

ANNALES DE L'I. H. P., SECTION B

J.-P. GEORGIN

Contrôle des chaînes de Markov sur des espaces arbitraires

Annales de l'I. H. P., section B, tome 14, n° 3 (1978), p. 255-277

http://www.numdam.org/item?id=AIHPB_1978__14_3_255_0

© Gauthier-Villars, 1978, tous droits réservés.

L'accès aux archives de la revue « Annales de l'I. H. P., section B » (<http://www.elsevier.com/locate/anihpb>) implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques

<http://www.numdam.org/>

Contrôle des chaînes de Markov sur des espaces arbitraires

par

J.-P. GEORGIN

Université Paris-Nord, CSP Avenue J.-B. Clément
93430 Villetaneuse

(appartenant à l'équipe de recherche associée au C. N. R. S. n° 532,
« statistiques appliquées ».)

RÉSUMÉ. — Nous donnons un cadre d'hypothèses qui nous permet d'affirmer l'existence d'une stratégie stationnaire optimale pour une chaîne de Markov contrôlée dans le cas du critère dévalué. (L'espace des actions est polonais et l'espace des états est un espace mesurable quelconque, ou un espace polonais muni de sa tribu borélienne.)

L'étude du problème dévalué permet de trouver la solution au problème en moyenne, moyennant la résolution d'une équation de Poisson dont nous donnons des conditions d'existence.

Nous appliquons enfin ces résultats dans le cas où la chaîne contrôlée est une chaîne récurrente positive au sens de Harris.

INTRODUCTION

Nous avons repris dans [4] les principaux résultats de la théorie du contrôle, sur des espaces non dénombrables, basée sur les travaux de Striebel [11] et Hinderer [5], qui nous conduisent ici au rappel du théorème fondamental sur la caractérisation d'une stratégie optimale dans le cas du critère dévalué.

L'intérêt de l'étude du critère dévalué est qu'il nous permet de passer sous nos hypothèses à l'étude du critère en moyenne.

L'étude particulière que nous proposons est celle où l'espace d'état est un espace mesurable quelconque, alors que S. M. Ross dans [10] considérait un espace d'état polonais; nous avons d'ailleurs dans [4] développé des conditions qui assurent la validité des hypothèses de S. M. Ross dans le cas d'espaces topologiques.

1. PRÉSENTATION DES OBJETS, NOTATIONS ET HYPOTHÈSES FONDAMENTALES

L'espace des temps est \mathbb{N} .

L'espace des épreuves est (Ω, \mathcal{A}) , un espace mesurable.

On étudie un système : processus défini sur (Ω, \mathcal{A}) à valeurs dans un espace mesurable des états $(\mathcal{Y}, \mathcal{Y})$: on note $Y = (Y_t)_{t \in \mathbb{N}}$ ce processus.

A chaque instant t , on peut entreprendre une action A_t , fonction mesurable de (Ω, \mathcal{A}) dans $(\mathcal{A}, \mathcal{A})$, $(\mathcal{A}, \mathcal{A})$ est l'espace mesurable des actions. Le processus $A = (A_t)_{t \in \mathbb{N}}$ est le processus des contrôles. On utilisera indifféremment les termes action ou contrôle.

On note pour toute suite $(a_t)_{t \in \mathbb{N}} : a^{(t)} = (a_0, \dots, a_t)$.

La tribu des événements observables jusqu'à l'instant t est \mathcal{B}_t :

$$\mathcal{B}_t = \sigma(Y^{(t)}, A^{(t-1)})$$

On devra à chaque instant t choisir l'action A_t au vu des observations antérieures à t . Autrement dit, on choisira une stratégie $\delta = (\delta_t)_{t \in \mathbb{N}}$, c'est-à-dire un processus de (Ω, \mathcal{A}) dans $(\mathcal{A}, \mathcal{A})$ adapté aux tribus \mathcal{B}_t . Seules seront acceptées certaines stratégies présentant des conditions de cohérence que l'on précisera. On notera \mathcal{D} l'ensemble des stratégies cohérentes.

On se donne une probabilité de transition π de $\mathcal{Y} \times \mathcal{A}$ dans \mathcal{Y} dont le sens intuitif est le suivant : si à l'instant t , $Y_t = y$, $A_t = a$; alors à l'instant $(t + 1)$, Y_{t+1} est dans Γ , $\Gamma \in \mathcal{Y}$, avec la probabilité $\pi(y, a, \Gamma)$. On donne enfin une probabilité λ_0 sur $(\mathcal{Y}, \mathcal{Y})$, la loi initiale du processus.

Les espaces produits seront munis des tribus produits associées. On ne mentionnera pas les tribus associées lorsqu'elles sont évidentes.

On se place dans le cadre suivant qui définit la notion de « chaîne de Markov contrôlée ».

Pour toute stratégie cohérente $d = (d_t)_{t \in \mathbb{N}}$, il existe une probabilité P^d sur (Ω, \mathcal{A}) telle que pour tout t , et tout $B \in \mathcal{Y}$

$$\begin{aligned} P^d[A_t = d_t \text{ pour tout } t] &= 1 \\ P^d[Y_0 \in B] &= \lambda_0(B) \\ P^d[Y_{t+1} \in B \mid \mathcal{B}_t] &= \pi(Y_t, d_t, B) \end{aligned}$$

Remarque. — La construction canonique d'un tel processus est possible (cf. [4]).

Définissons l'ensemble \mathcal{D} des stratégies cohérentes :

A chaque instant t , le domaine où peuvent être choisies les actions ne dépend que de Y_t : à chaque $y \in \mathcal{Y}$, on associe $D(y)$ partie non vide de \mathcal{A} tel que l'ensemble $\{(y, a) ; a \in D(y)\}$ soit élément de $\mathcal{Y} \otimes \mathcal{A}$.

Une stratégie $d = (d_t)_{t \in \mathbb{N}}$ est donc cohérente si et seulement si $d_t \in D(Y_t)$ pour tout t .

Le gain obtenu dans l'intervalle de temps $[t, t + 1]$ est \tilde{g}_t , variable aléatoire réelle (v. a. r.) mesurable et bornée.

On suppose que pour toute $d \in \mathcal{D}$ et tout $t \in \mathbb{N}$

$$E^d[\tilde{g}_t | \mathcal{B}_t] = g(Y_t, d_t) = g(Y_t, A_t) \quad [P^d \text{ p. s.}]$$

où g est une v. a. r. sur $\mathcal{Y} \times \mathcal{A}$.

Il existe une v. a. p sur $\mathcal{Y} \times \mathcal{A} \times \mathcal{Y}$ et une transition λ de \mathcal{Y} dans \mathcal{Y} telles que pour $y \in \mathcal{Y}, a \in \mathcal{A}, B \in \mathcal{Y}$

$$\pi(y, a ; B) = \int p(y, a, z) \lambda(y, dz) 1_B(z)$$

Il existe alors une probabilité ρ sur $(\mathcal{Y}, \mathcal{Y})$ telle que :

$$\rho(B) = 0 \Leftrightarrow [P^d[Y_t \in B] = 0 \text{ pour toute } d \text{ et tout } t \geq 0]$$

car les mesures $Y_t(P^d)$ sont toutes absolument continues par rapport à :

$$\sum_{t \geq 0} \frac{1}{2^t} \lambda_0(\lambda_t) \quad \text{où } \lambda_t \text{ est la } t\text{-ième itérée de } \lambda.$$

Notations : On notera :

$\mathcal{B}(\mathcal{Y})$ L'ensemble des v. a. bornées sur \mathcal{Y} muni de la norme définie par :

$$\|\phi\| = \sup_x |\phi(x)|$$

$\mathcal{M}(\mathcal{Y})$ L'espace vectoriel des mesures signées sur \mathcal{Y} muni de la norme

$$\|v\| = \sup_{\{\phi, \|\phi\| \leq 1\}} \left| \int \phi dv \right|$$

2. SOLUTION DU PROBLÈME DÉVALUÉ ET INTRODUCTION AU PROBLÈME EN MOYENNE

Sous les conditions précédentes, introduisons les notions de « critère dévalué » et de « critère en moyenne ».

Définissons le gain maximin en moyenne

$$\underline{\mathcal{M}} = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s$$

et le gain minimax en moyenne :

$$\bar{\mathcal{M}} = \overline{\lim}_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s$$

Le problème que nous étudions est : à quelles conditions sur π et sur g existe-t-il une stratégie d_0 pour laquelle :

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s = k \quad \mathbf{P}^{d_0} \text{ p. s.}$$

et pour toute autre stratégie d :

$$\bar{\mathcal{M}} \leq k \quad \mathbf{P}^d \text{ p. s.}$$

Un problème moins intéressant mais techniquement plus simple que celui posé précédemment est le cas du *problème dévalué* : on cherche à maximiser :

$$\mathbf{E}^d \left[\sum_{s=0}^{\infty} \beta^s \tilde{g}_s \right] \quad \text{ou} \quad \mathbf{E}^d \left[\sum_{s=0}^{\infty} \beta^s \tilde{g}_s \mid Y_0 \right]$$

où β , $0 < \beta < 1$, est le taux de dévaluation. Introduisons la définition :

Définition. — Une stratégie est dite *markovienne stationnaire* relativement à Y , si elle est de la forme :

$$\{ d(Y_t) \}_{t \in \mathbb{N}}$$

où d est une fonction mesurable de \mathcal{Y} dans \mathcal{A} telle que $d(y) \in D(y)$ pour tout $y \in \mathcal{Y}$.

On notera \mathcal{S} l'ensemble des stratégies markoviennes stationnaires asso-

ciées à Y . Si $d \in \mathcal{S}$, $Y = (Y_s)_{s \in \mathbb{N}}$ est une chaîne de Markov sur $(\Omega, \mathcal{A}, \mathbb{P}^d)$, dûment complété, de transition : $(y, B) \rightarrow \pi(y, d(y), B)$ notée Π^d .

On obtient alors le théorème que nous ne démontrons pas ici (cf. [4]).

THÉORÈME 1. — Dans le cas « stationnaire dévalué ».

1. La classe \mathcal{S} des stratégies markoviennes stationnaires associées à Y est complète. C'est-à-dire :

$$\text{ess sup}_{d \in \mathcal{D}} E^d \left[\sum_{s=0}^{\infty} \beta^s \tilde{g}_s \mid Y_0 \right] = \text{ess sup}_{d \in \mathcal{S}} E^d \left[\sum_{s=0}^{\infty} \beta^s \tilde{g}_s \mid Y_0 \right] = l_\beta(Y_0)$$

2.

$$l_\beta(y) = Tl_\beta(y) = \rho - \text{ess sup}_{d \in \mathcal{S}} \left[g(y, d(y)) + \beta \int \pi(y, d(y), dz) l_\beta(z) \right]$$

et la valeur de l_β peut-être obtenue par approximation dans $L^\infty(\rho)$ à l'aide de la suite des T^nh (n -ièmes itérées de T) pour $h \in L^\infty(\rho)$ quelconque.

3. Si d_β est telle que :

$$l_\beta(y) = g(y, d_\beta(y)) + \int \pi(y, d_\beta(y), dz) l_\beta(z) \quad [\rho \text{ p. s.}]$$

alors d_β est une stratégie optimale :

$$l_\beta(Y_0) = E^{d_\beta} \left[\sum_{s \geq 0} \beta^s \tilde{g}_s \mid Y_0 \right] \quad [\mathbb{P}^{d_\beta} \text{ p. s.}]$$

Quant au problème du critère en moyenne, nous recherchons la stratégie qui maximise l'espérance du gain moyen par unité de temps. Ce problème apparaît comme limite du précédent quand $\beta \rightarrow 1$.

Posons :

$$M_{\lambda_0}^d = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} E_{\lambda_0}^d [\tilde{g}_s]$$

Rappelons que si (u_k) est une suite de nombres réels et $0 < \beta < 1$ on a :

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} u_k \leq \lim_{\beta \rightarrow 1-} (1 - \beta) \sum_{k=0}^{\infty} \beta^k u_k$$

et aussi si $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} u_k$ existe, alors celle de $(1 - \beta) \sum_{k=0}^{\infty} \beta^k u_k$ existe et ce

sont les mêmes.

Donc :

$$\begin{aligned} \underline{M}_{\lambda_0}^d &= \underline{\lim}_t \frac{1}{t} \sum_{s=0}^{\infty} E_{\lambda_0}^d[\tilde{g}_s] \leq \underline{\lim}_{\beta \rightarrow 1-} (1 - \beta) \sum_{s=0}^{\infty} E_{\lambda_0}^d[\beta^s \tilde{g}_s] \\ &\leq \underline{\lim}_{\beta \rightarrow 1-} (1 - \beta) \sup_d E_{\lambda_0}^d \left[\sum_{s=0}^{\infty} \beta^s \tilde{g}_s \right] \end{aligned}$$

ainsi :

$$\underline{M}_{\lambda_0}^d \leq \underline{\lim}_{\beta \rightarrow 1-} (1 - \beta) \int l_{\beta} d\lambda_0 = k(\lambda_0)$$

Ce résultat donne une idée du lien entre le critère en moyenne et le critère dévalué.

Une autre approche de ce critère est l'exemple suivant :

Exemple. — Chaîne de Markov non contrôlée où seul le gain dépend de l'action.

Cet exemple conduit à étudier le cas où Y est une chaîne de Markov récurrente positive au sens de Harris non contrôlée. Nous supposons seulement que le gain g dépend de l'action : si à l'instant t , $Y_t = y$, $A_t = a$, on gagne $g(y, a)$.

On sait d'après un théorème ergodique que si μ est la probabilité invariante, alors pour toute $d \in \mathcal{S}$

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} g(Y_s, d(Y_s)) = \int \mu(dy) g(y, d(y)) = M^d \quad (\text{p. s.})$$

On notera :

$$g_d(\mu) = \int \mu(dy) g(y, d(y)) \quad \text{et} \quad g_d = g(\cdot, d(\cdot))$$

Supposons qu'il existe une fonction $d_0 : \mathcal{Y} \rightarrow \mathcal{A}$, mesurable, telle que $d_0(x) \in D(x)$ et :

$$g(x, d_0(x)) = \sup_{a \in D(x)} g(x, a)$$

donc d_0 est une telle stratégie optimale. Alors :

$$M^{d_0} = k_{d_0} = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} g(Y_s, d_0(Y_s)) \quad (\text{p. s.})$$

et pour toute autre stratégie d

$$\overline{\lim} \frac{1}{t} \sum_{s=0}^{t-1} g(Y_s, d_s) \leq M^d \quad (\text{p. s.})$$

Si de plus, on a récurrence au sens de Doeblin [cf. Revuz, Orey [9] [8]] alors, comme $(g_{d_0} - k_{d_0})$ est d'intégrale nulle par rapport à μ , il existe une v. a. l , bornée, unique à une constante près, telle que l soit solution de l'équation de Poisson $l - \pi l = g_{d_0} - k_{d_0}$ ou

$$l(x) + k_{d_0} = g(x, d_0(x)) + \int \pi(x, dy)l(y) = \sup_a \left[g(x, a) + \int \pi(x, dy)l(y) \right]$$

on obtient donc, sous des conditions de récurrence très fortes, dans ce cas simple, des équations analogues à celles que nous avons obtenues dans le cas dévalué. Nous chercherons à exhiber des fonctions k et l analogues dans le paragraphe sur le critère en moyenne.

Le problème qui se pose est celui des sélections mesurables. A quelles conditions la borne essentielle supérieure pour $d \in \mathcal{S}$ peut-elle être remplacée par une borne supérieure pour $a \in D(y)$ pour chaque y et existe-t-il une stratégie $d \in \mathcal{S}$ optimale ?

Envisageons maintenant des conditions qui répondent aux questions posées, que nous formulons dans les hypothèses suivantes :

HYPOTHÈSE 1. — *Conditions d'Evstigneev (3)*

$(\mathcal{Y}, \mathcal{Y})$ est un espace mesurable quelconque.

\mathcal{A} est un espace polonais (métrique, complet, séparable) muni de la tribu borélienne \mathcal{A} .

Pour tout $y \in D(y)$ est compact, et

$$\begin{cases} a \rightarrow g(y, a) & \text{est continue de } D(y) \text{ dans } \mathbb{R} \\ a \rightarrow \int \pi(y, a, dz)h(z) & \text{est continue de } D(y) \text{ dans } \mathbb{R} \end{cases}$$

quelle que soit la v. a. bornée h sur \mathcal{Y} .

On complète cette dernière hypothèse en prenant \mathcal{Y} un espace polonais.

HYPOTHÈSE 2. — *Conditions d'Evstigneev et de Brown-Purves (1)*

Hypothèse 1

et \mathcal{Y} est un espace polonais muni de la tribu borélienne \mathcal{Y} .

On a alors le théorème :

THÉORÈME 2. — *Dans le cadre de l'hypothèse 2 [resp. 1] on peut choisir une version de l_β telle que :*

$$l_\beta(y) = \sup_{a \in D(y)} \left[g(y, a) + \beta \int \pi(y, a, dz)l_\beta(z) \right]$$

et il existe une stratégie stationnaire markovienne optimale $d_\beta(\cdot)$ mesurable sur \mathcal{Y} telle que :

$$l_\beta(y) = g(y, d_\beta(y)) + \beta \int \pi(y, d_\beta(y); dz) l_\beta(z) \quad [\text{resp. } \rho \text{ p. s.}]$$

Démonstration :

1) *Sous l'hypothèse 1*

Rappelons un théorème de sélection dû à Evstigneev [3] que nous transposons ici avec nos notations :

THÉORÈME DE EVSTIGNEEV. — Soient $(\mathcal{Y}, \mathcal{Y})$ un espace mesurable quelconque muni d'une probabilité ρ , \mathcal{A} un espace polonais muni de sa tribu borélienne \mathcal{A} .

ϕ une v. a. r. sur $\mathcal{Y} \times \mathcal{A}$ telle que

$$\left\{ \begin{array}{l} \cdot \text{ pour tout } c, \text{ pour tout } y \in \mathcal{Y}, \{ a, \phi(y, a) \geq c \} \text{ est compact} \\ \cdot \phi(y, a) \leq q(y) \quad \text{où} \quad \int |q| d\rho < \infty \end{array} \right.$$

Alors il existe $d : \mathcal{Y} \rightarrow \mathcal{A}$ mesurable telle que

$$\phi(y, d(y)) = \sup_{a \in \mathcal{A}} \phi(y, a) \quad [\rho \text{ p. s.}]$$

Appliquons ce théorème à la fonction $\Phi(y, a)$

$$\begin{aligned} \Phi(y, a) &= g(y, a) + \beta \int \pi(y, a, dz) l_\beta(z) && \text{si } a \in D(y) \\ &= -\infty && \text{sinon} \end{aligned}$$

Les conditions du théorème de Evstigneev sont vérifiées sous l'hypothèse 1, en effet :

- $a \rightarrow \Phi(y, a)$ est continue sur $D(y)$, $a \rightarrow \Phi(y, a)$ est s. c. s. sur \mathcal{A} ,
- $\Phi(y, a) \leq \|g\| \frac{1}{1-\beta}$, pour tout c et pour tout y , l'ensemble $\{ a, \Phi(y, a) \geq c \}$ est fermé inclus dans le compact $D(y)$, il est compact.

Ainsi il existe une $d_\beta \in \mathcal{S}$ telle que, d_β mesurable de \mathcal{Y} dans \mathcal{A} et :

$$\Phi(y, d_\beta(y)) = \sup_{a \in D(y)} \Phi(y, a) \quad [\rho \text{ p. s.}]$$

On avait :

$$l_\beta(y) = \rho - \text{ess sup}_{d \in \mathcal{S}} \Phi(y, d(y))$$

Pour toute $d \in \mathcal{S}$

$$\Phi(y, d(y)) \leq \sup_{a \in D(y)} \Phi(y, a)$$

mais ici :

$$\Phi(y, d_\beta(y)) = \sup_{a \in D(y)} \Phi(y, a) \quad [\rho \text{ p. s.}]$$

Il en résulte que $\Phi(\cdot, d_\beta(\cdot)) = \text{ess sup}_{d \in \mathcal{S}} \Phi(\cdot, d(\cdot))$, en effet

$$\Phi(\cdot, d_\beta(\cdot)) \leq \text{ess sup}_{d \in \mathcal{S}} \Phi(\cdot, d(\cdot)) \quad \text{car } d_\beta \in \mathcal{S}$$

et

$$\Phi(\cdot, d_\beta(\cdot)) \geq \text{ess sup}_{d \in \mathcal{S}} \Phi(\cdot, d(\cdot)) \quad \text{car } \Phi(\cdot, d_\beta(\cdot)) \geq \Phi(\cdot, d(\cdot)) \quad (\rho \text{ p. s.})$$

pour toute fonction de la famille.

On choisit désormais cette version de l_β et ainsi :

$$\left\{ \begin{aligned} l_\beta(y) &= \sup_{a \in D(y)} \left[g(y, a) + \beta \int \pi(y, a, dz) l_\beta(z) \right] \\ &= g(y, d_\beta(y)) + \beta \int \pi(y, d_\beta(y), dz) l_\beta(z) \end{aligned} \right. \quad [\rho \text{ p. s.}]$$

2) *Sous l'hypothèse 2*

Le théorème de sélection est dû ici à Brown-Purves [1], que nous rappelons.

THÉORÈME DE BROWN-PURVES. — *Sous l'hypothèse 2, si ϕ est une v. a. r. sur $\mathcal{Y} \times \mathcal{A}$ telle que pour tout $y \in \mathcal{Y}$, l'application $a \rightarrow \phi(y, a)$ est continue de $D(y)$ dans \mathbb{R} , alors il existe une application $d : \mathcal{Y} \rightarrow \mathcal{A}$ mesurable telle que, pour tout y ,*

$$\phi(y, d(y)) = \sup_{a \in \mathcal{A}} \phi(y, a)$$

Le raisonnement est identique au précédent et on obtient une égalité sûre de l_β et de $\Phi(\cdot, d_\beta(\cdot))$.

Indépendance par rapport à la loi initiale

Il est intéressant de s'assurer de l'indépendance des résultats par rapport à la loi initiale.

Considérons l'opérateur T_β défini sur l'ensemble $\mathcal{B}(\mathcal{Y})$ des v. a. bornées sur \mathcal{Y} en posant pour $f \in \mathcal{B}(\mathcal{Y})$

$$T_\beta f(y) = \sup_{a \in D(y)} \left[g(y, a) + \beta \int \pi(y, a, dz) f(z) \right]$$

T_β est une contraction dont l_β est l'unique point fixe, donc l_β est indépendante de la loi initiale sous les hypothèses 1 et 2.

Dans le cadre de l'hypothèse 2, d_β est aussi indépendante de la loi initiale.

Dans le cadre de l'hypothèse 1, formulons l'hypothèse la suivante.

HYPOTHÈSE 1 a. — La transition λ est remplacée par une probabilité λ : pour $(y, a, B) \in \mathcal{Y} \times \mathcal{A} \times \mathcal{Y}$

$$\pi(y, a, B) = \int p(y, a, z) 1_B(z) \lambda(dz)$$

et

- $p(y, a, \cdot) \leq q(y, \cdot)$ λ intégrable,
- pour $(y, z) \in \mathcal{Y}^2$, $a \rightarrow p(y, a, z)$ est continue de $D(y)$ dans \mathbb{R} .

Conséquences

a) Si pour tout (y, z) , la fonction $a \rightarrow p(y, a, z)$ est continue de $D(y)$ dans \mathbb{R} , pour tout a_0 dans $D(y)$ on a :

$$\lim_{a \rightarrow a_0} \int |p(y, a, z) - p(y, a_0, z)| \lambda(dz) = 0$$

L'application $a \rightarrow \pi(y, a, \cdot)$ de $D(y)$ dans $\mathcal{M}(\mathcal{Y})$ est continue, $D(y)$ étant compact la continuité est uniforme. La dernière condition de l'hypothèse 1 est automatiquement réalisée, et en outre :

b) Soit $\mathcal{H} \subset \mathcal{B}(\mathcal{Y})$ une famille de fonctions de $\mathcal{B}(\mathcal{Y})$ uniformément bornées, alors à y donné, la famille

$$\left\{ a \rightarrow \int p(y, a, z) h(z) \lambda(dz) \right\}_{h \in \mathcal{H}} \text{ est équicontinue.}$$

Pour toute mesure initiale λ_0 , $\lambda_0 \ll \lambda$, on a $\rho \ll \lambda$. On peut appliquer le théorème d'Estigneev en remplaçant ρ par λ : on a λ p. s.

$$l_\beta(y) = g(y, d_\beta(y)) + \beta \int \pi(y, d_\beta(y); dz) l_\beta(z)$$

Donc d_β ne dépend pas de la loi initiale absolument continue par rapport à λ .

On notera N_β l'ensemble λ négligeable où la dernière relation n'est pas vraie.

Si $d \in \mathcal{D}$, on désignera par $P_{\lambda_0}^d$ la probabilité qui est désignée ci-dessus par P^d , si la mesure initiale est λ_0 . Lorsque λ_0 est la mesure de Dirac au point $x \in \mathcal{Y}$, on notera P_x^d . Les espérances sont notées $E_{\lambda_0}^d$ ou E_x^d . Lorsque

les résultats énoncés seront indépendants de la loi initiale, on utilisera les notations P^d et E^d sans préciser cette loi.

Remarquons que sous l'hypothèse 2 l'existence de la fonction d_0 , rencontrée lors de l'exemple, est assurée. D'après la formule obtenue dans le théorème 2, on a, pour un état 0 fixé, sous les hypothèses 1 et 2

$$(1 - \beta)l_\beta(0) + l_\beta(x) - l_\beta(0) = \sup_a \left[g(x, a) + \beta \int \pi(x, a, dy)(l_\beta(y) - l_\beta(0)) \right]$$

et on aimerait pouvoir passer à la limite quand $\beta \rightarrow 1$ afin de retrouver les fonctions k et l de l'équation de Poisson que nous avons rencontrées lors de l'exemple. Nous étudierons séparément les conditions de Evstigneev et les conditions de Brown-Evstigneev.

3. PROBLÈME EN MOYENNE

3.1. Conditions d'existence d'une équation de Poisson

Afin de simplifier l'exposé de ce qui suit, nous nous intéresserons aux deux hypothèses suivantes :

HYPOTHÈSE 1'. — L'hypothèse 1 et l'hypothèse 1 a.

. $\lambda_0 \ll \lambda$

. $L^1(\mathcal{Q})$ est un espace de Banach séparable

HYPOTHÈSE 2'. — L'hypothèse 2 et l'hypothèse 1 a.

3.1.1. Dans le cadre de l'hypothèse 1' :

On a

$$l_\beta(x) = \sup_{a \in D(x)} \left[g(x, a) + \beta \int \pi(x, a, dz)l_\beta(z) \right]$$

Considérons $(l_\beta(x) - l_\beta(0))$

$$l_\beta(x) - l_\beta(0) = \sup_{a \in D(x)} \left[g(x, a) + \beta \int \pi(x, a, dz)(l_\beta(z) - l_\beta(0)) + l_\beta(0)(\beta - 1) \right]$$

en posant :

$$f_\beta(x) = l_\beta(x) - l_\beta(0),$$

on a :

$$(1 - \beta)l_\beta(0) + f_\beta(x) = \sup_a \left\{ g(x, a) + \beta \int \pi(x, a, dz)f_\beta(z) \right\}$$

Supposons que la famille $\{f_\beta\}$ est uniformément bornée par un nombre M .

$L^\infty(\lambda)$ est le dual de l'espace de Banach séparable $L^1(\lambda)$ donc la boule unité de $L^\infty(\lambda)$ est séquentiellement relativement compacte pour la topologie faible $\sigma(L^\infty(\lambda), L^1(\lambda))$. Par suite, il existe une suite (β_n) qui croît vers 1 telle que :

$$\begin{cases} \lim_{n \rightarrow \infty} (1 - \beta_n)l_{\beta_n}(0) = k \\ \lim_{n \rightarrow \infty} f_{\beta_n}(x) = \hat{l}(x) \end{cases} \quad \text{dans } \sigma(L^\infty(\lambda), L^1(\lambda))$$

C'est-à-dire que pour tout (y, a) en utilisant l'hypothèse 1 a,

$$\lim_{n \rightarrow \infty} \int p(y, a, z) f_{\beta_n}(z) \lambda(dz) = \int p(y, a, z) \hat{l}(z) \lambda(dz)$$

D'après la remarque suivant l'hypothèse 1 a, on sait alors que, pour tout y , la famille $\left\{ a \rightarrow \int p(y, a, z) f_{\beta_n}(z) \lambda(dz) \right\}_n$ est équicontinue et donc d'après Ascoli

$$a \rightarrow \int p(y, a, z) \hat{l}(z) \lambda(dz)$$

est continue et la convergence est uniforme sur le compact $D(y)$.

Ainsi pour tout $\varepsilon > 0$, il existe n_0 tel que $|1 - \beta_{n_0}| \leq \varepsilon/M$ et tel que pour $n \geq n_0$ et pour tout $a \in D(y)$ on ait :

$$\begin{cases} \left| \int p(y, a, z) f_{\beta_n}(z) \lambda(dz) - \int p(y, a, z) \hat{l}(z) \lambda(dz) \right| \leq \varepsilon \\ |k - (1 - \beta_n)l_{\beta_n}(0)| \leq \varepsilon \end{cases}$$

Montrons que $\{f_{\beta_n}(y)\}$ est une suite de Cauchy.

Pour tout n , il existe d'après le théorème 2 et la remarque qui suit H_{1a} , un ensemble N_{β_n} , λ négligeable et une stratégie $d_{\beta_n} \in \mathcal{S}$ tels que pour tout $y \notin N_{\beta_n}$

$$f_{\beta_n}(y) + (1 - \beta_n)l_{\beta_n}(0) = g(y, d_{\beta_n}(y)) + \beta_n \int p(y, d_{\beta_n}(y), z) f_{\beta_n}(z) \lambda(dz)$$

Prenons $y \notin \dot{N} = \bigcup_n N_{\beta_n}$. Si n et m sont supérieurs à n_0

$$f_{\beta_m}(y) + (1 - \beta_m)l_{\beta_m}(0) \geq g(y, d_{\beta_n}(y)) + \beta_m \int p(y, d_{\beta_n}(y), z) f_{\beta_m}(z) \lambda(dz)$$

donc :

$$f_{\beta_n}(y) - f_{\beta_m}(y) + (1 - \beta_n)l_{\beta_n}(0) - (1 - \beta_m)l_{\beta_m}(0) \leq \beta_n \int p(y, d_{\beta_n}(y), z) f_{\beta_n}(z) \lambda(dz) - \beta_m \int p(y, d_{\beta_n}(y), z) f_{\beta_m}(z) \lambda(dz)$$

Ainsi :

$$f_{\beta_n}(y) - f_{\beta_m}(y) \leq (1 - \beta_m)l_{\beta_m}(0) - (1 - \beta_n)l_{\beta_n}(0) + \beta_n \int p(y, d_{\beta_n}(y), z) (f_{\beta_n}(z) - \hat{l}(z)) \lambda(dz) - \beta_m \int p(y, d_{\beta_n}(y), z) (f_{\beta_m}(z) - \hat{l}(z)) \lambda(dz) + (\beta_m - \beta_n) \int p(y, d_{\beta_n}(y), z) \hat{l}(z) \lambda(dz) \leq 4\varepsilon$$

Donc $\{ f_{\beta_n}(y) \}$ est une suite de Cauchy pour tout $y \notin \dot{N}$.

Soit

$$l(y) = \lim_{n \rightarrow \infty} f_{\beta_n}(y)$$

On a $l = \hat{l}$, λ p. s.

On peut alors prendre la limite quand $n \rightarrow \infty$ de

$$(1 - \beta_n)l_{\beta_n}(0) + f_{\beta_n}(y) = \sup_a \left\{ g(y, a) + \beta_n \int \pi(y, a, dz) f_{\beta_n}(z) \right\}$$

On obtient :

$$k + l(y) = \lim_{n \rightarrow \infty} \sup_a \left[g(y, a) + \beta_n \int \pi(y, a, dz) f_{\beta_n}(z) \right] \geq \sup_a \left[g(y, a) + \int \pi(y, a, dz) l(z) \right]$$

et on a aussi pour $y \notin \dot{N} = \bigcup_n N_{\beta_n}$

$$k + l(y) = \lim_{n \rightarrow \infty} \left[g(y, d_{\beta_n}(y)) + \beta_n \int \pi(y, d_{\beta_n}(y), dz) f_{\beta_n}(z) \right]$$

Pour $n \geq n_0$

$$g(y, d_{\beta_n}(y)) + \beta_n \int \pi(y, d_{\beta_n}(y), dz) f_{\beta_n}(z) \leq g(y, d_{\beta_n}(y)) + \beta_n \int \pi(y, d_{\beta_n}(y), dz) l(z) + \varepsilon$$

donc :

$$k + l(y) \leq \varepsilon + \lim_n \left\{ g(y, d_{\beta_n}(y)) + \beta_n \int \pi(y, d_{\beta_n}(y), dz) l(z) \right\} \leq \varepsilon + \sup_{a \in D(y)} \left\{ g(y, a) + \int \pi(y, a, dz) l(z) \right\}$$

et comme ceci est vrai pour tout ε

$$k + l(y) \leq \sup_{a \in D(y)} \left\{ g(y, a) + \int \pi(y, a, dz)l(z) \right\}$$

En conclusion :

$$k + l(y) \leq \sup_{a \in D(y)} \left\{ g(y, a) + \int \pi(y, a, dz)l(z) \right\}$$

3.1.2. Remarquons avant de nous résumer, qu'une des conditions imposées dans l'hypothèse 1' est de prendre $L^1(\mathcal{Y})$ un espace de Banach séparable, condition donnée dans l'hypothèse 2 puisque \mathcal{Y} est polonais (\mathcal{Y} espace métrique séparable, muni de sa tribu borélienne \mathcal{Y} , suffit).

Dans le cadre de l'hypothèse 2', l'ensemble N de la démonstration précédente est vide.

Résumons nos résultats dans la proposition suivante.

PROPOSITION 3. — *Dans le cadre de l'hypothèse 2' (respectivement dans le cadre de l'hypothèse 1'), si la condition suivante est réalisée.*

CONDITION 1. — *La famille $\{x \rightarrow l_{\beta_n}(x) - l_{\beta_n}(0)\}_{\beta}$ est bornée uniformément. Alors pour une suite $\{\beta_n\}$ croissant vers 1*

$$\begin{cases} \lim_{\beta_n \uparrow 1} (1 - \beta_n)l_{\beta_n}(0) = k \\ \lim_{\beta_n \uparrow 1} (l_{\beta_n}(\cdot) - l_{\beta_n}(0)) = l(\cdot) \quad [\text{resp. } \lambda \text{ p. s.}] \end{cases}$$

et les fonctions k et l vérifient l'équation

$$k + l(x) = \sup_{a \in D(x)} \left[g(x, a) + \int \pi(x, a, dz)l(z) \right]$$

3.2. Stratégies optimales en moyenne

Montrons maintenant un théorème qui caractérise les stratégies optimales. Ce théorème est adapté de théorèmes dus à Ross [10] et à Mandl [6]. Nous réunissons dans ce théorème les cas où l'hypothèse 2 ou l'hypothèse 1 sont vérifiées, on notera [respectivement] pour l'hypothèse 1.

THÉORÈME 4. — *Dans le cadre de l'hypothèse 2 [resp. 1].*

S'il existe une fonction aléatoire bornée l sur \mathcal{Y} et une constante k telles que pour tout $x \in \mathcal{X}$

$$k + l(x) = \sup_{a \in D(x)} \left[g(x, a) + \int \pi(x, a, dy)l(y) \right] \quad (1)$$

alors :

a) Il existe une stratégie stationnaire markovienne d telle que

$$k = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s \quad [\mathbb{P}^d \text{ p. s.}]$$

et pour toute $d \in \mathcal{D}$

$$\overline{\lim}_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s \leq k \quad [\mathbb{P}^d \text{ p. s.}]$$

b) Si

$$\phi(x, a) = g(x, a) + \int \pi(x, a, dy)l(y) - k - l(x)$$

On a :

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s = k \quad [\mathbb{P}^d \text{ p. s.}] \Leftrightarrow \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \phi(Y_s, A_s) = 0 \quad [\mathbb{P}^d \text{ p. s.}] \quad (3)$$

En particulier dans le cas des stratégies stationnaires markoviennes $d \in \mathcal{S}$, si elles existent, telles que $\phi(\cdot, d(\cdot)) = 0$ [resp. λ p. s.], c'est-à-dire :

$$k + l(\cdot) = g(\cdot, d(\cdot)) + \int \pi(\cdot, d(\cdot), dy)l(y) \quad [\text{resp. } \lambda \text{ p. s.}] \quad (4)$$

Démonstration. — La fonction ϕ introduite mesure la différence entre l'emploi d'un contrôle optimal et d'un contrôle quelconque.

On a posé :

$$\phi(x, a) = g(x, a) + \int \pi(x, a, dy)l(y) - k - l(x)$$

alors : $\phi \leq 0$.

Introduisons pour tout $t \in \mathbb{N}$

$$U_t = \tilde{g}_t - k + l(Y_{t+1}) - l(Y_t) - \phi(Y_t, A_t) \quad \text{et} \quad M_t = \sum_{s=0}^{t-1} U_s$$

Montrons que $(M_t)_{t \in \mathbb{N}}$ est une martingale sur $(\Omega, \mathcal{A}, \mathbb{P}^d)$ adaptée aux tribus $(\mathcal{B}_t)_{t \in \mathbb{N}}$ quelle que soit la stratégie d .

- . U_t est \mathcal{B}_{t+1} mesurable
- . $E^d[M_{t+1} - M_t | \mathcal{B}_t] = E^d[\tilde{g}_t | \mathcal{B}_t] + E^d[l(Y_{t+1}) | \mathcal{B}_t] - k - l(Y_t) - \phi(Y_t, A_t) = 0$.

On a donc d'après la loi des grands nombres pour les martingales si :

$$\sum_{t=1}^{\infty} \frac{1}{t^2} E^d[U_t^2] < \infty$$

(condition vérifiée car (U_t) est bornée uniformément en t)

$$\lim_{t \rightarrow \infty} \frac{M_t}{t} = 0 \quad [P^d \text{ p. s.}]$$

Or

$$\frac{M_t}{t} = \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s - k + \frac{l(Y_t) - l(Y_0)}{t} - \frac{1}{t} \sum_{s=0}^{t-1} \phi(Y_s, A_s)$$

donc

$$\lim_{t \rightarrow \infty} \left[\frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s - k - \frac{1}{t} \sum_{s=0}^{t-1} \phi(Y_s, A_s) \right] = 0 \quad [P^d \text{ p. s.}]$$

1° Il existe sous les hypothèses du théorème une stratégie d telle que (4) soit vérifiée, en effet en posant :

$$\psi(y, a) = g(y, a) + \int \pi(y, a, dz)l(z),$$

il existe $d \in \mathcal{S}$, mesurable de \mathcal{Y} dans \mathcal{A} telle que pour tout y

$$\psi(y, d(y)) = \sup_{a \in D(y)} \psi(y, a) \quad [\text{resp. } \lambda \text{ p. s.}]$$

et pour cette d :

$$\phi(y, d(y)) = 0 \quad [\text{resp. } \lambda \text{ p. s.}]$$

Dans le cadre de l'hypothèse 2 le résultat a) du théorème est alors immédiat. Pour l'hypothèse 1, il suffit de vérifier que si :

$$\phi(\cdot, d(\cdot)) = 0 \quad [\lambda \text{ p. s.}]$$

alors

$$\phi(Y_s, d(Y_s)) = 0 \quad [P^d \text{ p. s.}]$$

pour $s \geq 1$.

Soit $\Gamma = \{ \cdot, \phi(\cdot, d(\cdot)) \neq 0 \}$; on a $\lambda(\Gamma) = 0$.

Or $P_x^d[Y_s \in \Gamma] = \pi_s^d[x, \Gamma]$ et comme $\pi(x, d(x), \cdot) \ll \lambda(\cdot)$, $P_x^d[Y_s \in \Gamma]$ est nulle.

Alors d est telle que :

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s = k \quad [P^d \text{ p. s.}]$$

Et comme $\phi \leq 0$, pour toute autre stratégie d

$$\overline{\lim}_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \tilde{g}_s \leq k \quad [P^d \text{ p. s.}]$$

d est optimale.

2° Si

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \phi(Y_s, A_s) = 0 \quad [P^d \text{ p. s.}]$$

alors d est optimale au sens du 1°.

**3.3. Étude du cas où pour toute $d \in \mathcal{S}$,
la chaîne de transition π^d
est récurrente positive au sens de Harris**

Étude de $k(\cdot)$

Revenons à l'étude du comportement si β tend vers 1, de $(1 - \beta)l_\beta(\cdot)$.

Plaçons-nous dans le cadre de l'hypothèse 1'.

La famille $\{(1 - \beta)l_\beta(\cdot)\}_{0 < \beta < 1}$ est bornée par $\|g\|$ dans $L^\infty(\lambda)$, reprenons le raisonnement vu plus haut. $L^\infty(\lambda)$ est le dual de l'espace de Banach séparable $L^1(\lambda)$, donc la boule unité de $L^\infty(\lambda)$ est séquentiellement relativement compacte pour la topologie faible $\sigma(L^\infty(\lambda), L^1(\lambda))$, donc il existe une sous-suite $\{\beta_n\}$ qui tend vers 1 telle que :

$$\lim_{n \rightarrow \infty} (1 - \beta_n)l_{\beta_n} = \hat{k} \quad \text{dans } \sigma(L^\infty(\lambda), L^1(\lambda))$$

c'est-à-dire que pour tout (y, a)

$$\lim_{n \rightarrow \infty} \int p(y, a, z)(1 - \beta_n)l_{\beta_n}(z)\lambda(dz) = \int p(y, a, z)\hat{k}(z)\lambda(dz)$$

Or, d'après la conséquence b) de l'hypothèse 1 a, pour tout $y \in \mathcal{Y}$

$$\left\{ a \rightarrow \int p(y, a, z)(1 - \beta_n)l_{\beta_n}(z)\lambda(dz) \right\}_{\beta_n}$$

est une famille équicontinue, et de plus, d'après Ascoli la convergence sur le compact $D(y)$ de $\int p(y, a, z)l_{\beta_n}(z)(1 - \beta_n)\lambda(dz)$ vers $\int p(y, a, z)\hat{k}(z)\lambda(dz)$ est uniforme en a .

Ainsi pour tout y dans le complémentaire de :

$$N = \bigcup_n N_{\beta_n}$$

[N_{β} étant l'ensemble de λ -mesure nulle où il n'y avait pas l'égalité dans le théorème 2, N est vide sous l'hypothèse 2] on a :

$$l_{\beta_n}(y) = \sup_{a \in D(y)} \left\{ g(y, a) + \beta_n \int p(y, a, z) l_{\beta_n}(z) \lambda(dz) \right\}$$

donc :

$$\lim_{n \rightarrow \infty} (1 - \beta_n) l_{\beta_n}(y) = \sup_{a \in D(y)} \int p(y, a, z) \hat{k}(z) \lambda(dz) = \sup_{a \in D(y)} \int \pi(y, a, dz) \hat{k}(z)$$

Donc sur N^c :

$$\hat{k}(\cdot) = \sup_{a \in D(\cdot)} \int \pi(\cdot, a, dz) \hat{k}(z).$$

Remarquons que la borne supérieure étant prise sur $D(y)$, compact, et

$$a \rightarrow \int \pi(y, a, dz) \hat{k}(z)$$

étant continue, il existe une stratégie $d \in \mathcal{S}$ telle que

$$\pi^d \hat{k} = \sup_{a \in D(\cdot)} \int \pi(\cdot, a, dz) \hat{k}(z) \quad [\text{resp. } \lambda \text{ p. s.}]$$

Posons :

$$k = \pi^d \hat{k}$$

$$k = \begin{cases} \pi^d k \\ \sup_{a \in D(\cdot)} \int \pi(\cdot, a, dz) k(z) \end{cases} \quad [\text{resp. } \lambda \text{ p. s.}]$$

Ainsi dans le cadre de l'hypothèse 1' (en particulier sous l'hypothèse 2') :

$$k(x) = \int \pi(x, d(x), dz) k(z)$$

Supposons, pour la suite de ce paragraphe, que pour toute $d \in \mathcal{S}$, la chaîne Y est récurrente Harris positive.

Ainsi comme k est harmonique pour une certaine stratégie stationnaire, k est constante. Pour tout y

$$\lim_{n \rightarrow \infty} (1 - \beta_n) \int \pi(y, a, dz) l_{\beta_n}(z) = k$$

et la convergence est uniforme en $a \in D(y)$.

Donc :

$$\lim_{n \rightarrow \infty} (1 - \beta_n) \int \pi(y, d_{\beta_n}(y), dz) l_{\beta_n}(z) = k$$

Soit $\mu_{d_{\beta_n}}$ la probabilité invariante par $\pi^{d_{\beta_n}}$: $\mu_{d_{\beta_n}} \pi^{d_{\beta_n}} = \mu_{d_{\beta_n}}$, $\mu_{d_{\beta_n}} \ll \lambda$.
Comme

$$\lim_{n \rightarrow \infty} (1 - \beta_n) \pi^{d_{\beta_n}} l_{\beta_n} = k \quad [\lambda \text{ p. s.}]$$

$$\lim_{n \rightarrow \infty} (1 - \beta_n) \mu_{d_{\beta_n}} \pi^{d_{\beta_n}} l_{\beta_n} = \lim_{n \rightarrow \infty} \mu_{d_{\beta_n}} k = k \quad .$$

$$\lim_{n \rightarrow \infty} (1 - \beta_n) l_{\beta_n}(\mu_{d_{\beta_n}}) = k$$

en notant

$$l(\mu) = \int l(x) \mu(dx) .$$

Or pour $d_{\beta_n} \in \mathcal{S}$, en considérant le critère de l'espérance du gain moyen

$$M^{d_{\beta_n}} = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} g(Y_s, A_s) = \int g(x, d_{\beta_n}(x)) d\mu_{d_{\beta_n}}(x) \quad [P^{d_{\beta_n}} \text{ p. s.}]$$

Montrons que :

$$\lim_{n \rightarrow \infty} M^{d_{\beta_n}} = k$$

Par définition

$$l_{\beta_n}(Y_0) = E^{d_{\beta_n}} \left\{ \sum_{s=0}^{\infty} \beta_n^s g(Y_s, A_s) \mid Y_0 \right\}$$

donc :

$$\int l_{\beta_n}(x) \mu_{d_{\beta_n}}(dx) = \int \mu_{d_{\beta_n}}(dx) \sum_{s=0}^{\infty} \beta_n^s \left(\int \pi(x, d_{\beta_n}(x), dy) g(y, d_{\beta_n}(y)) \right)$$

Comme $\mu_{d_{\beta_n}}$ est $\pi^{d_{\beta_n}}$ invariante

$$\int l_{\beta_n}(x) \mu_{d_{\beta_n}}(dx) = \sum_{s=0}^{\infty} \beta_n^s \int g(x, d_{\beta_n}(x)) \mu_{d_{\beta_n}}(dx) = \sum_{s=0}^{\infty} \beta_n^s M^{d_{\beta_n}} = \frac{M^{d_{\beta_n}}}{1 - \beta_n}$$

ainsi $(1 - \beta_n) \int l_{\beta_n}(x) \mu_{d_{\beta_n}}(dx) = M^{d_{\beta_n}}$ et $\lim_{n \rightarrow \infty} M^{d_{\beta_n}} = k$.

Donc pour toute suite (β_n) qui tend vers 1 telle que $(1 - \beta_n) l_{\beta_n}$ converge, la limite est $k = \sup_d M^d$. Par suite $k = \lim_{\beta \rightarrow 1} (1 - \beta) l_{\beta}$. S'il existe une $d_0 \in \mathcal{S}$ telle que $k = M^{d_0}$ et si π^{d_0} est récurrente au sens de Doebelin comme $(g - k)$ est d'intégrale nulle par rapport à μ_{d_0} ($\mu_{d_0}(g) = M^{d_0} = k$), l'équa-

tion de Poisson est vérifiée, et on sait qu'il existe une v. a. l bornée unique à une constante près telle que :

$$l - \pi^{d_0}(l) = g - k$$

On vient donc de généraliser le résultat de l'exemple d'une chaîne de Markov récurrente au sens de Harris, non contrôlée.

Résumons les résultats dans le théorème suivant :

THÉORÈME 5. — *Dans le cadre de l'hypothèse 1'.*

Si pour toute $d \in \mathcal{S}$, la chaîne Y est récurrente Harris positive, alors,

1. *il existe une constante k telle que :*

$$k = \sup_{d \in \mathcal{S}} M^d = \lim_{\beta \rightarrow 1} (1 - \beta)l_\beta(x)$$

Si de plus il existe une $d_0 \in \mathcal{S}$, telle que $M^{d_0} = k$ et si π^{d_0} est récurrente au sens de Doeblin, alors il existe une v. a. l bornée, unique à une constante près telle que : $l - \pi^{d_0}l = g_{d_0} - k$.

2. *Si*

$$\phi(x, a) = g(x, a) - k - l(x) + \int \pi(x, a, dz)l(z)$$

et si

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \phi(Y_s, A_s) = 0 \quad [\mathbb{P}^d \text{ p. s.}]$$

alors

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} g(Y_s, A_s) = k \quad [\mathbb{P}^d \text{ p. s.}]$$

En particulier si pour $d \in \mathcal{S}$ on a $\phi(\cdot, d(\cdot)) = 0$, alors d est optimale.

Remarque. — Des résultats récents de Numelin [7] assurent l'existence d'une fonction l , μ_{d_0} intégrable, solution de l'équation de Poisson $l - \pi^{d_0}l = g_{d_0} - k$ en supposant seulement que π^{d_0} est récurrente positive. Mais l peut ne pas être μ_d intégrable pour $d \in \mathcal{S}$ quelconque et l'application du théorème 4 pose des problèmes avec un tel l .

Exemple d'application

Nous nous plaçons dans le cadre où l'hypothèse 1' est réalisée, et nous supposons que la condition 2 suivante est réalisée :

CONDITION 2. — Il existe une constante k' , $k' < 1$ telle que si :

$$K = \{ (x, a) ; a \in D(x) \}$$

alors

$$\sup_{((x,a), (y,b)) \in K \times K} \|\pi(x, a, \cdot) - \pi(y, b, \cdot)\| = 2k'$$

Ce qui est réalisé en particulier sous la condition 2' suivante :

CONDITION 2'. — 1. Il existe un point 0 tel que pour tout $(x, a) \in K$

$$\pi(x, a, \{0\}) \geq \alpha > 0$$

ou

2. Il existe une mesure ν sur $(\mathcal{Y}, \mathcal{Y})$ non nulle telle que pour tout $(x, a) \in K$

$$\pi(x, a, \cdot) \geq \nu(\cdot)$$

En effet, pour la condition 2'.2 (qui implique 2'.1), il existe deux ensembles disjoints \mathcal{Y}^+ et \mathcal{Y}^- de \mathcal{Y} tels que $\mathcal{Y}^+ \cup \mathcal{Y}^- = \mathcal{Y}$

$$\begin{aligned} \|\pi(x, a, \cdot) - \pi(y, b, \cdot)\| &= \pi(x, a, \mathcal{Y}^+) - \pi(y, b, \mathcal{Y}^+) \\ &\quad + \pi(y, b, \mathcal{Y}^-) - \pi(x, a, \mathcal{Y}^-) \\ &\leq 1 - \nu(\mathcal{Y}^+) + 1 - \nu(\mathcal{Y}^-) \leq 2 - \|\nu\| \end{aligned}$$

Nous allons montrer que nous sommes dans les conditions du théorème 5. Nous montrons tout d'abord qu'il existe $d_0 \in \mathcal{S}$, telle que $M^{d_0} = k$, ce qui revient à montrer suivant la proposition 3 que $\{l_\beta(\cdot) - l_\beta(0)\}_\beta$ est uniformément bornée. Nous montrerons ensuite que π^{d_0} est récurrente au sens de Doeblin.

Rappelons alors un *théorème de Ueno* [12].

Si P est une probabilité de transition, on note P_s ses itérées et alors :

$$\|P_s(x, \cdot) - P_s(y, \cdot)\| \leq \frac{1}{2^{s-1}} \sup_{x,y} \|P(x, \cdot) - P(y, \cdot)\|^s$$

Appliquons ceci aux transitions π

$$\begin{aligned} \|\pi_s^d(x, \cdot) - \pi_s^d(y, \cdot)\| &\leq \frac{1}{2^{s-1}} (2k')^s, & k' < 1 \\ &\leq 2k'^s & k' < 1 \end{aligned}$$

et dans le cadre de nos hypothèses il existe une stratégie d_β qui vérifie :

$$l_\beta(x) = E_x^{d_\beta} \left[\sum_s \beta^s g(Y_s, d_\beta(Y_s)) \right]$$

Donc :

$$l_\beta(x) = \sum_s \beta^s \int g(z, d_\beta(z)) \pi_s[x, d_\beta(x), dz]$$

que nous notons :

$$\begin{aligned}
 l_\beta(x) &= \sum_s \beta^s \pi_s^{d_\beta}(x, g_{d_\beta}) \\
 |l_\beta(x) - l_\beta(o)| &= \left| \sum_s \beta^s (\pi_s^{d_\beta}(x, g_{d_\beta}) - \pi_s^{d_\beta}(0, g_{d_\beta})) \right| \\
 &\leq \sum_s \beta^s \|g\| \|\pi_s^{d_\beta}(x, \cdot) - \pi_s^{d_\beta}(0, \cdot)\| \\
 &\leq \sum_s \beta^s \|g\| 2k'^s \quad \text{avec } k' < 1 \\
 &\leq \frac{2 \|g\|}{1 - \beta k'} \leq \frac{2 \|g\|}{1 - k'}
 \end{aligned}$$

en considérant évidemment g uniformément borné, on vient de montrer que $\{l_\beta(\cdot) - l_\beta(0)\}_\beta$, sous les conditions précédentes, est uniformément bornée.

Montrons maintenant que π^d est une chaîne de Markov récurrente au sens de Doeblin, pour tout $d \in \mathcal{S}$.

On vient de voir que :

$$\|\pi_s^d(x, \cdot) - \pi_s^d(y, \cdot)\| \leq 2k'^s$$

Soit $\Gamma \in \mathcal{Y}$,

$$\left| \int \pi_t^d(x, dy) [\pi_s^d(x, \Gamma) - \pi_s^d(y, \Gamma)] \right| \leq 2k'^s$$

donc :

$$\|\pi_s^d(x, \cdot) - \pi_{t+s}^d(x, \cdot)\| \leq 2k'^s$$

Comme $\mathcal{M}(\mathcal{Y})$ est un espace de Banach, on en déduit que

$$\lim_{s \rightarrow \infty} \pi_s^d(x, \cdot) = \mu_d(\cdot) \quad \text{dans } \mathcal{M}(\mathcal{Y})$$

et

$$\|\pi_s^d(x, \cdot) - \mu_d(\cdot)\| \leq 2k'^s$$

et on est alors dans le cadre du théorème 5.

D'où le corollaire

COROLLAIRE 6. — *On suppose que :*

$$\sup_{[(x,a), (y,b)] \in \mathbb{K} \times \mathbb{K}} \|\pi(x, a, \cdot) - \pi(y, b, \cdot)\| = 2k', \quad k' < 1$$

Sous l'hypothèse 1', il existe une constante k et une v. a. bornée l sur \mathcal{Y} , et une stratégie $d_0 \in \mathcal{S}$, telles que :

$$\left\{ \begin{array}{l} k = \sup_{d \in \mathcal{S}} M^d = \lim_{\beta \rightarrow 1} (1 - \beta) l_\beta(\cdot) = M^{d_0} \\ k = l(\cdot) = \sup_a \left[g(\cdot, a) + \int \pi(\cdot, a, dy) l(y) \right] \\ k + l(\cdot) = g(\cdot, d_0(\cdot)) + \pi^{d_0} l(\cdot) \quad [\text{resp. } \lambda \text{ p. s.}] \end{array} \right.$$

Pour toute $d \in \mathcal{S}$, π^d est récurrente au sens de Doeblin.

BIBLIOGRAPHIE

- [1] BROWN-PURVES, Measurable selection of extrema. *Annals of statistics*, vol. 1, n° 5, 1973, p. 902-912.
- [2] DERMANN, *Finite state Markovian decision process*. Academic Press, 1970.
- [3] EŠTIGNEEV, Measurable selection and dynamic programming. *Mathematics of operations research*, vol. 1, n° 3, 1976, p. 267-272.
- [4] GEORGIN, Contrôle des chaînes de Markov sur des espaces arbitraires. Estimation et contrôle optimal dans le cadre adaptatif. *Thèse de 3^e Cycle*. Université Paris-Sud, Centre d'Orsay, 1977.
- [5] K. HINDERER, Foundations of non-stationary dynamic programming with discrete time parameters. *Lecture notes in operations research and Mathematical systems*. Springer-Verlag, Berlin.
- [6] MANDL, Estimation and control in Markov chains. *Adv. Appl. Prob.*, t. 6, 1974, p. 40-60.
- [7] NUMELIN, On the Poisson equation for ϕ recurrent Markov chains. *A paraître 1977*.
- [8] OREY, *Limit theorems for Markov chains transitions probabilities*. Van Nostrand, 1971.
- [9] D. REVUZ, *Markov chains*. North Holland, 1975.
- [10] ROSS, Arbitrary state markovian decision processes. *The annals of Mathematical statistics*, vol. 39, n° 6, 1968.
- [11] STRIEBEL, a) Optimal control of discrete time stochastic systems. *Lectures notes in operations research and mathematical systems*, Springer (19).
b) Martingale conditions for the optimal control of continuous time stochastic systems. International workshop of stochastic-filtering and control, Los Angeles, Californy.
- [12] UENO, Some limit theorems for temporally discrete Markov processes. *J. Fac. Sciences*, Université Tokyo (I), 1957, p. 7.

(Manuscrit reçu le 1^{er} février 1978)