

## ΜΕΡΙΚΑ ΠΑΡΑΤΗΡΗΣΙΜΕΣ ΑΛΥΣΙΔΕΣ ΜΑΡΚΟΒΙΑΝΗΣ ΑΠΟΦΑΣΗΣ ΜΕ ΟΜΟΙΟΜΟΡΦΗ ΚΑΤΑΝΟΜΗ ΜΗΝΥΜΑΤΩΝ

Υπό

*Γκουλιώνη Ιωάννη*

Τμήμα Στατιστικής και Ασφαλιστικής Επιστήμης, Πανεπιστήμιο Πειραιώς

### Abstract

#### PARTIALLY OBSERVABLE MARKOV DECISION PROCESSES WITH UNIFORMLY DISTRIBUTED SIGNAL PROCESSES

A Partially observed Markov decision process (P.O.M.D.P.) is a sequential decision problem where information concerning parameters of interest is incomplete, and possible actions include sampling, surveying, or otherwise collecting additional information. Such problems can theoretically be solved as dynamics programs, but the relevant state space is infinite which inhibits algorithmic solution. We formulate a (P.O.M.D.P.) with a continuous signal space and a method to convert a problem with uniformly distributed signal processes. We discussed how to solve (P.O.M.D.P.) problems with continuous signal processes. However, in order to obtain a value function which is close to the optimal value function, we might need to construct a step function with large number of signals.

Keywords: Maintenance, dynamic-programming, P.O.M.D.P.

### 1. Εισαγωγή

Μια P.O.M.D.P. (Μερικά παρατηρήσιμη Μαρκοβιανή αλυσίδα) είναι μια γενικευμένη αλυσίδα Μαρκοβιανής απόφασης, που επιτρέπει μιαν ατελή πληροφόρηση του συστήματος των καταστάσεων. Η γενίκευση αυτή είναι σημαντική σε προβλήματα, όπου η αβεβαιότητα ως προς την κατάσταση είναι το κεντρικό και ουσιώδες.

Επειδή λοιπόν σε ένα πρόβλημα P.O.M.D.P. δεν ξέρουμε ακριβώς την κατάσταση του συστήματος σε κάθε χρόνο, μέσω ενός μηχανισμού ελέγχου λαμβάνουμε ένα μήνυμα που συνδέεται στοχαστικά με την πραγματική

κατάσταση. Για παράδειγμα, όταν εξετάζουμε την κατάσταση ενός ασθενούς που πάσχει από στεφανιαία νόσο, το αποτέλεσμα του λεγόμενου τεστ κόπωσης, το επίπεδο ισχαιμίας, καθώς και ο πόνος στο στήθος είναι μηνύματα που συνυφαίνονται με την ζωτική κατάσταση του ασθενούς. Αρκετοί ερευνητές ασχολήθηκαν με το πρόβλημα εύρεσης βέλτιστων στρατηγικών σε πεπερασμένο και σε άπειρο χρονικόν ορίζοντα με διακριτό χώρο μηνυμάτων όπως οι Sondik [6,7,8], Denardo [2], Littman [5], Lovejoy [4], Howard [3], Bertsekas [1].

Στην εργασία αυτή επεκτείνουμε τα αποτελέσματα, στην περίπτωση όπου ο χώρος μηνυμάτων είναι συνεχής, με ομοιόμορφη κατανομή μηνυμάτων.

Η εργασία οργανώνεται ως εξής: Στην παράγραφο 2 περιγράφουμε το μοντέλο P.O.M.D.P. με διακριτό χώρο μηνυμάτων. Στην παράγραφο 3 δίνουμε τα ήδη γνωστά αποτελέσματα, που αφορούν την ύπαρξη βέλτιστης στάσιμης στρατηγικής για το πρόβλημα του άπειρου χρονικού ορίζοντα με διακριτό αριθμό μηνυμάτων. Οι ιδιότητες της συνάρτησης τιμών σε πεπερασμένο χρονικόν ορίζοντα είναι ουσιαστικές, διότι επιτρέπουν την επανάληψη στον συνεχή χώρο των καταστάσεων μιας P.O.M.D.P.

Στην παράγραφο 4 επεκτείνουμε τις τεχνικές που ισχύουν για το πρόβλημα P.O.M.D.P. με διακριτό χώρο μηνυμάτων. Κλειδί για την παραπάνω επέκταση είναι η διαμέριση του χώρου μηνυμάτων, με βάση ορισμένες υποθέσεις, σε τρόπον ώστε κάθε κελλί της διαμέρισης να μπορεί να θεωρηθεί σαν διακριτό μήνυμα.

Τέλος δίνουμε μια μέθοδο που μετατρέπει το πρόβλημα με συνεχή χώρο μηνυμάτων σε ένα πρόβλημα με διακριτό χώρο μηνυμάτων, καθώς και μια εφαρμογή για να γίνουν κατανοητές οι παρακάτω σκέψεις.

## **2. Περιγραφή του μοντέλου των μερικά παρατηρήσιμων Μαρκοβιανών αλυσίδων**

Μια P.O.M.D.P. περιγράφει μια στοχαστική διαδικασία που τυπικά αντιστοιχεί στην εξάδα (S, A, Θ, P, R,q).

- Το  $S = \{1,2,3...N\}$  είναι το σύνολο των δυνατών καταστάσεων του συστήματος.

- Το A είναι το σύνολο των αποφάσεων, που υποτίθεται ορισμένο. Η απόφαση που εκλέγεται στον χρόνο t ορίζεται σαν  $Y_t$ .

Στον χρόνο  $t$ , αν το σύστημα βρίσκεται στην κατάσταση  $i$ , και επιλέξουμε την απόφαση  $a$ , τότε είναι γνωστό ένα σύνολο από πιθανότητες  $\{p_{ij}^a(t) \text{ για } j=1,2,\dots,N\}$ , όπου  $(\sum_j p_{ij}^a(t)=1 \text{ για κάθε κατάσταση } i \text{ του συστήματος } t)$ .

$$P_{ij}^a(t) = \Pr[x_{t+1} = j | x_t = i, Y_t = a] \quad (1)$$

Υποθέτουμε ότι οι πιθανότητες μεταφοράς είναι ανεξάρτητες του χρόνου. Υποτίθεται ότι γνωρίζουμε τον πίνακα μεταφοράς  $P^a$

• Ένα ορισμένο όφελος  $q^a(i)$  προκύπτει, όταν το σύστημα βρίσκεται στην κατάσταση  $i$  και εκλέγεται η απόφαση  $a$ . Το  $q^a$  είναι ένα  $N$ -διάστατο διάνυσμα στήλη της μορφής  $q^a = (q^a(1), q^a(2), \dots, q^a(N))$ .

Αν γνωρίζουμε την τρέχουσα κατάσταση του συστήματος, πριν εκλέξουμε μια απόφαση τότε το πρόβλημα είναι ένα πρόβλημα πλήρως παρατηρήσιμης M.D.P. Μπορούμε τότε να εκλέξουμε μια απόφαση, που βασίζεται στην τρέχουσα κατάσταση του συστήματος, μεγιστοποιώντας το ολικό προσδοκώμενο αποπληθωρισμένο όφελος.

Οι POMDPS έχουν το χαρακτηριστικό ότι η κατάσταση του συστήματος δεν μπορεί να παρατηρηθεί απευθείας. Αντί για την κατάσταση, μέσω ενός μηχανισμού ελέγχου αυτός/ή που αποφασίζει λαμβάνει ένα τυχαίο μήνυμα  $\theta$ .

$$r_{i\theta}^a = \Pr[z_t = \theta | x_t = i, Y_{t-1} = a] \forall t \quad (2)$$

Μολονότι το μήνυμα είναι τυχαίο, γνωρίζουμε την παραπάνω υπό συνθήκη πιθανότητα να λάβουμε το μήνυμα  $\theta$ , όταν η κατάσταση του συστήματος, καθώς και η απόφαση είναι γνωστές. Υποθέτοντας ότι οι υπό συνθήκη πιθανότητες δεν μεταβάλλονται με τον χρόνο μπορούμε να θεωρήσουμε τον πίνακα  $R^a = [r_{i\theta}^a]$

Υποθέτοντας ότι το αντικείμενο του μοντέλου απόφασης είναι να μεγιστοποιήσει το ολικό προσδοκώμενο εκπίπτον όφελος, αυτός/ή που αποφασίζει, αφού δεν γνωρίζει την πραγματική κατάσταση του συστήματος παίρνει μια απόφαση βασιζόμενος σε όλη την ιστορία της εξέλιξης. Κριτήρια για την επιλογή βέλτιστων στρατηγικών σε πεπερασμένο και άπειρο χρονικών ορίζοντα, είναι η μεγιστοποίηση του ολικού προσδοκώμενου εκπίπτοντος

οφέλους, που υλοποιείται μέσω των παρακάτω σχέσεων (3), (4) για πεπερασμένο και άπειρο χρονικόν ορίζοντα αντίστοιχα:

$$E_{\pi(0)} \left\{ \sum_{t=0}^{T-1} \beta^t \cdot q[s(t), \alpha(t)] + \beta^T q[s(T)] \right\} \quad (3)$$

$$E_{\pi(0)} \left\{ \sum_{t=0}^{\infty} \beta^t \cdot q[s(t), \alpha(t)] \right\} \quad (4)$$

### 3. Η βέλτιστη συνάρτηση τιμών

Η ιστορία της εξέλιξης στον χρόνο  $t$ , δηλώνεται σαν  $I_t$ , όπου

$$I_0 = [\pi(0), Z_0]$$

$$I_t = [\pi(0), Z_0, Y_0, Z_1, Y_1, \dots, Z_{t-1}, Y_{t-1}, Z_t] \quad \text{και}$$

$$I_{t+1} = [I_t, Y_t, Z_{t+1}]$$

Σημειώνουμε ότι  $\{I_t, t=0,1,2,\dots\}$  είναι μια αλυσίδα Μαρκοβιανής απόφασης σύντομα, (M.D.P).

Αν τώρα,  $v_t(\cdot)$ , δηλώνει το βέλτιστο ολικό προσδοκώμενο όφελος από τον χρόνο  $t$  μέχρι το τέλος του σχεδιασθέντος ορίζοντα αποφάσεων.

τότε το δυναμικό πρόγραμμα (επαναληπτική εξίσωση D.P.) μπορεί να γράφει σαν:

$$\begin{aligned} v_t(I_t) &= \max_{Y_t \in A} E \left\{ q^{Y_t}(x_t) + \beta \cdot v_{t+1}(I_{t+1}) \mid I_t, Y_t \right\} \\ &= \max_{Y_t \in A} E \left\{ E(q^{Y_t}(x_t) \mid I_t, Y_t) + \beta \cdot E(v_{t+1}(I_{t+1}) \mid I_t, Y_t) \right\} \\ &= \max_{Y_t \in A} E \left\{ \sum_{i=1}^N \Pr(x_t = i \mid I_t) \cdot q^{Y_t}(i) + \beta \cdot \sum_{\theta \in \Theta} \Pr(z_{t+1} = \theta \mid I_t, Y_t) \cdot v_{t+1}([I_t, Y_t, Z_{t+1} = \theta]) \right\} \end{aligned} \quad (5)$$

όπου  $\beta$  είναι ένας παράγοντας έκπτωσης (discount factor).

Υποθέτουμε ότι  $0 < \beta < 1$  για τα προβλήματα άπειρου χρονικού ορίζοντα, ενώ  $0 < \beta$  για τα προβλήματα του πεπερασμένου χρονικού ορίζοντα.

Για να λύσουμε την επαναληπτική εξίσωση (5), οι συναρτήσεις τιμών στον χρόνο  $t$  πρέπει να υπολογισθούν για κάθε πιθανή ιστορία  $I_t$ . Όταν υπάρχει ένας μεγάλος αριθμός από πιθανές αποφάσεις, ένας μεγάλος αριθμός μηνυμάτων, και ο χρονικός ορίζοντας είναι μεγάλος, τότε ο αριθμός των πιθανών ιστοριών  $I_t$  είναι μεγάλος (ή άπειρος) και αυξάνει γραμμικά με το  $t$ . Οι υπολογιστικές απαιτήσεις ενός τέτοιου DP (dynamic - programming) αλγόριθμου είναι αληθινά τεράστιες (Οι λεπτομέρειες αυτού του προβλήματος σχολιάζονται στον Bertsekas (1976).

Επομένως, μπορούμε να θεωρήσουμε όλες τις πιθανές μεταβλητές του μοντέλου. Ας είναι  $\pi_i(t) = \Pr(x_t = i | I_t)$  και  $\pi(t) = [\pi_1(t), \pi_2(t), \dots, \pi_N(t)]$ ,

$$\text{όπου } \sum_{i=1}^N \pi_i(t) = 1 \quad 0 \leq \pi_i(t) \leq 1 \text{ για } i = 1, 2, \dots, N.$$

$$\text{Είναι εύκολο να δείξουμε } \Pr(x_{t+1} = i | I_t, Y_t, Z_{t+1}) = \Pr(x_{t+1} = i | \pi(t), Y_t, Z_{t+1})$$

$$\text{και } \Pr(z_{t+1} = \theta | I_t, Y_t) = \Pr(z_{t+1} = \theta | \pi(t), Y_t)$$

Επομένως, αν δοθούν  $I_t, Y_t = \alpha, z_{t+1} = \theta$  από τον κανόνα Bayes'

$$\pi_j(t+1) = \Pr(x_{t+1} = j | I_t, Y_t = \alpha, z_{t+1} = \theta) = \Pr(x_{t+1} = j | \pi(t), Y_t = \alpha, z_{t+1} = \theta)$$

$$= \frac{\Pr(x_{t+1} = j, z_{t+1} = \theta | \pi(t), Y_t = \alpha)}{\Pr(z_{t+1} = \theta | \pi(t), Y_t = \alpha)} = \frac{\sum_{i=1}^N \pi_i(t) \cdot P_{ij}^a \cdot r_{j,\theta}^a}{\sum_{k=1}^N \sum_{i=1}^N \pi_i(t) \cdot P_{ik}^a \cdot r_{k,\theta}^a} \quad (6)$$

ή σε μορφή πινάκων

$$T(\pi(t), \theta, \alpha) \equiv \Pr(x_{t+1} | I_t, Y_t = \alpha, z_{t+1} = \theta) = \Pr(x_{t+1} | \pi(t), Y_t = \alpha, z_{t+1} = \theta) = \frac{\pi(t) \cdot P^a \cdot R_\theta^a}{\pi(t) \cdot P^a \cdot R_\theta^a \cdot 1}$$

όπου  $R^a_0$  είναι ένας διαγώνιος πίνακας που έχει  $r^a_{j_0}$  σαν διαγώνια στοιχεία, και 1 είναι ένας πίνακας στήλη με όλα τα στοιχεία του 1.

Είναι καλά γνωστό, ότι  $\pi(t)$  είναι μια επαρκής στατιστική για την ιστορία της εξέλιξης  $I_t$ , όταν εκλέγουμε μια απόφαση σε χρόνο  $t$ . Συγκεκριμένα  $\pi(t)$  περιέχει όλη την αναγκαία πληροφορία, όσον αφορά την ιστορία της εξέλιξης, προκειμένου να εκλεγεί μια απόφαση στον χρόνο  $t$  (*Bertsekas 1976, Monahan 1982, Sondik 1971, Striebel 1965*).

Ας είναι  $\Pi \equiv \{\pi \in R^N : \sum_{i=1}^N \pi_i = 1 \text{ και } \pi_i \geq 0 \text{ για } i=1, 2, \dots, N\}$ . Το ακόλουθο αποτέλεσμα οφείλεται στον (Sondik 1971): Για κάθε ακολουθία αποφάσεων  $Y_1, Y_2, \dots, Y_t \in A$ , η ακολουθία των πιθανοτήτων  $\{\pi(k)\}$  όπου  $k = 0, 1, \dots, t$  είναι μια Μαρκοβιανή αλυσίδα, που σημαίνει ότι, αν  $\Gamma \subset \Pi$ , τότε:

$$\Pr(\pi(t+1) \in \Gamma \mid \pi(0), \pi(1), \dots, \pi(t), Y_t) = \Pr(\pi(t+1) \in \Gamma \mid \pi(t), Y_t)$$

Επομένως,  $\forall \pi \in \Pi$ ,

$$\begin{aligned} v_t(\pi) &= \max_{a \in A} E \{q^a(x_t) + \beta \cdot v_{t+1}(T(\pi, a, \theta)) \mid \pi, a\} \\ &= \max_{\alpha \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^a(i) + \beta \cdot \sum_{\theta \in \Theta} \Pr(z_{t+1} = \theta \mid \pi, a) \cdot v_{t+1}(T(\pi, a, \theta)) \right\} \end{aligned} \quad (7)$$

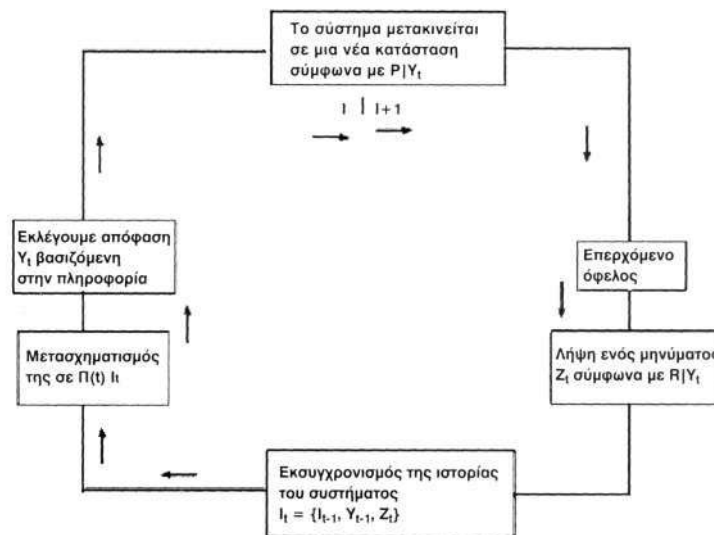
Σημειώνουμε ότι η εξίσωση (7) είναι μια επαναληπτική εξίσωση, που θεωρητικά τουλάχιστον μπορεί να βρει μια βέλτιστη στρατηγική, με βάση το κριτήριο του μέσου αποπληθωρισμένου οφέλους. Ωστόσο επειδή ο χώρος  $\Pi$  είναι συνεχής, δεν επιτρέπονται οι επαναλήψεις, οπότε το να βρούμε μια βέλτιστη στρατηγική δεν είναι εύκολο έργο. Ευτυχώς όμως οι POMDPs έχουν κάποιες θαυμάσιες ιδιότητες, που μας βοηθούν να αναπτύξουμε υπολογιστικούς αλγόριθμους για την εύρεση μιας βέλτιστης στρατηγικής. Αυτές τις ιδιότητες παραθέτουμε στην επόμενη παράγραφο.

#### 4. Συμβολισμοί και τελεστές

Στην παράγραφο αυτή, κάποιοι συμβολισμοί και τελεστές εισάγονται. Εισάγουμε, την νόρμα  $\|\cdot\|$  του Supremum. Έστω  $B(\pi)$  το σύνολο των φραγμένων συναρτήσεων (πραγματικών τιμών) στο  $\Pi$ .

ΣΧΗΜΑ 1

Διάγραμμα απόφασης για τις μερικά παρατηρήσιμες  
Μαρκοβιανές αλυσίδες



Τότε  $\|v\| = \sup \{v(\pi) : \pi \in \Pi\}$

Έστω τώρα συνάρτηση πραγματικών τιμών  $h: \Pi \times A \times B(\Pi) \rightarrow R$

$$h(\pi, \alpha, v) = \pi \cdot q^\alpha + \beta \cdot \sum_{\theta \in \Theta} \Pr(\theta | \pi, \alpha) \cdot v(T(\pi, \alpha, \theta)) \quad (8)$$

Τότε μέσω της  $h$  η εξίσωση (7) μπορεί να ξαναγραφεί σαν:

$$v_t(\pi) = \max_{\alpha \in A} h(\pi, \alpha, v_{t+1}).$$

Αν,  $\delta$ , είναι τώρα μια στρατηγική,  $\delta: \Pi \rightarrow A$ , και  $\Delta$  το σύνολο όλων των στάσιμων στρατηγικών, μπορούμε να γενικεύσουμε μέσω ενός τελεστού  $H_\delta$  στο  $B(\Pi)$ ,  $\forall \delta \in \Delta$ .

$$[H_\delta(v)](\pi) = h(\pi, \delta(\pi), v).$$

Αν  $\delta(\pi) = \alpha$  για  $\forall \pi \in \Pi$ , τότε χρησιμοποιούμε την έκφραση  $H_\alpha$  αντί για την  $H_\delta$ .

Τελικά ορίζεται ένας τελεστής βελτιστοποίησης σαν:  $Hv = \max_{\delta \in \Delta} [H_\delta v]$

### Ιδιότητες τελεστών $H_\delta$ , $H$

α) Είναι φραγμένοι (*Whitt 1978*)

β) Έχουν την ιδιότητα της μονοτονίας

γ) Έχουν την ιδιότητα της συστολής για το πρόβλημα του άπειρου χρονικού ορίζοντα.

Στις ιδιότητες αυτές των παραπάνω τελεστών θα πρέπει να προσθέσουμε, ότι η βέλτιστη συνάρτηση τιμών, σε πεπερασμένο χρονικό ορίζοντα  $T$ , είναι κατά τμήματα γραμμική και μάλιστα κυρτή (*Sondik 1971*).

Μια συνάρτηση  $v$  καλείται κατά τμήματα γραμμική και κυρτή αν υπάρχει ένα ορισμένο σύνολο από  $N$ -διάστατα "gradients",  $\Gamma$ , ώστε:

$$v(\pi) = \max_{i=1,2,3,\dots,x} \{ \pi \cdot \gamma^i : \text{όπου } \gamma^i \in \Gamma \}$$

Για  $\pi \in \Pi$ ,  $\alpha \in A$  και  $\theta \in \Theta$ , ορίζουμε τό σύνολο:

$$\Gamma_{\pi,\alpha,\theta} = \{ \bar{\gamma} \in \Gamma : T(\pi, \alpha, \theta) \cdot \bar{\gamma} \geq T(\pi, \alpha, \theta) \cdot \gamma \quad \forall \gamma \in \Gamma \}$$

Τότε  $v(T(\pi, \alpha, \theta)) = T(\pi, \alpha, \theta) \cdot \bar{\gamma} \quad \forall \bar{\gamma} \in \Gamma_{\pi,\alpha,\theta}$ .

Αν  $\gamma_{\pi,\alpha,\theta}$  είναι ένα "gradient" στο σύνολο  $\Gamma_{\pi,\alpha,\theta}$ , τότε,  $Hv(\pi)$  μπορεί να γραφεί σαν:

$$\begin{aligned} H_\alpha v(\pi) &= \pi \cdot q^\alpha + \beta \sum_{\theta \in \Theta} \Pr(\theta | \pi, \alpha) \cdot v(T(\pi, \alpha, \theta)) \\ &= \pi \cdot q^\alpha + \beta \sum_{\theta \in \Theta} \left( \pi \cdot P^\alpha \cdot R_\theta^\alpha \cdot 1 \right) \frac{\pi \cdot P^\alpha \cdot R_\theta^\alpha}{\pi \cdot P^\alpha \cdot R_\theta^\alpha \cdot 1} \gamma_{\pi,\alpha,\theta} = \pi \cdot \left[ q^\alpha + \beta \cdot \sum_{\theta \in \Theta} P^\alpha \cdot R_\theta^\alpha \cdot \gamma_{\pi,\alpha,\theta} \right] \quad (9) \end{aligned}$$

Επειδή  $q^\alpha + \beta \sum_{\theta \in \Theta} P^\alpha \cdot R_\theta^\alpha \cdot \gamma_{\pi,\alpha,\theta}$  είναι ένα  $N$ -διάστασης "gradient" η (9) μπορεί να απλοποιηθεί σαν:  $H_\alpha v(\pi) = \pi \cdot \gamma_{\pi,\alpha}$

$$\text{όπου } \gamma_{\pi,\alpha} = q^\alpha + \beta \cdot P^\alpha \sum_{\theta \in \Theta} R_\theta^\alpha \cdot \gamma_{\pi,\alpha,\theta}$$



Επιπλέον,  $Hv(\pi) = \max_{\alpha \in A} H_{\alpha} v(\pi)$ , τότε:

$$Hv(\pi) = \max_{\alpha \in A} [H_{\alpha} v(\pi)] = \max_{\alpha \in A} \{ \pi \cdot \gamma_{\pi, \alpha} \} = \pi \cdot \gamma_{\pi}$$

$$\text{όπου } \gamma_{\pi} = q^{\alpha} + \beta \sum_{\theta \in \Theta} P_{\theta}^{\alpha^*} R_{\theta}^{\alpha^*} \cdot \gamma_{\pi, \alpha^*, \theta} \quad (10)$$

Όπου,  $\alpha^*$ , είναι μια βέλτιστη απόφαση για το  $\pi$ .

Επειδή μάλιστα οι τελεστές  $H_{\delta}$ ,  $H$  είναι συστολές, υπάρχουν μοναδικά "fixed-point" στο  $B(\pi)$ , ας τα πούμε  $v_{\delta}$ , και  $v^*$  αντίστοιχα ώστε:

$$H_{\delta} v_{\delta} = v_{\delta}$$

και

$$H v^* = v^* \quad (\text{Elgolic 1964})$$

Μια στρατηγική  $\delta \in \Delta$  για την οποία ισχύει  $H_{\delta} v^* = H v^*$  είναι μια βέλτιστη στρατηγική. Οι Sawaragi και Yoshikawa (1970) έδειξαν, ότι αν ο χώρος  $A$  των αποφάσεων είναι ορισμένος, τότε υπάρχει μια βέλτιστη στάσιμη στρατηγική. Αυτό εκφράζει ότι, αν έχουμε ένα αρχικό (i.v) (information-vector),  $\pi$ , και η στρατηγική  $\delta \in \Delta$  (στάσιμη) χρησιμοποιείται για το πρόβλημα του άπειρου χρονικού ορίζοντα, τότε το ολικό προσδοκώμενο εκπίπτον όφελος θα είναι  $v_{\delta}(\pi)$ .

Από το ορισμό τώρα του τελεστή  $H$ ,

$$v^*(\pi) = \max_{\delta \in \Delta} \{ v_{\delta}(\pi) \} \quad \text{για όλα τα } \pi \in \Pi.$$

Άρα,  $v^*$  είναι η βέλτιστη συνάρτηση τιμών.

Επειδή μια P.O.M.D.P. έχει έναν συνεχή χώρο καταστάσεων, είναι πολύ δύσκολο να βρούμε την  $v_{\delta}$  για  $\delta \in \Delta$  ή το  $v^*$ . Επομένως, φράγματα για τις τιμές των παραπάνω συναρτήσεων αναζητούνται σε αρκετούς ερευνητές.

Συνοπτικά μπορούμε να πούμε, ότι οι αλγόριθμοι που αφορούν το πρόβλημα του πεπερασμένου χρονικού ορίζοντα, και βασίζονται σε μια επαναληπτική διαδικασία "value-iteration" χρησιμοποιούνται επίσης στο βήμα "policy-improvement" για το πρόβλημα του άπειρου χρονικού ορίζοντα.

Επομένως ουσιαστικό βήμα στην κατασκευή αλγορίθμων για το πρόβλημα του άπειρου χρονικού ορίζοντα είναι ο υπολογισμός  $H_u$  για μια δεδομένη κατά τμήματα γραμμική και κυρτή συνάρτηση  $u$ .

### 5. Μερικά παρατηρήσιμες Μαρκοβιανές αλυσίδες με συνεχή χώρο μηνυμάτων

Το σχέδιο αυτής της παραγράφου έχει ως ακολούθως: υποθέσεις, συμβολισμοί και διατύπωση μιας P.O.M.D.P. με συνεχή χώρο μηνυμάτων, μια μέθοδο που να μετατρέπει ένα πρόβλημα P.O.M.D.P. με συνεχή χώρο μηνυμάτων σε μια P.O.M.D.P. με ορισμένο αριθμό μηνυμάτων.

Οι βασικές υποθέσεις είναι ίδιες με αυτές που αναπτύχθηκαν στην παράγραφο (3). Η μόνη διαφορά είναι στην φύση των μηνυμάτων. Στην παράγραφο αυτή υποθέτουμε, ότι οι κατανομές μηνυμάτων έχουν συναρτήσεις πυκνότητας, ενώ στην παράγραφο (3) υποθέσαμε ότι είναι διακριτές. Οι παράμετροι αυτών των συναρτήσεων πυκνότητας εξαρτώνται από την κατάσταση του συστήματος, καθώς επίσης και από την απόφαση που λαμβάνεται στον συγκεκριμένο χρόνο. Για να γίνουμε πιο κατανοητοί, για κάθε κατάσταση του συστήματος  $i \in S$ , κάθε απόφαση  $a \in A$ , και κάθε χρόνο  $t=0, 1, 2, \dots$  υπάρχει μια συνάρτηση πυκνότητας πιθανότητας  $f_{i,u}^a(\cdot)$  στο σύνολο των μηνυμάτων.

Για το πρόβλημα του άπειρου χρονικού ορίζοντα, υποθέτουμε ότι η συνάρτηση πυκνότητας πιθανότητας είναι ανεξάρτητη από τον χρόνο και επομένως η εξάρτηση από τον χρόνο παραλείπεται στον συμβολισμό.

Ας είναι  $\pi \in \Pi$ , όπως ακριβώς ορίστηκε στην παράγραφο (3). Δεδομένου ότι το τρέχον (i.v) είναι  $\pi$  και η απόφαση  $a$  χρησιμοποιείται, η υπό συνθήκη πυκνότητα πιθανότητας είναι  $\bar{f}(\theta|\pi, a) = \sum_{k=1}^N \sum_{j=1}^N \pi_j \cdot P_{j,k}^a \cdot f_k^a(\theta)$  ή σε μορφή πινάκων

$$\bar{f}(\theta|\pi, a) = \pi \cdot P^a \bar{R}_\theta^a \cdot \mathbf{1} \quad (11)$$

όπου  $\bar{R}_\theta^a$  είναι ένας διαγώνιος πίνακας που έχει διαγώνια στοιχεία τα  $f_i^a(\theta)$  και  $\mathbf{1}$  είναι ένα N-διάστασης διάνυσμα στήλη με όλα τα στοιχεία αποτελούμενα από 1.

Ανάλογα με τον ορισμό του  $T(\pi, \alpha, \theta)$ , ορίζουμε  $\bar{T}(\pi, \alpha, \theta)$  την κατανομή πιθανότητας, που αφορά την κατάσταση του συστήματος, στην επόμενη χρονική περίοδο, δεδομένου ότι η κατανομή πιθανότητας για την τρέχουσα κατάσταση του συστήματος είναι  $\pi$ , λαμβάνεται η απόφαση  $\alpha$ , και παρατηρείται το μήνυμα  $\theta$  στην αρχή της επόμενης περιόδου.

$$\bar{T}_i(\pi, \alpha, \theta) = \Pr(x_{t+1} = i | \pi, Y_t = \alpha, Z_{t+1} = \theta),$$

και

$$\bar{T}_i(\pi, \alpha, \theta) = \Pr(x_{t+1} = i | \pi, Y_{t+1} = \alpha, Z_{t+1} = \theta) = \frac{\sum_{j=1}^N \pi_j \cdot P_{j,i}^\alpha \cdot f_i^\alpha(\theta)}{\sum_{\kappa=1}^N \sum_{j=1}^N \pi_j \cdot P_{j,\kappa}^\alpha \cdot f_\kappa^\alpha(\theta)} \quad (12)$$

$$\text{ή, in matrix form, } \bar{T}(\pi, \alpha, \theta) = \frac{\pi \cdot P^\alpha \cdot \bar{R}_\theta^\alpha}{\bar{f}(\theta | \pi, \alpha)} = \frac{\pi \cdot P^\alpha \cdot \bar{R}_\theta^\alpha}{\pi \cdot P^\alpha \cdot \bar{R}_\theta^\alpha \cdot \mathbf{1}} \quad (13)$$

Αναλογικά με την παράγραφο (3),  $\alpha$ s είναι  $v_t$  και  $v_{t+1}$  κυρτές, φραγμένες και συνεχείς συναρτήσεις τιμών στις χρονικές περιόδους  $t$  και  $t+1$ , αντίστοιχα. Τότε, για  $\pi \in \Pi$ ,

$$v_t(\pi) = \max_{a \in A} E \{ q^a(x_t) + \beta \cdot v_{t+1}(\bar{T}(\pi, \alpha, \theta)) | \pi, \alpha \}.$$

$$= \max_{a \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^a(i) + \beta \cdot \int_{\theta \in \Theta} \bar{f}(\theta | \pi, \alpha) \cdot v_{t+1}(\bar{T}(\pi, \alpha, \theta)) d\theta \right\} \quad (14)$$

Αν  $v_{t+1}$  είναι κατά τμήματα γραμμική συνάρτηση, με πεδίο ορισμού το  $\Pi$ , τότε ο τύπος (14) μπορεί να ξαναγραφεί σαν:

$$v_t(\pi) = \max_{\alpha \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^\alpha(i) + \beta \cdot \int_{\theta \in \Theta} \pi \cdot P^\alpha \cdot \bar{R}_\theta^\alpha \cdot \gamma^{l(\pi, \alpha, \theta)} \cdot d\theta \right\} \quad (15)$$

όπου  $\gamma^{l(\pi, \alpha, \theta)}$  είναι το “gradient” του  $v_{t+1}$ , ώστε:

$$T(\pi, \alpha, \theta) \cdot \gamma^{l(\pi, \alpha, \theta)} \geq \bar{T}(\pi, \alpha, \theta) \cdot \gamma^*$$

για όλα τα  $\gamma^*$  που υποστηρίζουν την  $v_{t+1}$ .

Σημειώνουμε ότι ο τύπος (15) είναι παρόμοιος με τον αντίστοιχο τύπο που αφορά την περίπτωση που έχουμε διακριτό αριθμό μηνυμάτων. Η μόνη διαφορά είναι ότι το άθροισμα αντικαθίσταται με ένα ολοκλήρωμα.

Οι White and Harrington (1980) έδειξαν, ότι αν  $v_i+1$  είναι μια κυρτή συνάρτηση, το ίδιο ισχύει και για την  $v_i$ . Εν τούτοις, σε αντιδιαστολή με την περίπτωση που έχουμε ορισμένο αριθμό μηνυμάτων,  $v_i$  δεν είναι κατά τμήματα γραμμική ακόμη και αν  $v_i+1$  είναι κατά τμήματα γραμμική. Αυτή η ιδιότητα καθιστά δυσκολότερο το πρόβλημα όπου έχουμε συνεχή κατανομή μηνυμάτων, σε σύγκριση με εκείνο, όπου έχουμε διακριτή κατανομή μηνυμάτων.

Μια ομοιόμορφη κατανομή χρησιμοποιείται σε μοντέλα όπου δεν έχουμε αρκετή πληροφόρηση. Σχολιάζουμε αυτή την κατανομή, διότι μας δίνει την δυνατότητα να αντιμετωπίσουμε το πρόβλημα όπου έχουμε συνεχή χώρο μηνυμάτων, εφαρμόζοντας τις μεθόδους που αφορούν τον διακριτό χώρο μηνυμάτων.

Ας είναι  $\Theta$  ο χώρος μηνυμάτων. Επίσης στο χρόνο  $t$ , ας είναι  $\Theta(t, i, \alpha)$  ο χώρος μηνυμάτων για την εξέλιξη, δεδομένου ότι η κατάσταση του συστήματος είναι  $i$  και η απόφαση που λαμβάνεται στο χρόνο  $t$  είναι  $\alpha$ . Η συνάρτηση πυκνότητας πιθανότητας για το μήνυμα στον χώρο μηνυμάτων υποθέτουμε ότι έχει την ομοιόμορφη κατανομή.

Μια τετριμμένη περίπτωση έχουμε, όταν  $\Theta(t, i, \alpha) = \Theta(t, j, \alpha) = 0$  για όλα τα ζεύγη των καταστάσεων  $i$  και  $j$ . Τότε μπορούμε ξεκάθαρα να διατυπώσουμε το πρόβλημα, στα πλαίσια μιας πλήρως παρατηρήσιμης MDP. Για παράδειγμα, υποθέτουμε ότι έχουμε ένα πρόβλημα με δύο καταστάσεις, σε τρόπο ώστε μια μηχανή να βρίσκεται στην καλή κατάσταση όταν έχει λειτουργικό κόστος από 100-250 ευρώ την ημέρα, ενώ βρίσκεται στην κακή κατάσταση, όταν το λειτουργικό κόστος είναι από 300-450 ευρώ την ημέρα. Έτσι αν το τρέχον λειτουργικό κόστος είναι 350 ευρώ την ημέρα, τότε είναι προφανές ότι το σύστημα βρίσκεται στην κακή κατάσταση.

Βέβαια, η παραπάνω τεχνική αποτυγχάνει, όταν έχουμε αλληλοεπικάλυψη στις κατανομές μηνυμάτων. Εν τούτοις, αν οι κατανομές μηνυμάτων έχουν την ομοιόμορφη κατανομή, τότε το πρόβλημα μπορεί να αναδιατυπωθεί στα πλαίσια μιας POMDP με ορισμένο αριθμό μηνυμάτων.

$$\text{Ας είναι } \Theta'(t, i, \alpha) = \Theta - \Theta(t, i, \alpha) \quad \text{και} \quad \bar{\Theta}(t, i, \alpha) = \{\Theta(t, i, \alpha), \Theta'(t, i, \alpha)\}$$

Τότε  $\bar{\Theta}(t, i, \alpha)$  είναι μια διαμέριση του χώρου των μηνυμάτων  $\Theta$ . Ας είναι τώρα

$$\Theta_i = \{ \hat{\Theta}_{(t,1)}, \hat{\Theta}_{(t,2)} \dots \hat{\Theta}_{(t,k)} \}$$

το γινόμενο των διαμερίσεων των  $\bar{\Theta}(t, i, \alpha)$  για όλες τις καταστάσεις του συστήματος  $i$  και όλες τις αποφάσεις  $\alpha$ .

Με την παραπάνω μέθοδο χωρίζουμε τον παραμετρικό χώρο  $\Theta$  που είναι συνεχής σε  $k$  το πλήθος κελιά. Κάθε κελί μπορεί να θεωρηθεί σαν ένα διακριτό μήνυμα. Τα μηνύματα είναι όμως ομοιόμορφα κατανεμημένα στον παραμετρικό χώρο  $\Theta$ , οπότε οι πιθανότητες για καθένα από τα μηνύματα που ορίστηκαν με τον νέο αυτό τρόπο, μπορούν να υπολογισθούν σαν το ολοκλήρωμα της συνάρτησης πυκνότητας σε κάθε περιοχή του  $\Theta_i$ .

Επειδή υπάρχει μόνον ορισμένος αριθμός  $N$ , από καταστάσεις, το σύνολο  $\Theta$  έχει ορισμένο αριθμό από στοιχεία. Το κλειδί για την μετατροπή ενός προβλήματος P.O.M.D.P. με συνεχή χώρο μηνυμάτων, που όμως οι κατανομές μηνυμάτων (Signal distributions) είναι ομοιόμορφα καταναμεμημένες, σε ένα πρόβλημα P.O.M.D.P. με διακριτό αριθμό μηνυμάτων, είναι το γεγονός, ότι η μόνη πληροφορία που παρέχεται από το μήνυμα είναι το κελί της διαμέρισης στην οποίαν αυτό ανήκει.

Κάθε στοιχείο του  $\Theta_i$ , μπορεί να θεωρηθεί σαν ένα μήνυμα, οπότε πλέον έχουμε ένα πρόβλημα P.O.M.D.P. με διακριτό αριθμό μηνυμάτων. Το κλειδί για την μετατροπή μιας P.O.M.D.P. με ομοιόμορφη κατανομή μηνυμάτων, σε ένα πρόβλημα P.O.M.D.P. με ορισμένο διακριτό αριθμό μηνυμάτων είναι ότι το  $T(\pi, \alpha, \cdot)$  είναι σταθερό πάνω σε κάθε στοιχείο του  $\Theta_i$ . Ας δείξουμε αυτό το αποτέλεσμα στα λήμματα που ακολουθούν.

Για να γίνουν κατανοητά τα παραπάνω, θεωρούμε μια μηχανή με λειτουργικό κόστος που εξαρτάται από την κατάστασή της. Αν η μηχανή είναι σε άριστη κατάσταση, το λειτουργικό κόστος είναι ομοιόμορφα κατανεμημένο από 100-250 ευρώ ανά ημέρα. Αν τώρα η μηχανή είναι σε κατάσταση επισκευής, το λειτουργικό κόστος είναι ομοιόμορφα κατανεμημένο από 200-350 ευρώ ανά ημέρα. Τέλος αν η μηχανή είναι σε άθλια κατάσταση, το λειτουργικό κόστος είναι ομοιόμορφα κατανεμημένο από 300-450 ευρώ ανά ημέρα. Άρα το κόστος μπορεί να διαιρεθεί σε πέντε περιοχές.

100-200

200-250

250-300

300-350

350-450

Κάθε περιοχή μπορούμε να την δοῦμε σαν διακριτό (ξεχωριστό μήνυμα). Από το γεγονός ότι τα μηνύματα είναι ομοιόμορφα κατανεμημένα στον χώρο των μηνυμάτων, οι πιθανότητες του καθενός από τα νέα μηνύματα όπως τα ορίσαμε, μπορούν να ορισθούν με ολοκλήρωμα της συνάρτησης πυκνότητας πιθανότητας σε κάθε μία από τις περιοχές του  $\Theta_i$ . Για παράδειγμα, υπάρχει 66,67% πιθανότητα το λειτουργικό κόστος να είναι μεταξύ 100-200 ευρώ και 33,33% το λειτουργικό κόστος να είναι μεταξύ 200-250 ευρώ, αν η μηχανή βρίσκεται σε εξαιρετική κατάσταση. Όμοια η πιθανότητα για κάθε διάστημα κόστους μπορεί να υπολογισθεί για την κατάσταση επισκευής, καθώς και για την χειρότερη κατάσταση.

Το κλειδί όμως για την μετατροπή μιας P.O.M.D.P. με ομοιόμορφη κατανομή μηνυμάτων, σε μία P.O.M.D.P. με ορισμένο διακριτό αριθμό μηνυμάτων βασίζεται στο αποτέλεσμα ότι η  $\bar{T}$  ( $\pi, \alpha, .$ ) είναι σταθερή πάνω σε κάθε στοιχείο του  $\Theta_i$ . Προτείνουμε λοιπόν το ακόλουθο λήμμα και θεώρημα που αποδεικνύουν το αποτέλεσμα.

**Λήμμα 3.1:** Για κάθε κατάσταση του συστήματος  $i \in S, f_i^a$  είναι σταθερή για κάθε στοιχείο του  $\Theta_i$ .

**Απόδειξη:** Ας είναι  $\theta_1$  και  $\theta_2$  δύο αυθαίρετα μηνύματα σε ένα οποιοδήποτε  $\hat{\Theta}(t, j) \in \Theta_i$ . Σύμφωνα με την μέθοδο που αφορά την γενίκευση των στοιχείων του συνόλου  $\Theta_i$  θα ισχύουν:  $\hat{\Theta}(t, j) \cap \Theta(t, i, \alpha) = \emptyset$  είτε  $\hat{\Theta}(t, j) \subseteq \Theta(t, i, \alpha)$ . Αυτό μάλιστα θα ισχύει για όλες τις καταστάσεις  $i$ .

Αν  $\hat{\Theta}(t, j) \cap \Theta(t, i, \alpha) = \emptyset$  τότε  $f_i^a(\theta_1) = f_i^a(\theta_2) = 0$ . Αν  $\hat{\Theta}(t, j) \subseteq \Theta(t, i, \alpha)$  τότε, από την υπόθεση της ομοιόμορφης κατανομής μηνυμάτων  $f_i^a(\theta_1) = f_i^a(\theta_2)$ .

Το λήμμα είναι ουσιαστικό, διότι κάνει εμφανή την μέθοδο με την οποία μία P.O.M.D.P. με ομοιόμορφη κατανομή μηνυμάτων, αναδιαμορφώνεται σε

μία P.O.M.D.P. με ορισμένο διακριτό αριθμό μηνυμάτων. Επιπλέον, στην απόδειξη του λήμματος απαιτούμε ότι  $\hat{\Theta}(t, j) \subseteq \Theta(t, i, \alpha)$  ή  $\hat{\Theta}(t, j) \cap \Theta(t, i, \alpha) = \emptyset$  για όλες τις καταστάσεις του συστήματος  $i$ , και επίσης ότι οι στοχαστικές διαδικασίες μηνυμάτων είναι ομοιόμορφα κατανεμημένες στο  $\hat{\Theta}(t, j) \in \Theta_i$ . Στο σημείο αυτό επισημαίνουμε, ότι μπορούμε να επεκτείνουμε αυτό το αποτέλεσμα σε περισσότερες γενικές περιπτώσεις, όπου οι στοχαστικές διαδικασίες (signal-processes) είναι step-functions.

Από το δεδομένο ότι  $f_i^a(\cdot)$  είναι σταθερή για κάθε στοιχείο του  $\Theta_i$ ,  $\bar{R}_\theta^a$  είναι το ίδιο για όλα τα  $\theta$  που ανήκουν στο ίδιο στοιχείο (κελί) της διαμέρισης του  $\Theta_i$ . Παράλληλα,  $\bar{f}(\theta | \pi, \alpha) = \pi \cdot P^\alpha \cdot R_\theta^\alpha \cdot 1$ , επομένως,  $\bar{f}(\cdot | \pi, \alpha)$  είναι σταθερό σε κάθε στοιχείο του  $\Theta_i$ . Τώρα θα αποδείξουμε ότι είναι σταθερό σε κάθε στοιχείο του  $\Theta_i$ .

**Θεώρημα 3.2:**  $\bar{T}(\pi, \alpha, \cdot)$  είναι σταθερό σε κάθε στοιχείο του  $\Theta_i$ .

**Απόδειξη:** Ας είναι  $\hat{\Theta}(t, j) \in \Theta_i$ , ένα αυθαίρετο στοιχείο του  $\Theta_i$  και  $\theta_1, \theta_2$  δύο αυθαίρετα μηνύματα στο  $\hat{\Theta}(t, j)$ . Αφού  $\bar{R}_{\theta_1}^a = \bar{R}_{\theta_2}^a$  και  $\bar{f}(\theta_1 | \pi, \alpha) = \bar{f}(\theta_2 | \pi, \alpha)$  από την (11). Οπότε πλέον θά έχουμε ότι  $\bar{T}(\pi, \alpha, \theta_1) = \bar{T}(\pi, \alpha, \theta_2)$  λόγω της (13). Το αποτέλεσμα λοιπόν αποδείχθηκε.

Τώρα ας δείξουμε, ότι αν οι στοχαστικές διαδικασίες μηνυμάτων είναι ομοιόμορφα κατανεμημένες, τότε η εξίσωσή (14) μπορεί να ξαναγραφεί μορφολογικά εντελώς όμοια με την (7).

Ας θεωρήσουμε από τη σχέση (15) πρώτα το τμήμα του ολοκληρώματος. Ας είναι  $\theta_i$  ένα αυθαίρετο μήνυμα ( $\theta_i \in \hat{\Theta}(t, i)$ )

$$\begin{aligned}
& \text{Από το θεώρημα } \int_{\theta \in \hat{\Theta}(t,i)} \bar{f}(\theta|\pi, \alpha) \cdot v_{t+1}(\bar{T}(\pi, \alpha, \theta)) d\theta \\
& = v_{t+1}(\bar{T}(\pi, \alpha, \theta)) \cdot \int_{\theta \in \hat{\Theta}(t,i)} \bar{f}(\theta|\pi, \alpha) \cdot d\theta = v_{t+1}(\bar{T}(\pi, \alpha, \theta_i)) \cdot \Pr(\theta \in \hat{\Theta}(t, i)|\pi, \alpha) \\
& = \Pr(\theta \in \hat{\Theta}(t, i)|\pi, \alpha) \cdot v_{t+1}(\bar{T}(\pi, \alpha, \theta_i))
\end{aligned}$$

Από το θεώρημα (3.2), κάθε στοιχείο του  $\Theta_i$  μπορούμε να το δούμε σαν ένα ξεχωριστό μήνυμα. Ορίζοντας τώρα  $T(\pi, \alpha, \hat{\Theta}(t, i)) = \bar{T}(\pi, \alpha, \theta_i)$  για  $(\theta_i \in \hat{\Theta}(t, i))$  και λαμβάνοντας υπόψη ότι έχουμε μονάχα ορισμένον αριθμό από στοιχεία του  $\Theta$ , η εξίσωση (15) μπορεί να ξαναγραφεί σαν:

$$\begin{aligned}
v_t(\pi) &= \max_{\alpha \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^\alpha(i) + \beta \cdot \int_{\theta \in \Theta_i} \bar{f}(\theta|\pi, \alpha) \cdot v_{t+1}(\bar{T}(\pi, \alpha, \theta)) d\theta \right\} \\
&= \max_{\alpha \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^\alpha(i) + \beta \cdot \sum_{\Theta(t,i) \in \Theta, \theta \in \Theta(t,i)} \int \bar{f}(\theta|\pi, \alpha) \cdot v_{t+1}(\bar{T}(\pi, \alpha, \theta)) d\theta \right\} \\
&= \max_{\alpha \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^\alpha(i) + \beta \cdot \sum_{\hat{\Theta}(t,i) \in \Theta_i} \Pr(\hat{\Theta}(t,i)|\pi, \alpha) \cdot v_{t+1}(T(\pi, \alpha, \Theta(t, i))) \right\}
\end{aligned}$$

όπου

$$\Pr(\hat{\Theta}(t,i)|\pi, \alpha) = \int_{\theta \in \hat{\Theta}(t,i)} \bar{f}(\theta|\pi, \alpha) d\theta$$

Η τελευταία εξίσωση είναι ίδια μορφολογικά με την εξίσωση (7). Επομένως η (15) μπορεί να ξαναγραφεί όπως ακριβώς η (7), που πρακτικά σημαίνει ότι μία P.O.M.D.P. με ομοιόμορφα καταναμεμημένες τις στοχαστικές διαδικασίες μηνυμάτων, μπορεί να αναδιατυπωθεί σαν μία P.O.M.D.P. με ορισμένο διακριτό αριθμό μηνυμάτων.



## 6. Μέθοδοι για την λύση των P.O.M.D.P.S. με συνεχή χώρο μηνυμάτων

Στην παράγραφο αυτή θα υπολογίσουμε την τιμή και τα gradients του  $Hv(\pi)$  για αυθαίρετο  $\pi \in \Pi$ , αν  $v$  είναι κατά τμήματα γραμμική συνάρτηση (υπολογισμός μέσω της (15)). Ας είναι  $\Gamma$  ένα σύνολο από “gradients” του  $v$ . Ας είναι τώρα  $\bar{\gamma}$  ένα gradient, ( $\bar{\gamma} \in \Gamma$ ), και  $a$  μία απόφαση ( $a \in A$ ), τότε ορίζουμε:

$$\text{Ορισμός } \Theta_{\pi, a, \bar{\gamma}} = \{\theta \in \Theta : \pi \cdot P^a \cdot \bar{R}_\theta^a \cdot \bar{\gamma} \geq \pi \cdot P^a \cdot \bar{R}_\theta^a \cdot \gamma, \forall \gamma \in \Gamma\} \quad (16)$$

Τότε η (15) λόγω της (16) μπορεί να ξαναγραφεί σαν:

$$\begin{aligned} Hv(\pi) &= \max_{a \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^a(i) + \beta \int_{\theta \in \Theta} \pi \cdot P^a \cdot \bar{R}_\theta^a \cdot \gamma^{l(\pi(a, \theta))} \cdot d\theta \right\} = \\ &= \max_{a \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^a(i) + \beta \cdot \sum_{\bar{\gamma} \in \Gamma} \int_{\theta \in \Theta_{\pi, a, \bar{\gamma}}} \pi \cdot P^a \cdot \bar{R}_\theta^a \cdot \bar{\gamma} \cdot d\theta \right\} = \\ &= \max_{a \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^a(i) + \beta \cdot \sum_{\bar{\gamma} \in \Gamma} \pi \cdot P^a \cdot \int_{\theta \in \Theta_{\pi, a, \bar{\gamma}}} \bar{R}_\theta^a \cdot \bar{\gamma} \cdot d\theta \right\} \end{aligned} \quad (17)$$

Προκειμένου να απλοποιήσουμε τον υπολογισμό του  $\int_{\theta \in \Theta_{\pi, a, \bar{\gamma}}} \pi \cdot P^a \cdot \bar{R}_\theta^a \cdot \bar{\gamma} \cdot d\theta$  ας υποθέσουμε ότι για ένα δεδομένο  $\pi \in \Pi$ ,  $a \in A$ , και  $\bar{\gamma} \in \Gamma$ , υπάρχει μονάχα ένας ορισμένος αριθμός από επικοινωνούντα σύνολα στο  $\Theta_{\pi, a, \bar{\gamma}}$

Επιπλέον, μπορούμε να υποθέσουμε ότι τα σύνορα του  $\Theta_{\pi, a, \bar{\gamma}}$  έχουν μέτρο 0.

Αφού  $\bar{R}_\theta^a$  είναι ένας διαγώνιος πίνακας που έχει σαν ι-διαγώνιο στοιχείο το  $f_i^a(\theta)$  και  $\bar{\gamma}$  είναι ένα N-διάστατο διάνυσμα, που έχει  $\bar{\gamma}_i$  το i-στοιχείο του. Παραστατικά  $\bar{\gamma} = (\bar{\gamma}_1, \bar{\gamma}_2, \dots, \bar{\gamma}_N)$ .

Τότε με σκοπό την απλοποίηση των παραστάσεων, και για να φανεί η μορφολογική ομοιότητα θεωρούμε

$$\xi^\alpha(\pi, \bar{\gamma}) = \begin{pmatrix} \int_{\theta \in \Theta_{\pi, \alpha, \bar{\gamma}}} f_1^\alpha(\theta) \cdot \bar{\gamma}_1 \cdot d\theta \\ \int_{\theta \in \Theta_{\pi, \alpha, \bar{\gamma}}} f_2^\alpha(\theta) \cdot \bar{\gamma}_2 \cdot d\theta \\ \vdots \\ \int_{\theta \in \Theta_{\pi, \alpha, \bar{\gamma}}} f_N^\alpha(\theta) \cdot \bar{\gamma}_N \cdot d\theta \end{pmatrix} \quad (18)$$

Επομένως, η (17) μπορεί να ξαναγραφεί σαν

$$\begin{aligned} \text{Hu}(\pi) &= \max_{\alpha \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^\alpha(i) + \beta \sum_{\Theta_{\pi, \alpha, \bar{\gamma}}} \pi P^\alpha \cdot \int_{\theta \in \Theta_{\pi, \alpha, \bar{\gamma}}} \bar{R}_\theta^\alpha \cdot \bar{\gamma} \cdot d\theta \right\} \\ &= \max_{\alpha \in A} \left\{ \sum_{i=1}^N \pi_i \cdot q^\alpha(i) + \beta \cdot \sum_{\Theta_{\pi, \alpha, \bar{\gamma}}} \pi \cdot P^\alpha \cdot \xi^\alpha(\pi, \bar{\gamma}) \right\} \\ &= \max_{\alpha \in A} \left\{ \pi \cdot \left[ q^\alpha + \beta \cdot \sum_{\Theta_{\pi, \alpha, \bar{\gamma}}} P^\alpha \cdot \xi^\alpha(\pi, \bar{\gamma}) \right] \right\} \end{aligned} \quad (19)$$

Παρατηρούμε ότι η παράσταση στις αγκύλες  $q^\alpha + \beta \sum_{\Theta_{\pi, \alpha, \bar{\gamma}}} P^\alpha \cdot \xi^\alpha(\pi, \bar{\gamma})$  είναι ένα N-διάστατο διάνυσμα (gradient) και ένα “gradient” επομένως του τελεστού Hu στο  $\pi$ .

Μολονότι οι υποθέσεις σχετικά με το  $\Theta_{\pi, \alpha, \bar{\gamma}}$  είναι ισχυρές, αυτές οι υποθέσεις είναι στην πραγματικότητα αληθείς για πολλές κατανομές που χρησιμοποιούνται συνήθως. Πριν ωστόσο αρχίσουμε τους υπολογισμούς, θεωρείται απαραίτητο να επιβεβαιώσουμε αν οι δεδομένες κατανομές μηνυμάτων

έχουν έναν ορισμένο αριθμό από επικοινωνούντα σύνολα στο  $\Theta_{\pi, \alpha, \bar{\gamma}}$ .

Για παράδειγμα, αν οι στοχαστικές διαδικασίες μηνυμάτων είναι εκθετικές και υπάρχουν μόνον δυο καταστάσεις στο σύστημά μας, μπορούμε να δείξουμε ότι οι υποθέσεις ισχύουν. Τώρα ας χρησιμοποιήσουμε εκθετική κατανομή στο παράδειγμα που ακολουθεί.

### Παράδειγμα

Για να διευκρινίσουμε την μέθοδο παραθέτουμε το παράδειγμα

$$P^1 = \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix} \quad q^1 = \begin{bmatrix} -4 \\ 4 \end{bmatrix}$$

$$P^2 = \begin{bmatrix} 0.5 & 0.5 \\ 0.4 & 0.6 \end{bmatrix} \quad q^2 = \begin{bmatrix} 0 \\ 3 \end{bmatrix}$$

Οι στοχαστικές διαδικασίες μηνυμάτων έχουν τις ακόλουθες συναρτήσεις πυκνότητας πιθανότητας: Για  $\theta > 0$

$$f_1^1(\theta) = e^{-\theta}, \quad f_2^1(\theta) = 10 \cdot e^{-10\theta}$$

$$f_1^2(\theta) = 3 \cdot e^{-3\theta}, \quad f_2^2(\theta) = 2 \cdot e^{-2\theta}$$

και  $f_1^1(\theta) = f_2^1(\theta) = f_1^2(\theta) = f_2^2(\theta) = 0$  για  $\theta \leq 0$

Τώρα ας υποθέσουμε ότι  $\beta = 0.9$ ,  $\pi = (0, 1)$  και  $\Gamma = \{\gamma^1, \gamma^2\}$  όπου  $\gamma^1 = [-4, 4]^T$   $\gamma^2 = [0, 3]^T$ . Υπολογίζουμε πρώτα το  $\Theta_{\pi, 1, \gamma^1}$ , έχουμε

$$[0, 1] \cdot \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix} \cdot \begin{bmatrix} e^{-\theta} & 0 \\ 0 & 10 \cdot e^{-10\theta} \end{bmatrix} \cdot \begin{bmatrix} -4 \\ 4 \end{bmatrix} \geq [0, 1] \cdot \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix} \cdot \begin{bmatrix} e^{-\theta} & 0 \\ 0 & 10 \cdot e^{-10\theta} \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 3 \end{bmatrix}$$

Λύνοντας την ανισότητα έχουμε:  $\Theta_{\pi, 1, \gamma^1} = \{0 < \theta \leq 0.1\}$ . Άρα  $\Theta_{\pi, 1, \gamma^2} = \{0 < \theta < 0.2\}$ . Όμοια βρίσκουμε  $\Theta_{\pi, 2, \gamma^1} = \{\theta < -177\}$ , που είναι έξω από το πεδίο ορισμού των συναρτήσεων πυκνότητας, οπότε  $\Theta_{\pi, 2, \gamma^1} = \emptyset$  και  $\Theta_{\pi, 2, \gamma^2} = \{\theta > 0\}$ . Κατόπιν χρησιμοποιώντας τα  $\Theta_{\pi, 1, \gamma^1}$ ,  $\Theta_{\pi, 1, \gamma^2}$ ,  $\Theta_{\pi, 2, \gamma^1}$ ,  $\Theta_{\pi, 2, \gamma^2}$  στην (19), πετυχαίνουμε:

$$\begin{aligned}
Hv([0,1]) &= \max_{[0,1]} \left\{ \begin{bmatrix} -4 \\ 4 \end{bmatrix} + 0.9 \cdot \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix} \cdot \left( \begin{bmatrix} \int_0^{0.1} -4 \cdot e^{-\theta} \cdot d\theta \\ \int_0^{0.1} 4 \cdot 10 \cdot e^{-10\theta} \cdot d\theta \end{bmatrix} + \begin{bmatrix} 0 \\ \int_{0.1}^{\infty} 3 \cdot 10 \cdot e^{-10\theta} \cdot d\theta \end{bmatrix} \right) \right\}; \\
&= \max_{[0,1]} \left\{ \begin{bmatrix} 0 \\ 3 \end{bmatrix} + 0.9 \cdot \begin{bmatrix} 0.5 & 0.5 \\ 0.4 & 0.6 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ \int_0^{\infty} 3 \cdot 2 \cdot e^{-2\theta} \cdot d\theta \end{bmatrix} \right\} \\
&= \max \left\{ [0,1] \cdot \begin{bmatrix} -3.62 \\ 5.46 \end{bmatrix}; [0,1] \cdot \begin{bmatrix} 1.35 \\ 4.62 \end{bmatrix} \right\} = 5.46
\end{aligned}$$

Επομένως  $Hv$  στο  $[0, 1]$  είναι 5.46 και το "gradient" που αντιστοιχεί στην κατάσταση αυτή είναι  $[-3.62, 5.46]^T$ .

### Συμπεράσματα

Αναπτύξαμε μια μέθοδο που μπορεί να λύσει ένα πρόβλημα μερικά παρατηρήσιμων Μαρκοβιανών αλυσίδων με συνεχή χώρο μηνυμάτων. Ακολουθήσαμε την ίδια πορεία, με εκείνη που ακολούθησαν αρκετοί ερευνητές και αφορούσε το πρόβλημα P.O.M.D.P. με διακριτό αριθμό μηνυμάτων. Επεκτείναμε τα αποτελέσματα για την περίπτωση που ο χώρος μηνυμάτων είναι συνεχής. Κλειδί για την επέκταση αυτή είναι το λήμμα (3.1) και το θεώρημα (3.2) που ουσιαστικά διαμερίζουν τον χώρο μηνυμάτων,  $\Theta$ , σε τρόπον ώστε κάθε κελλί της διαμέρισης που προτείνεται, να μπορεί να θεωρηθεί σαν ένα ξεχωριστό διακριτό μήνυμα.

Επομένως μπορούμε να αναδιατυπώσουμε το πρόβλημα P.O.M.D.P. με συνεχή χώρο μηνυμάτων, ώστε να αντιμετωπισθεί με μια μέθοδο επίλυσης ανάλογη με εκείνη που λύνει το πρόβλημα με διακριτό αριθμό μηνυμάτων.

Η παραπάνω επέκταση κρίθηκε αναγκαία, διότι σε αρκετά προβλήματα της καθημερινότητας ο χώρος μηνυμάτων είναι συνεχής. Για παράδειγμα η θερμοκρασία μιας μηχανής είναι συνεχής μεταβλητή.

## References

1. Bertsekas, D.P. (1975) Dynamic programming and stochastic control. Academic Press New York.
2. Denardo, E.V. (1967) Contraction mapping in theory underlying dynamic programming. SIAM Reviews 9, (165-177).
3. Howard, R.A. (1960) Dynamic programming and Markov processes. MIT.
4. Lovejoy, W.S. (1991) Computationally feasible bounds for P.O.M.D.P. Operation research, vol 39, No 1, 1991 (162-176).
5. Littman, M.L. Algorithms for sequential decision making PhD thesis, Department of computer science, Brown University (1996).
6. Sondik E.J. (1978) Optimal control of P.O.M.D.P. over the infinite horizon: discount costs, Operation research, vol 2, March 1978 (282-303).
7. Sondik E.J. (1973) Optimal control of P.O.M.D.P. over the infinite horizon: discount costs, Operation research, 21 (1071-1088).
8. Sondik E.J. (1971) Optimal control of P.O.M.D.P.. PhD thesis. Department of electrical engineering, Stanford University 1971.
9. White C.C. and Harrington D.P. (1980) Application of Jensen's Inequality to adaptive Suboptimal Design, Journal of optimization theory and applications 32, (89-99).
10. Whitt, W. (1978). Approximations of dynamics programs, I Mathematics of operation research 3, (231-243).
11. Abraham Grosfeld-Nir (1996) A Two-state partially observable Markov decision process with uniformly distributed observations. Operation research vol 44, no 3, (458-462).