pmatrix},$$

le calcul de la commande

$$u = - \left(\frac{k_{12}b_0}{r} \quad \frac{k_{22}b_0}{r} \right) \underline{x}$$

montre que seuls les coefficients k_{12} et k_{22} sont à déterminer. En développant l'équation de Riccati, nous obtenons :

$$\begin{bmatrix} -2k_{12}a_0 - \frac{k_{12}^2b_0^2}{r} & k_{11} - k_{12}a_1 - k_{22}a_0 - \frac{k_{12}b_0^2k_{22}}{r} \\ k_{11} - k_{12}a_1 - k_{22}a_0 - \frac{k_{12}b_0^2k_{22}}{r} & 2k_{12} - 2k_{22}a_1 - \frac{k_{22}^2b_0^2}{r} + f_v \end{bmatrix} = 0$$

Nous en déduisons que :

$$\begin{aligned} k_{12} &= 0 \\ k_{22} &= 1/2 \frac{-2a_1r + 2\sqrt{a_1^2r^2 + b_0^2f_vr}}{b_0^2} \\ \text{ou} \\ k_{22} &= 1/2 \frac{-2a_1r - 2\sqrt{a_1^2r^2 + b_0^2f_vr}}{b_0^2} \end{aligned}$$

comme \mathbf{K} est définie négative, la solution est :

$$\begin{aligned} k_{12} &= 0 \\ k_{22} &= 1/2 \frac{-2a_1r - 2\sqrt{a_1^2r^2 + b_0^2f_vr}}{b_0^2} \end{aligned}$$

donc la commande est :

$$u = - \left(0 \quad \frac{b_0}{2r} \frac{-2a_1r - 2\sqrt{a_1^2r^2 + b_0^2f_vr}}{b_0^2} \right) \underline{x}$$

Chapitre 7

Synthèse d'observateurs d'état

7.1 Introduction au problème de la reconstruction d'état

Dans tout ce qui précède, nous sommes partis du principe que nous avons accès à toutes les composantes du vecteur d'état. Nous avons donc supposé que le système est complètement instrumenté. En réalité, les systèmes physiques sont très peu instrumentés, les raisons sont :

- le coût,
- la difficulté d'accéder à certaines variables
- la fiabilité,
- l'encombrement, ...

Nous allons donc voir comment "reconstruire" l'état à partir des commandes appliquées au système réel et les quelques mesures du vecteur d'état effectuées.

Soit le système,

$$\begin{cases} \dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}\underline{u} \\ \underline{y} = \mathbf{C}\underline{x} + \mathbf{D}\underline{u} \end{cases} \quad \text{ou} \quad \begin{cases} \dot{\underline{x}}_{k+1} = \mathbf{\Phi}\underline{x}_k + \mathbf{\Gamma}\underline{u}_k \\ \underline{y}_k = \mathbf{C}\underline{x}_k + \mathbf{D}\underline{u}_k \end{cases} \quad (7.1)$$

Comment obtenir \underline{x} ?

7.1.1 Par calcul direct

Avec cette méthode \underline{x} est obtenu par la formule suivante :

$$\underline{x} \stackrel{?}{=} \mathbf{C}^{-1}(\underline{y} - \mathbf{D}\underline{u})$$

Ce calcul est impossible car en général \mathbf{C} n'est pas inversible, dans les systèmes SISO, \mathbf{C} est un vecteur, donc jamais inversible.

7.1.2 Par simulation du processus

L'idée consiste à dire que puisque nous possédons un modèle du système, il suffit de simuler le processus dans le ordinateur. Ainsi, les différentes variables comme le vecteur d'état sont parfaitement accessibles. Cette méthode est illustrée par le schéma 7.1.

Les indices r (réelle) et m (mesurée) sont là pour rappeler qu'il existe toujours une différence entre la réalité et le modèle utilisé. Par la suite, comme dans tout ce qui précède, nous ne ferons aucune différence entre les deux mais ayez toujours à l'esprit qu'elle existe.

En fait cette méthode ne fonctionne pas du tout car le modèle, pour précis qu'il soit, est toujours faux (imprécisions de mesure des coefficients des matrices, non linéarités du système, ...). Aussi après quelques temps de simulation le simulateur donne n'importe quoi. Néanmoins dans quelques cas, ce sera la seule solution, pour obtenir une variable non mesurée, on parle alors de reconstruction par simulation. La robustesse du processus corrigé est alors faible et la synthèse de la commande doit prendre en compte cet aspect (commande robuste).

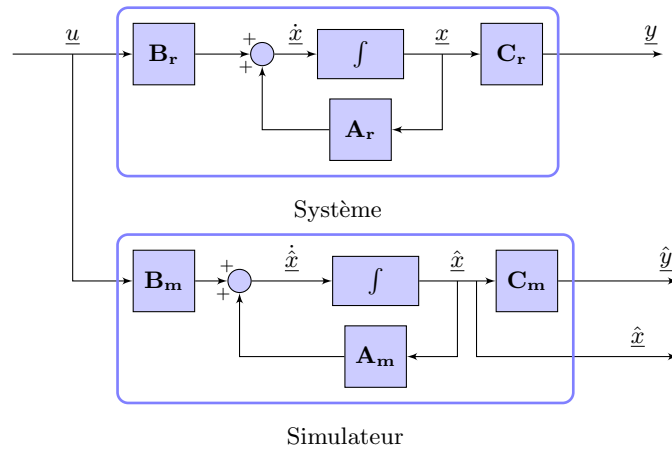


FIGURE 7.1 – Schéma général d'un simulateur

7.1.3 Par simulation du processus et asservissement sur les parties connues du vecteur d'état

L'idée consiste comme précédemment à simuler le système, mais à l'asservir de façon à ce que les sorties mesurées concordent avec les sorties simulées, en injectant à l'entrée l'erreur d'estimation pondérée par un gain.

C'est exactement ce que l'on appelle un observateur de Luenberger.

7.2 Observateurs de Luenberger

Définition

On appelle observateur (de Luenberger) du système,

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}\underline{u} \quad (7.2)$$

$$\underline{y} = \mathbf{C}\underline{x} \quad (7.3)$$

un système qui est décrit par le schéma fonctionnel 7.2,

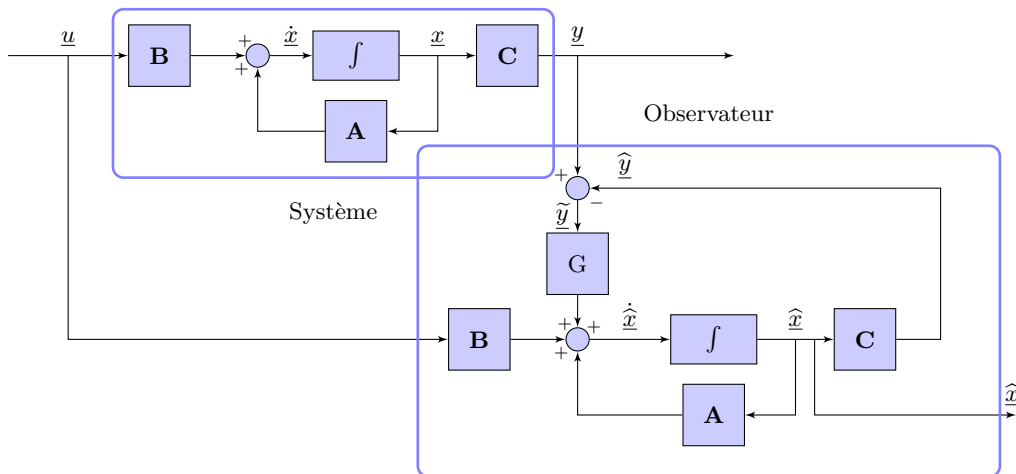


FIGURE 7.2 – Schéma général d'un observateur de Luenberger complet.

l'équation différentielle d'état de l'observateur est :

$$\begin{aligned}
 \dot{\hat{\underline{x}}} &= \mathbf{A}\hat{\underline{x}} + \mathbf{B}\underline{u} + \mathbf{G}\tilde{\underline{y}} \\
 &= \mathbf{A}\hat{\underline{x}} + \mathbf{B}\underline{u} + \mathbf{G}\underline{y} - \mathbf{G}\hat{\underline{y}} \\
 &= \mathbf{A}\hat{\underline{x}} + \mathbf{B}\underline{u} + \mathbf{G}\underline{y} - \mathbf{G}\mathbf{C}\hat{\underline{x}} \\
 &= (\mathbf{A} - \mathbf{G}\mathbf{C})\hat{\underline{x}} + \mathbf{B}\underline{u} + \mathbf{G}\underline{y}
 \end{aligned}$$

En définissant

$$\mathbf{A} - \mathbf{G}\mathbf{C} = \mathbf{F}$$

$$\dot{\hat{\underline{x}}} = \mathbf{F}\hat{\underline{x}} + \mathbf{G}\underline{y} + \mathbf{B}\underline{u} \quad (7.4)$$

où les valeurs propres de la matrice \mathbf{F} sont stables.

$\hat{\underline{x}}(t)$ est une estimation de $\underline{x}(t)$,

on définit $\tilde{\underline{x}} = \underline{x} - \hat{\underline{x}}$ l'erreur d'estimation et $\tilde{\underline{x}} \rightarrow 0$ pour $t \rightarrow \infty$.

Il ne reste plus qu'à définir la matrice de "retour" \mathbf{G} .

Dualité

$$\begin{aligned}
\text{synthèse d'un régulateur d'état} &\longrightarrow \text{synthèse d'un observateur} \\
\mathbf{A} &\longrightarrow \mathbf{A}^T \\
\mathbf{B} &\longrightarrow \mathbf{C}^T \\
\mathbf{L} &\longrightarrow \mathbf{G}^T \\
\det [p\mathbf{I} - (\mathbf{A} - \mathbf{B}\mathbf{R})] &\longrightarrow \det [p\mathbf{I} - (\mathbf{A}^T - \mathbf{C}^T\mathbf{G}^T)] \\
[\mathbf{B}, \mathbf{A}\mathbf{B}, \mathbf{A}^2\mathbf{B}, \dots, \mathbf{A}^{n-1}\mathbf{B}] &\longrightarrow \begin{bmatrix} \mathbf{C} \\ \mathbf{C}\mathbf{A} \\ \mathbf{C}\mathbf{A}^2 \\ \vdots \\ \mathbf{C}\mathbf{A}^{n-1} \end{bmatrix}^T
\end{aligned}$$

Toutes les méthodes de calcul d'un régulateur sont valables pour le calcul de la matrice \mathbf{G}^T .

7.3 Observateurs d'ordre réduit

Il est parfois intéressant de n'observer que la partie de l'état qui n'est pas accessible afin de minimiser le temps de calcul. On définit alors un observateur d'ordre réduit.

7.3.1 Hypothèses

soit le système :

$$\begin{aligned}
\dot{\underline{x}} &= \mathbf{A}\underline{x} + \mathbf{B}u \\
y &= \mathbf{C}\underline{x}
\end{aligned} \tag{7.5}$$

avec :

$$\mathbf{C} = \left[\begin{array}{cccccc} 0 & \dots & 0 & 1 & \dots & 0 \\ \vdots & & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & \dots & 1 \end{array} \right] \left. \vphantom{\begin{array}{cccccc} 0 & \dots & 0 & 1 & \dots & 0 \\ \vdots & & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & \dots & 1 \end{array}} \right\} q = [0 \ \mathbf{I}_q] \tag{7.6}$$

$$\underbrace{\hspace{10em}}_{n-q} \quad \underbrace{\hspace{10em}}_q \tag{7.7}$$

Bien entendu, il est probable qu'il faille procéder à une transformation de coordonnées pour mettre \mathbf{C} sous la forme ci-dessus.

Partitionnement

$$\begin{aligned}
\underline{x} &= \begin{bmatrix} x_1 \\ \vdots \\ x_{n-q} \\ \dots \\ x_{n-q+1} \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \underline{v} \\ \dots \\ \underline{y} \end{bmatrix} \left. \vphantom{\begin{bmatrix} x_1 \\ \vdots \\ x_{n-q} \\ \dots \\ x_{n-q+1} \\ \vdots \\ x_n \end{bmatrix}} \right\} \begin{array}{l} n-q \\ q \end{array} \\
\dot{\underline{x}} &= \begin{bmatrix} \dot{\underline{v}} \\ \dots \\ \dot{\underline{y}} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \dots & \dots \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \underline{v} \\ \dots \\ \underline{y} \end{bmatrix} + \begin{bmatrix} \mathbf{B}_1 \\ \dots \\ \mathbf{B}_2 \end{bmatrix} u \left. \vphantom{\begin{bmatrix} \dot{\underline{v}} \\ \dots \\ \dot{\underline{y}} \end{bmatrix}} \right\} \begin{array}{l} n-q \\ q \end{array}
\end{aligned}$$

soit :

$$\dot{\underline{v}} = \mathbf{A}_{11}\underline{v} + \mathbf{A}_{12}\underline{y} + \mathbf{B}_1\underline{u} \quad (7.8)$$

$$\dot{\underline{y}} = \mathbf{A}_{21}\underline{v} + \mathbf{A}_{22}\underline{y} + \mathbf{B}_2\underline{u} \quad (7.9)$$

Il reste à estimer \underline{v} , ($n - q$ composantes) uniquement.

Idée de base : Considérer ces deux équations comme l'équation différentielle d'état et l'équation de sortie d'un système qui aurait

- \underline{v} : comme vecteur d'état
- $\mathbf{A}_{12}\underline{y} + \mathbf{B}_1\underline{u}$: comme vecteur d'entrée
- $\dot{\underline{y}} - \mathbf{A}_{22}\underline{y} - \mathbf{B}_2\underline{u}$: comme vecteur de sortie

$$\underbrace{\dot{\underline{v}}}_{\dot{X}} = \underbrace{\mathbf{A}_{11}}_A \underbrace{\underline{v}}_X + \underbrace{\mathbf{A}_{12}\underline{y} + \mathbf{B}_1\underline{u}}_{BU} \quad (7.10)$$

$$\underbrace{\dot{\underline{y}} - \mathbf{A}_{22}\underline{y} - \mathbf{B}_2\underline{u}}_Y = \underbrace{\mathbf{A}_{21}}_C \underbrace{\underline{v}}_X \quad (7.11)$$

Développons un observateur complet (d'ordre $n - q$) pour ce système

$$\dot{\hat{\underline{v}}} = (\mathbf{A}_{11} - \mathbf{G}\mathbf{A}_{21})\hat{\underline{v}} + \mathbf{A}_{12}\underline{y} + \mathbf{B}_1\underline{u} + \mathbf{G}(\dot{\underline{y}} - \mathbf{A}_{22}\underline{y} - \mathbf{B}_2\underline{u})$$

Posons : $\underline{z} = \hat{\underline{v}} - \mathbf{G}\underline{y}$

D'où les équations d'état de l'observateur :

$$\dot{\underline{z}} = (\mathbf{A}_{11} - \mathbf{G}\mathbf{A}_{21})\underline{z} + [(\mathbf{A}_{11} - \mathbf{G}\mathbf{A}_{21})\mathbf{G}\mathbf{A}_{12} - \mathbf{G}\mathbf{A}_{22}]\underline{y}(\mathbf{B}_1 - \mathbf{G}\mathbf{B}_2)\underline{u} \quad (7.12)$$

$$\hat{\underline{v}} = \underline{z} + \mathbf{G}\underline{y} \quad (7.13)$$

dont le schéma fonctionnel est donné en figure 7.3

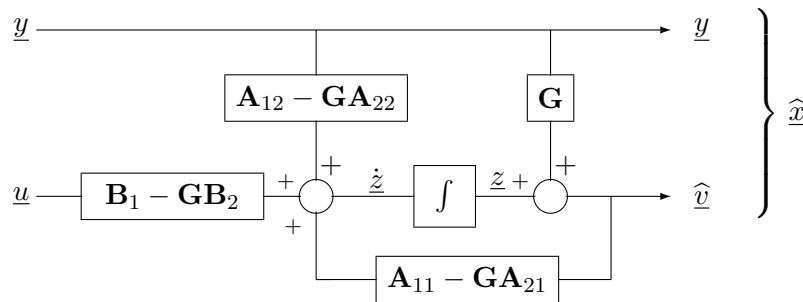


FIGURE 7.3 – Schéma fonctionnel d'un observateur de Luenberger d'ordre réduit

Remarque : La construction précédente fournit un observateur de ce type, avec la matrice de système

$$\mathbf{F} = \mathbf{A}_{11} - \mathbf{G}\mathbf{A}_{21}$$

7.4 Observateur généralisé

Pour observer un système linéaire invariant, on introduit un second système linéaire invariant, d'ordre r et ayant l'équation d'état suivante :

$$\dot{\underline{z}}(t) = \mathbf{F}\underline{z}(t) + \mathbf{G}\underline{y}(t) + \mathbf{E}\underline{u}(t)$$

On dit qu'un tel système est un observateur du premier système si, étant donné une matrice de transformation arbitraire \mathbf{K} de dimension $(r \times n)$, son vecteur d'état représente une estimation de $\mathbf{K}\underline{x}$:

$$\underline{z} = \mathbf{K}\hat{\underline{x}}$$

c'est-à-dire que $\underline{z} \rightarrow \mathbf{K}\hat{\underline{x}}$ pour $t \rightarrow \infty$

\underline{z} doit donc être solution d'une équation homogène, de la forme

$$(\underline{z} - \mathbf{K}\underline{x})' = \mathbf{F}(\underline{z} - \mathbf{K}\underline{x})$$

$$\dot{\underline{z}} = \mathbf{F}\underline{z} + \mathbf{K}\dot{\underline{x}} - \mathbf{F}\mathbf{K}\underline{x}$$

et

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}\underline{u}$$

donc

$$\dot{\underline{z}} = \mathbf{F}\underline{z} + (\mathbf{K}\mathbf{A} - \mathbf{F}\mathbf{K})\hat{\underline{x}} - \mathbf{K}\mathbf{B}\underline{u}$$

d'où l'on en déduit que

$$\mathbf{E} = \mathbf{K}\mathbf{B}$$

Cherchons à introduire $\underline{y} = \mathbf{C}\underline{x}$, on cherche à trouver une matrice \mathbf{G} telle que

$$(\mathbf{K}\mathbf{A} - \mathbf{F}\mathbf{K})\underline{x} = \mathbf{G}\underline{y} = \mathbf{G}\mathbf{C}\underline{x}$$

donc

$$(\mathbf{K}\mathbf{A} - \mathbf{F}\mathbf{K}) = \mathbf{G}\mathbf{C}$$

d'où l'équation d'état de l'observateur généralisé

$$\dot{\underline{z}} = \mathbf{F}\underline{z} + \mathbf{G}\underline{y} + \mathbf{K}\mathbf{B}\underline{u}$$

7.5 Equation d'état d'un système asservi avec observateur

Loi de commande

$$\underline{u} = -\mathbf{L}\hat{\underline{x}}$$

à $\underline{w} = 0$, étude de la stabilité

Vu du point de vue de l'observateur généralisé, \underline{u} doit dépendre de la grandeur de sortie \underline{y} du procédé et du vecteur d'état \underline{z} de l'observateur :

$$\underline{u} = \mathbf{D}\hat{\underline{y}} + \mathbf{E}\hat{\underline{z}}$$

avec $\underline{z} = -\mathbf{K}\hat{\underline{x}}$ et $\underline{y} = \mathbf{C}\underline{x}$ il vient

$$-\mathbf{L}\hat{\underline{x}} = \mathbf{D}\mathbf{C}\underline{x} + \mathbf{E}\mathbf{K}\hat{\underline{x}}$$

cette expression doit être valable quel que soit t , donc aussi pour $t \rightarrow \infty$ ($\hat{\underline{x}} = \underline{x}$) donc,

$$-\mathbf{L} = \mathbf{D}\mathbf{C} + \mathbf{E}\mathbf{K}$$

donc, l'équation d'état de l'ensemble observateur+ régulateur est

$$\begin{aligned} \dot{\underline{z}} &= \mathbf{F}\underline{z} + \mathbf{G}\underline{y} + \mathbf{K}\mathbf{B}\underline{u} \\ \underline{u} &= \mathbf{D}\underline{y} + \mathbf{E}\underline{z} \end{aligned} \quad (7.14)$$

schéma fonctionnel du système asservi

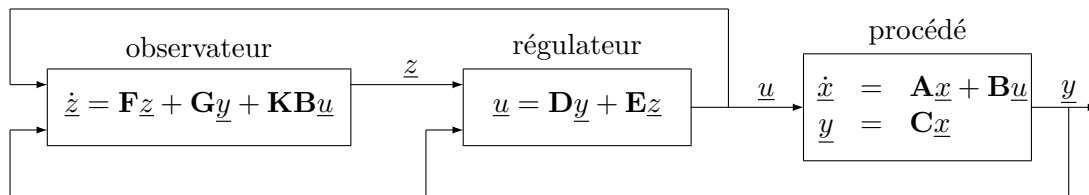


FIGURE 7.4 – Schéma fonctionnel d'un système asservi via un observateur

vecteur d'état de ce système composite : $\begin{bmatrix} \underline{x} \\ \underline{z} \end{bmatrix}$

7.5.1 Théorème de séparation

posons : $\underline{\xi} = \underline{z} - \mathbf{K}\underline{x}$

pour l'observateur :

$$\begin{aligned} \dot{\underline{\xi}} &= \dot{\underline{z}} - \mathbf{K}\dot{\underline{x}} \\ &= \mathbf{F}\underline{z} + \mathbf{G}\underline{y} + \mathbf{K}\mathbf{B}\underline{u} - \mathbf{K}\mathbf{A}\underline{x} - \mathbf{K}\mathbf{B}\underline{u} \\ &= \mathbf{F}(\underline{z} - \mathbf{K}\underline{x}) \quad \text{car } (\mathbf{K}\mathbf{A} - \mathbf{F}\mathbf{K}) = \mathbf{G}\mathbf{C} \\ &= \mathbf{K}\underline{\xi} \end{aligned} \quad (7.15)$$

pour le procédé :

$$\begin{aligned} \dot{\underline{x}} &= \mathbf{A}\underline{x} + \mathbf{B}\underline{u} \\ &= \mathbf{A}\underline{x} + \mathbf{B}(\mathbf{D}\underline{y} + \mathbf{E}\underline{z}) \\ &= \mathbf{A}\underline{x} + \mathbf{B}\mathbf{D}\mathbf{C}\underline{x} + \mathbf{B}\mathbf{E}\underline{z} \\ &= \mathbf{A}\underline{x} - \mathbf{B}\mathbf{L}\underline{x} - \mathbf{B}\mathbf{E}\mathbf{K}\underline{x} + \mathbf{B}\mathbf{E}\underline{z} \quad \text{car } (\mathbf{D}\mathbf{C} - \mathbf{E}\mathbf{K}) = -\mathbf{L} \\ &= (\mathbf{A} - \mathbf{B}\mathbf{L})\underline{x} + \mathbf{B}\mathbf{E}\underline{\xi} \end{aligned} \quad (7.16)$$

Des deux développements précédents nous en déduisons l'équation d'état du système composite observateur + procédé.

$$\begin{pmatrix} \dot{\underline{x}} \\ \dot{\underline{\xi}} \end{pmatrix} = \begin{pmatrix} \mathbf{A} - \mathbf{BL} & \mathbf{BE} \\ 0 & \mathbf{F} \end{pmatrix} \begin{pmatrix} \underline{x} \\ \underline{\xi} \end{pmatrix}$$

dont l'équation caractéristique est :

$$\begin{vmatrix} \lambda \mathbf{I}_n - (\mathbf{A} - \mathbf{BL}) & -\mathbf{BE} \\ 0 & \lambda \mathbf{I}_r - \mathbf{F} \end{vmatrix} = 0$$

c'est-à-dire :

$$|\lambda \mathbf{I}_n - (\mathbf{A} - \mathbf{BL})| \cdot |\lambda \mathbf{I}_r - \mathbf{F}| = 0$$

d'où :

théorème de séparation :

Les valeurs propres d'un système commandé par retour d'état et comportant un observateur dans sa boucle se composent de la réunion des valeurs propres du système bouclé sans observateur et de celles de l'observateur.

7.6 Filtrage de Kalman

Soit le système,

$$\begin{cases} \dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}\underline{u} + \mathbf{K}\underline{v} \\ \underline{y} = \mathbf{C}\underline{x} + \mathbf{D}\underline{u} + \underline{w} \end{cases} \quad \text{ou} \quad \begin{cases} \dot{\underline{x}}_{k+1} = \mathbf{\Phi}\underline{x}_k + \mathbf{\Gamma}\underline{u}_k + \mathbf{K}\underline{v}_k \\ \underline{y}_k = \mathbf{C}\underline{x}_k + \mathbf{D}\underline{u}_k + \underline{w}_k \end{cases} \quad (7.17)$$

avec :

$\underline{v}(t)$ ou \underline{v}_k vecteur aléatoire superposé à l'état (bruit d'état)

$\underline{w}(t)$ ou \underline{w}_k vecteur aléatoire superposé à la sortie (bruit de mesure)

Problème du filtrage linéaire :

A partir des données du problème (représentation d'état du problème) et des données statistiques concernant les bruits \underline{v} et \underline{w} (distribution statistique, moyenne, variance), trouver un système causal, à l'entrée duquel sont appliquées les signaux accessibles à la mesure, u et y et qui fournit à sa sortie $\hat{\underline{x}}$ aussi proche que possible de l'état \underline{x} inconnu.

La solution optimale de ce problème, au sens de la variance de $\hat{\underline{x}} - \underline{x}$ est le filtrage optimal de Kalman.

Chapitre 8

Représentation d'état des systèmes linéaires échantillonnés

8.1 Système discret

Le schéma de la figure 8.1 représente la forme générale d'un système discret.

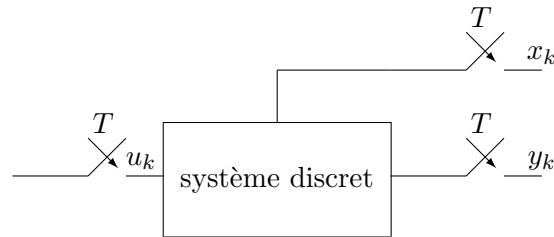


FIGURE 8.1 – Représentation générale d'un système discret

On définit :

$$\underline{x}_k = \begin{pmatrix} x_1(kT_e) \\ x_2(kT_e) \\ \vdots \\ x_n(kT_e) \end{pmatrix}, \quad \underline{y}_k = \begin{pmatrix} y_1(kT_e) \\ y_2(kT_e) \\ \vdots \\ y_q(kT_e) \end{pmatrix}, \quad \underline{u}_k = \begin{pmatrix} u_1(kT_e) \\ u_2(kT_e) \\ \vdots \\ u_p(kT_e) \end{pmatrix}$$

La résolution de l'équation d'état donne (voir (1.63, page 14)) :

$$\underline{x}(t) = e^{\mathbf{A}t} \left[\int_0^t e^{-\mathbf{A}\tau} \mathbf{B} \underline{u}(\tau) d\tau \right] + e^{\mathbf{A}(t-t_0)} \underline{x}(t_0) \quad (8.1)$$

Appliquons (8.1) à l'intervalle $kT_e \leq t < (k+1)T_e$ en y faisant $t_0 = kt$, nous obtenons :

$$\underline{x}_{k+1} = e^{\mathbf{A}T_e} \underline{x}_k + \int_{kT_e}^{(k+1)T_e} e^{\mathbf{A}((k+1)T_e-\tau)} \mathbf{B} \underline{u}(\tau) d\tau$$

Supposons que $\underline{u}(t) = \underline{u}_k = cte$ pour $kT_e \leq t < (k+1)T_e$ alors,

$$\underline{x}_{k+1} = e^{\mathbf{A}T_e} \underline{x}_k + \int_{kT_e}^{(k+1)T_e} e^{\mathbf{A}((k+1)T_e-\tau)} d\tau \cdot \mathbf{B} \underline{u}_k \quad (8.2)$$

En procédant au changement de variable : $\tau' = (k+1)T_e - \tau$ alors :

$$\int_{kT_e}^{(k+1)T_e} e^{\mathbf{A}((k+1)T_e-\tau)} d\tau = - \int_{T_e}^0 e^{\mathbf{A}\tau'} d\tau' = \int_0^{T_e} e^{\mathbf{A}\tau'} d\tau'$$

l'équation 8.2 devient :

$$\underline{x}_{k+1} = e^{\mathbf{A}T_e} \underline{x}_k + \int_0^{T_e} e^{\mathbf{A}\tau'} d\tau' \cdot \mathbf{B} \underline{u}_k. \quad (8.3)$$

On voit que x_{k+1} peut se mettre sous la forme :

$$\underline{x}_{k+1} = \mathbf{\Phi} \underline{x}_k + \mathbf{\Gamma} \underline{u}_k. \quad (8.4)$$

avec :

$$\mathbf{\Phi} = e^{\mathbf{A}T_e} \quad (8.5)$$

$$\mathbf{\Gamma} = \int_0^{T_e} e^{\mathbf{A}\tau'} d\tau' \cdot \mathbf{B} \quad (8.6)$$

De même on a :

$$\underline{y}_k = \mathbf{C} \underline{x}_k + \mathbf{D} \underline{u}_k. \quad (8.7)$$

Cas particulier : si \mathbf{A} est inversible,

$$\mathbf{\Gamma} = \mathbf{A}^{-1} [e^{\mathbf{A}\tau}]_0^{T_e} \cdot \mathbf{B} = \mathbf{A}^{-1} (e^{\mathbf{A}T_e} - \mathbf{I}) \cdot \mathbf{B} \quad (8.8)$$

soit :

$$\mathbf{\Gamma} = \mathbf{A}^{-1} (\mathbf{\Phi} - \mathbf{I}) \mathbf{B} \quad (8.9)$$

8.2 Résolution des équations dans le domaine du temps

En supposant que le vecteur d'état initial soit x_0 , l'application répétée de (8.4) donne :

$$\underline{x}_k = \mathbf{\Phi}^k \underline{x}_0 + \sum_{i=0}^{k-1} \mathbf{\Phi}^{k-1-i} \mathbf{\Gamma} \underline{u}_i, \quad (8.10)$$

la résolution se fait alors sur ordinateur.

8.3 Application de la transformée en z

Théorème de l'avance :

$$\mathcal{Z}[\underline{x}_{k+1}] = z \underline{X}(z) - z \underline{x}_0,$$

donc la transformée en z de 8.4 est :

$$z \underline{X}(z) - z \underline{x}_0 = \mathbf{\Phi} \underline{X}(z) + \mathbf{\Gamma} \underline{U}(z),$$

soit :

$$\underline{X}(z) = (z\mathbf{I} - \mathbf{\Phi})^{-1} z \underline{x}_0 + (z\mathbf{I} - \mathbf{\Phi})^{-1} \mathbf{\Gamma} \underline{U}(z), \quad (8.11)$$

et,

$$\underline{Y}(z) = \mathbf{C} \underline{X}(z) + \mathbf{D} \underline{U}(z) \quad (8.12)$$

8.4 Matrice de transfert

Dans le cas SISO : Fonction de transfert.

En supposant $x_0 = 0$, des équations (8.11) et (8.12) on déduit :

$$\underline{Y}(z) = \mathbf{C}(z\mathbf{I} - \mathbf{\Phi})^{-1}\mathbf{\Gamma}\underline{U}(z) + \mathbf{D}\underline{U}(z).$$

soit

$$\underline{Y}(z) = G(z)\underline{U}(z),$$

avec :

$$G(z) = \mathbf{C}(z\mathbf{I} - \mathbf{\Phi})^{-1}\mathbf{\Gamma} + \mathbf{D}.$$

$G(z)$ est la fonction (matrice dans le cas multivariable) de transfert du système échantillonné.

A comparer avec le résultat obtenu dans le cas continu :

$$\mathbf{H}(p) = \mathbf{C}(p\mathbf{I} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$$

8.5 Obtention d'un modèle d'état à partir de la fonction de transfert en z

Le modèle d'état s'obtient de la même façon qu'en continu mais avec la substitution

intégrateur pur \rightarrow retard pur

$$\frac{1}{p} \rightarrow z^{-1}$$

8.6 Résolution de l'équation d'état dans le domaine de z

La réponse libre du système ($u_k = 0$) s'obtient en utilisant (8.10) et (8.11).

$$\begin{aligned} \text{(8.10) devient :} \quad \underline{x}_k &= \mathbf{\Phi}^k \underline{x}_0 \\ \text{sa transformée en } z \text{ est :} \quad \underline{X}(z) &= \mathcal{Z}[\mathbf{\Phi}^k] \underline{x}_0 \\ \text{(8.11) devient :} \quad \underline{X}(z) &= (z\mathbf{I} - \mathbf{\Phi})^{-1} z \underline{x}_0 \end{aligned}$$

en identifiant terme à terme :

$$\mathcal{Z}[\mathbf{\Phi}^k] = (z\mathbf{I} - \mathbf{\Phi})^{-1} z,$$

d'où,

$$\mathbf{\Phi}(kT_e) = \mathbf{\Phi}^k = \mathcal{Z}^{-1} \left[(z\mathbf{I} - \mathbf{\Phi})^{-1} z \right],$$

à comparer avec : $\mathcal{L}^{-1}(p\mathbf{I} - \mathbf{A})^{-1}$.

8.7 Commandabilité et observabilité

8.7.1 Commandabilité d'un système échantillonné

Un système échantillonné représenté par une équation d'état aux différences pour une entrée scalaire u_k ,

$$\underline{x}_{k+1} = \mathbf{\Phi}\underline{x}_k + \mathbf{\Gamma}u_k. \quad (8.13)$$

$$\underline{y}_k = \mathbf{C}\underline{x}_k \quad (8.14)$$

$x(0)$ étant donné, le vecteur d'état solution de l'équation (8.13) s'exprime par : (en écrivant les termes successifs) :

$$\underline{x}_k = \Phi^k \underline{x}(0) + \Gamma \underline{u}_{k-1} + \Phi \Gamma \underline{u}_{k-2} + \Phi^2 \Gamma \underline{u}_{k-3} + \dots + \Phi^{k-1} \Gamma \underline{u}_0 \quad (8.15)$$

$$\underline{x}_k = \Phi^k \underline{x}(0) + \sum_{i=1}^k \Phi^{i-1} \Gamma \underline{u}_{k-i} \quad (8.16)$$

que l'on peut réécrire sous la forme :

$$\underline{x}_k = \Phi^k \underline{x}(0) + \mathbf{Q}_S \underline{U}^T \quad (8.17)$$

avec :

$$\mathbf{Q}_S = \begin{bmatrix} \Gamma, \Phi \Gamma, \Phi^2 \Gamma, \dots, \Phi^{k-1} \Gamma \end{bmatrix}$$

$$\underline{U} = \begin{bmatrix} \underline{u}_{k-1}, \underline{u}_{k-2}, \dots, \underline{u}_1, \underline{u}_0 \end{bmatrix}$$

Pour que l'état du système passe de l'état \underline{x}_0 à l'état \underline{x}_k , il faut que la séquence de commande \underline{U}^T respecte

$$\mathbf{Q}_S \underline{U}^T = \underline{x}_k - \Phi^k \underline{x}(0). \quad (8.18)$$

Cette séquence existe si \mathbf{Q}_S est inversible donc :

$$\underline{U}^T = \mathbf{Q}_S^{-1} \left[\underline{x}_k - \Phi^k \underline{x}(0) \right]. \quad (8.19)$$

Les conditions de commandabilité sont donc que :

1. \mathbf{Q}_S soit carrée,
2. $\det(\mathbf{Q}_S) \neq 0$.

Ainsi la condition de commandabilité devient (n est le nombre d'états) :

$$\det \mathbf{Q}_S = \det \left[\Gamma, \Phi \Gamma, \Phi^2 \Gamma, \dots, \Phi^{n-1} \Gamma \right] \neq 0. \quad (8.20)$$

8.7.2 Observabilité d'un système échantillonné

La sortie d'un système échantillonné peut être écrite sous la forme :

$$\underline{y}_k = \mathbf{C} \Phi^k \underline{x}_0 + \sum_{i=1}^k \mathbf{C} \Phi^{i-1} \Gamma \underline{u}_{k-i} \quad (8.21)$$

$$\begin{aligned} \underline{y}_0 &= \mathbf{C} \underline{x}_0 \\ \underline{y}_1 &= \mathbf{C} \underline{x}_1 = \mathbf{C} \Phi \underline{x}_0 + \mathbf{C} \Gamma \underline{u}_0 \\ \underline{y}_2 &= \mathbf{C} \underline{x}_2 = \mathbf{C} \Phi \underline{x}_1 + \mathbf{C} \Gamma \underline{u}_1 = \mathbf{C} \Phi^2 \underline{x}_0 + \mathbf{C} \Phi \Gamma \underline{u}_0 + \mathbf{C} \Gamma \underline{u}_1 \\ &\vdots \\ \underline{y}_{k-1} &= \mathbf{C} \underline{x}_{k-1} = \mathbf{C} \Phi \underline{x}_{k-2} + \mathbf{C} \Gamma \underline{u}_{k-2} = \mathbf{C} \Phi^{k-1} \underline{x}_0 + \mathbf{C} \Phi^{k-2} \Gamma \underline{u}_0 + \dots + \mathbf{C} \Gamma \underline{u}_{k-2} \end{aligned}$$

que l'on peut réécrire sous la forme :

$$\underline{Y}^T = \mathbf{Q}_B \underline{x}_0 + \mathbf{H} \underline{U}_1 \quad (8.22)$$

avec :

$$\mathbf{Q}_B^T = \left[\mathbf{C}, \mathbf{C} \Phi, \mathbf{C} \Phi^2, \dots, \mathbf{C} \Phi^{k-1} \right]$$

$$\underline{U}_1 = \left[0, \underline{u}_0, \underline{u}_1, \dots, \underline{u}_{k-2} \right]$$

donc

$$\mathbf{Q}_B \underline{x}_0 = \underline{Y}^T - \mathbf{H} \underline{U}_1 \quad (8.23)$$

$$\underline{x}_0 = \mathbf{Q}_B^{-1} (\underline{Y}^T - \mathbf{H} \underline{U}_1) \quad (8.24)$$

Remonter à l'état initial \underline{x}_0 en mesurant \underline{Y}^T il faut que :

1. \mathbf{Q}_B soit carrée,
2. $\det(\mathbf{Q}_B) \neq 0$.

Un système est observable si et seulement si

$$\det \mathbf{Q}_B = \det \begin{bmatrix} \mathbf{C} \\ \mathbf{C}\Phi \\ \mathbf{C}\Phi^2 \\ \vdots \\ \mathbf{C}\Phi^{n-1} \end{bmatrix} \neq 0. \quad (8.25)$$

8.8 Stabilité des systèmes échantillonnés

La stabilité des systèmes échantillonnés est identique au cas continu sauf qu'un pôle stable en z est un pôle dont le module est inférieur à 1.

$$\text{stable} \Rightarrow \|z\| < 1$$

Bien entendu le critère de Routh appliqué sur la dernière ligne de la matrice sous la forme canonique de commandabilité devient le critère de Jury. Une autre solution consiste à appliquer le critère de Routh sur cette dernière ligne après lui avoir appliqué la transformée en w .

Rappel : la transformée en w consiste à remplacer z par $\frac{w-1}{w+1}$

8.9 Commandes des systèmes échantillonnés

La commande des systèmes échantillonnés se fait de la même façon que pour les systèmes continus. Les éventuelles transformations sont identiques.

Le choix des pôles en z dans le cadre d'une méthode de placement de pôles est par contre différent (voir fig. 8.2).

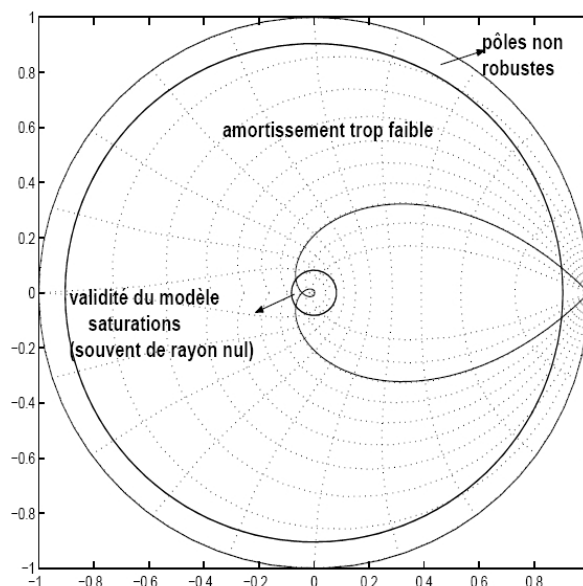


FIGURE 8.2 – Choix des pôles dans le plan en z

8.9.1 Calcul de la matrice de préfiltre

Le modèle du système en boucle fermée est :

$$\begin{aligned}\underline{x}_{k+1} &= (\Phi - \Gamma L)\underline{x}_k + \Gamma S \underline{y}_{c_k} \\ \underline{y}_k &= C \underline{x}_k\end{aligned}$$

Si le système est stable, alors $\lim_{k \rightarrow \infty} \underline{x}_{k+1} = \underline{x}_k$ donc,

$$\lim_{k \rightarrow \infty} \underline{y}_k = -C(\Phi - \Gamma L - I)^{-1} \Gamma S \underline{y}_{c_k}$$

or, on désire que $\lim_{k \rightarrow \infty} \underline{y}_k = \underline{y}_{c_k}$, donc,

$$S^{-1} = -C(\Phi - \Gamma L - I)^{-1} \Gamma$$

$$S = -\left(C(\Phi - \Gamma L - I)^{-1} \Gamma\right)^{-1}$$

à comparer avec le résultat obtenu dans le cas continu

$$S = -\left(C(A - BL)^{-1} B\right)^{-1}$$

8.9.2 Commande optimale dans le cas discret

$$J = \sum_{n=1}^{\infty} \underline{x}_n^T Q \underline{x}_n + \underline{u}_n^T R \underline{u}_n \quad (8.26)$$

avec : Q = matrice symétrique définie positive

R = matrice symétrique non négative

la commande u est alors définie par l'équation suivante :

$$u_{k-1} = -\left(R + \Gamma^T K \Gamma\right)^{-1} \Gamma K \Phi \underline{x}_{k-1}$$

où K est une matrice symétrique, solution définie négative de l'équation de Riccati :

$$K = \Phi^T (K - K \Gamma (R + \Gamma^T K \Gamma)^{-1} \Gamma^T K) \Phi + Q$$

Chapitre 9

Annales d'examens

Examen final d'automatique

Durée : 2 heures.

Aucun document autorisé.

Le sujet de cet examen est la modélisation à l'aide d'une représentation d'état d'une micro pince. Dans un deuxième temps, nous développerons une commande par retour d'état du système fondée sur un observateur complet de Luenberger. Enfin nous développerons un observateur de l'effort exercé par la pince sur l'objet.

Quelques conseils

- La justification des choix est aussi importante que les calculs effectués.
- Ne restez pas bloqué sur une question, revenez-y par la suite.

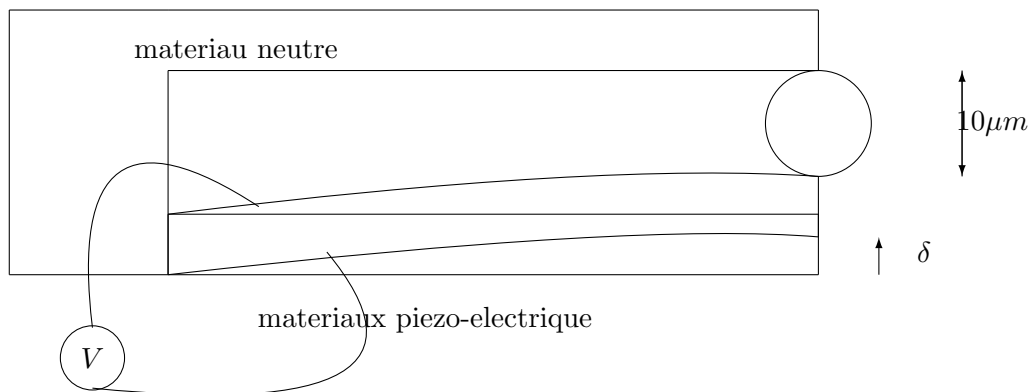


FIGURE 9.1 – Présentation de la pince étudiée.

Cette pince est en matériau piézo-électrique, la tension V appliqué sur les électrodes provoque le même effet qu'une force αV appliquée au bout de la mâchoire de la pince.

En première approximation le comportement dynamique de la pince est le même que celui du modèle mécanique présenté figure 9.2 :

La masse M est supposée ponctuelle, toutes les forces sont appliquées à une longueur L du centre de rotation.

Modélisation

1. En appliquant la relation fondamentale de la dynamique sur modèle mécanique précédent, déterminer l'équation différentielle reliant le mouvement de rotation aux forces appliquées.
2. En faisant l'hypothèse que $\theta \ll 1$ et que les force de gravité sont négligeables vis-à-vis des autres forces, déterminez l'équation différentielle linéarisée reliant δ aux autres paramètres du système.
3. Montrez que l'on obtient une équation différentielle de la forme :

$$m\ddot{\delta} + f\dot{\delta} + k\delta = -F_{ext} + \alpha V$$

Dans tout ce qui suit, l'équation différentielle du mouvement utilisée sera celle-ci!

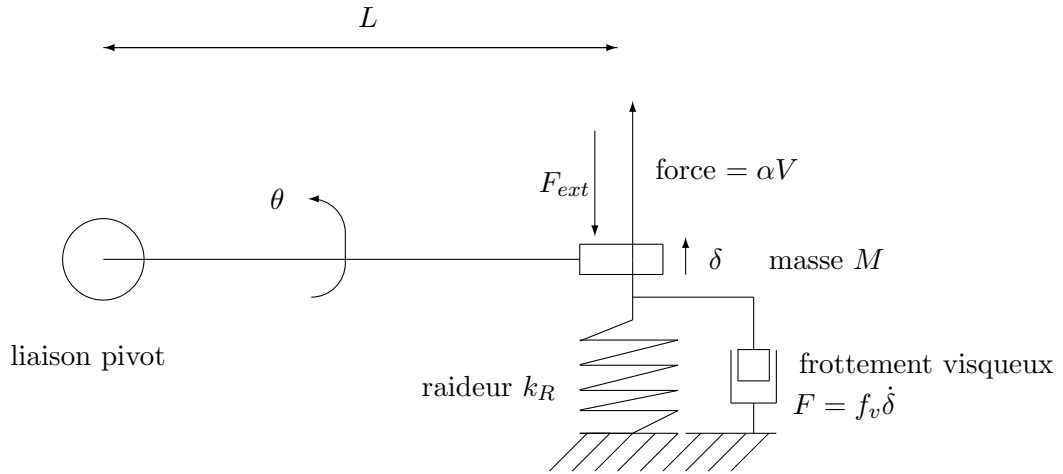


FIGURE 9.2 – Modèle mécanique équivalent.

4. En supposant que la force F_{ext} exercée par l'objet sur la pince est nulle. Déduisez de ce qui précède une représentation d'état de la forme :

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}u \tag{9.1}$$

$$\underline{y} = \mathbf{C}\underline{x} + \mathbf{D}u \tag{9.2}$$

Applications numériques :

$$\alpha = 0.002, \quad m = 10^{-3}, \quad f = 11, \quad k = 1^5$$

Analyse

1. Le système est-il stable en boucle ouverte ?
2. Le système est-il commandable ?
3. Le système est-il observable ?

Commande

1. Donnez la valeur numérique des pôles en boucle ouverte.
2. Au vu des pôles en boucle ouverte, de votre expérience et du système proposez des pôles en boucle fermée.
3. Justifiez-les en 4 lignes maximum (+ un dessin si besoin).
4. Calculer un régulateur d'état qui permet d'obtenir la dynamique que vous vous êtes fixée en boucle fermée.
5. Par le calcul de la fonction de transfert, vérifiez que la dynamique souhaitée est bien obtenue.
6. A l'aide de cette fonction de transfert, calculer la matrice (ici un réel) de préfiltre.
7. Proposez une autre méthode pour obtenir une erreur statique nulle.

Observation

1. Donnez le schéma-bloc du système sous sa forme de représentation d'état avec son observateur de Luenberger complet.
2. Proposez et justifiez les pôles de l'observateur.
3. Calculer, par la méthode de votre choix, la matrice de "retour" \mathbf{G}

Observation de la perturbation

Dans cette deuxième partie, la force F_{ext} exercée par l'objet sur la pince est non nulle, par contre on supposera de faibles variations de force, sa dérivé \dot{F}_{ext} est nulle.

1. A partir des équations établies lors de la modélisation et avec les hypothèses précédentes, donnez une représentation d'état du système incluant la force F_{ext} exercée par l'objet sur la pince.

$$\dot{\underline{z}} = \mathbf{A}_a \underline{z} + \mathbf{B}_a u \quad (9.3)$$

$$\delta = \mathbf{C}_a \underline{z} + \mathbf{D}_a u \quad (9.4)$$

avec : $\underline{z} = \begin{bmatrix} \underline{x} \\ F_{ext} \end{bmatrix}$

2. Donnez le schéma bloc du système complet incluant la pince, son observateur augmenté, le point d'application de la perturbation, le retour d'état utilisant le vecteur estimé, la force estimée \hat{F}_{ext} .
3. Montrez que lorsque $t \rightarrow \infty$ et en l'absence d'erreurs de modèle, $\hat{F}_{ext} \rightarrow F_{ext}$
4. Calculez la matrice de retour \mathbf{G} pour ce nouvel observateur.

Examen de septembre d'automatique

Durée : 2 heures.

Aucun document autorisé.

Le sujet de cet examen est la modélisation à l'aide d'une représentation d'état d'un hélicoptère. Dans un deuxième temps, nous développerons une commande par retour d'état du système fondée sur un observateur complet de Luenberger. Enfin nous développerons un nouvel asservissement permettant de garantir une erreur statique nulle

Quelques conseils

- La justification des choix est aussi importante que les calculs effectués.
- Ne restez pas bloqué sur une question, revenez-y par la suite.

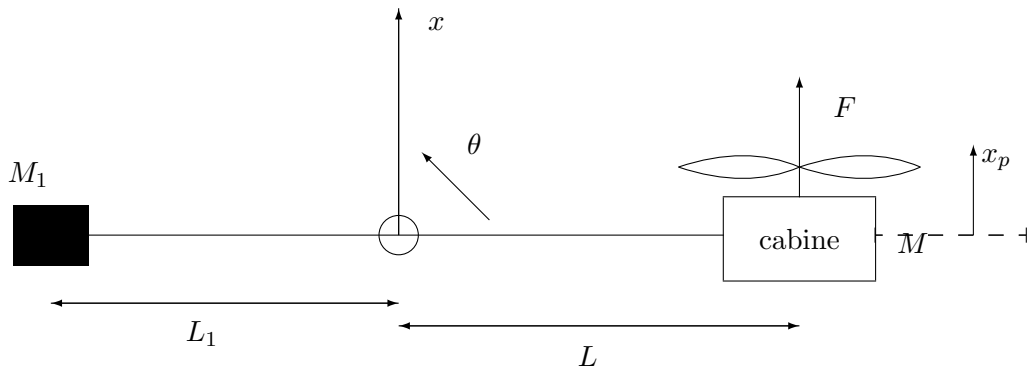


FIGURE 9.3 – Présentation de l'hélicoptère ($\theta = \pi/2$).

Nous supposons que le problème est plan. La position instantanée de l'hélicoptère peut donc être décrite sans ambiguïté par l'angle du bras avec la verticale θ .

La cabine de masse M est surmontée d'une hélice, le rotor, qui provoque une force F qui tend à faire monter la cabine. Un contre-poids de masse M_1 est situé à une distance L_1 du centre de rotation. La liaison pivot présente un couple de frottement visqueux de la forme $f_v \dot{\theta}$ ainsi qu'un couple de rappel de la forme $k_R \theta$, ce couple est nul à la position d'équilibre en $\theta = \pi/2 rad$.

Modélisation

1. Du modèle mécanique précédent, déterminer l'équation différentielle reliant le mouvement de l'hélicoptère à la force appliquée par le rotor F .
2. En faisant l'hypothèse que $\theta \ll 1$, déterminez l'équation différentielle linéarisée reliant \ddot{x}_p aux autres paramètres du système.
3. Montrez que l'on obtient une équation différentielle de la forme :

$$m\ddot{x}_p + f\dot{x}_p + kx_p = F$$

Dans tout ce qui suit, l'équation différentielle du mouvement utilisée sera celle-ci !

4. Dédisez de ce qui précède une représentation d'état de la forme :

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}u \quad (9.5)$$

$$\underline{y} = \mathbf{C}\underline{x} + \mathbf{D}u \quad (9.6)$$

avec $y = x_p$

Applications numériques :

$$m = 2, \quad f = 9, \quad k = 9$$

Analyse

1. Le système est-il stable en boucle ouverte ?
2. Le système est-il commandable ?
3. Le système est-il observable ?

Commande

1. Donnez la valeur numérique des pôles en boucle ouverte.
2. Au vu des pôles en boucle ouverte, de votre expérience et du système proposez des pôles en boucle fermée.
3. Justifiez-les en 4 lignes maximum (+ un dessin si besoin).
4. Calculer un régulateur d'état qui permet d'obtenir la dynamique que vous vous êtes fixée en boucle fermée.
5. Par le calcul de la fonction de transfert, vérifiez que la dynamique souhaitée est bien obtenue.
6. A l'aide de cette fonction de transfert, calculer la matrice (ici un réel) de préfiltre.
7. Proposez une autre méthode pour obtenir une erreur statique nulle.

Observation

1. Donnez le schéma-bloc du système sous sa forme de représentation d'état avec son observateur de Luenberger complet.
2. Proposez et justifiez les pôles de l'observateur.
3. Calculer, par la méthode de votre choix, la matrice de "retour" \mathbf{G}

Obtention de l'erreur statique nulle

1. Donnez le schéma complet du système incluant :
 - le système lui-même,
 - un intégrateur pur pour obtenir une erreur statique nulle.
 On considérera que la matrice d'anticipation est nulle.
2. Dédisez de ce qui précède une représentation d'état de la forme :

$$\dot{\underline{z}} = \mathbf{A}_a \underline{z} + \mathbf{B}_a y_c \quad (9.7)$$

$$y = \mathbf{C}_a \underline{z} + \mathbf{D}_a y_c \quad (9.8)$$

$$\text{avec : } \underline{z} = \begin{bmatrix} x \\ \nu \end{bmatrix}$$

où ν représente la sortie de l'intégrateur pur.

3. calculez un nouveau correcteur tel que ce nouveau système présente à peu près les mêmes caractéristiques que celle que vous avez précédemment choisies.

Examen final d'automatique avancée

Durée : 2 heures.

Tous documents autorisés.

1 Premier problème : Commande des systèmes non commandables

Un système linéaire continu est donné par sa représentation d'état :

$$\dot{\underline{x}} = \mathbf{A}\underline{x} + \mathbf{B}u \quad (9.9)$$

$$\underline{y} = \mathbf{C}\underline{x} + \mathbf{D}u \quad (9.10)$$

$$\text{avec : } \mathbf{A} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 2 & 1 & -2 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}, \quad \mathbf{C} = [1 \ 0 \ 0], \quad \mathbf{D} = 0$$

1.1 Analyse

1. Montrer que le système est instable en boucle ouverte.

Instable : Critère de Routh

2. Montrer que le système est non commandable.

$$Q_s = \begin{bmatrix} 0 & 1 & -1 \\ 1 & -1 & 3 \\ -1 & 3 & -5 \end{bmatrix} \quad \det Q_s = 0 \implies \text{Non commandable}$$

3. Le système est-il observable?

$$Q_b = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \det Q_b = 1 \implies \text{Observable}$$

1.2 Réduction du modèle

Dans cette partie on "enlève" les parties non commandables et/ou non observables

1. Déterminer la fonction de transfert du système $H(p) = \frac{Y(p)}{U(p)}$.

$$H(p) = \frac{2p + 3}{(p - 1)(2 + p)}$$

2. A partir de la fonction de transfert $H(p)$ déduisez-en une représentation d'état sous forme canonique de commandabilité.

$$\dot{\underline{x}}_r = \mathbf{A}_r \underline{x}_r + \mathbf{B}_r u \quad (9.11)$$

$$\underline{y} = \mathbf{C}_r \underline{x}_r + \mathbf{D}_r u \quad (9.12)$$

avec :

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 2 & -1 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad \mathbf{C} = [1 \ 0], \quad \mathbf{D} = 0$$

3. Ce nouveau système est-il commandable et observable?
oui, car il provient de la fonction de transfert. Pôles en -1 et -2. Commandable

$$Q_s = \begin{bmatrix} 0 & 1 \\ 1 & -1 \end{bmatrix} \quad \det Q_s = 0 \implies \text{Commandable}$$

Observable?

$$Q_b = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \det Q_s = 1 \implies \text{Observable}$$

4. Vérifier que le pôle instable est toujours présent.
 Pôles en 1 et -2

1.3 Commande du modèle réduit

1. Calculer le régulateur \mathbf{L} tel que les pôles du système en boucle fermée soient :

$$p_{s1,2} = -3 \pm 3i$$

$$Q_s = \begin{bmatrix} 0 & 1 \\ 1 & -1 \end{bmatrix} \quad \text{donc} \quad Q_s^{-1} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \quad \text{et} \quad q_s^T = [1, 0]$$

$$P(p) = (p + 3 - 3\sqrt{-1})(p + 3 + 3\sqrt{-1}) = p^2 + 6p + 18$$

$$\mathbf{L} = [20, 5]$$

2. Expliquez brièvement le comportement en boucle fermée que l'on aura avec ces pôles.
Sans doute pas celui voulu à cause du zéro de la fonction de transfert.
3. Calculer la matrice (ici un réel) de préfiltre \mathbf{S} .

$$\mathbf{S} = -(\mathbf{C}(\mathbf{A} - \mathbf{BL})^{-1}\mathbf{B})^{-1} = \frac{-1}{18}$$

1.4 Observation

1. Donnez le schéma-bloc du système sous sa forme de représentation d'état avec son observateur de Luenberger complet.
2. Calculer, par la méthode de votre choix, la matrice de "retour" \mathbf{G} , telle que les pôles de l'observateur soient placés en

$$p_{o1,2} = -10 \pm 10i$$

$$Q_b = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{donc} \quad Q_b^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{et} \quad q_b = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$P(p) = (p + 10 - 10\sqrt{-1})(p + 10 + 10\sqrt{-1}) = p^2 + 20p + 200$$

$$\mathbf{G} = \begin{bmatrix} 38 \\ 183 \end{bmatrix}$$

1.5 Erreur statique nulle

On se propose d'ajouter un intégrateur pur tel que l'erreur statique devienne nulle, y compris en présence de perturbations ou d'erreurs de modélisation.

1. Le schéma-bloc du système est donné en figure 9.4. Déterminer la représentation d'état du système augmenté :

$$\dot{\underline{x}}_a = \mathbf{A}_a \underline{x}_a + \mathbf{B}_a \underline{y}_c \tag{9.13}$$

$$\underline{y} = \mathbf{C}_a \underline{x}_a + \mathbf{D}_a \underline{y}_c \tag{9.14}$$

avec $\underline{x}_a = \begin{bmatrix} x \\ \eta \end{bmatrix}$

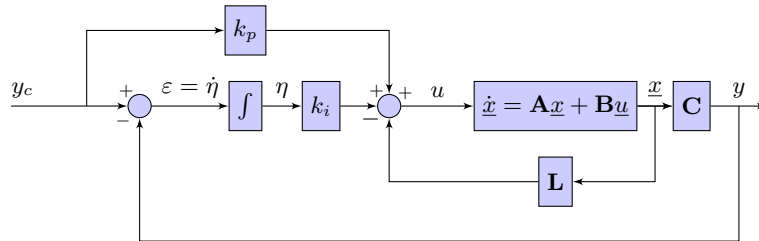


FIGURE 9.4 – Schéma d'une commande intégrale.

$$\mathbf{A}_a = \begin{bmatrix} 0 & 1 & 0 \\ 2 - l_0 & -1 - l_1 & k_i \\ -1 & 0 & 1 \end{bmatrix}, \quad \mathbf{B}_a = \begin{bmatrix} 0 \\ k_p \\ 1 \end{bmatrix}, \quad \mathbf{C}_a = [1 \ 0 \ 0], \quad \mathbf{D}_a = 0$$

2. En posant $\mathbf{L} = [l_0 \ l_1]$, calculer l_0 , l_1 et k_i tels que les pôles du système soient placés en :

$$p_{s_{1,2,3}} = -3 \quad \text{et} \quad -3 \pm 3i$$

par identification :

$$\det \begin{bmatrix} p & -1 & 0 \\ -2 + l_0 & 1 + l_1 + p & -k_i \\ 1 & 0 & p - 1 \end{bmatrix} = (p + 3)(p + 3 - 3i)(p + 3 + 3i)$$

$$p^3 + l_1 p^2 + (l_0 - 3 - l_1)p + 2 - l_0 + k_i = p^3 + 9p^2 + 36p + 54$$

$$\implies \{k_i = 100, l_1 = 9, l_0 = 48\}$$

1.6 Observateur numérique

L'observateur du système précédent est en fait implanté sous forme numérique.

1. transformer le système

$$\dot{\underline{x}}_r = \mathbf{A}_r \underline{x}_r + \mathbf{B}_r u \tag{9.15}$$

$$\underline{y} = \mathbf{C}_r \underline{x}_r + \mathbf{D}_r u \tag{9.16}$$

en un système échantillonné de la forme :

$$\underline{x}_{k+1} = \Phi \underline{x}_k + \Gamma u_k \tag{9.17}$$

$$\underline{y}_k = \mathbf{C} \underline{x}_k \tag{9.18}$$

sachant que en prenant $T_e = 1\text{s}$, $\Phi = \begin{bmatrix} 1.86 & 0.86 \\ 1.72 & 0.99 \end{bmatrix}$

$$\Gamma = \mathbf{A}^{-1}(\Phi - \mathbf{I})\mathbf{B} = \begin{bmatrix} 0.43 \\ 0.86 \end{bmatrix}$$

2. Déterminer un observateur complet pour ce système numérique. Les pôles de cet observateur complet seront placés en 0.

$$Q_b = \begin{bmatrix} 1 & 0 \\ 1.86 & 0.86 \end{bmatrix} \quad \text{donc} \quad Q_b^{-1} = \begin{bmatrix} 1.0 & 0 \\ -2.16 & 1.16 \end{bmatrix} \quad \text{et} \quad q_b = \begin{bmatrix} 0 \\ 1.16 \end{bmatrix}$$

$$P(z) = z^2$$

$$\mathbf{G} = \begin{bmatrix} -2.32 \\ 3.48 \end{bmatrix}$$

Chapitre 10

Travaux dirigés

1 TD1 - Représentation d'état

1.1 A partir d'un système

Soit un moteur à courant continu commandé par l'inducteur.

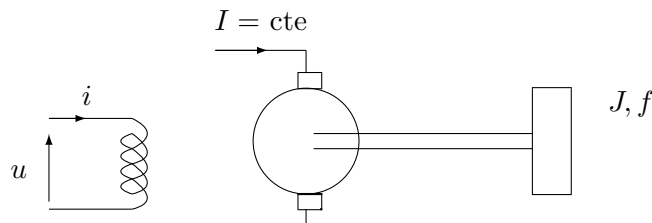


FIGURE 10.1 – Moteur à courant continu commandé par l'inducteur

- commande : u
- sortie : ω

Les équations fondamentales de ce système sont données ci-après :

$$U = Ri + L \frac{di}{dt} \quad (10.1)$$

$$J \frac{d\omega}{dt} = -f\omega + \gamma \quad (10.2)$$

$$\gamma = ki \quad (10.3)$$

$$(10.4)$$

1. Déduisez-en l'équation différentielle qui régit la vitesse du rotor en fonction de la tension. Déduisez-en la fonction de transfert du système, puis la réponse impulsionnelle du système.
2. Déduisez-en un schéma-bloc du système.
3. En posant :

$$x_1 = \omega, \quad x_2 = \frac{d\omega}{dt} = \dot{\omega}, \quad (10.5)$$

donnez l'équation différentielle du système sous la forme matricielle suivante :

Equation d'état :

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} u \quad (10.6)$$

Equation de sortie :

$$\omega = \begin{bmatrix} \cdot & \cdot \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (10.7)$$

1.2 A partir d'un schéma bloc

Soit le système représenté figure 10.2

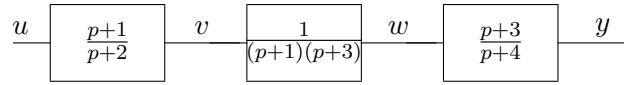


FIGURE 10.2 – Système étudié

Montrez que ce système peut être mis sous la forme suivante (fig. 10.3)

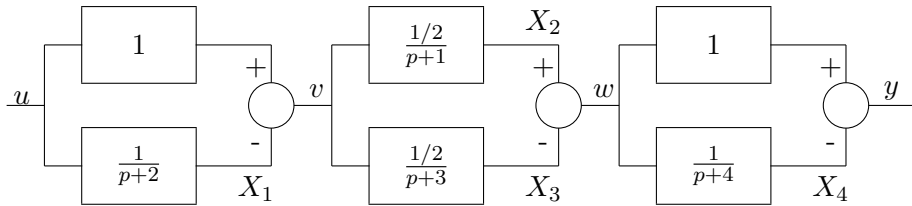


FIGURE 10.3 – Nouvelle forme du système

Mise en équation du système

Mettre le système en équation avec les trois méthodes proposées.

1. la représentation d'état
2. la fonction de transfert
3. l'équation différentielle

Comparez l'information présente dans les trois mises en équation.

1.3 Pour l'entraînement

Etudiez le système suivant (figure 10.4)

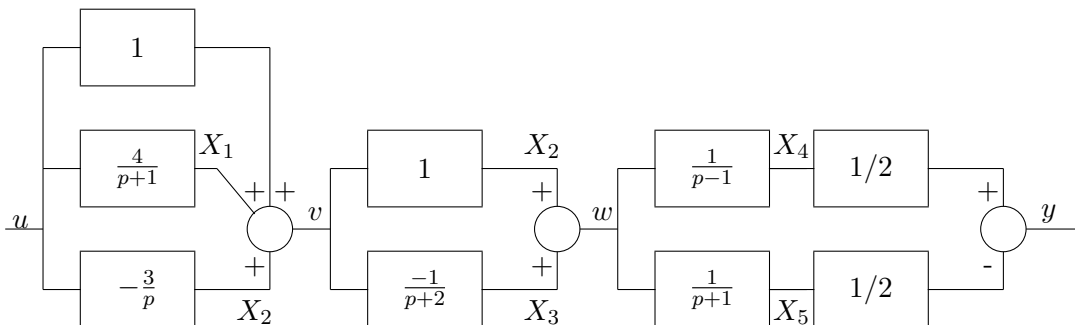


FIGURE 10.4 – Système étudié

2 Commandabilité et observabilité

soit le système représenté figure 10.5

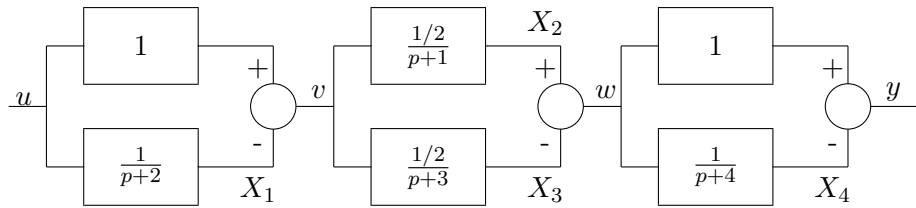


FIGURE 10.5 – Système étudié

2.1 Mise en équation du système

Mettre le système en équation avec les trois méthodes proposées.

1. la représentation d'état
2. la fonction de transfert
3. l'équation différentielle

Comparez l'information présente dans les trois mises en équation.

2.2 Mise sous la forme canonique de Jordan

1. Résoudre l'équation

$$\det(p\mathbf{I} - \mathbf{A}) = 0$$

2. Calculez le vecteur propre associé à chaque valeur propre
3. Déterminez la matrice de transformation \mathbf{T}
4. Vérifiez que $\mathbf{T}^{-1}\mathbf{A}\mathbf{T} = \mathbf{\Lambda}$
5. Calculez les vecteurs $\hat{\mathbf{C}} = \mathbf{C}\mathbf{T}$ et $\hat{\mathbf{B}} = \mathbf{T}^{-1}\mathbf{B}$
6. Donnez la représentation d'état complète du système sous la forme canonique de Jordan.

2.3 Commandabilité et observabilité

1. Vérifiez la commandabilité et l'observabilité du système sur la représentation initiale.
2. Vérifiez la commandabilité et l'observabilité sur la représentation sous la forme canonique de Jordan.
3. Tracez le graphe de fluence du système sous la forme canonique de Jordan.

3 TD asservissement dans l'espace d'état d'un pont roulant

Le rôle d'un pont roulant est de s'emparer d'une charge à un endroit déterminé, de la déplacer sur une distance déterminée, puis de la déposer à un autre endroit donné. Le problème qui surgit lors de l'automatisation de ce processus est que la charge entre en oscillation à la mise en marche et au freinage du chariot. Ces oscillations sont nuisibles au fonctionnement puisqu'elles retardent la prise et le dépôt de la charge.

Etant donné que les performances d'un pont roulant résident principalement dans la vitesse à laquelle il est capable d'exécuter un cycle de travail, il est nécessaire pour améliorer ses performances d'empêcher la naissance de ces oscillations ou, tout du moins de les réduire à une valeur acceptable.

Nous supposons que le problème est plan. La position instantanée de la benne peut donc être décrite sans ambiguïté à l'aide de deux coordonnées, qui peuvent être, par exemple, l'abscisse du chariot x_c et soit l'abscisse de la benne x_b soit l'angle du filin avec la verticale θ .

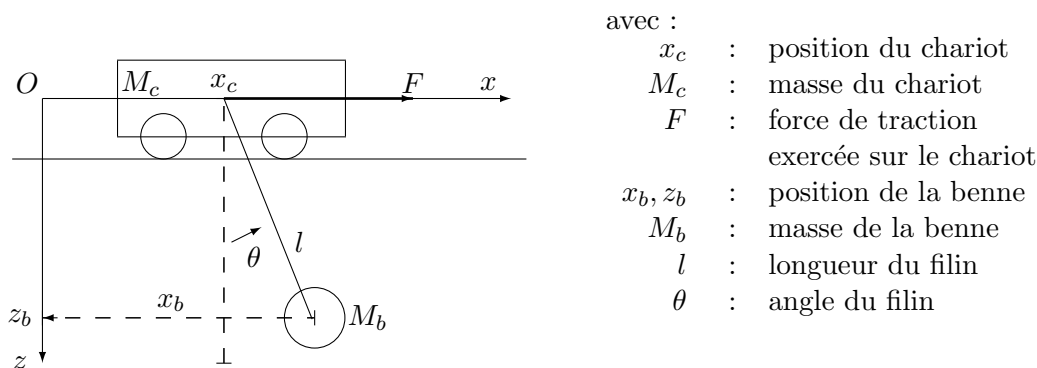


FIGURE 10.6 – Schéma du chariot et de la benne

3.1 Mise en équation

Mettez en équations ce problème. Linéarisez ensuite dans l'hypothèse, généralement vérifiée dans la réalité, d'excursions angulaires faibles de la benne (θ petit).

3.2 Représentation d'état

Etablir les équations d'état de ce système en choisissant les variables suivantes :

x_1 : position du chariot

x_2 : vitesse du chariot

x_3 : angle du filin

x_4 : vitesse angulaire du filin

L'équation de sortie sera, pour l'instant, ignorée.

Note : Je vous suggère d'utiliser les notations suivantes :

$$\alpha = \frac{M_b}{M_c}g, \quad \beta = -\left(1 + \frac{M_b}{M_c}\right)\frac{g}{l}, \quad b_2 = \frac{1}{M_c}, \quad b_4 = -\frac{1}{M_c l}$$

3.3 Simulation Matlab

Simulez le système sur un ordinateur. Relevez et interprétez les courbes de x_c , \dot{x}_c , θ en réponse à un échelon de force appliquée $F = 1\text{kN}$.

Valeurs numériques :

$$M_c = 1000\text{kg}, \quad M_b = 4000\text{kg}, \quad l = 10\text{m}, \quad g = 9,81\text{ms}^{-2}$$

3.4 Commandabilité et observabilité

Déterminez la commandabilité en calculant la matrice de commandabilité puis interprétez physiquement le résultat obtenu.

Déterminez l'observabilité du système en se plaçant dans les deux hypothèses suivantes :

1. mesure de la position du chariot seule,
2. mesure de l'angle du filin seul.

3.5 Commande par retour d'état

Effectuez la synthèse de l'asservissement de ce système par retour d'état par la méthode du placement de pôles en respectant le cahier des charges suivant :

1. en régime transitoire,
 - (a) assez bien amorti, dépassement faible,
 - (b) rapide, le temps d'établissement à 2% ne doit pas dépasser $t_{2\%} = 20\text{s}$,
2. en régime permanent, erreur statique nulle

Indications :

1. Commencez par déterminer les pôles du système à asservir,
2. Choisir les pôles du système bouclé en fonction du cahier des charges de la façon suivante :
 - (a) 1 paire de pôles complexes dominants,
 - (b) 1 pôle réel unique d'ordre ?? en $p = -1$.
3. Calculez les paramètres du régulateur d'état.
4. Calculez la matrice de préfiltre.
5. Tracez le schéma fonctionnel du système bouclé, en représentant le système à asservir par un seul bloc.

3.6 Observateur complet

Déterminez un observateur de Luenberger complet, dans le cas de la mesure de la position du chariot, ayant une valeur propre unique d'ordre ??, en $p = -2$. Justifiez ce choix. Calculez les paramètres de cet observateur numériquement et tracez son schéma fonctionnel.

3.7 Observateur réduit

Déterminez un observateur de Luenberger réduit, dans le cas de la mesure de la position du chariot, ayant une valeur propre unique d'ordre ??, en $p = -2$. Justifiez ce choix. Calculez les paramètres de cet observateur numériquement et tracez son schéma fonctionnel.

4 Commande du tangage d'un avion

Les équations du mouvement d'un avion sont particulièrement complexes puis qu'il s'agit de 6 équations différentielles non linéaires couplées. Moyennant certaines hypothèses, celles-ci peuvent être réduites à des équations différentielles découplées et linéarisées. Dans cet exercice on se propose d'étudier un pilote automatique contrôlant le tangage de l'avion. Attention en cas d'erreur, les conséquences seraient très graves!

La complexité des calculs est telle que l'utilisation du logiciel Matlab s'impose.

4.1 Modélisation

Les équation du mouvement de tangage sont les suivantes :

$$\dot{\alpha} = \mu\Omega[-(C_L + C_D)\alpha + (1/\mu - C_L)q - (C_w \sin(\gamma_e))\theta + C_L] \quad (10.8)$$

$$\dot{q} = \frac{\mu\Omega}{2i_{yy}} \{ [C_M - \nu(C_L + C_D)]\alpha + [C_M + \sigma C_M(1 - \mu C_L)]q + (\nu C_W \sin(\gamma_e))\delta_e \} \quad (10.9)$$

$$\dot{\theta} = \Omega q \quad (10.10)$$

Définition de quelques variables :

α = angle d'attaque

q = variation d'angle de tangage

θ = angle de tangage

δ_e = angle des volets de l'avion

ρ_e = densité de l'air ambiant

S = surface de l'aile

m = masse de l'avion

U = vitesse de l'avion

Déduisez des équations précédentes la forme générale de la représentation d'état du système.

Pour la suite, les valeurs numériques suivantes seront utilisées.

$$\dot{\alpha} = -0,313\alpha + 56,7q + 0.232\delta_e \quad (10.11)$$

$$\dot{q} = -0,0139\alpha - 0,426q + 0.0203\delta_e \quad (10.12)$$

$$\dot{\theta} = 56,7q \quad (10.13)$$

4.2 Cahier des charges

Les critères de confort des passagers font apparaître le cahier des charges suivant :

dépassement de moins de 10%

temps de montée de moins de 2 secondes

régime permanent en moins de 10 secondes

erreur statique de moins de 2%

4.3 Etude en boucle ouverte - fonction de transfert

1. Déterminez la fonction de transfert du système.
entrée : δ_e , sortie : θ
2. le système est-il stable en boucle ouverte.
3. le système est-il contrôlable, est-il observable ?

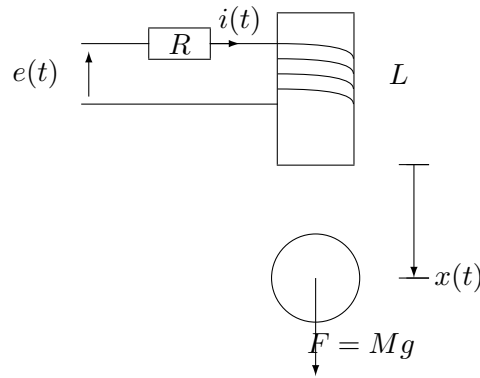


FIGURE 10.7 – Suspension magnétique

4.4 Commande du système

Placement de pôles

1. Déterminez les pôles complexes conjugués répondant au cahier des charges. Les autres pôles seront choisis tous identiques et négligeables devant les pôles dominants.
2. Mettez en place ce correcteur et vérifiez que le cahier des charges est bien respecté.

Commande optimale

1. Calculez le correcteur optimal avec les matrices de pondération suivantes : $\mathbf{R} = 1$ et $\mathbf{Q} =$

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Contrôle numérique

Après avoir déterminé la représentation d'état échantillonnée du système, reprenez les deux méthodes de correction du système précédentes et appliquez-les au cas discret.

4.5 Synthèse d'observateurs

Synthétisez un observateur complet de Luenberger sur le système précédent.

5 Asservissement d'une suspension magnétique

Principe : L'électro-aimant d'inductance L variable est parcouru par un courant $i(t)$ délivré par un générateur de tension $e(t)$ et de résistance interne R . Cet électro-aimant exerce une force d'attraction sur la bille, celle-ci se rapproche de l'électro-aimant et donc l'inductance augmente. La loi de variation de L est :

$$L(x) = \frac{L_0}{x}$$

5.1 Modélisation

1. A l'aide des équations de Lagrange, modélisez le système. Les coordonnées généralisées choisies seront :
 - q_1 = la charge électrique, \dot{q}_1 = le courant dans l'inductance,
 - q_2 = la position de la bille, \dot{q}_2 = la vitesse de la bille

2. Ce système étant non linéaire, la première étape consiste donc à le linéariser autour d'un point de fonctionnement par un développement limité au premier ordre ($a = a_0 + \hat{a}$).

Ainsi, si :

$$\ddot{q}_1 = F_1(q_1, q_2, \dot{q}_1, \dot{q}_2, e)$$

alors, au point de fonctionnement

$$\ddot{\hat{q}}_1 = \left. \frac{\partial F_1}{\partial q_1} \right|_E \hat{q}_1 + \left. \frac{\partial F_1}{\partial q_2} \right|_E \hat{q}_2 + \left. \frac{\partial F_1}{\partial \dot{q}_1} \right|_E \dot{\hat{q}}_1 + \left. \frac{\partial F_1}{\partial \dot{q}_2} \right|_E \dot{\hat{q}}_2 + \left. \frac{\partial F_1}{\partial e} \right|_E \hat{e}$$

3. En supposant qu'au point de fonctionnement la vitesse de la bille est nulle et que la tension moyenne appliquée est $e_0 = Ri_0$, montrez que la représentation d'état obtenue est :

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & -\frac{Rq_{20}}{L} & 0 & \frac{\dot{q}_{10}}{q_{20}} \\ 0 & 0 & 0 & 1 \\ 0 & -\frac{L\dot{q}_{10}}{q_{20}^2 M} & \frac{L\dot{q}_{10}^2}{q_{20}^3 M} & 0 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 0 \\ \frac{q_{20}}{L} \\ 0 \\ 0 \end{bmatrix}$$

4. Déterminez l'équation différentielle reliant la position et ses dérivées successives à la tension d'entrée.
5. Déterminez la fonction de transfert du système.

5.2 Analyse du système

1. Le système est-il stable ?
2. Le système est-il observable ? Est-il commandable ?

5.3 Correction du système

1. Donnez la représentation d'état du système sous forme diagonale, sous la forme canonique de commandabilité, sous la forme canonique d'observabilité.
2. Afin de stabiliser le système, proposez un choix de pôles en boucle fermée. Calculez le correcteur et la matrice de préfiltre.

5.4 Observation du système

1. Le système étant non observable, le correcteur précédent est inapplicable, proposez alors une autre solution.

5.5 Réduction du système

1. A partir de la fonction de transfert, déterminer une représentation d'état sous la forme canonique de commandabilité du système.
2. Le modèle ainsi obtenu est-il stable, commandable, observable ?

5.6 Synthèse de la commande

Applications numériques : $M = 0.1kg$, $i_0 = 1A$, $x_0 = 0,03m$, ; $L = 1H$, $R = 0.1\Omega$.

1. Avec les applications numériques précédentes, le système présente les pôles suivants :

$$+10.35, -5.18 - 8.97 * I, -5.18 + 8.97 * I$$

2. Choisissez les pôles en boucle fermée, calculez le correcteur.

5.7 Synthèse de l'observateur

1. Synthétisez un observateur complet pour ce système.

6 Commande numérique d'un AMF

Le système se compose d'un fil en Alliage à Mémoire de Forme (AMF) qui présente la propriété de se contracter lorsqu'il est chauffé. Ce fil de $150\mu\text{m}$ de diamètre est relié mécaniquement à un ressort de traction. Le principe de chauffage est l'effet Joule.

Les modèles, simplifiés à l'extrême sont donnés ci-après.

$$mC_p\dot{T} + hS(T - T_a) = Ri^2$$

$$\dot{x} + bx = aT$$

	$T_a = 300 \text{ K}$	température ambiante
	$h = 38 \text{ J m}^{-2}$	coefficient de convection
	$\phi = 150.10^{-6} \text{ m}$	diamètre du fil
	$L = 0.10 \text{ m}$	longueur du fil
avec :	$R = 5\Omega$	résistance électrique
	$\rho = 6450 \text{ kg/m}^{-3}$	densité volumique
	$C_p = 1886 \text{ K.kg}^{-1}$	chaleur spécifique
	$a = 5.10^{-6} \text{ m.K}^{-1}$	
	$b = 0.1 \text{ s}^{-1}$	

6.1 Modélisation et analyse

1. Donnez la représentation d'état du système.
2. Déduisez-en la représentation d'état échantillonnée du système.
3. Le système est-il stable, commandable, observable ?

6.2 Première commande

1. Au vu des pôles en boucle ouverte du système, proposez des pôles en boucle fermée.
2. Calculez le correcteur pour obtenir le comportement dynamique souhaité, ainsi que la matrice de préfiltre.
3. Dans quelles conditions, la matrice de préfiltre garantie-t-elle une erreur statique nulle ?

6.3 Amélioration de la commande

Afin de garantir une erreur statique nulle, y compris en cas de variation de la température ambiante, introduisons un intégrateur pur dans le système en amont du point d'application des perturbations.

1. Déterminez la nouvelle représentation d'état en z du système.
2. Calculez un correcteur pour ce nouveau système, présentant approximativement les mêmes caractéristiques que le précédent.

6.4 Synthèse de l'observateur

1. Calculez un observateur de Luenberger complet.
2. Après avoir déterminé les variables non accessibles à la mesure, déduisez-en un observateur d'ordre réduit.

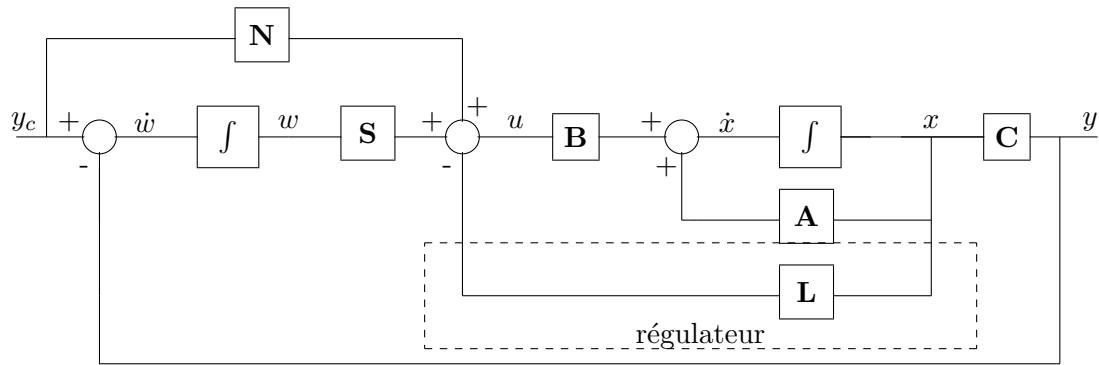


FIGURE 10.8 – Introduction d'un intégrateur pur.

6.5 Création du code C

1. Donnez la fonction de commande du système précédent dont la définition est donnée ci-après. Nous supposons que la mise à jour des tableaux est faite par ailleurs. Ainsi, à l'instant k , la sortie du système de l'instant $k - 1$ est stockée dans `sortie[0]`, $y(k - 2)$ est dans `sortie[1]` ...

```
float regulateur(float* commande, float* sortie)
{

return(commande)
}
```

6.6 Prise en compte d'un retard pur

1. La mise en place du code précédent montre que le temps de calcul est grand par rapport à la période d'échantillonnage, l'hypothèse de simultanéité temporelle de l'échantillonnage et de l'envoi de la commande est donc fautive. Afin de prendre en compte ce phénomène, introduisez un retard pur (z^{-1}) dans le système.
2. Reprenez la conception du correcteur et de l'observateur d'ordre réduit.

7 Bille sur un rail

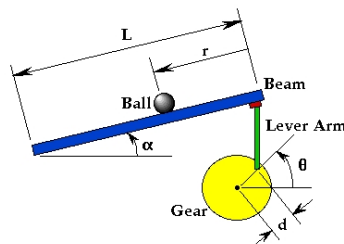


FIGURE 10.9 – Bille sur un rail

7.1 Modélisation

M	masse de la bille	0.11kg
R	rayon de la bille	0.015 m
d	rayon du réducteur	0.03 m
g	accélération terrestre	9.8 m/s ²
L	longueur du rail	1.0 m
J	moment d'inertie de la bille	1e - 5 kg.m ²
r	position de la bille	m
α	angle du rail	rad.
θ	angle du servomoteur	rad.

1. Avec les hypothèses suivantes :

- la commande u du système commande directement l'accélération du servomoteur $u = \ddot{\alpha}$,
 - l'angle α est faible, on considère que $\sin \alpha = 0$,
- montrez que la représentation d'état est de la forme :

$$\begin{pmatrix} \dot{r} \\ \ddot{r} \\ \dot{\alpha} \\ \ddot{\alpha} \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & \frac{-mg}{R^2+m} & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} r \\ \dot{r} \\ \alpha \\ \dot{\alpha} \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix} u \quad (10.14)$$

$$y = (1, 0, 0, 0) \begin{pmatrix} r \\ \dot{r} \\ \alpha \\ \dot{\alpha} \end{pmatrix} \quad (10.15)$$

2. De la représentation d'état continue, déduisez la représentation d'état discrète.

7.2 Analyse

1. Le système est-il stable ?
2. Le système est-il observable ? Est-il commandable ?
3. Que faire ? Dans quelle mesure les hypothèses formulées sont-elles sources de problèmes.

7.3 Commande

1. Calculer un régulateur numérique pour le système précédent.

7.4 Observation

1. Calculer un observateur numérique pour le système précédent.

8 Pendule inverse : Traité avec Matlab

8.1 Modélisation

Description et valeurs numériques :

M	masse du chariot	0.5 kg
m	masse du pendule	0.5kg
b	coefficient de frottement du chariot	0.1N.m ⁻¹ .s
l	longueur du pendule	0.3 m
I	inertie du pendule	0.006 kg*m ²
F	force appliquée sur le chariot	N
x	position du chariot	m
θ	angle du pendule	rad.

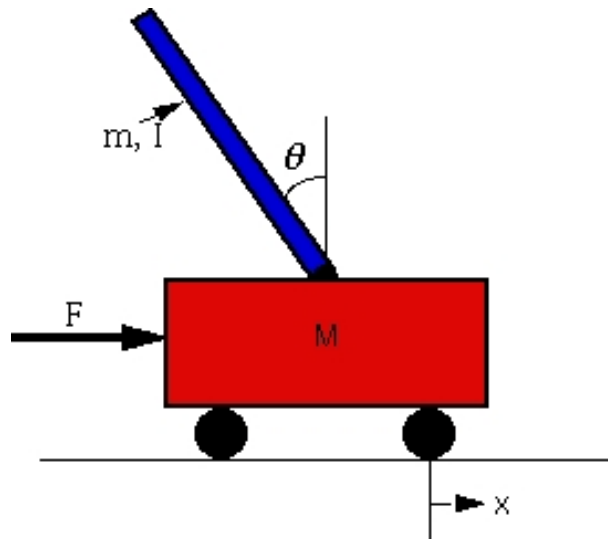


FIGURE 10.10 – Pendule inverse

1. Déterminez les équations différentielles qui régissent le mouvement des deux corps.

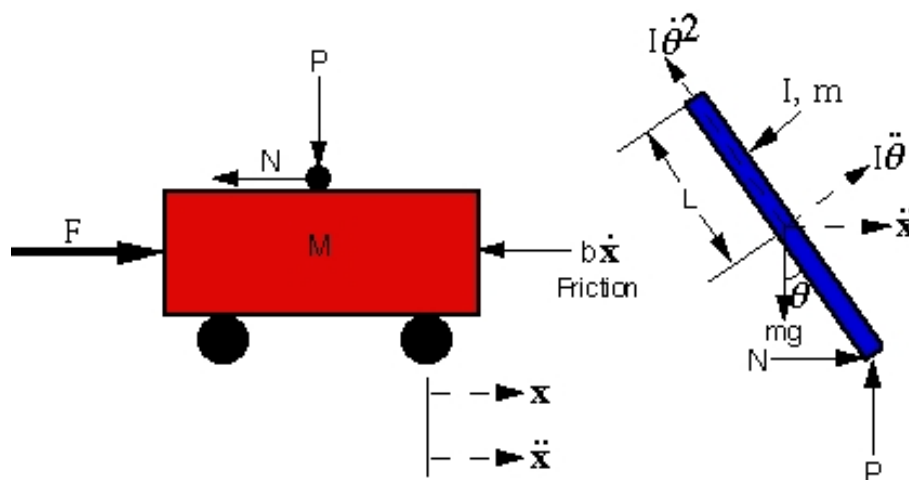


FIGURE 10.11 – Mise en équation pendule inverse

2. Après linéarisation des équations précédentes ($\theta \simeq \pi + \phi$ avec ϕ petit), déterminez la fonction de transfert du système.
3. Calculer la représentation d'état du système.
4. Vérifier que la représentation obtenue est :

$$\begin{bmatrix} \dot{x} \\ \dot{\phi} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ \phi \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u$$

$$\begin{bmatrix} \ddot{x} \\ \ddot{\phi} \end{bmatrix} = \begin{bmatrix} \frac{1}{I(M+m) + Mml^2} & \frac{0}{I(M+m) + Mml^2} \\ \frac{-(I+ml^2)b}{I(M+m) + Mml^2} & \frac{m^2gl^2}{I(M+m) + Mml^2} \end{bmatrix} \begin{bmatrix} \dot{x} \\ \dot{\phi} \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u$$

$$y = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \\ \phi \\ \dot{\phi} \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} u$$

FIGURE 10.12 – Représentation d'état du pendule inverse

8.2 Analyse

1. Le système est-il stable? Est-il observable?
2. Quelles variables faut-il mesurer pour que le système soit observable.

8.3 Commande

1. Calculer un régulateur par la méthode du placement de pôles qui respecte le cahier des charges suivant :
 - Temps de réponse sur θ inférieur à 5 s.
 - Temps de montée de x inférieur à 1 s.
 - Dépassement de moins de 0.035 radians.
 - Erreur statique nulle.
2. Comparer les réponses temporelles obtenues en boucle fermée par rapport à celle obtenues par une commande optimale obtenue avec :

$$\mathbf{R} = 5000, \quad \mathbf{Q} = \begin{pmatrix} 100 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 100 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

9 Exercices divers

9.1 Placement de pôles

Soit la fonction de transfert :

$$F(p) = \frac{1 + 2p}{(p - 1)((p + 1)^2 + 2)}$$

1. Déterminez directement la représentation d'état de cette fonction de transfert sous la forme de commandabilité.
2. On désire que le système en boucle fermée présente un pôle triple en -1, calculez le correcteur.
3. Par ailleurs, on désire que le système présente une erreur nulle pour une entrée en échelon, déterminer la matrice de préfiltre.

9.2 Passage de l'équation récurrente - représentation d'état - fonction de transfert

soit le système décrit par les équation récurrentes suivantes :

$$x_1(k + 1) = x_2(k), \tag{10.16}$$

$$x_2(k + 1) = x_3(k), \tag{10.17}$$

$$x_3(k + 1) = -4x_2(k) - 5x_3(k) + u(k), \tag{10.18}$$

$$y(k) = 4x_1(k) + 4x_2(k) + x_3(k). \tag{10.19}$$

1. Déterminez directement la représentation d'état du système décrit par les équation récurrentes précédentes.
2. Calculez la fonction de transfert du système.

Première partie

Annexes

Annexe A

Quelques publications originales

A New Approach to Linear Filtering and Prediction Problems¹

R. E. KALMAN

Research Institute for Advanced Study,²
Baltimore, Md.

The classical filtering and prediction problem is re-examined using the Bode-Shannon representation of random processes and the "state transition" method of analysis of dynamic systems. New results are:

(1) *The formulation and methods of solution of the problem apply without modification to stationary and nonstationary statistics and to growing-memory and infinite-memory filters.*

(2) *A nonlinear difference (or differential) equation is derived for the covariance matrix of the optimal estimation error. From the solution of this equation the coefficients of the difference (or differential) equation of the optimal linear filter are obtained without further calculations.*

(3) *The filtering problem is shown to be the dual of the noise-free regulator problem. The new method developed here is applied to two well-known problems, confirming and extending earlier results.*

The discussion is largely self-contained and proceeds from first principles; basic concepts of the theory of random processes are reviewed in the Appendix.

Introduction

AN IMPORTANT class of theoretical and practical problems in communication and control is of a statistical nature. Such problems are: (i) Prediction of random signals; (ii) separation of random signals from random noise; (iii) detection of signals of known form (pulses, sinusoids) in the presence of random noise.

In his pioneering work, Wiener [1]³ showed that problems (i) and (ii) lead to the so-called Wiener-Hopf integral equation; he also gave a method (spectral factorization) for the solution of this integral equation in the practically important special case of stationary statistics and rational spectra.

Many extensions and generalizations followed Wiener's basic work. Zadeh and Ragazzini solved the finite-memory case [2]. Concurrently and independently of Bode and Shannon [3], they also gave a simplified method [2] of solution. Booton discussed the nonstationary Wiener-Hopf equation [4]. These results are now in standard texts [5-6]. A somewhat different approach along these main lines has been given recently by Darlington [7]. For extensions to sampled signals, see, e.g., Franklin [8], Lees [9]. Another approach based on the eigenfunctions of the Wiener-Hopf equation (which applies also to nonstationary problems whereas the preceding methods in general don't), has been pioneered by Davis [10] and applied by many others, e.g., Shinbrot [11], Blum [12], Pugachev [13], Solodovnikov [14].

In all these works, the objective is to obtain the specification of a linear dynamic system (Wiener filter) which accomplishes the prediction, separation, or detection of a random signal.⁴

¹ This research was supported in part by the U. S. Air Force Office of Scientific Research under Contract AF 49 (638)-382.

² 7212 Bellona Ave.

³ Numbers in brackets designate References at end of paper.

⁴ Of course, in general these tasks may be done better by nonlinear filters. At present, however, little or nothing is known about how to obtain (both theoretically and practically) these nonlinear filters.

Contributed by the Instruments and Regulators Division and presented at the Instruments and Regulators Conference, March 29-April 2, 1959, of THE AMERICAN SOCIETY OF MECHANICAL ENGINEERS.

NOTE: Statements and opinions advanced in papers are to be understood as individual expressions of their authors and not those of the Society. Manuscript received at ASME Headquarters, February 24, 1959. Paper No. 59-IRD-11.

Present methods for solving the Wiener problem are subject to a number of limitations which seriously curtail their practical usefulness:

(1) The optimal filter is specified by its impulse response. It is not a simple task to synthesize the filter from such data.

(2) Numerical determination of the optimal impulse response is often quite involved and poorly suited to machine computation. The situation gets rapidly worse with increasing complexity of the problem.

(3) Important generalizations (e.g., growing-memory filters, nonstationary prediction) require new derivations, frequently of considerable difficulty to the nonspecialist.

(4) The mathematics of the derivations are not transparent. Fundamental assumptions and their consequences tend to be obscured.

This paper introduces a new look at this whole assemblage of problems, sidestepping the difficulties just mentioned. The following are the highlights of the paper:

(5) *Optimal Estimates and Orthogonal Projections.* The Wiener problem is approached from the point of view of conditional distributions and expectations. In this way, basic facts of the Wiener theory are quickly obtained; the scope of the results and the fundamental assumptions appear clearly. It is seen that all statistical calculations and results are based on first and second order averages; no other statistical data are needed. Thus difficulty (4) is eliminated. This method is well known in probability theory (see pp. 75-78 and 148-155 of Doob [15] and pp. 455-464 of Loève [16]) but has not yet been used extensively in engineering.

(6) *Models for Random Processes.* Following, in particular, Bode and Shannon [3], arbitrary random signals are represented (up to second order average statistical properties) as the output of a linear dynamic system excited by independent or uncorrelated random signals ("white noise"). This is a standard trick in the engineering applications of the Wiener theory [2-7]. The approach taken here differs from the conventional one only in the way in which linear dynamic systems are described. We shall emphasize the concepts of *state* and *state transition*; in other words, linear systems will be specified by systems of first-order difference (or differential) equations. This point of view is

natural and also necessary in order to take advantage of the simplifications mentioned under (5).

(7) *Solution of the Wiener Problem.* With the state-transition method, a single derivation covers a large variety of problems: growing and infinite memory filters, stationary and nonstationary statistics, etc.; difficulty (3) disappears. Having guessed the "state" of the estimation (i.e., filtering or prediction) problem correctly, one is led to a nonlinear difference (or differential) equation for the covariance matrix of the optimal estimation error. This is vaguely analogous to the Wiener-Hopf equation. Solution of the equation for the covariance matrix starts at the time t_0 when the first observation is taken; at each later time t the solution of the equation represents the covariance of the optimal prediction error given observations in the interval (t_0, t) . From the covariance matrix at time t we obtain at once, without further calculations, the coefficients (in general, time-varying) characterizing the optimal linear filter.

(8) *The Dual Problem.* The new formulation of the Wiener problem brings it into contact with the growing new theory of control systems based on the "state" point of view [17–24]. It turns out, *surprisingly*, that the Wiener problem is the *dual* of the noise-free optimal regulator problem, which has been solved previously by the author, using the state-transition method to great advantage [18, 23, 24]. The mathematical background of the two problems is identical—this has been suspected all along, but until now the analogies have never been made explicit.

(9) *Applications.* The power of the new method is most apparent in theoretical investigations and in numerical answers to complex practical problems. In the latter case, it is best to resort to machine computation. Examples of this type will be discussed later. To provide some feel for applications, two standard examples from nonstationary prediction are included; in these cases the solution of the nonlinear difference equation mentioned under (7) above can be obtained even in closed form.

For easy reference, the main results are displayed in the form of theorems. Only Theorems 3 and 4 are original. The next section and the Appendix serve mainly to review well-known material in a form suitable for the present purposes.

Notation Conventions

Throughout the paper, we shall deal mainly with *discrete* (or *sampled*) dynamic systems; in other words, signals will be observed at equally spaced points in time (*sampling instants*). By suitable choice of the time scale, the constant intervals between successive sampling instants (*sampling periods*) may be chosen as unity. Thus variables referring to time, such as t, t_0, τ, T will always be integers. The restriction to discrete dynamic systems is not at all essential (at least from the engineering point of view); by using the discreteness, however, we can keep the mathematics rigorous and yet elementary. Vectors will be denoted by small bold-face letters: $\mathbf{a}, \mathbf{b}, \dots, \mathbf{u}, \mathbf{x}, \mathbf{y}, \dots$. A *vector* or more precisely an *n-vector* is a set of n numbers x_1, \dots, x_n ; the x_i are the *co-ordinates* or *components* of the vector \mathbf{x} .

Matrices will be denoted by capital bold-face letters: $\mathbf{A}, \mathbf{B}, \mathbf{Q}, \Phi, \Psi, \dots$; they are $m \times n$ arrays of elements $a_{ij}, b_{ij}, q_{ij}, \dots$. The *transpose* (interchanging rows and columns) of a matrix will be denoted by the prime. In manipulating formulas, it will be convenient to regard a vector as a matrix with a single column.

Using the conventional definition of matrix multiplication, we write the *scalar product* of two n -vectors \mathbf{x}, \mathbf{y} as

$$\mathbf{x}'\mathbf{y} = \sum_{i=1}^n x_i y_i = \mathbf{y}'\mathbf{x}$$

The scalar product is clearly a scalar, i.e., not a vector, quantity.

Similarly, the quadratic form associated with the $n \times n$ matrix \mathbf{Q} is,

$$\mathbf{x}'\mathbf{Q}\mathbf{x} = \sum_{i,j=1}^n x_i q_{ij} x_j$$

We define the expression $\mathbf{x}\mathbf{y}'$ where \mathbf{x}' is an m -vector and \mathbf{y} is an n -vector to be the $m \times n$ matrix with elements $x_i y_j$.

We write $E(\mathbf{x}) = E\mathbf{x}$ for the expected value of the random vector \mathbf{x} (see Appendix). It is usually convenient to omit the brackets after E . This does not result in confusion in simple cases since constants and the operator E commute. Thus $E\mathbf{x}\mathbf{y}' =$ matrix with elements $E(x_i y_j)$; $E\mathbf{x}E\mathbf{y}' =$ matrix with elements $E(x_i)E(y_j)$.

For ease of reference, a list of the principal symbols used is given below.

Optimal Estimates

t	time in general, present time.
t_0	time at which observations start.
$x_1(t), x_2(t)$	basic random variables.
$y(t)$	observed random variable.
$x_1^*(t_1 t)$	optimal estimate of $x_1(t_1)$ given $y(t_0), \dots, y(t)$.
L	loss function (non random function of its argument).
ε	estimation error (random variable).

Orthogonal Projections

$\mathcal{Y}(t)$	linear manifold generated by the random variables $y(t_0), \dots, y(t)$.
$\bar{x}(t_1 t)$	orthogonal projection of $x(t_1)$ on $\mathcal{Y}(t)$.
$\tilde{x}(t_1 t)$	component of $x(t_1)$ orthogonal to $\mathcal{Y}(t)$.

Models for Random Processes

$\Phi(t+1; t)$	transition matrix
$\mathbf{Q}(t)$	covariance of random excitation

Solution of the Wiener Problem

$\mathbf{x}(t)$	basic random variable.
$\mathbf{y}(t)$	observed random variable.
$\mathcal{Y}(t)$	linear manifold generated by $\mathbf{y}(t_0), \dots, \mathbf{y}(t)$.
$\mathcal{Z}(t)$	linear manifold generated by $\mathbf{y}(t t-1)$.
$\mathbf{x}^*(t_1 t)$	optimal estimate of $\mathbf{x}(t_1)$ given $\mathcal{Y}(t)$.
$\tilde{\mathbf{x}}(t_1 t)$	error in optimal estimate of $\mathbf{x}(t_1)$ given $\mathcal{Y}(t)$.

Optimal Estimates

To have a concrete description or the type of problems to be studied, consider the following situation. We are given signal $x_1(t)$ and noise $x_2(t)$. Only the sum $y(t) = x_1(t) + x_2(t)$ can be observed. Suppose we have observed and know exactly the values of $y(t_0), \dots, y(t)$. What can we infer from this knowledge in regard to the (unobservable) value of the signal at $t = t_1$, where t_1 may be less than, equal to, or greater than t ? If $t_1 < t$, this is a *data-smoothing* (*interpolation*) problem. If $t_1 = t$, this is called *filtering*. If $t_1 > t$, we have a *prediction* problem. Since our treatment will be general enough to include these and similar problems, we shall use hereafter the collective term *estimation*.

As was pointed out by Wiener [1], the natural setting of the estimation problem belongs to the realm of probability theory and statistics. Thus signal, noise, and their sum will be random variables, and consequently they may be regarded as random processes. From the probabilistic description of the random processes we can determine the probability with which a particular sample of the signal and noise will occur. For any given set of measured values $\eta(t_0), \dots, \eta(t)$ of the random variable $y(t)$ one can then also determine, in principle, the probability of simultaneous occurrence of various values $\xi_1(t)$ of the random variable $x_1(t_1)$. This is the conditional probability distribution function

$$Pr[x_1(t_1) \leq \xi_1 | y(t_0) = \eta(t_0), \dots, y(t) = \eta(t)] = F(\xi_1) \quad (1)$$

Evidently, $F(\xi_1)$ represents all the information which the measurement of the random variables $y(t_0), \dots, y(t)$ has conveyed about the random variable $x_1(t_1)$. Any statistical estimate of the random variable $x_1(t_1)$ will be some function of this distribution and therefore a (nonrandom) function of the random variables $y(t_0), \dots, y(t)$. This statistical estimate is denoted by $X_1(t_1|t)$, or by just $X_1(t_1)$ or X_1 when the set of observed random variables or the time at which the estimate is required are clear from context.

Suppose now that X_1 is given as a fixed function of the random variables $y(t_0), \dots, y(t)$. Then X_1 is itself a random variable and its actual value is known whenever the actual values of $y(t_0), \dots, y(t)$ are known. In general, the actual value of $X_1(t_1)$ will be different from the (unknown) actual value of $x_1(t_1)$. To arrive at a rational way of determining X_1 , it is natural to assign a *penalty* or *loss* for incorrect estimates. Clearly, the loss should be a (i) positive, (ii) nondecreasing function of the *estimation error* $\varepsilon = x_1(t_1) - X_1(t_1)$. Thus we define a *loss function* by

$$\begin{aligned} L(0) &= 0 \\ L(\varepsilon_2) \geq L(\varepsilon_1) &\geq 0 \quad \text{when} \quad \varepsilon_2 \geq \varepsilon_1 \geq 0 \\ L(\varepsilon) &= L(-\varepsilon) \end{aligned} \quad (2)$$

Some common examples of loss functions are: $L(\varepsilon) = a\varepsilon^2$, $a\varepsilon^4$, $a|\varepsilon|$, $a[1 - \exp(-\varepsilon^2)]$, etc., where a is a positive constant.

One (but by no means the only) natural way of choosing the random variable X_1 is to require that this choice should minimize the average loss or risk

$$E\{L[x_1(t_1) - X_1(t_1)]\} = E\{E\{L[x_1(t_1) - X_1(t_1)] | y(t_0), \dots, y(t)\}\} \quad (3)$$

Since the first expectation on the right-hand side of (3) does not depend on the choice of X_1 but only on $y(t_0), \dots, y(t)$, it is clear that minimizing (3) is equivalent to minimizing

$$E\{L[x_1(t_1) - X_1(t_1)] | y(t_0), \dots, y(t)\} \quad (4)$$

Under just slight additional assumptions, optimal estimates can be characterized in a simple way.

Theorem 1. Assume that L is of type (2) and that the conditional distribution function $F(\xi)$ defined by (1) is:

(A) symmetric about the mean $\bar{\xi}$:

$$F(\xi - \bar{\xi}) = 1 - F(\bar{\xi} - \xi)$$

(B) convex for $\xi \leq \bar{\xi}$:

$$F(\lambda\xi_1 + (1-\lambda)\xi_2) \leq \lambda F(\xi_1) + (1-\lambda)F(\xi_2)$$

for all $\xi_1, \xi_2 \leq \bar{\xi}$ and $0 \leq \lambda \leq 1$

Then the random variable $x_1^*(t_1|t)$ which minimizes the average loss (3) is the conditional expectation

$$x_1^*(t_1|t) = E[x_1(t_1) | y(t_0), \dots, y(t)] \quad (5)$$

Proof: As pointed out recently by Sherman [25], this theorem follows immediately from a well-known lemma in probability theory.

Corollary. If the random processes $\{x_1(t)\}$, $\{x_2(t)\}$, and $\{y(t)\}$ are gaussian, Theorem 1 holds.

Proof: By Theorem 5, (A) (see Appendix), conditional distributions on a gaussian random process are gaussian. Hence the requirements of Theorem 1 are always satisfied.

In the control system literature, this theorem appears sometimes in a form which is more restrictive in one way and more general in another way:

Theorem 1-a. If $L(\varepsilon) = \varepsilon^2$, then Theorem 1 is true without assumptions (A) and (B).

Proof: Expand the conditional expectation (4):

$$E[x_1^2(t_1) | y(t_0), \dots, y(t)] - 2X_1(t_1)E[x_1(t_1) | y(t_0), \dots, y(t)] + X_1^2(t_1)$$

and differentiate with respect to $X_1(t_1)$. This is not a completely rigorous argument; for a simple rigorous proof see Doob [15], pp. 77–78.

Remarks. (a) As far as the author is aware, it is not known what is the most general class of random processes $\{x_1(t)\}$, $\{x_2(t)\}$ for which the conditional distribution function satisfies the requirements of Theorem 1.

(b) Aside from the note of Sherman, Theorem 1 apparently has never been stated explicitly in the control systems literature. In fact, one finds many statements to the effect that loss functions of the general type (2) cannot be conveniently handled mathematically.

(c) In the sequel, we shall be dealing mainly with vector-valued random variables. In that case, the estimation problem is stated as: Given a vector-valued random process $\{\mathbf{x}(t)\}$ and observed random variables $\mathbf{y}(t_0), \dots, \mathbf{y}(t)$, where $\mathbf{y}(t) = \mathbf{M}\mathbf{x}(t)$ (\mathbf{M} being a singular matrix; in other words, not all co-ordinates of $\mathbf{x}(t)$ can be observed), find an estimate $\mathbf{X}(t_1)$ which minimizes the expected loss $E\{L(\|\mathbf{x}(t_1) - \mathbf{X}(t_1)\|)\}$, $\|\cdot\|$ being the norm of a vector.

Theorem 1 remains true in the vector case also, provided we require that the conditional distribution function of the n co-ordinates of the vector $\mathbf{x}(t_1)$,

$$Pr\{x_1(t_1) \leq \xi_1, \dots, x_n(t_1) \leq \xi_n | \mathbf{y}(t_0), \dots, \mathbf{y}(t)\} = F(\xi_1, \dots, \xi_n)$$

be symmetric with respect to the n variables $\xi_1 - \bar{\xi}_1, \dots, \xi_n - \bar{\xi}_n$ and convex in the region where all of these variables are negative.

Orthogonal Projections

The explicit calculation of the optimal estimate as a function of the observed variables is, in general, impossible. There is an important exception: The processes $\{x_1(t)\}$, $\{x_2(t)\}$ are gaussian.

On the other hand, if we attempt to get an optimal estimate under the restriction $L(\varepsilon) = \varepsilon^2$ and the additional requirement that the estimate be a linear function of the observed random variables, we get an estimate which is identical with the optimal estimate in the gaussian case, without the assumption of linearity or quadratic loss function. This shows that results obtainable by linear estimation can be bettered by nonlinear estimation only when (i) the random processes are nongaussian and even then (in view of Theorem 5, (C)) only (ii) by considering at least third-order probability distribution functions.

In the special cases just mentioned, the explicit solution of the estimation problem is most easily understood with the help of a geometric picture. This is the subject of the present section.

Consider the (real-valued) random variables $y(t_0), \dots, y(t)$. The set of all linear combinations of these random variables with real coefficients

$$\sum_{i=t_0}^t a_i y(i) \quad (6)$$

forms a *vector space* (*linear manifold*) which we denote by $\mathcal{Y}(t)$. We regard, abstractly, any expression of the form (6) as “point” or “vector” in $\mathcal{Y}(t)$; this use of the word “vector” should not be confused, of course, with “vector-valued” random variables, etc. Since we do not want to fix the value of t (i.e., the total number of possible observations), $\mathcal{Y}(t)$ should be regarded as a finite-dimensional subspace of the space of all possible observations.

Given any two vectors u, v in $\mathcal{Y}(t)$ (i.e., random variables expressible in the form (6)), we say that u and v are *orthogonal* if $Euv = 0$. Using the Schmidt orthogonalization procedure, as described for instance by Doob [15], p. 151, or by Loève [16], p. 459, it is easy to select an *orthonormal basis* in $\mathcal{Y}(t)$. By this is meant a set of vectors e_{t_0}, \dots, e_t in $\mathcal{Y}(t)$ such that any vector in $\mathcal{Y}(t)$ can be expressed as a unique linear combination of e_{t_0}, \dots, e_t and

$$Ee_i e_j = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (i, j = t_0, \dots, t) \quad (7)$$

Thus any vector \bar{x} in $\mathcal{Y}(t)$ is given by

$$\bar{x} = \sum_{i=t_0}^t a_i e_i$$

and so the coefficients a_i can be immediately determined with the aid of (7):

$$E\bar{x}e_j = E\left(\sum_{i=t_0}^t a_i e_i\right)e_j = \sum_{i=t_0}^t a_i Ee_i e_j = \sum_{i=t_0}^t a_i \delta_{ij} = a_j \quad (8)$$

It follows further that any random variable x (not necessarily in $\mathcal{Y}(t)$) can be uniquely decomposed into two parts: a part \bar{x} in $\mathcal{Y}(t)$ and a part \tilde{x} orthogonal to $\mathcal{Y}(t)$ (i.e., orthogonal to every vector in $\mathcal{Y}(t)$). In fact, we can write

$$x = \bar{x} + \tilde{x} = \sum_{i=t_0}^t (Exe_i)e_i + \tilde{x} \quad (9)$$

Thus \bar{x} is uniquely determined by equation (9) and is obviously a vector in $\mathcal{Y}(t)$. Therefore \tilde{x} is also uniquely determined; it remains to check that it is orthogonal to $\mathcal{Y}(t)$:

$$E\tilde{x}e_i = E(x - \bar{x})e_i = Exe_i - E\bar{x}e_i$$

Now the co-ordinates of \bar{x} with respect to the basis e_{t_0}, \dots, e_t are given either in the form $E\bar{x}e_i$ (as in (8)) or in the form Exe_i (as in (9)). Since the co-ordinates are unique, $Exe_i = E\bar{x}e_i$ ($i = t_0, \dots, t$); hence $E\tilde{x}e_i = 0$ and \tilde{x} is orthogonal to every base vector e_i ; and therefore to $\mathcal{Y}(t)$. We call \bar{x} the *orthogonal projection* of x on $\mathcal{Y}(t)$.

There is another way in which the orthogonal projection can be characterized: \bar{x} is that vector in $\mathcal{Y}(t)$ (i.e., that linear function of the random variables $y(t_0), \dots, y(t)$) which minimizes the quadratic loss function. In fact, if \bar{w} is any other vector in $\mathcal{Y}(t)$, we have

$$E(x - \bar{w})^2 = E(\tilde{x} + \bar{x} - \bar{w})^2 = E[(x - \bar{x}) + (\bar{x} - \bar{w})]^2$$

Since \tilde{x} is orthogonal to every vector in $\mathcal{Y}(t)$ and in particular to $\bar{x} - \bar{w}$ we have

$$E(x - \bar{w})^2 = E(x - \bar{x})^2 + E(\bar{x} - \bar{w})^2 \geq E(x - \bar{x})^2 \quad (10)$$

This shows that, if \bar{w} also minimizes the quadratic loss, we must have $E(\bar{x} - \bar{w})^2 = 0$ which means that the random variables \bar{x} and \bar{w} are equal (except possibly for a set of events whose probability is zero).

These results may be summarized as follows:

Theorem 2. Let $\{x(t)\}, \{y(t)\}$ random processes with zero mean (i.e., $Ex(t) = Ey(t) = 0$ for all t). We observe $y(t_0), \dots, y(t)$. If either

- (A) the random processes $\{x(t)\}, \{y(t)\}$ are gaussian; or
- (B) the optimal estimate is restricted to be a linear function of the observed random variables and $L(\varepsilon) = \varepsilon^2$;

then

$$x^*(t_1|t) = \text{optimal estimate of } x(t_1) \text{ given } y(t_0), \dots, y(t) \\ = \text{orthogonal projection } \bar{x}(t_1|t) \text{ of } x(t_1) \text{ on } \mathcal{Y}(t). \quad (11)$$

These results are well-known though not easily accessible in the control systems literature. See Doob [15], pp. 75–78, or Pugachev [26]. It is sometimes convenient to denote the orthogonal projection by

$$\bar{x}(t_1|t) \equiv x^*(t_1|t) = \hat{E}[x(t_1)|\mathcal{Y}(t)]$$

The notation \hat{E} is motivated by part (b) of the theorem: If the stochastic processes in question are gaussian, then orthogonal projection is actually identical with conditional expectation.

Proof. (A) This is a direct consequence of the remarks in connection with (10).

(B) Since $x(t), y(t)$ are random variables with zero mean, it is clear from formula (9) that the orthogonal part $\tilde{x}(t_1|t)$ of $x(t_1)$ with respect to the linear manifold $\mathcal{Y}(t)$ is also a random variable with zero mean. Orthogonal random variables with zero mean are uncorrelated; if they are also gaussian then (by Theorem 5 (B)) they are independent. Thus

$$0 = E\tilde{x}(t_1|t) = E[\tilde{x}(t_1|t)|y(t_0), \dots, y(t)] \\ = E[x(t_1) - \bar{x}(t_1|t)|y(t_0), \dots, y(t)] \\ = E[x(t_1)|y(t_0), \dots, y(t)] - \bar{x}(t_1|t) = 0$$

Remarks. (d) A rigorous formulation of the contents of this section as $t \rightarrow \infty$ requires some elementary notions from the theory of Hilbert space. See Doob [15] and Loève [16].

(e) The physical interpretation of Theorem 2 is largely a matter of taste. If we are not worried about the assumption of gaussianity, part (A) shows that the orthogonal projection is the optimal estimate for all reasonable loss functions. If we do worry about gaussianity, even if we are resigned to consider only linear estimates, we know that orthogonal projections are *not* the optimal estimate for many reasonable loss functions. Since in practice it is difficult to ascertain to what degree of approximation a random process of physical origin is gaussian, it is hard to decide whether Theorem 2 has very broad or very limited significance.

(f) Theorem 2 is immediately generalized for the case of vector-valued random variables. In fact, we define the linear manifold $\mathcal{Y}(t)$ generated by $\mathbf{y}(t_0), \dots, \mathbf{y}(t)$ to be the set of all linear combinations

$$\sum_{i=t_0}^t \sum_{j=1}^m a_{ij} y_j(i)$$

of all m co-ordinates of each of the random vectors $\mathbf{y}(t_0), \dots, \mathbf{y}(t)$. The rest of the story proceeds as before.

(g) Theorem 2 states in effect that the optimal estimate under conditions (A) or (B) is a linear combination of all previous observations. In other words, the optimal estimate can be regarded as the output of a linear filter, with the input being the actually occurring values of the observable random variables; Theorem 2 gives a way of computing the impulse response of the optimal filter. As pointed out before, knowledge of this impulse response is not a complete solution of the problem; for this reason, no explicit formulas for the calculation of the impulse response will be given.

Models for Random Processes

In dealing with physical phenomena, it is not sufficient to give an empirical description but one must have also some idea of the underlying causes. Without being able to separate in some sense causes and effects, i.e., without the assumption of causality, one can hardly hope for useful results.

It is a fairly generally accepted fact that primary macroscopic sources of random phenomena are independent gaussian processes.⁵ A well-known example is the noise voltage produced in a resistor due to thermal agitation. In most cases, *observed* random phenomena are not describable by independent random variables. The statistical dependence (correlation) between random signals observed at different times is usually explained by the presence of a dynamic system between the primary random source and the observer. Thus a random function of time may be thought of as the output of a dynamic system excited by an independent gaussian random process.

An important property of gaussian random signals is that they remain gaussian after passing through a linear system (Theorem 5 (A)). Assuming independent gaussian primary random sources, if the observed random signal is also gaussian, we may assume that the dynamic system between the observer and the primary source is *linear*. This conclusion may be forced on us also because of lack of detailed knowledge of the statistical properties of the observed random signal: Given any random process with known first and second-order averages, we can find a gaussian random process with the same properties (Theorem 5 (C)). Thus gaussian distributions and linear dynamics are natural, mutually plausible assumptions particularly when the statistical data are scant.

How is a dynamic system (linear or nonlinear) described? The fundamental concept is the notion of the *state*. By this is meant, intuitively, some quantitative information (a set of numbers, a function, etc.) which is the least amount of data one has to know about the past behavior of the system in order to predict its future behavior. The dynamics is then described in terms of *state transitions*, i.e., one must specify how one state is transformed into another as time passes.

A linear dynamic system may be described in general by the vector differential equation

$$\left. \begin{aligned} \frac{d\mathbf{x}}{dt} &= \mathbf{F}(t)\mathbf{x} + \mathbf{D}(t)\mathbf{u}(t) \\ \mathbf{y}(t) &= \mathbf{M}(t)\mathbf{x}(t) \end{aligned} \right\} \quad (12)$$

where \mathbf{x} is an n -vector, the *state* of the system (the components x_i of \mathbf{x} are called *state variables*); $\mathbf{u}(t)$ is an m -vector ($m \leq n$) representing the *inputs* to the system; $\mathbf{F}(t)$ and $\mathbf{D}(t)$ are $n \times n$, respectively, $n \times m$ matrices. If all coefficients of $\mathbf{F}(t)$, $\mathbf{D}(t)$, $\mathbf{M}(t)$ are constants, we say that the dynamic system (12) is *time-invariant* or *stationary*. Finally, $\mathbf{y}(t)$ is a p -vector denoting the outputs of the system; $\mathbf{M}(t)$ is an $n \times p$ matrix; $p \leq n$

The physical interpretation of (12) has been discussed in detail elsewhere [18, 20, 23]. A look at the block diagram in Fig. 1 may be helpful. This is not an ordinary but a matrix block diagram (as revealed by the fat lines indicating signal flow). The integrator in

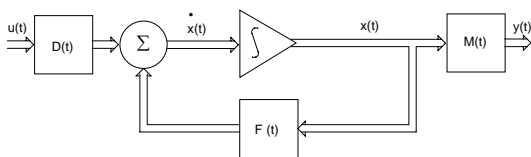


Fig. 1. Matrix block diagram of the general linear continuous-dynamic system

⁵ The probability distributions will be gaussian because macroscopic random effects may be thought of as the superposition of very many microscopic random effects; under very general conditions, such aggregate effects tend to be gaussian, regardless of the statistical properties of the microscopic effects. The assumption of independence in this context is motivated by the fact that microscopic phenomena tend to take place much more rapidly than macroscopic phenomena; thus primary random sources would appear to be independent on a macroscopic time scale.

Fig. 1 actually stands for n integrators such that the output of each is a state variable; $\mathbf{F}(t)$ indicates how the outputs of the integrators are fed back to the inputs of the integrators. Thus $f_{ij}(t)$ is the coefficient with which the output of the j th integrator is fed back to the input of the i th integrator. It is not hard to relate this formalism to more conventional methods of linear system analysis.

If we assume that the system (12) is stationary and that $\mathbf{u}(t)$ is constant during each sampling period, that is

$$\mathbf{u}(t + \tau) = \mathbf{u}(t); \quad 0 \leq \tau < 1, \quad t = 0, 1, \dots \quad (13)$$

then (12) can be readily transformed into the more convenient discrete form.

$$\mathbf{x}(t + 1) = \mathbf{\Phi}(1)\mathbf{x}(t) + \mathbf{\Delta}(1)\mathbf{u}(t); \quad t = 0, 1, \dots$$

where [18, 20]

$$\mathbf{\Phi}(1) = \exp \mathbf{F} = \sum_{i=0}^{\infty} \mathbf{F}^i / i! \quad (\mathbf{F}^0 = \text{unit matrix})$$

and

$$\mathbf{\Delta}(1) = \left(\int_0^1 \exp \mathbf{F} \tau d\tau \right) \mathbf{D}$$

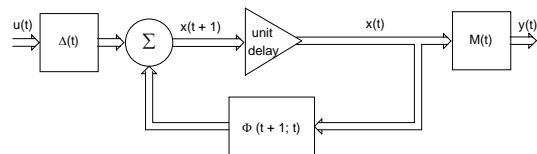


Fig. 2. Matrix block diagram of the general linear discrete-dynamic system

See Fig. 2. One could also express $\exp \mathbf{F} \tau$ in closed form using Laplace transform methods [18, 20, 22, 24]. If $\mathbf{u}(t)$ satisfies (13) but the system (12) is nonstationary, we can write analogously

$$\left. \begin{aligned} \mathbf{x}(t + 1) &= \mathbf{\Phi}(t + 1; t) + \mathbf{\Delta}(t)\mathbf{u}(t) \\ \mathbf{y}(t) &= \mathbf{M}(t)\mathbf{x}(t) \end{aligned} \right\} \quad t = 0, 1, \dots \quad (14)$$

but of course now $\mathbf{\Phi}(t + 1; t)$, $\mathbf{\Delta}(t)$ cannot be expressed in general in closed form. Equations of type (14) are encountered frequently also in the study of complicated sampled-data systems [22]. See Fig. 2

$\mathbf{\Phi}(t + 1; t)$ is the *transition matrix* of the system (12) or (14). The notation $\mathbf{\Phi}(t_2; t_1)$ ($t_2, t_1 = \text{integers}$) indicates transition from time t_1 to time t_2 . Evidently $\mathbf{\Phi}(t; t) = \mathbf{I}$ = unit matrix. If the system (12) is stationary then $\mathbf{\Phi}(t + 1; t) = \mathbf{\Phi}(t + 1 - t) = \mathbf{\Phi}(1) = \text{const}$. Note also the product rule: $\mathbf{\Phi}(t; s)\mathbf{\Phi}(s; r) = \mathbf{\Phi}(t; r)$ and the inverse rule $\mathbf{\Phi}^{-1}(t; s) = \mathbf{\Phi}(s; t)$, where t, s, r are integers. In a stationary system, $\mathbf{\Phi}(t; \tau) = \exp \mathbf{F}(t - \tau)$.

As a result of the preceding discussion, we shall represent random phenomena by the model

$$\mathbf{x}(t + 1) = \mathbf{\Phi}(t + 1; t)\mathbf{x}(t) + \mathbf{u}(t) \quad (15)$$

where $\{\mathbf{u}(t)\}$ is a vector-valued, independent, gaussian random process, with zero mean, which is completely described by (in view of Theorem 5 (C))

$$E\mathbf{u}(t) = \mathbf{0} \quad \text{for all } t;$$

$$E\mathbf{u}(t)\mathbf{u}^1(s) = \mathbf{0} \quad \text{if } t \neq s$$

$$E\mathbf{u}(t)\mathbf{u}^1(t) = \mathbf{G}(t).$$

Of course (Theorem 5 (A)), $\mathbf{x}(t)$ is then also a gaussian random process with zero mean, but it is no longer independent. In fact, if we consider (15) in the steady state (assuming it is a stable system), in other words, if we neglect the initial state $\mathbf{x}(t_0)$, then

$$\mathbf{x}(t) = \sum_{r=-\infty}^{t-1} \Phi(t; r+1)\mathbf{u}(r).$$

Therefore if $t \geq s$ we have

$$E\mathbf{x}(t)\mathbf{x}'(s) = \sum_{r=-\infty}^{s-1} \Phi(t; r+1)\mathbf{Q}(r)\Phi'(s; r+1).$$

Thus if we assume a linear dynamic model and know the statistical properties of the gaussian random excitation, it is easy to find the corresponding statistical properties of the gaussian random process $\{\mathbf{x}(t)\}$.

In real life, however, the situation is usually reversed. One is given the covariance matrix $E\mathbf{x}(t)\mathbf{x}'(s)$ (or rather, one attempts to estimate the matrix from limited statistical data) and the problem is to get (15) and the statistical properties of $\mathbf{u}(t)$. This is a subtle and presently largely unsolved problem in experimentation and data reduction. As in the vast majority of the engineering literature on the Wiener problem, we shall find it convenient to start with the model (15) and regard the problem of obtaining the model itself as a separate question. To be sure, the two problems *should* be optimized jointly if possible; the author is not aware, however, of any study of the *joint* optimization problem.

In summary, the following assumptions are made about random processes:

Physical random phenomena may be thought of as due to primary random sources exciting dynamic systems. The primary sources are assumed to be independent gaussian random processes with zero mean; the dynamic systems will be linear. The random processes are therefore described by models such as (15). The question of how the numbers specifying the model are obtained from experimental data will not be considered.

Solution of the Wiener problem

Let us now define the principal problem of the paper.

Problem I. Consider the dynamic model

$$\mathbf{x}(t+1) = \Phi(t+1; t)\mathbf{x}(t) + \mathbf{u}(t) \quad (16)$$

$$\mathbf{y}(t) = \mathbf{M}(t)\mathbf{x}(t) \quad (17)$$

where $\mathbf{u}(t)$ is an independent gaussian random process of n -vectors with zero mean, $\mathbf{x}(t)$ is an n -vector, $\mathbf{y}(t)$ is a p -vector ($p \leq n$), $\Phi(t+1; t)$, $\mathbf{M}(t)$ are $n \times n$, resp. $p \times n$, matrices whose elements are nonrandom functions of time.

Given the observed values of $\mathbf{y}(t_0), \dots, \mathbf{y}(t)$ find an estimate $\mathbf{x}^*(t_1|t)$ of $\mathbf{x}(t_1)$ which minimizes the expected loss. (See Fig. 2, where $\Delta(t) = \mathbf{I}$.)

This problem includes as a special case the problems of filtering, prediction, and data smoothing mentioned earlier. It includes also the problem of reconstructing all the state variables of a linear dynamic system from noisy observations of some of the state variables ($p < n$!).

From Theorem 2-a we know that the solution of Problem I is simply the orthogonal projection of $\mathbf{x}(t_1)$ on the linear manifold $\mathcal{Y}(t)$ generated by the observed random variables. As remarked in the Introduction, this is to be accomplished by means of a linear (not necessarily stationary!) dynamic system of the general form (14). With this in mind, we proceed as follows.

Assume that $\mathbf{y}(t_0), \dots, \mathbf{y}(t-1)$ have been measured, i.e., that $\mathcal{Y}(t-1)$ is known. Next, at time t , the random variable $\mathbf{y}(t)$ is measured. As before let $\tilde{\mathbf{y}}(t|t-1)$ be the component of $\mathbf{y}(t)$ orthogonal to $\mathcal{Y}(t-1)$. If $\tilde{\mathbf{y}}(t|t-1) \equiv 0$, which means that the values of all components of this random vector are zero for almost every possible event, then $\mathcal{Y}(t)$ is obviously the same as $\mathcal{Y}(t-1)$ and therefore the measurement of $\mathbf{y}(t)$ does not convey any additional information. This is not likely to happen in a physically meaningful situation. In any case, $\tilde{\mathbf{y}}(t|t-1)$ generates a linear

manifold (possibly 0) which we denote by $Z(t)$. By definition, $\mathcal{Y}(t-1)$ and $Z(t)$ taken together are the same manifold as $\mathcal{Y}(t)$, and every vector in $Z(t)$ is orthogonal to every vector in $\mathcal{Y}(t-1)$.

Assuming by induction that $\mathbf{x}^*(t_1-1|t-1)$ is known, we can write:

$$\begin{aligned} \mathbf{x}^*(t_1|t) &= \hat{E}[\mathbf{x}(t_1)|\mathcal{Y}(t)] = \hat{E}[\mathbf{x}(t_1)|\mathcal{Y}(t-1)] + \hat{E}[\mathbf{x}(t_1)|Z(t)] \\ &= \Phi(t+1; t)\mathbf{x}^*(t_1-1|t-1) + \hat{E}[\mathbf{u}(t_1-1)|\mathcal{Y}(t-1)] \\ &\quad + \hat{E}[\mathbf{x}(t_1)|Z(t)] \quad (18) \end{aligned}$$

where the last line is obtained using (16).

Let $t_1 = t + s$, where s is any integer. If $s \geq 0$, then $\mathbf{u}(t_1-1)$ is independent of $\mathcal{Y}(t-1)$. This is because $\mathbf{u}(t_1-1) = \mathbf{u}(t+s-1)$ is then independent of $\mathbf{u}(t-2), \mathbf{u}(t-3), \dots$ and therefore by (16-17), independent of $\mathbf{y}(t_0), \dots, \mathbf{y}(t-1)$, hence independent of $\mathcal{Y}(t-1)$. Since, for all t , $\mathbf{u}(t_0)$ has zero mean by assumption, it follows that $\mathbf{u}(t_1-1)$ ($s \geq 0$) is orthogonal to $\mathcal{Y}(t-1)$. Thus if $s \geq 0$, the second term on the right-hand side of (18) vanishes; if $s < 0$, considerable complications result in evaluating this term. We shall consider only the case $t_1 \geq t$. Furthermore, it will suffice to consider in detail only the case $t_1 = t+1$ since the other cases can be easily reduced to this one.

The last term in (18) must be a linear operation on the random variable $\tilde{\mathbf{y}}(t|t-1)$:

$$\hat{E}[\mathbf{x}(t+1)|Z(t)] = \Delta^*(t)\tilde{\mathbf{y}}(t|t-1) \quad (19)$$

where $\Delta^*(t)$ is an $n \times p$ matrix, and the star refers to "optimal filtering".

The component of $\mathbf{y}(t)$ lying in $\mathcal{Y}(t-1)$ is $\bar{\mathbf{y}}(t|t-1) = \mathbf{M}(t)\mathbf{x}^*(t|t-1)$. Hence

$$\tilde{\mathbf{y}}(t|t-1) = \mathbf{y}(t) - \bar{\mathbf{y}}(t|t-1) = \mathbf{y}(t) - \mathbf{M}(t)\mathbf{x}^*(t|t-1). \quad (20)$$

Combining (18-20) (see Fig. 3) we obtain

$$\mathbf{x}^*(t+1|t) = \Phi^*(t+1; t)\mathbf{x}^*(t|t-1) + \Delta^*(t)\mathbf{y}(t) \quad (21)$$

where

$$\Phi^*(t+1; t) = \Phi(t+1; t) - \Delta^*(t)\mathbf{M}(t) \quad (22)$$

Thus optimal estimation is performed by a linear dynamic system of the same form as (14). The state of the estimator is the previous estimate, the input is the last measured value of the observable random variable $\mathbf{y}(t)$, the transition matrix is given by (22). Notice that physical realization of the optimal filter requires only (i) the model of the random process (ii) the operator $\Delta^*(t)$.

The estimation error is also governed by a linear dynamic system. In fact,

$$\begin{aligned} \tilde{\mathbf{x}}(t+1|t) &= \mathbf{x}(t+1) - \mathbf{x}^*(t+1|t) \\ &= \Phi(t+1; t)\mathbf{x}(t) + \mathbf{u}(t) - \Phi^*(t+1; t)\mathbf{x}^*(t|t-1) \\ &\quad - \Delta^*(t)\mathbf{M}(t)\mathbf{x}(t) \end{aligned}$$

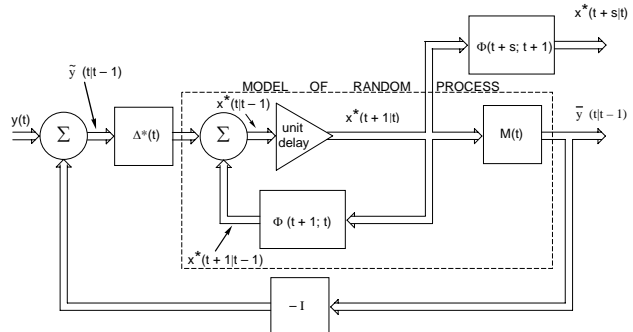


Fig. 3 Matrix block diagram of optimal filter

$$= \Phi^*(t+1; t) \tilde{\mathbf{x}}(t|t-1) + \mathbf{u}(t) \quad (23)$$

Thus Φ^* is also the transition matrix of the linear dynamic system governing the error.

From (23) we obtain at once a recursion relation for the covariance matrix $\mathbf{P}^*(t)$ of the optimal error $\tilde{\mathbf{x}}(t|t-1)$. Noting that $\mathbf{u}(t)$ is independent of $\mathbf{x}(t)$ and therefore of $\tilde{\mathbf{x}}(t|t-1)$ we get

$$\begin{aligned} \mathbf{P}^*(t+1) &= E \tilde{\mathbf{x}}(t+1|t) \tilde{\mathbf{x}}'(t+1|t) \\ &= \Phi^*(t+1; t) E \tilde{\mathbf{x}}(t|t-1) \tilde{\mathbf{x}}'(t|t-1) \Phi^{*\prime}(t+1; t) + \mathbf{Q}(t) \\ &= \Phi^*(t+1; t) E \tilde{\mathbf{x}}(t|t-1) \tilde{\mathbf{x}}'(t|t-1) \Phi'(t+1; t) + \mathbf{Q}(t) \\ &= \Phi^*(t+1; t) \mathbf{P}^*(t) \Phi'(t+1; t) + \mathbf{Q}(t) \end{aligned} \quad (24)$$

where $\mathbf{Q}(t) = E\mathbf{u}(t)\mathbf{u}'(t)$.

There remains the problem of obtaining an explicit formula for Δ^* (and thus also for Φ^*). Since,

$$\tilde{\mathbf{x}}(t+1|Z(t)) = \mathbf{x}(t+1) - \hat{E}[\mathbf{x}(t+1)|Z(t)]$$

is orthogonal to $\tilde{\mathbf{y}}(t|t-1)$, it follows that by (19) that

$$\begin{aligned} 0 &= E[\mathbf{x}(t+1) - \Delta^*(t) \tilde{\mathbf{y}}(t|t-1)] \tilde{\mathbf{y}}'(t|t-1) \\ &= E\mathbf{x}(t+1) \tilde{\mathbf{y}}'(t|t-1) - \Delta^*(t) E \tilde{\mathbf{y}}(t|t-1) \tilde{\mathbf{y}}'(t|t-1). \end{aligned}$$

Noting that $\tilde{\mathbf{x}}(t+1|t-1)$ is orthogonal to $Z(t)$, the definition of $\mathbf{P}(t)$ given earlier, and (17), it follows further

$$\begin{aligned} 0 &= E \tilde{\mathbf{x}}(t+1|t-1) \tilde{\mathbf{y}}'(t|t-1) - \Delta^*(t) \mathbf{M}(t) \mathbf{P}^*(t) \mathbf{M}'(t) \\ &= E[\Phi(t+1; t) \tilde{\mathbf{x}}(t|t-1) + \mathbf{u}(t|t-1)] \tilde{\mathbf{x}}'(t|t-1) \mathbf{M}'(t) \\ &\quad - \Delta^*(t) \mathbf{M}(t) \mathbf{P}^*(t) \mathbf{M}'(t). \end{aligned}$$

Finally, since $\mathbf{u}(t)$ is independent of $\mathbf{x}(t)$,

$$0 = \Phi(t+1; t) \mathbf{P}^*(t) \mathbf{M}'(t) - \Delta^*(t) \mathbf{M}(t) \mathbf{P}^*(t) \mathbf{M}'(t).$$

Now the matrix $\mathbf{M}(t) \mathbf{P}^*(t) \mathbf{M}'(t)$ will be positive definite and hence invertible whenever $\mathbf{P}^*(t)$ is positive definite, provided that none of the rows of $\mathbf{M}(t)$ are linearly dependent at any time, in other words, that none of the observed scalar random variables $y_1(t), \dots, y_m(t)$, is a linear combination of the others. Under these circumstances we get finally:

$$\Delta^*(t) = \Phi(t+1; t) \mathbf{P}^*(t) \mathbf{M}'(t) [\mathbf{M}(t) \mathbf{P}^*(t) \mathbf{M}'(t)]^{-1} \quad (25)$$

Since observations start at t_0 , $\tilde{\mathbf{x}}(t_0|t_0-1) = \mathbf{x}(t_0)$; to begin the iterative evaluation of $\mathbf{P}^*(t)$ by means of equation (24), we must obviously specify $\mathbf{P}^*(t_0) = E\mathbf{x}(t_0)\mathbf{x}'(t_0)$. Assuming this matrix is positive definite, equation (25) then yields $\Delta^*(t_0)$; equation (22) $\Phi^*(t_0+1; t_0)$, and equation (24) $\mathbf{P}^*(t_0+1)$, completing the cycle. If now $\mathbf{Q}(t)$ is positive definite, then all the $\mathbf{P}^*(t)$ will be positive definite and the requirements in deriving (25) will be satisfied at each step.

Now we remove the restriction that $t_1 = t+1$. Since $\mathbf{u}(t)$ is orthogonal to $\mathbf{y}(t)$, we have

$$\mathbf{x}^*(t+1|t) = \hat{E}[\Phi(t+1; t)\mathbf{x}(t) + \mathbf{u}(t)|\mathbf{y}(t)] = \Phi(t+1; t)\mathbf{x}^*(t|t)$$

Hence if $\Phi(t+1; t)$ has an inverse $\Phi(t; t+1)$ (which is always the case when Φ is the transition matrix of a dynamic system describable by a differential equation) we have

$$\mathbf{x}^*(t|t) = \Phi(t; t+1)\mathbf{x}^*(t+1|t)$$

If $t_1 \geq t+1$, we first observe by repeated application of (16) that

$$\begin{aligned} \mathbf{x}(t+s) &= \Phi(t+s; t+1)\mathbf{x}(t+1) \\ &\quad + \sum_{r=1}^{s-1} \Phi(t+s; t+r)\mathbf{u}(t+r) \end{aligned} \quad (s \geq 1)$$

Since $\mathbf{u}(t+s-1), \dots, \mathbf{u}(t+1)$ are all orthogonal to $\mathbf{y}(t)$,

$$\begin{aligned} \mathbf{x}^*(t+s|t) &= \hat{E}[\mathbf{x}(t+s)|\mathbf{y}(t)] \\ &= \hat{E}[\Phi(t+s; t+1)\mathbf{x}(t+1)|\mathbf{y}(t)] \\ &= \Phi(t+s; t+1)\mathbf{x}^*(t+1|t) \quad (s \geq 1) \end{aligned}$$

If $s < 0$, the results are similar, but $\mathbf{x}^*(t-s|t)$ will have $(1-s)(n-p)$ co-ordinates.

The results of this section may be summarized as follows:

Theorem 3. (Solution of the Wiener Problem)

Consider Problem I. The optimal estimate $\mathbf{x}^*(t+1|t)$ of $\mathbf{x}(t+1)$ given $\mathbf{y}(t_0), \dots, \mathbf{y}(t)$ is generated by the linear dynamic system

$$\mathbf{x}^*(t+1|t) = \Phi^*(t+1; t)\mathbf{x}^*(t|t-1) + \Delta^*(t)\mathbf{y}(t) \quad (21)$$

The estimation error is given by

$$\tilde{\mathbf{x}}(t+1|t) = \Phi^*(t+1; t) \tilde{\mathbf{x}}(t|t-1) + \mathbf{u}(t) \quad (23)$$

The covariance matrix of the estimation error is

$$\text{cov } \tilde{\mathbf{x}}(t|t-1) = E \tilde{\mathbf{x}}(t|t-1) \tilde{\mathbf{x}}'(t|t-1) = \mathbf{P}^*(t) \quad (26)$$

The expected quadratic loss is

$$\sum_{i=1}^n E\tilde{x}_i^2(t|t-1) = \text{trace } \mathbf{P}^*(t) \quad (27)$$

The matrices $\Delta^*(t)$, $\Phi^*(t+1; t)$, $\mathbf{P}^*(t)$ are generated by the recursion relations

$$\Delta^*(t) = \Phi(t+1; t) \mathbf{P}^*(t) \mathbf{M}'(t) [\mathbf{M}(t) \mathbf{P}^*(t) \mathbf{M}'(t)]^{-1} \quad (28)$$

$$\Phi^*(t+1; t) = \Phi(t+1; t) - \Delta^*(t) \mathbf{M}(t) \quad (29)$$

$$\mathbf{P}^*(t+1) = \Phi^*(t+1; t) \mathbf{P}^*(t) \Phi'(t+1; t) + \mathbf{Q}(t) \quad (30)$$

In order to carry out the iterations, one must specify the covariance $\mathbf{P}^*(t_0)$ of $\mathbf{x}(t_0)$ and the covariance $\mathbf{Q}(t)$ of $\mathbf{u}(t)$. Finally, for any $s \geq 0$, if Φ is invertible

$$\begin{aligned} \mathbf{x}^*(t+s|t) &= \Phi(t+s; t+1)\mathbf{x}^*(t+1|t) \\ &= \Phi(t+s; t+1)\Phi^*(t+1; t)\Phi(t; t+s-1) \\ &\quad \times \mathbf{x}^*(t+s-1|t-1) \\ &\quad + \Phi(t+s; t+1)\Delta^*(t)\mathbf{y}(t) \end{aligned} \quad (31)$$

so that the estimate $\mathbf{x}^*(t+s|t)$ ($s \geq 0$) is also given by a linear dynamic system of the type (21).

Remarks. (h) Eliminating Δ^* and Φ^* from (28–30), a nonlinear difference equation is obtained for $\mathbf{P}^*(t)$:

$$\mathbf{P}^*(t+1) = \Phi(t+1; t) \{ \mathbf{P}^*(t) - \mathbf{P}^*(t) \mathbf{M}'(t) [\mathbf{M}(t) \mathbf{P}^*(t) \mathbf{M}'(t)]^{-1} \times \mathbf{P}^*(t) \mathbf{M}(t) \} \Phi'(t+1; t) + \mathbf{Q}(t) \quad t \geq t_0 \quad (32)$$

This equation is linear only if $\mathbf{M}(t)$ is invertible but then the problem is trivial since all components of the random vector $\mathbf{x}(t)$ are observable $\mathbf{P}^*(t+1) = \mathbf{Q}(t)$. Observe that equation (32) plays a role in the present theory analogous to that of the Wiener-Hopf equation in the conventional theory.

Once $\mathbf{P}^*(t)$ has been computed via (32) starting at $t = t_0$, the explicit specification of the optimal linear filter is immediately available from formulas (29–30). Of course, the solution of Equation (32), or of its differential-equation equivalent, is a much simpler task than solution of the Wiener-Hopf equation.

(i) The results stated in Theorem 3 do not resolve completely Problem I. Little has been said, for instance, about the physical significance of the assumptions needed to obtain equation (25), the convergence and stability of the nonlinear difference equation (32), the stability of the optimal filter (21), etc. This can actually be done in a completely satisfactory way, but must be left to a future paper. In this connection, the principal guide and

tool turns out to be the duality theorem mentioned briefly in the next section. See [29].

(j) By letting the sampling period (equal to one so far) approach zero, the method can be used to obtain the specification of a differential equation for the optimal filter. To do this, i.e., to pass from equation (14) to equation (12), requires computing the logarithm \mathbf{F}^* of the matrix Φ^* . But this can be done only if Φ^* is nonsingular—which is easily seen *not* to be the case. This is because it is sufficient for the optimal filter to have $n - p$ state variables, rather than n , as the formalism of equation (22) would seem to imply. By appropriate modifications, therefore, equation (22) can be reduced to an equivalent set of only $n - p$ equations whose transition matrix is nonsingular. Details of this type will be covered in later publications.

(k) The dynamic system (21) is, in general, nonstationary. This is due to two things: (1) The time dependence of $\Phi(t + 1; t)$ and $\mathbf{M}(t)$; (2) the fact that the estimation starts at $t = t_0$ and improves as more data are accumulated. If Φ, \mathbf{M} are constants, it can be shown that (21) becomes a stationary dynamic system in the limit $t \rightarrow \infty$. This is the case treated by the classical Wiener theory.

(l) It is noteworthy that the derivations given are not affected by the nonstationarity of the model for $\mathbf{x}(t)$ or the finiteness of available data. In fact, as far as the author is aware, the only explicit recursion relations given before for the growing-memory filter are due to Blum [12]. However, his results are much more complicated than ours.

(m) By inspection of Fig. 3 we see that the optimal filter is a feedback system, and that the signal after the first summer is white noise since $\tilde{\mathbf{y}}(t|t - 1)$ is obviously an orthogonal random process. This corresponds to some well-known results in Wiener filtering, see, e.g., Smith [28], Chapter 6, Fig. 6-4. However, this is apparently the first *rigorous* proof that every Wiener filter is realizable by means of a feedback system. Moreover, it will be shown in another paper that such a filter is always *stable*, under very mild assumptions on the model (16-17). See [29].

The Dual Problem

Let us now consider another problem which is conceptually very different from optimal estimation, namely, the noise-free regulator problem. In the simplest cases, this is:

Problem II. Consider the dynamic system

$$\mathbf{x}(t + 1) = \hat{\Phi}(t + 1; t)\mathbf{x}(t) + \hat{\mathbf{M}}(t)\mathbf{u}(t) \quad (33)$$

where $\mathbf{x}(t)$ is an n -vector, $\mathbf{u}(t)$ is an m -vector ($m \leq n$), $\hat{\Phi}, \hat{\mathbf{M}}$ are $n \times n$ resp. $n \times m$ matrices whose elements are nonrandom functions of time. Given any state $\mathbf{x}(t)$ at time t , we are to find a sequence $\mathbf{u}(t), \dots, \mathbf{u}(T)$ of control vectors which minimizes the performance index

$$V[\mathbf{x}(t)] = \sum_{\tau=t}^{T+1} \mathbf{x}'(\tau)\mathbf{Q}(\tau)\mathbf{x}(\tau)$$

Where $\hat{\mathbf{Q}}(t)$ is a positive definite matrix whose elements are nonrandom functions of time. See Fig. 2, where $\hat{\Delta} = \hat{\mathbf{M}}$ and $\mathbf{M} = \mathbf{I}$.

Probabilistic considerations play no part in Problem II; it is implicitly assumed that every state variable can be measured exactly at each instant $t, t + 1, \dots, T$. It is customary to call $T \geq t$ the *terminal time* (it may be infinity).

The first general solution of the noise-free regulator problem is due to the author [18]. The main result is that the optimal control vectors $\mathbf{u}^*(t)$ are nonstationary linear functions of $\mathbf{x}(t)$. After a change in notation, the formulas of the Appendix, Reference [18] (see also Reference [23]) are as follows:

$$\mathbf{u}^*(t) = -\hat{\Delta}^*(t)\mathbf{x}(t) \quad (34)$$

Under optimal control as given by (34), the "closed-loop" equations for the system are (see Fig. 4)

$$\mathbf{x}(t + 1) = \hat{\Phi}^*(t + 1; t)\mathbf{x}(t)$$

and the minimum performance index at time t is given by

$$V^*[\mathbf{x}(t)] = \mathbf{x}'(t)\mathbf{P}^*(t - 1)\mathbf{x}(t)$$

The matrices $\hat{\Delta}^*(t), \hat{\Phi}^*(t + 1; t), \hat{\mathbf{P}}^*(t)$ are determined by the recursion relations:

$$\hat{\Delta}^*(t) = [\hat{\mathbf{M}}'(t) \hat{\mathbf{P}}^*(t) \hat{\mathbf{M}}(t)]^{-1} \hat{\mathbf{M}}'(t) \hat{\mathbf{P}}^*(t) \hat{\Phi}(t + 1; t) \quad (35)$$

$$\hat{\Phi}^*(t + 1; t) = \hat{\Phi}(t + 1; t) - \hat{\mathbf{M}}(t) \hat{\Delta}^*(t) \quad (36)$$

$$\hat{\mathbf{P}}^*(t - 1) = \hat{\Phi}'(t + 1; t) \hat{\mathbf{P}}^*(t) \hat{\Phi}^*(t + 1; t) + \hat{\mathbf{Q}}(t) \quad (37)$$

Initially we must set $\hat{\mathbf{P}}^*(T) = \hat{\mathbf{Q}}(T + 1)$.

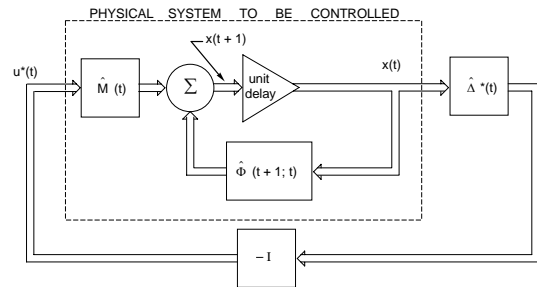


Fig. 4 Matrix block diagram of optimal controller

Comparing equations (35-37) with (28-30) and Fig. 3 with Fig. 4 we notice some interesting things which are expressed precisely by

Theorem 4. (Duality Theorem) Problem I and Problem II are duals of each other in the following sense:

Let $\tau \geq 0$. Replace every matrix $\mathbf{X}(t) = \mathbf{X}(t_0 + \tau)$ in (28-30) by $\hat{\mathbf{X}}'(t) = \hat{\mathbf{X}}'(T - \tau)$. Then One has (35-37). Conversely, replace every matrix $\hat{\mathbf{X}}(T - \tau)$ in (35-37) by $\mathbf{X}'(t_0 + \tau)$. Then one has (28-30).

Proof. Carry out the substitutions. For ease of reference, the dualities between the two problems are given in detail in Table 1.

	Problem I	Problem II
1	$\mathbf{x}(t)$ (unobservable) state variables of random process.	$\mathbf{x}(t)$ (observable) state variables of plant to be regulated.
2	$\mathbf{y}(t)$ observed random variables.	$\mathbf{u}(t)$ control variables
3	t_0 first observation.	T last control action.
4	$\Phi(t_0 + \tau + 1; t_0 + \tau)$ transition matrix.	$\hat{\Phi}(T - \tau + 1; T - \tau)$ transition matrix.
5	$\mathbf{P}^*(t_0 + \tau)$ covariance of optimized estimation error.	$\hat{\mathbf{P}}^*(T - \tau)$ matrix of quadratic form for performance index under optimal regulation.
6	$\Delta^*(t_0 + \tau)$ weighting of observation for optimal estimation.	$\hat{\Delta}^*(T - \tau)$ weighting of state for optimal control.
7	$\Phi^*(t_0 + \tau + 1; t_0 + \tau)$ transition matrix for optimal estimation error.	$\hat{\Phi}^*(T - \tau + 1; T - \tau)$ transition matrix under optimal regulation.
8	$\mathbf{M}(t_0 + \tau)$ effect of state on observation.	$\hat{\mathbf{M}}(T - \tau)$ effect of control vectors on state.
9	$\mathbf{Q}(t_0 + \tau)$ covariance of random excitation.	$\hat{\mathbf{Q}}(T - \tau)$ matrix of quadratic form defining error criterion.

Remarks. (n) The *mathematical* significance of the duality between Problem I and Problem II is that both problems reduce to the solution of the Wiener-Hopf-like equation (32).

(o) The *physical* significance of the duality is intriguing. Why are observations and control dual quantities?

Recent research [29] has shown that the essence of the Duality Theorem lies in the duality of constraints at the output (represented by the matrix $\hat{\mathbf{M}}(t)$ in Problem I) and constraints at the input (represented by the matrix $\hat{\mathbf{M}}(t)$ in Problem II).

(p) Applications of Wiener's methods to the solution of noise-free regulator problem have been known for a long time; see the recent textbook of Newton, Gould, and Kaiser [27]. However, the connections between the two problems, and in particular the duality, have apparently never been stated precisely before.

(q) The duality theorem offers a powerful tool for developing more deeply the theory (as opposed to the computation) of Wiener filters, as mentioned in Remark (i). This will be published elsewhere [29].

Applications

The power of the new approach to the Wiener problem, as expressed by Theorem 3, is most obvious when the data of the problem are given in numerical form. In that case, one simply performs the numerical computations required by (28–30). Results of such calculations, in some cases of practical engineering interest, will be published elsewhere.

When the answers are desired in closed analytic form, the iterations (28–30) may lead to very unwieldy expressions. In a few cases, Δ^* and Φ^* can be put into "closed form." Without discussing here how (if at all) such closed forms can be obtained, we now give two examples indicative of the type of results to be expected.

Example 1. Consider the problem mentioned under "Optimal Estimates." Let $x_1(t)$ be the signal and $x_2(t)$ the noise. We assume the model:

$$x_1(t+1) = \phi_{11}(t+1; t)x_1(t) + u_1(t)$$

$$x_2(t+1) = u_2(t)$$

$$y_1(t) = x_1(t) + x_2(t)$$

The specific data for which we desire a solution of the estimation problem are as follows:

- 1 $t_1 = t+1; t_0 = 0$
- 2 $Ex_1^2(0) = 0$, i.e., $x_1(0) = 0$
- 3 $Eu_1^2(t) = a^2, Eu_2^2(t) = b^2, Eu_1(t)u_2(t) = 0$ (for all t)
- 4 $\phi_{11}(t+1; t) = \phi_{11} = \text{const.}$

A simple calculation shows that the following matrices satisfy the difference equations (28–30), for all $t \geq t_0$:

$$\Delta^*(t) = \begin{bmatrix} \phi_{11}C(t) \\ 0 \end{bmatrix}$$

$$\Phi^*(t+1; t) = \begin{bmatrix} \phi_{11}[1-C(t)] & 0 \\ 0 & 0 \end{bmatrix}$$

$$\mathbf{P}^*(t+1) = \begin{bmatrix} a^2 + \phi_{11}^2 b^2 C(t) & 0 \\ 0 & b^2 \end{bmatrix}$$

$$\text{where } C(t+1) = 1 - \frac{b^2}{a^2 + b^2 + \phi_{11}^2 b^2 C(t)} \quad t \geq 0 \quad (38)$$

Since it was assumed that $x_1(0) = 0$, neither $x_1(1)$ nor $x_2(1)$ can be predicted from the measurement of $y_1(0)$. Hence the measurement at time $t = 0$ is useless, which shows that we should set $C(0) = 0$. This fact, with the iterations (38), completely determines the function $C(t)$. The nonlinear difference equation (38) plays the role of the Wiener-Hopf equation.

If $b^2/a^2 \ll 1$, then $C(t) \approx 1$ which is essentially pure prediction. If $b^2/a^2 \gg 1$, then $C(t) \approx 0$, and we depend mainly on $x_1^*(t|t-1)$ for the estimation of $x_1^*(t+1|t)$ and assign only very small weight

to the measurement $y_1(t)$; this is what one would expect when the measured data are very noisy.

In any case, $x_2^*(t|t-1) = 0$ at all times; one cannot predict independent noise! This means that ϕ^*_{12} can be set equal to zero. The optimal predictor is a first-order dynamic system. See Remark (j).

To find the stationary Wiener filter, let $t = \infty$ on both sides of (38), solve the resulting quadratic equation in $C(\infty)$, etc.

Example 2. A number of particles leave the origin at time $t_0 = 0$ with random velocities; after $t = 0$, each particle moves with a constant (unknown) velocity. Suppose that the position of one of these particles is measured, the data being contaminated by stationary, additive, correlated noise. What is the optimal estimate of the position and velocity of the particle at the time of the last measurement?

Let $x_1(t)$ be the position and $x_2(t)$ the velocity of the particle; $x_3(t)$ is the noise. The problem is then represented by the model,

$$x_1(t+1) = x_1(t) + x_2(t)$$

$$x_2(t+1) = x_2(t)$$

$$x_3(t+1) = \phi_{33}(t+1; t)x_3(t) + u_3(t)$$

$$y_1(t) = x_1(t) + x_3(t)$$

and the additional conditions

- 1 $t_1 = t; t_0 = 0$
- 2 $Ex_1^2(0) = Ex_2(0) = 0, Ex_2^2(0) = a^2 > 0;$
- 3 $Eu_3(t) = 0, Eu_3^2(t) = b^2.$
- 4 $\phi_{33}(t+1; t) = \phi_{33} = \text{const.}$

According to Theorem 3, $\mathbf{x}^*(t|t)$ is calculated using the dynamic system (31).

First we solve the problem of predicting the position and velocity of the particle one step ahead. Simple considerations show that

$$\mathbf{P}^*(1) = \begin{bmatrix} a^2 & a^2 & 0 \\ a^2 & a^2 & 0 \\ 0 & 0 & b^2 \end{bmatrix} \quad \text{and} \quad \Delta^*(0) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

It is then easy to check by substitution into equations (28–30) that

$$\mathbf{P}^*(t) = \frac{b^2}{C_1(t-1)} \times \begin{bmatrix} t^2 & t & -\phi_{33}t(t-1) \\ t & 1 & -\phi_{33}(t-1) \\ -\phi_{33}t(t-1) & -\phi_{33}(t-1) & \phi_{33}^2(t-1)^2 + C_1(t-1) \end{bmatrix}$$

is the correct expression for the covariance matrix of the prediction error $\tilde{\mathbf{x}}(t|t-1)$ for all $t \geq 1$, provided that we define

$$C_1(0) = b^2/a^2$$

$$C_1(t) = C_1(t-1) + [t - \phi_{33}(t-1)]^2, \quad t \geq 1$$

It is interesting to note that the results just obtained are valid also when ϕ_{33} depends on t . This is true also in Example 1. In conventional treatments of such problems there *seems* to be an essential difference between the cases of stationary and nonstationary noise. This misleading impression created by the conventional theory is due to the very special *methods* used in solving the Wiener-Hopf equation.

Introducing the abbreviation

$$C_2(0) = 0$$

$$C_2(t) = t - \phi_{33}(t-1), \quad t \geq 1$$

and observing that

$$\text{cov } \tilde{\mathbf{x}}(t+1|t) = \mathbf{P}^*(t+1)$$

$$= \Phi(t+1; t)[\text{cov } \tilde{\mathbf{x}}(t|t)]\Phi'(t+1; t) + \mathbf{Q}(t)$$

the matrices occurring in equation (31) and the covariance matrix of $\tilde{\mathbf{x}}(t)$ are found after simple calculations. We have, for all $t \geq 0$,

$$\Phi(t; t+1)\Delta^*(t) = \frac{1}{C_1(t)} \begin{bmatrix} tC_2(t) \\ C_2(t) \\ C_1(t) - tC_2(t) \end{bmatrix}$$

$\Phi(t; t+1)\Phi^*(t+1; t)\Phi(t+1; t)$

$$= \frac{1}{C_1(t)} \begin{bmatrix} C_1(t) - tC_2(t) & C_1(t) - tC_3(t) & -\phi_{33}tC_2(t) \\ -C_2(t) & C_1(t) - C_2(t) & -\phi_{33}C_2(t) \\ -C_1(t) + tC_2(t) & -C_1(t) + tC_2(t) & +\phi_{33}tC_2(t) \end{bmatrix}$$

and

$$\text{cov } \tilde{\mathbf{x}}(t|t) = E \tilde{\mathbf{x}}(t|t) \tilde{\mathbf{x}}^*(t|t) = \frac{b^2}{C_1(t)} \begin{bmatrix} t^2 & t & -t^2 \\ t & 1 & -t \\ -t^2 & -t & t^2 \end{bmatrix}$$

To gain some insight into the behavior of this system, let us examine the limiting case $t \rightarrow \infty$ of a large number of observations. Then $C_1(t)$ obeys approximately the differential equation

$$dC_1(t)/dt \approx C_2^2(t) \quad (t \gg 1)$$

from which we find

$$C_1(t) \approx (1 - \phi_{33})^2 t^3/3 + \phi_{33}(1 - \phi_{33})t^2 + \phi_{33}^2 t + b^2/a^2 \quad (t \gg 1) \quad (39)$$

Using (39), we get further,

$$\Phi^{-1}\Phi^*\Phi \approx \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ -1 & -1 & 0 \end{bmatrix} \quad \text{and} \quad \Phi^{-1}\Delta^* \approx \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad (t \gg 1)$$

Thus as the number of observations becomes large, we depend almost exclusively on $x_1^*(t|t)$ and $x_2^*(t|t)$ to estimate $x_1^*(t+1|t+1)$ and $x_2^*(t+1|t+1)$. Current observations are used almost exclusively to estimate the noise

$$x_3^*(t|t) \approx y_1^*(t) - x_1^*(t|t) \quad (t \gg 1)$$

One would of course expect something like this since the problem is analogous to fitting a straight line to an increasing number of points.

As a second check on the reasonableness of the results given, observe that the case $t \gg 1$ is essentially the same as prediction based on continuous observations. Setting $\phi_{33} = 0$, we have

$$E \tilde{x}_1^2(t|t) \approx \frac{a^2 b^2 t^2}{b^2 + a^2 t^3/3} \quad (t \gg 1; \phi_{33} = 0)$$

which is identical with the result obtained by Shinbrot [11], Example 1, and Solodovnikov [14], Example 2, in their treatment of the Wiener problem in the finite-length, continuous-data case, using an approach entirely different from ours.

Conclusions

This paper formulates and solves the Wiener problem from the "state" point of view. On the one hand, this leads to a very general treatment including cases which cause difficulties when attacked by other methods. On the other hand, the Wiener problem is shown to be closely connected with other problems in the theory of control. Much remains to be done to exploit these connections.

References

- 1 N. Wiener, "The Extrapolation, Interpolation and Smoothing of Stationary Time Series," John Wiley & Sons, Inc., New York, N.Y., 1949.
- 2 L. A. Zadeh and J. R. Ragazzini, "An Extension of Wiener's Theory of Prediction," *Journal of Applied Physics*, vol. 21, 1950, pp. 645-655.
- 3 H. W. Bode and C. E. Shannon, "A Simplified Derivation of Linear Least-Squares Smoothing and Prediction Theory," *Proceedings IRE*, vol. 38, 1950, pp. 417-425.
- 4 R. C. Booton, "An Optimization Theory for Time-Varying Linear Systems With Nonstationary Statistical Inputs," *Proceedings IRE*, vol. 40, 1952, pp. 977-981.
- 5 J. H. Laning and R. H. Battin, "Random Processes in Automatic Control," McGraw-Hill Book Company, Inc., New York, N. Y., 1956.
- 6 W. B. Davenport, Jr., and W. L. Root, "An Introduction to the Theory of Random Signals and Noise," McGraw-Hill Book Company, Inc., New York, N. Y., 1958.
- 7 S. Darlington, "Linear Least-Squares Smoothing and Prediction, With Applications," *Bell System Tech. Journal*, vol. 37, 1958, pp. 1221-1294.
- 8 G. Franklin, "The Optimum Synthesis of Sampled-Data Systems" Doctoral dissertation, Dept. of Elect. Engr., Columbia University, 1955.
- 9 A. B. Lees, "Interpolation and Extrapolation of Sampled Data," *Trans. IRE Prof. Group on Information Theory*, IT-2, 1956, pp. 173-175.
- 10 R. C. Davis, "On the Theory of Prediction of Nonstationary Stochastic Processes," *Journal of Applied Physics*, vol. 23, 1952, pp. 1047-1053.
- 11 M. Shinbrot, "Optimization of Time-Varying Linear Systems With Nonstationary Inputs," *TRANS. ASME*, vol. 80, 1958, pp. 457-462.
- 12 M. Blum, "Recursion Formulas for Growing Memory Digital Filters," *Trans. IRE Prof. Group on Information Theory*, IT-4, 1958, pp. 24-30.
- 13 V. S. Pugachev, "The Use of Canonical Expansions of Random Functions in Determining an Optimum Linear System," *Automatics and Remote Control (USSR)*, vol. 17, 1956, pp. 489-499; translation pp. 545-556.
- 14 V. V. Solodovnikov and A. M. Batkov, "On the Theory of Self-Optimizing Systems (in German and Russian)," *Proc. Heidelberg Conference on Automatic Control*, 1956, pp. 308-323.
- 15 J. L. Doob, "Stochastic Processes," John Wiley & Sons, Inc., New York, N. Y., 1955.
- 16 M. Loève, "Probability Theory," Van Nostrand Company, Inc., New York, N. Y., 1955.
- 17 R. E. Bellman, I. Glicksberg, and O. A. Gross, "Some Aspects of the Mathematical Theory of Control Processes," *RAND Report R-313*, 1958, 244 pp.
- 18 R. E. Kalman and R. W. Koepcke, "Optimal Synthesis of Linear Sampling Control Systems Using Generalized Performance Indexes," *TRANS. ASME*, vol. 80, 1958, pp. 1820-1826.
- 19 J. E. Bertram, "Effect of Quantization in Sampled-Feedback Systems," *Trans. AIEE*, vol. 77, II, 1958, pp. 177-182.
- 20 R. E. Kalman and J. E. Bertram, "General Synthesis Procedure for Computer Control of Single and Multi-Loop Linear Systems" *Trans. AIEE*, vol. 77, II, 1958, pp. 602-609.
- 21 C. W. Merriam, III, "A Class of Optimum Control Systems," *Journal of the Franklin Institute*, vol. 267, 1959, pp. 267-281.
- 22 R. E. Kalman and J. E. Bertram, "A Unified Approach to the Theory of Sampling Systems," *Journal of the Franklin Institute*, vol. 267, 1959, pp. 405-436.
- 23 R. E. Kalman and R. W. Koepcke, "The Role of Digital Computers in the Dynamic Optimization of Chemical Reactors," *Proc. Western Joint Computer Conference*, 1959, pp. 107-116.
- 24 R. E. Kalman, "Dynamic Optimization of Linear Control Systems, I. Theory," to appear.
- 25 S. Sherman, "Non-Mean-Square Error Criteria," *Trans. IRE Prof. Group on Information Theory*, IT-4, 1958, pp. 125-126.
- 26 V. S. Pugachev, "On a Possible General Solution of the Problem of Determining Optimum Dynamic Systems," *Automatics and Remote Control (USSR)*, vol. 17, 1956, pp. 585-589.
- 27 G. C. Newton, Jr., L. A. Gould, and J. F. Kaiser, "Analytical Design of Linear Feedback Controls," John Wiley & Sons, Inc., New York, N. Y., 1957.
- 28 O. J. M. Smith, "Feedback Control Systems," McGraw-Hill Book Company, Inc., New York, N. Y., 1958.
- 29 R. E. Kalman, "On the General Theory of Control Systems," *Proceedings First International Conference on Automatic Control*, Moscow, USSR, 1960.

APPENDIX RANDOM PROCESSES: BASIC CONCEPTS

For convenience of the reader, we review here some elementary definitions and facts about probability and random processes. Everything is presented with the utmost possible simplicity; for greater depth and breadth, consult Laning and Battin [5] or Doob [15].

A *random variable* is a function whose values depend on the outcome of a chance event. The *values* of a random variable may be any convenient mathematical entities; real or complex numbers, vectors, etc. For simplicity, we shall consider here only real-valued random variables, but this is no real restriction. Random variables will be denoted by x, y, \dots and their values by ξ, η, \dots . Sums, products, and functions of random variables are also random variables.

A random variable x can be explicitly defined by stating the probability that x is less than or equal to some real constant ξ . This is expressed symbolically by writing

$$Pr(x \leq \xi) = F_x(\xi); F_x(-\infty) = 0, F_x(+\infty) = 1$$

$F_x(\xi)$ is called the *probability distribution function* of the random variable x . When $F_x(\xi)$ is differentiable with respect to ξ , then $f_x(\xi) = dF_x(\xi)/d\xi$ is called the *probability density function* of x .

The *expected value* (*mathematical expectation*, *statistical average*, *ensemble average*, *mean*, etc., are commonly used synonyms) of any nonrandom function $g(x)$ of a random variable x is defined by

$$Eg(x) = E[g(x)] = \int_{-\infty}^{\infty} g(\xi) dF_x(\xi) = \int_{-\infty}^{\infty} g(\xi) f_x(\xi) d\xi \quad (40)$$

As indicated, it is often convenient to omit the brackets after the symbol E . A sequence of random variables (finite or infinite)

$$\{x(t)\} = \dots, x(-1), x(0), x(1), \dots \quad (41)$$

is called a *discrete* (or *discrete-parameter*) *random* (or *stochastic*) *process*. One particular set of observed values of the random process (41)

$$\dots, \xi(-1), \xi(0), \xi(1), \dots$$

is called a *realization* (or a *sample function*) of the process. Intuitively, a random process is simply a set of random variables which are indexed in such a way as to bring the notion of time into the picture.

A random process is *uncorrelated* if

$$Ex(t)x(s) = Ex(t)Ex(s) \quad (t \neq s)$$

If, furthermore,

$$Ex(t)x(s) = 0 \quad (t \neq s)$$

then the random process is *orthogonal*. Any uncorrelated random process can be changed into orthogonal random process by replacing $x(t)$ by $x'(t) = x(t) - Ex(t)$ since then

$$Ex'(t)x'(s) = E[x(t) - Ex(t)][x(s) - Ex(s)] \\ = Ex(t)x(s) - Ex(t)Ex(s) = 0$$

It is useful to remember that, if a random process is orthogonal, then

$$E[x(t_1) + x(t_2) + \dots]^2 = Ex^2(t_1) + Ex^2(t_2) + \dots \quad (t_1 \neq t_2 \neq \dots)$$

If \mathbf{x} is a vector-valued random variable with components x_1, \dots, x_n (which are of course random variables), the matrix

$$[E(x_i - Ex_i)(x_j - Ex_j)] = E(\mathbf{x} - E\mathbf{x})(\mathbf{x}' - E\mathbf{x}') \\ = \text{cov } \mathbf{x} \quad (42)$$

is called the *covariance matrix* of x .

A random process may be specified explicitly by stating the probability of simultaneous occurrence of any finite number of events of the type

$$x(t_1) \leq \xi_1, \dots, x(t_n) \leq \xi_n; (t_1 \neq \dots \neq t_n), \text{ i.e.,} \\ Pr[x(t_1) \leq \xi_1, \dots, x(t_n) \leq \xi_n] = F_{x(t_1), \dots, x(t_n)}(\xi_1, \dots, \xi_n) \quad (43)$$

where $F_{x(t_1), \dots, x(t_n)}$ is called the *joint probability distribution function* of the random variables $x(t_1), \dots, x(t_n)$. The *joint probability density function* is then

$$f_{x(t_1), \dots, x(t_n)}(\xi_1, \dots, \xi_n) = \partial^n F_{x(t_1), \dots, x(t_n)} / \partial \xi_1 \dots \partial \xi_n$$

provided the required derivatives exist. The expected value $Eg[x(t_1), \dots, x(t_n)]$ of any nonrandom function of n random variables is defined by an n -fold integral analogous to (40).

A random process is *independent* if for any finite $t_1 \neq \dots \neq t_n$, (43) is equal to the product of the first-order distributions

$$Pr[x(t_1) \leq \xi_1] \dots Pr[x(t_n) \leq \xi_n]$$

If a set of random variables is independent, then they are obviously also uncorrelated. The converse is not true in general. For a set of more than 2 random variables to be independent, it is not sufficient that any pair of random variables be independent.

Frequently it is of interest to consider the probability distribution of a random variable $x(t_{n+1})$ of a random process given the actual values $\xi(t_1), \dots, \xi(t_n)$ with which the random variables $x(t_1), \dots, x(t_n)$ have occurred. This is denoted by

$$Pr[x(t_{n+1}) \leq \xi_{n+1} | x(t_1) = \xi_1, \dots, x(t_n) = \xi_n] \\ = \frac{\int_{-\infty}^{\xi_{n+1}} f_{x(t_1), \dots, x(t_{n+1})}(\xi_1, \dots, \xi_{n+1}) d\xi_{n+1}}{f_{x(t_1), \dots, x(t_n)}(\xi_1, \dots, \xi_n)} \quad (44)$$

which is called the *conditional probability distribution function* of $x(t_{n+1})$ given $x(t_1), \dots, x(t_n)$. The *conditional expectation*

$$E\{g[x(t_{n+1})] | x(t_1), \dots, x(t_n)\}$$

is defined analogously to (40). The conditional expectation is a random variable; it follows that

$$E[E\{g[x(t_{n+1})] | x(t_1), \dots, x(t_n)\}] = E\{g[x(t_{n+1})]\}$$

In all cases of interest in this paper, integrals of the type (40) or (44) need never be evaluated explicitly, only the *concept* of the expected value is needed.

A random variable x is *gaussian* (or *normally distributed*) if

$$f_x(\xi) = \frac{1}{[2\pi E(x - Ex)^2]^{1/2}} \exp \left[-\frac{1}{2} \frac{(\xi - Ex)^2}{E(x - Ex)^2} \right]$$

which is the well-known bell-shaped curve. Similarly, a random vector \mathbf{x} is *gaussian* if

$$f_x(\xi) = \frac{1}{(2\pi)^{n/2} (\det \mathbf{C})^{1/2}} \exp \left[-\frac{1}{2} (\xi - E\mathbf{x})' \mathbf{C}^{-1} (\xi - E\mathbf{x}) \right]$$

where \mathbf{C}^{-1} is the inverse of the covariance matrix (42) of \mathbf{x} . A *gaussian random process* is defined similarly.

The importance of gaussian random variables and processes is largely due to the following facts:

Theorem 5. (A) *Linear functions (and therefore conditional expectations) on a gaussian random process are gaussian random variables.*

(B) *Orthogonal gaussian random variables are independent.*

(C) *Given any random process with means $Ex(t)$ and covariances $Ex(t)x(s)$, there exists a unique gaussian random process with the same means and covariances.*

Explanation of this transcription, John Lukesh, 20 January 2002.

Using a photo copy of R. E. Kalman's 1960 paper from an original of the ASME "Journal of Basic Engineering", March 1960 issue, I did my best to make an accurate version of this rather significant piece, in an up-to-date computer file format. For this I was able to choose page formatting and type font spacings that resulted in a document that is a close match to the original. (All pages start and stop at about the same point, for example; even most individual lines of text do.) I used a recent version of Word for Windows and a recent Hewlett Packard scanner with OCR (optical character recognition) software. The OCR software is very good on plain text, even distinguishing between italic versus regular characters quite reliably, but it does not do well with subscripts, superscripts, and special fonts, which were quite prevalent in the original paper. And I found there was no point in trying to work from the OCR results for equations. A lot of manual labor was involved.

Since I wanted to make a faithful reproduction of the original, I did not make any changes to correct (what I believed were) mistakes in it. For example, equation (32) has a $\mathbf{P}^*(t)\mathbf{M}(t)$ product that should be reversed, I think. I left this, and some other things that I thought were mistakes in the original, as is. (I didn't find very many other problems with the original.) There may, *of course*, be problems with my transcription. The plain text OCR results, which didn't require much editing, are pretty accurate I think. But the subscripts etc and the equations which I copied essentially manually, are suspect. I've reviewed the resulting document quite carefully, several times finding mistakes in what I did each time. The last time there were five, four cosmetic and one fairly inconsequential. There are probably more. I would be very pleased to know about it if any reader of this finds some of them; jlukesh@deltanet.com.