



**ΑΡΙΣΤΟΤΕΛΕΙΟ ΠΑΝΕΠΙΣΤΗΜΙΟ**

**ΘΕΣΣΑΛΟΝΙΚΗΣ**

**ΤΜΗΜΑ ΜΑΘΗΜΑΤΙΚΩΝ**

**ΠΡΟΓΡΑΜΜΑ ΜΕΤΑΠΤΥΧΙΑΚΩΝ ΣΠΟΥΔΩΝ**

**«ΣΤΑΤΙΣΤΙΚΗ ΚΑΙ ΜΟΝΤΕΛΟΠΟΙΗΣΗ»**

**ΜΕΤΑΠΤΥΧΙΑΚΗ ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ**

Επίλυση στοχαστικών παιγνίων με μεθόδους γραμμικού προγραμματισμού

Κωνσταντίνος Ζυγογιάννης

**ΕΠΙΒΛΕΠΟΥΣΑ:** Αλεξάνδρα Παπαδοπούλου

Αναπ. Καθηγ. Α.Π.Θ.

**Θεσσαλονίκη, Δεκέμβριος 2022**





ΑΡΙΣΤΟΤΕΛΕΙΟ ΠΑΝΕΠΙΣΤΗΜΙΟ

ΘΕΣΣΑΛΟΝΙΚΗΣ

ΤΜΗΜΑ ΜΑΘΗΜΑΤΙΚΩΝ

ΠΡΟΓΡΑΜΜΑ ΜΕΤΑΠΤΥΧΙΑΚΩΝ ΣΠΟΥΔΩΝ

«ΣΤΑΤΙΣΤΙΚΗ ΚΑΙ ΜΟΝΤΕΛΟΠΟΙΗΣΗ»

**ΜΕΤΑΠΤΥΧΙΑΚΗ ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ**

Επίλυση στοχαστικών παιγνίων με μεθόδους γραμμικού προγραμματισμού

*Κωνσταντίνος Ζυγογιάννης*

**ΕΠΙΒΛΕΠΟΥΣΑ:** Αλεξάνδρα Παπαδοπούλου

Αναπ. Καθηγ. Α.Π.Θ.

Εγκρίθηκε από την Τριμελή Εξεταστική Επιτροπή:

.....

Γ.Κωνσταντινίδης

Διδάκτορας Α.Π.Θ.

.....

Α.Παπαδοπούλου

Αναπ. Καθηγήτρια Α.Π.Θ.

.....

Χ. Πελέκης

Επικ. Καθηγητής ΑΠΘ

**Θεσσαλονίκη, Δεκέμβριος 2022**

.....

Κωνσταντίνος Ζυγογιάννης

Πτυχιούχος Μαθηματικός Α.Π.Θ.

Copyright © Κωνσταντίνος Ζυγογιάννης, 2022

Με επιφύλαξη παντός δικαιώματος. All rights reserved.

Απαγορεύεται η αντιγραφή, αποθήκευση και διανομή της παρούσας εργασίας, εξ ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, αποθήκευση και διανομή για σκοπό μη κερδοσκοπικό, ερευνητικής ή εκπαιδευτικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν κείμενο. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα.

Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και δεν πρέπει να ερμηνευθεί ότι εκφράζουν τις επίσημες θέσεις του Α.Π.Θ.

## ΠΡΟΛΟΓΟΣ

Η παρούσα διπλωματική εργασία εκπονήθηκε στο πλαίσιο του Μεταπτυχιακού Προγράμματος Σπουδών «Στατιστική και Μοντελοποίηση» του τμήματος Μαθηματικών του Αριστοτελείου Πανεπιστημίου Θεσσαλονίκης. Θα ήθελα να ευχαριστήσω θερμά την κυρία Αλεξάνδρα Παπαδοπούλου για τις πολύ χρήσιμες συμβουλές, την άψογη συνεργασία και καθοδήγησή της καθ' όλη τη διάρκεια προετοιμασίας και υλοποίησης της εργασίας αυτής. Ακόμη, θα ήθελα να ευχαριστήσω την οικογένεια και τους φίλους μου, για την αδιάκοπη και ανιδιοτελή στήριξή τους όλα αυτά τα χρόνια.



## ΠΕΡΙΛΗΨΗ

Η παρούσα εργασία περιέχει κατ' αρχήν μια ανασκόπηση με κεντρικό αντικείμενο τα Στοχαστικά Παίγνια. Το πρώτο κεφάλαιο αποτελεί μια εισαγωγή στις Μαρκοβιανές διαδικασίες απόφασης , ένα πολύ ισχυρό εργαλείο για την μελέτη των πιο «πολύπλοκων» στοχαστικών παιγνίων, τα οποία αποτελούν αντικείμενο μελέτης του δεύτερου κεφαλαίου. Σημαντικό μέρος της θεωρίας που αναπτύσσεται στα δύο πρώτα κεφάλαια αντλήθηκε από το βιβλίο των Jerzy Filar και Koos Vrieze με τίτλο *Competitive Markov Decision Processes* (1997). Το τρίτο κεφάλαιο βασίστηκε στην εργασία των S. Sinha και K. G. Bakshi με τίτλο *Zero-Sum Two Person Perfect Information Semi-Markov Games: A Reduction* (2022). Στην εργασία αναδεικνύεται ο ρόλος του Γραμμικού Προγραμματισμού, στην ανάλυση και επίλυση προβλημάτων που μοντελοποιούνται είτε ως διαδικασίες απόφασης είτε ως παίγνια. Όπως θα γίνει αντιληπτό, μελετώντας συγκεκριμένες κατηγορίες παιγνίων ως προς τη δομή τους, εν τέλει, ο υπολογισμός μιας βέλτιστης πολιτικής η οποία θα αποφέρει το μέγιστο αναμενόμενο κέρδος (ή την ελάχιστη ζημία) αποτελεί πρόβλημα που ισοδυναμεί με την επίλυση κατάλληλα κατασκευασμένων γραμμικών προγραμμάτων. Στο τρίτο κεφάλαιο γενικεύονται τα παίγνια που προαναφέρθηκαν, τα οποία πλέον θα καλούνται Ημιμαρκοβιανά Παίγνια. Πραγματοποιείται ανάλυση μιας συγκεκριμένης κλάσης αυτής της ιδιαίτερα ευρείας κατηγορίας παιγνίων, τα Ημιμαρκοβιανά Παίγνια Τέλειας Πληροφόρησης. Τέλος, γίνεται υλοποίηση του σχετικού αλγορίθμου επίλυσης τέτοιων παιγνίων με το λογισμικό MATLAB και παρατίθενται τα σχετικά αποτελέσματα.





## ***ABSTRACT***

The present paper initially includes a review with Stochastic Games as the main subject. The first chapter is an introduction to Markov Decision Processes, a significant tool in order to study the more complicated case of Stochastic Games, which are going to be the subject of the next chapter. A big part of the theory developed in the first two chapters is due to the book Competitive Markov Decision Processes (1997), written by Jerzy Filar και Koos Vrieze. The third chapter is based on a research paper by S. Sinha και K. G. Bakshi, titled as Zero-Sum Two Person Perfect Information Semi-Markov Games: A Reduction (2022). Through this paper, the role of Linear Programming stands out in analyzing and solving problems which can be modelled as either decision processes or stochastic games. When analyzing specific subclasses of games regarding their structure, it is concluded that finding the policy that maximizes the expected profit (or minimizes the expected damage) is a problem equivalent to solving suitably constructed linear programs. In the third chapter, semi-Markov games, a generalization of the games mentioned so far, are being introduced. The main subject of this chapter is a specific category of semi-Markov games, called Perfect Information semi-Markov games. Finally, a specific algorithm regarding the solution of such games is implemented through the MATLAB software and its results are briefly discussed.



## Περιεχόμενα

<b>1 Μαρκοβιανές Διαδικασίες Απόφασης.....</b>	<b>12</b>
1.0 Εισαγωγή.....	12
1.1 Η διαδικασία με συντελεστή προεξόφλησης και η τερματιζόμενη διαδικασία .....	12
1.1.1 Μαρκοβιανό μοντέλο απόφασης $G_R$ με συντελεστή προεξόφλησης $\beta$ .....	13
1.1.2 Τερματιζόμενο μαρκοβιανό μοντέλο απόφασης $G_T$ .....	16
1.2 Η μαρκοβιανή διαδικασία απόφασης πεπερασμένου ορίζοντα .....	16
1.3 Γραμμικός προγραμματισμός και αθροίστιμες μαρκοβιανές διαδικασίες απόφασης .....	23
1.4 Η αδιαχώριστη διαδικασία οριακού μέσου .....	31
1.5 Συμπεριφορικές και μαρκοβιανές στρατηγικές.....	35
<b>2 Στοχαστικά παίγνια και γραμμικός προγραμματισμός .....</b>	<b>40</b>
2.0 Εισαγωγή.....	40
2.1 Στοχαστικά παίγνια με συντελεστή προεξόφλησης.....	40
2.2 Γραμμικός προγραμματισμός και στοχαστικά παίγνια με συντελεστή προεξόφλησης.....	48
2.2.1 Παίγνια με προεξόφληση με έναν ελεγκτή .....	49
2.2.2. Παίγνια με προεξόφληση τύπου SER-SIT .....	54
2.2.3. Παίγνια με προεξόφληση εναλλασσόμενου ελεγκτή .....	57
2.3 Στοχαστικά παίγνια οριακού μέσου .....	59
2.4 Στοχαστικά παίγνια οριακού μέσου με έναν ελεγκτή .....	63
<b>3 Ημιμαρκοβιανά παίγνια τέλει πληροφόρησης.....</b>	<b>76</b>
3.0 Εισαγωγή .....	76
3.1 Πεπερασμένα ημιμαρκοβιανά παίγνια δύο παικτών και μηδενικού αθροίσματος .....	76
3.2 Ημιμαρκοβιανά παίγνια με κριτήριο πληρωμής τον οριακό λόγο των μέσων.....	78
3.3 Πεπερασμένες ημιμαρκοβιανές διαδικασίες απόφασης.....	81
3.4 Αλγόριθμος επίλυσης ημιμαρκοβιανών παιγνίων τέλει πληροφόρησης.....	86
3.5 Επίλυση μέσω γραμμικού προγραμματισμού και μέσω πλήρους απαρίθμησης .....	91
<b>Παράρτημα .....</b>	<b>102</b>
<b>Βιβλιογραφία .....</b>	<b>108</b>

## Κεφάλαιο 1

### Μαρκοβιανές Διαδικασίες Απόφασης

#### 1.0 Εισαγωγή

Το 1<sup>ο</sup> κεφάλαιο της εργασίας αφορά στις Μαρκοβιανές Διαδικασίες Απόφασης και, συγκεκριμένα, πώς οι διαδικασίες αυτές μπορούν να μελετηθούν ως Στοχαστικά παίγνια, τα οποία αποτελούν αντικείμενο μελέτης του επόμενου κεφαλαίου, όπου υφίσταται μόνον ένας παίκτης-ελεγκτής. Τόσο στο παρόν κεφάλαιο όσο και γενικότερα στα πλαίσια της εργασίας αυτής, θα μας απασχολήσουν διαδικασίες απόφασης όπου ο χώρος καταστάσεων και ενεργειών του/των παίκτη/παικτών είναι πεπερασμένος. Γίνεται αναφορά σε συγκεκριμένες κατηγορίες διαδικασιών απόφασης, οι οποίες διαφέρουν ως προς περιορισμούς που πιθανόν να συνδέονται με τη δομή αυτών ή/και το κριτήριο πληρωμής του ελεγκτή στην εκάστοτε διαδικασία. Παράλληλα, θα γίνει αντιληπτός ο τρόπος που μπορούν να επιλυθούν τέτοιου τύπου μοντέλα απόφασης με μεθόδους που προκύπτουν από το αντικείμενο του Γραμμικού Προγραμματισμού.

#### 1.1 Η διαδικασία με συντελεστή προεξόφλησης και η τερματιζόμενη διαδικασία

Θα συμβολίζουμε με  $S_t$  την τυχαία μεταβλητή που παίρνει τιμές από το σύνολο  $S = \{1, 2, \dots, N\}$ , το οποίο καλούμε χώρο καταστάσεων. Όταν λέμε ότι η διαδικασία βρίσκεται στην κατάσταση  $s \in S$  στον χρόνο  $t$ , θα εννοούμε ότι πραγματοποιείται το γεγονός  $\{S_t = s\}$ .

Η διαδικασία απόφασης ελέγχεται από έναν «decision-maker» (d-m), ο οποίος, δεδομένου ότι το σύστημα βρίσκεται στην κατάσταση  $s \in S$ , επιλέγει μια ενέργεια  $a \in A(s) = \{1, 2, \dots, m(s)\}$ . Η επιλογή της  $a \in A(s)$  είναι συνδεδεμένη με μία άμεσα παραγόμενη αμοιβή  $r(s, a)$ , που θα καλούμε αμοιβή ενός βήματος, την οποία λαμβάνει ο decision-maker και έπειτα η διαδικασία μεταπηδά σε μία νέα κατάσταση  $s' \in S$  και εξελίσσεται με τον ίδιο τρόπο από την νέα κατάσταση  $s'$ .

Ορίζουμε:

$$p(s'|s, a) = \mathbb{P}\{S_{t+1} = s' | S_t = s, A_t = a\}$$

να είναι η πιθανότητα μετάβασης στην κατάσταση  $s'$ , δεδομένου ότι βρισκόμαστε στην κατάσταση  $s$  και επιλέγεται η ενέργεια  $a \in A(s)$ . Όπως είναι φανερό, έχουμε ομοιογένεια ως προς το χρόνο, δηλαδή οι πιθανότητες μετάβασης είναι ανεξάρτητες του  $t$ .

Μία στρατηγική  $\mathbf{f} = [f(1), f(2), \dots, f(N)]$  είναι ένα διάνυσμα-γραμμή, όπου το  $s$ -οστό στοιχείο του είναι το μη αρνητικό διάνυσμα γραμμής:

$\mathbf{f}(s) = [f(s, 1), f(s, 2), \dots, f(s, m(s))]$  για το οποίο ισχύει ότι :

$$\sum_{a=1}^{m(s)} f(s, a) = 1$$

$f(s, a) \doteq$  Η πιθανότητα να επιλεγεί η ενέργεια  $a \in A(s)$ , δεδομένου ότι η διαδικασία βρίσκεται στην κατάσταση  $s$ .

Μία στρατηγική  $\mathbf{f}$  θα καλείται *καθαρή* (ή *ντετερμινιστική*) αν

$f(s, a) \in \{0, 1\}$  , για όλα τα  $a \in A(s), s \in S$

Οπότε μία στρατηγική  $\mathbf{f}$  ορίζει έναν  $N \times N$  πίνακα πιθανοτήτων μετάβασης

$$\mathbf{P}(\mathbf{f}) = (p(s'|s, \mathbf{f}))_{s, s' \in S}$$

τα στοιχεία του οποίου ορίζονται ως :

$$p(s'|s, \mathbf{f}) = \sum_{a=1}^{m(s)} p(s'|s, a) f(s, a)$$

Φυσικά, για συγκεκριμένο  $s \in S$ , για όλα τα  $a \in A(s)$ , θα έχουμε ότι :

$$\sum_{s'=1}^N p(s'|s, a) = 1$$

Δηλαδή, ο πίνακας  $\mathbf{P}(\mathbf{f})$  είναι στοχαστικός, που σημαίνει ότι όλες του οι γραμμές αθροίζουν στη μονάδα, και συνεπώς, ορίζει μια μοναδική Μαρκοβιανή Αλυσίδα. Φυσικά, μια διαφορετική στρατηγική  $\mathbf{f}'$  ορίζει διαφορετική Μαρκοβιανή Αλυσίδα μέσω του  $\mathbf{P}(\mathbf{f}')$ .

### 1.1.1 Μαρκοβιανό μοντέλο απόφασης $\Gamma_\beta$ με συντελεστή προεξόφλησης $\beta$

Έστω  $\{R_t\}_{t=0}^\infty$  η ακολουθία των αμοιβών, όπου  $R_t$  θα είναι η αμοιβή για τον d-m που αντιστοιχεί στην χρονική περίοδο  $[t, t + 1)$ . Εφόσον οριστεί η αρχική κατάσταση  $s$  και μία στρατηγική  $\mathbf{f}$ , τότε γνωρίζουμε την κατανομή πιθανότητας της  $R_t$  και η μέση τιμή αυτής είναι καλά ορισμένη. Θα συμβολίζουμε τη μέση τιμή με :

$$E_{sf}[R_t] \doteq E_f[R_t | S_0 = s]$$

Η συνολική αναμενόμενη τιμή με προεξόφληση για την  $S_0 = s$  μέσω της  $\mathbf{f}$  θα είναι :

$$v_\beta(s, \mathbf{f}) = \sum_{t=0}^{\infty} \beta^t E_{sf}[R_t], \quad \text{όπου } \beta \in [0,1) \text{ είναι ο συντελεστής προεξόφλησης} \quad (1.1.1)$$

Η φυσική ερμηνεία αυτού είναι πως: Η απόδοση μιας μονάδας στο χρόνο  $t + 1$  'αξίζει' μόνο  $\beta$  φορές απ'ότι άξιζε στο χρόνο  $t$  ( $\beta < 1$ ).

Στη συνέχεια, ορίζουμε το διάνυσμα αναμενόμενης αμοιβής ενός βήματος ως:

$$\mathbf{r}(\mathbf{f}) = [r(1, \mathbf{f}), r(2, \mathbf{f}), \dots, r(N, \mathbf{f})]^T$$

όπου, για κάθε  $s \in S$

$$r(s, \mathbf{f}) \doteq \sum_{a \in A(s)} r(s, a) f(s, a)$$

Για αυθαίρετο  $s \in S$ , υπολογίζουμε:

$$\begin{aligned} E_{sf}[R_0] &= r(s, \mathbf{f}) = [\mathbf{r}(\mathbf{f})]_s \\ E_{sf}[R_1] &= \sum_{s'=1}^N p(s'|s, \mathbf{f}) r(s', \mathbf{f}) = [P(\mathbf{f})\mathbf{r}(\mathbf{f})]_s \\ &\dots \\ E_{sf}[R_t] &= \sum_{s'=1}^N p_t(s'|s, \mathbf{f}) r(s', \mathbf{f}) = [P^t(\mathbf{f})\mathbf{r}(\mathbf{f})]_s \end{aligned}$$

όπου με  $[\mathbf{a}]_s$  συμβολίζουμε το  $s$ -οστό στοιχείο του διανύσματος  $\mathbf{a}$  και με  $p_t(s'|s, \mathbf{f})$  την πιθανότητα μετάβασης  $t$ -βήματος από την κατάσταση  $s$  στην  $s'$  για την αλυσίδα που ορίζει η  $\mathbf{f}$ .

Από τη θεωρία γνωρίζουμε ότι η  $t$ -δύναμη του πίνακα  $P(\mathbf{f})$  περιγράφει ακριβώς την πιθανότητα μετάβασης  $t$ -βήματος. Δηλαδή :

$$P^t(\mathbf{f}) = (p_t(s'|s, \mathbf{f}))_{s,s'=1}^N$$

Αντικαθιστώντας στην (1.1.1) , για  $\mathbf{v}_\beta(\mathbf{f}) = (v(1, \mathbf{f}), v(2, \mathbf{f}), \dots, v(N, \mathbf{f}))^T$  παίρνουμε ότι:

$$\mathbf{v}_\beta(\mathbf{f}) = \sum_{t=0}^{\infty} \beta^t P^t(\mathbf{f}) \mathbf{r}(\mathbf{f}) \quad (1.1.2)$$

όπου ορίζουμε  $P^0(\mathbf{f}) \doteq I_N$  (με  $I_N$  συμβολίζουμε τον μοναδιαίο πίνακα διάστασης  $N \times N$ )

Ο πίνακας  $[I - \beta P(\mathbf{f})]$  είναι αντιστρέψιμος και

$$[I - \beta P(f)]^{-1} = I + \beta P(f) + \beta^2 P^2(f) + \dots$$

### Απόδειξη

Αρχικά παρατηρούμε ότι ο πίνακας  $\beta P(f)$  είναι υποστοχαστικός, εφόσον ο  $P(f)$  είναι στοχαστικός και  $0 < \beta < 1$ . Οπότε θα ισχύει ότι :

$$\lim_{n \rightarrow \infty} (\beta P(f))^n = 0$$

$$\text{Ισχύει ότι: } I^n - (\beta P(f))^n = (I - \beta P(f)) \left( I + \beta P(f) + (\beta P(f))^2 + \dots + (\beta P(f))^{n-1} \right)$$

Παίρνοντας όρια στην παραπάνω σχέση, λαμβάνουμε:

$$\lim_{n \rightarrow \infty} (I^n - (\beta P(f))^n) = \lim_{n \rightarrow \infty} (I - \beta P(f)) \left( I + \beta P(f) + (\beta P(f))^2 + \dots + (\beta P(f))^{n-1} \right)$$

$$\Rightarrow I = (I - \beta P(f)) \sum_{n=0}^{\infty} (\beta P(f))^n \quad (1.1.3)$$

Παίρνοντας ορίζουσες στην (1.1.3) , έχουμε ότι :

$$1 = \det(I - \beta P(f)) \det \left( \sum_{n=0}^{\infty} (\beta P(f))^n \right)$$

Από την τελευταία σχέση συμπεραίνουμε πως  $\det(I - \beta P(f)) \neq 0$ , δηλαδή ορίζεται ο πίνακας  $(I - \beta P(f))^{-1}$ .

Οπότε, από την (1.1.3), πολλαπλασιάζοντας επί  $(I - \beta P(f))^{-1}$  :

$$(I - \beta P(f))^{-1} = \sum_{n=0}^{\infty} (\beta P(f))^n$$

δηλαδή, αποδείχθηκε το ζητούμενο. ■

Τέλος, αντικαθιστώντας στην (1.1.2) , λαμβάνουμε ότι :

$$v_\beta(f) = [I - \beta P(f)]^{-1} r(f)$$

ή, ισοδύναμα, ότι:

$$v_\beta(f) = r(f) + \beta P(f) v_\beta(f) \quad (1.1.4)$$

### 1.1.2 Τερματιζόμενο Μαρκοβιανό μοντέλο Απόφασης $\Gamma_T$

Σε αυτή τη μοντελοποίηση της διαδικασίας απόφασης, κάνουμε την υπόθεση πως :

$$\sum_{s'=1}^N p(s'|s, a) < 1, \quad \text{για όλα τα } a \in A(s), s \in S$$

Σύμφωνα με την υπόθεση αυτή, για κάθε ενέργεια  $a \in A(s)$  και για κάθε κατάσταση  $s$ , υπάρχει μια θετική πιθανότητα απορρόφησης :

$$p(0|s, a) \doteq 1 - \sum_{s'=1}^N p(s'|s, a) > 0$$

η οποία σηματοδοτεί τον τερματισμό της διαδικασίας (ή την απορρόφησή της από μία ‘τεχνητή’ κατάσταση  $s_T = 0$ )

Ανάλογα με την περίπτωση της διαδικασίας με προεξόφληση, για το διάνυσμα αξίας έχουμε ότι:

$$v_T(f) = \sum_{t=0}^{\infty} P^t(f) r(f) = [I - P(f)]^{-1} r(f)$$

### 1.2 Η Μαρκοβιανή Διαδικασία Απόφασης Πεπερασμένου Ορίζοντα

Στην προηγούμενη παράγραφο θεωρήσαμε πως ο χρονικός ορίζοντας των μοντέλων ήταν ,εν γένει, άπειρος. Στο σημείο αυτό θα θεωρήσουμε πεπερασμένο χρονικό ορίζοντα, έχοντας ότι

$$t \in \{0, 1, 2, \dots, T\}$$

Οι στρατηγικές που έχουμε ορίσει μέχρι στιγμής χαρακτηρίζονται ως *στάσιμες*, με την έννοια ότι δεν μας ενδιαφέρει η χρονική στιγμή  $t$  κατά την οποία επισκεπτόμαστε την κατάσταση  $s \in S$ . Θα θέλαμε να επεκταθούμε σε μια γενικότερη κλάση στρατηγικών. Ο λόγος για αυτή την γενίκευση είναι πως η επιλογή μιας ενέργειας ενδέχεται να είναι μη ευνοϊκή στα πρώτα στάδια της διαδικασίας, επειδή θα μπορούσε να οδηγήσει σε μία ‘κακή’ κατάσταση, ενώ θα μπορούσε να είναι ευνοϊκή σε μια μεταγενέστερη χρονικά στιγμή, εφόσον δεν θα υπάρχει αρκετός χρόνος προτού τερματίσει η διαδικασία ώστε να επισκεφτούμε αυτή την ‘κακή’ κατάσταση. Αυτό σημαίνει πως η αξία μιας ενέργειας είναι πλέον συνάρτηση του χρόνου που απομένει ως τον τερματισμό, και έτσι, εικάζουμε πως και μια βέλτιστη στρατηγική οφείλει και αυτή να εξαρτάται από τον υπολειπόμενο χρόνο έως τον τερματισμό της διαδικασίας.

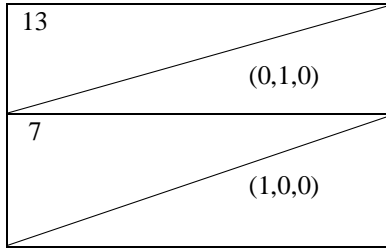


### Παράδειγμα 1.2.1

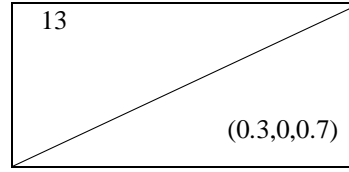
Έστω  $S = \{1,2,3\}$ ,  $A(1) = \{1,2\}$ ,  $A(2) = \{1\} = A(3)$ ,  $T = 2$  (δηλαδή  $t \in \{0,1,2\}$ )

Ενώ οι αμοιβές και οι πιθανότητες μετάβασης δίνονται παρακάτω:

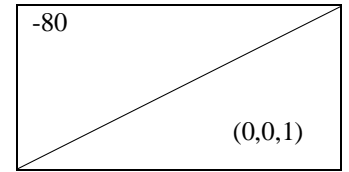
(Σε κάθε ορθογώνιο, πάνω αριστερά αναγράφονται οι αμοιβές ενός βήματος και μέσα στην παρένθεση αναγράφονται οι πιθανότητες μετάβασης σε όλες τις καταστάσεις)



Κατάσταση 1



Κατάσταση 2



Κατάσταση 3

Υποθέτουμε ότι ο ελεγκτής προσπαθεί να μεγιστοποιήσει το άθροισμα των αναμενόμενων αμοιβών στους χρόνους  $t = 0,1,2$ . Εφόσον ο ελεγκτής έχει περισσότερες από μία ενέργειες μόνο στην κατάσταση 1, κάθε στάσιμη στρατηγική θα είναι της μορφής :

$$f_p = ((p, 1-p), (1), (1)), \text{ με } p \in (0,1)$$

Στο συγκεκριμένο παράδειγμα, το σύνολο των καθαρών στάσιμων στρατηγικών θα είναι το  $F_D = \{f_0, f_1\}$ . Θέλουμε τώρα να υπολογίσουμε το άθροισμα των αμοιβών για τα στάδια 0,1,2, ξεκινώντας από την κατάσταση 1 και ακολουθώντας την στρατηγική  $f_p$ .

$$\begin{aligned} v_2(1, f_p) &= [13p + 7(1-p)] + \{(1-p)[13p + 7(1-p)] + 13p\} \\ &\quad + \{((1-p)^2 + 0.3p)(13p + 7(1-p)) + 13p(1-p) - 56p\} \\ &= 21 - 30.9p - 22.2p^2 + 6p^3 \end{aligned}$$

Η παραπάνω συνάρτηση του  $p$  λαμβάνει ακρότατο (μέγιστο) στο  $p = 0$ , με  $v_2(1, f_0) = 21$ . Συνεπώς, το καλύτερο που μπορεί να πετύχει ο ελεγκτής, στο σύνολο των στάσιμων στρατηγικών, είναι να επιλέξει την στρατηγική  $f_0$ .

Εάν, ωστόσο, θεωρήσουμε τη νέα στρατηγική  $\pi = (f_0, f_1, f_1)$ , που σημαίνει ότι θα χρησιμοποιήσουμε την  $f_0$  στο  $t = 0$  και την  $f_1$  για  $t = 1,2$ , τότε προκύπτει πως :

$$v_2(1, \pi) = 7 + 13 + 13 = 33$$

Μπορούμε να ερμηνεύσουμε το παραπάνω ως εξής:

Ο ελεγκτής μπορεί να έχει μεγαλύτερο κέρδος εάν επιλέξει την ενέργεια 1 στην κατάσταση 1 στον χρόνο  $t = 1$ , παρά στον χρόνο  $t = 0$ , εφόσον δεν θα χρειάζεται να ανησυχεί για το κόστος 80 μονάδων στην κατάσταση 3, αφού δεν θα υπάρχει αρκετός χρόνος ώστε η διαδικασία να φτάσει στην κατάσταση 3.

Το παραπάνω παράδειγμα μας οδηγεί να ορίσουμε μία στρατηγική ως μία πεπερασμένη ακολουθία

$$\pi = (f_0, f_1, f_2, \dots, f_T)$$

όπου όλα τα  $f_t$  είναι στάσιμες στρατηγικές ( $f_t \in F_s$ ). Θα συμβολίζουμε το σύνολο αυτών των στρατηγικών με  $F_M^T$  και θα τις καλούμε *Μαρκοβιανές στρατηγικές* της διαδικασίας απόφασης πεπερασμένου ορίζοντα.

Για κάθε  $\pi \in F_M^T$ , ορίζουμε:  $E_{s\pi}[R_t] \doteq E_\pi[R_t | S_0 = s]$

Η μέση τιμή είναι καλά ορισμένη για κάθε  $t = 0, 1, \dots, T$  και η αξία T-σταδίου για την  $\pi$  θα είναι

$$v_T(s, \pi) = \sum_{t=0}^T E_{s\pi}[R_t], \quad \text{για κάθε } s \in S$$

και το αντίστοιχο διάνυσμα για την  $\pi$ , για όλες τις καταστάσεις, θα είναι :

$$v_T(\pi) = (v_T(1, \pi), v_T(2, \pi), \dots, v_T(N, \pi))$$

Ένας από τους πιο διαδεδομένους τρόπους προκειμένου να υπολογίσουμε την αξία αλλά και την βέλτιστη στρατηγική της διαδικασίας, συνοψίζεται στον αλγόριθμο που θα ακολουθήσει, γνωστό και ως «Προς τα πίσω αναδρομή του Δυναμικού Προγραμματισμού». Η βασική ιδέα στην οποία στηρίζεται ο αλγόριθμος είναι η αρχή της βελτιστοποίησης, σύμφωνα με την οποία:

‘Εάν υποθέσουμε πως η μέγιστη αμοιβή είναι γνωστή όταν απομένουν  $(n - 1)$  στάδια για τον τερματισμό της διαδικασίας, τότε για να βρούμε τη μέγιστη αμοιβή ενώ απομένουν  $n$  στάδια για τον τερματισμό, αρκεί να μεγιστοποιήσουμε το άθροισμα της αμοιβής του τωρινού σταδίου με την μέγιστη αναμενόμενη πληρωμή για τα υπόλοιπα  $(n - 1)$  στάδια που απομένουν.’

Θα συμβολίζουμε με  $\underset{z \in Z}{\operatorname{argmax}} h(z)$  την τιμή του  $z$  για την οποία η πραγματική συνάρτηση  $h(z)$  λαμβάνει το μέγιστο στο σύνολο  $Z$ .

### Αλγόριθμος 1.2.1 – Προς τα πίσω αναδρομή του Δυναμικού Προγραμματισμού

#### Βήμα 1 (Αρχικοποίηση)

Θέτουμε  $V_{-1}(s) = 0$  για όλα τα  $s \in S$  και ορίζουμε:

$$f_T^*(s) \doteq a_s^T = \underset{a \in A(s)}{\operatorname{argmax}} \left\{ r(s, a) + \sum_{s'=1}^N p(s'|s, a) V_{-1}(s') \right\}$$

και

$$V_0(s) \doteq r(s, a_s^T) = \max_{a \in A(s)} \{r(s, a) + 0\}$$

### Βήμα 2 (Αναδρομή)

Για κάθε  $n = 1, 2, \dots, T$  υπολογίζουμε, για κάθε  $s \in S$  :

$$f_{T-n}^* = a_s^{T-n} = \operatorname{argmax}_{a \in A(s)} \left\{ r(s, a) + \sum_{s'=1}^N p(s'|s, a) V_{n-1}(s') \right\}$$

και

$$V_n(s) \doteq r(s, a_s^{T-n}) + \sum_{s'=1}^N p(s'|s, a_s^{T-n}) V_{n-1}(s')$$

### Βήμα 3

Κατασκευάζουμε την στρατηγική  $\boldsymbol{\pi}^* = (f_0^*, f_1^*, \dots, f_T^*) \in F_M^T$

### Σχόλια

- Η φράση ‘προς τα πίσω αναδρομή’ οφείλεται στα επαναληπτικά βήματα 2 και 3, όπου η  $\boldsymbol{\pi}^*$  κατασκευάζεται αναδρομικά από τις  $f_T^*, f_{T-1}^*, \dots, f_0^*$ .
- Για κάθε  $s \in S$ , για κάθε  $\boldsymbol{\pi} \in F_M^T$  και  $n = 0, 1, \dots, T$ , ορίζουμε:

$$V_n(s, \boldsymbol{\pi}) = \sum_{t=T-n}^T E_{\pi}[R_t | S_{T-n} = s]$$

οπότε το  $V_n(s, \boldsymbol{\pi})$  αναπαριστά την αναμενόμενη αμοιβή για τα τελευταία  $n$ -στάδια, δεδομένου ότι στο χρόνο  $T - n$  βρισκόμαστε στην  $s$ . Για  $n = T$ , λαμβάνουμε τη συνολική αναμενόμενη αμοιβή, δηλαδή :

$$V_T(s, \boldsymbol{\pi}) = v_T(s, \boldsymbol{\pi}), \quad \text{για όλα τα } s \in S, \boldsymbol{\pi} \in F_M^T$$

### Θεώρημα 1.2.1

Θεωρούμε την Μαρκοβιανή διαδικασία απόφασης  $T$ -ορίζοντα  $\Gamma_T$  και έστω  $\boldsymbol{\pi}^* \in F_M^T$  μια στρατηγική κατασκευασμένη όπως υποδεικνύεται από τον Αλγόριθμο 2.2.1. Τότε, η  $\boldsymbol{\pi}^*$  είναι βέλτιστη στρατηγική και, για όλα τα  $n = 0, 1, \dots, T, s \in S$  :

$$V_n(s) = \max_{a \in A(s)} \left\{ r(s, a) + \sum_{s'=1}^N p(s'|s, a) V_{n-1}(s') \right\} \quad (1.2.1)$$

### Απόδειξη

Η εξίσωση (1.2.1) είναι συνέπεια του ορισμού του  $V_n(s)$  από το Βήμα 2 του παραπάνω Αλγορίθμου. Μένει να δείξουμε ότι η  $\pi^*$  είναι βέλτιστη.

Για  $n = 0$  και μία αυθαίρετη στρατηγική  $\pi = (f_0, f_1, \dots, f_T) \in F_M^T$ , από το Βήμα 1 του αλγορίθμου, για κάθε  $s' \in S$ , λαμβάνουμε ότι:

$$V_0(s', \pi) = E_\pi[R_T | S_T = s'] \leq \max_{a \in A(s')} \{r(s', a) = E_{\pi^*}[R_T | S_T = s'] = V_0(s')\}$$

Επαγωγικά, υποθέτουμε ότι το ζητούμενο ισχύει για  $n - 1$ , δηλαδή ότι:

$$\begin{aligned} V_{n-1}(s', \pi) &= \sum_{t=T-n+1}^T E_\pi[R_t | S_{T-n+1} = s'] \leq \sum_{t=T-n+1}^T E_{\pi^*}[R_t | S_{T-n+1} = s'] = V_{n-1}(s', \pi^*) \\ &= V_{n-1}(s') \end{aligned} \quad (1.2.2)$$

για όλα τα  $s' \in S$ . Ακόμη, έχουμε ότι:

$$\begin{aligned} V_n(s'', \pi) &= E_\pi \left[ \sum_{t=T-n}^T (R_t | S_{T-n} = s'') \right] \\ &= \sum_{a \in A(s'')} E_\pi \left[ \sum_{t=T-n}^T (R_t | S_{T-n} = s'', A_{T-n} = a) \right] f_{T-n}(s'', a) \end{aligned} \quad (1.2.3)$$

Τώρα, θα γράψουμε τον όρο με τη μέση τιμή της (1.2.3) ως:

$$\begin{aligned} E_\pi \left[ \sum_{t=T-n}^T (R_t | S_{T-n} = s'', A_{T-n} = a) \right] &= \\ &= r(s'', a) + \sum_{s'=1}^N p(s' | s'', a) E_\pi \left[ \sum_{t=T-n+1}^T (R_t | S_{T-n} = s'', A_{T-n} = a, S_{T-n+1} = s') \right] \end{aligned}$$

και επειδή η κατανομή του  $\sum_{t=T-n+1}^T R_t$ , δεδομένου του γεγονότος  $\{S_{T-n+1} = s'\}$  και της  $\pi$ , είναι ανεξάρτητη από όσα έχουν προηγηθεί στη διαδικασία πριν από το  $(T - n + 1)$  - οστό στάδιο,

$$= r(s'', a) + \sum_{s'=1}^N p(s' | s'', a) E_\pi \left[ \sum_{t=T-n+1}^T (R_t | S_{T-n+1} = s') \right]$$

$$= r(s'', a) + \sum_{s'=1}^N p(s'|s'', a) V_{n-1}(s', \boldsymbol{\pi})$$

Αντικαθιστώντας την παραπάνω στην (1.2.3) , παίρνουμε:

$$V_n(s'', \boldsymbol{\pi}) = \sum_{a \in A(s'')} f_{T-n}(s'', a) \left( r(s'', a) + \sum_{s'=1}^N p(s'|s'', a) V_{n-1}(s', \boldsymbol{\pi}) \right)$$

Τώρα θα χρησιμοποιήσουμε την επαγωγική υπόθεση, οπότε:

$$\begin{aligned} V_n(s'', \boldsymbol{\pi}) &\leq \sum_{a \in A(s'')} f_{T-n}(s'', a) \left( r(s'', a) + \sum_{s'=1}^N p(s'|s'', a) V_{n-1}(s', \boldsymbol{\pi}) \right) \\ &\leq \sum_{a \in A(s'')} f_{T-n}(s'', a) \left\{ r(s'', a_{s''}^{T-n}) + \sum_{s'=1}^N p(s'|s'', a_{s''}^{T-n}) V_{n-1}(s') \right\} \end{aligned}$$

(λόγω του βήματος 2 του αλγορίθμου)

$$= \sum_{a \in A(s'')} f_{T-n}(s'', a) V_n(s'') = V_n(s'')$$

για όλα τα  $s'' \in S$  και  $\boldsymbol{\pi} \in F_M^T$  , το οποίο ολοκληρώνει την απόδειξη. ■

### Παράδειγμα 1.2.2

Σε αυτό το παράδειγμα θα υπολογίσουμε μια βέλτιστη στρατηγική για το Παράδειγμα 1.2.1. Θα δούμε ότι, χρησιμοποιώντας τον Αλγόριθμο 1.2.1, θα είμαστε σε θέση να επιβεβαιώσουμε ότι η στρατηγική  $\boldsymbol{\pi} = (\mathbf{f}_0, \mathbf{f}_1, \mathbf{f}_2)$  που χρησιμοποιήσαμε στο παράδειγμα, είναι πράγματι μια ‘καλή’ στρατηγική.

Αρχικοποιούμε τον αλγόριθμο με τις τιμές:

$$\mathbf{V}_0 = (V_0(1), V_0(2), V_0(3))^T = (13, 13, -80)^T, \quad f_2^*(1) = f_2^*(2) = f_2^*(3) = 1$$

Για  $n = 1$ ,

$$\begin{aligned} f_1^*(1) &= a_1^{2-1} = \operatorname{argmax} \left( r(1,1) + \sum_{s'} p(s'|1,1) V_0(s'), r(1,2) + \sum_{s'} p(s'|1,2) V_0(s') \right) \\ &= \operatorname{argmax} (13 + 13, 7 + 13) = \operatorname{argmax} (26) = 1 \end{aligned}$$

$$V_1(1) = r(1,1) + \sum_{s'} p(s'|1,1)V_0(s') = 13 + 13 = 26$$

$$\begin{aligned} f_1^*(2) &= a_2^{2-1} = \operatorname{argmax} \left( r(2,1) + \sum_{s'} p(s'|2,1)V_0(s') \right) \\ &= \operatorname{argmax} (13 + 13(0.3) - 80(0.7)) = \operatorname{argmax} (-39.1) = 1 \end{aligned}$$

$$V_1(2) = r(2,1) + \sum_{s'} p(s'|2,1)V_0(s') = 13 + 3.9 - 56 = -39.1$$

$$\begin{aligned} f_1^*(3) &= a_3^{2-1} = \operatorname{argmax} \left( r(3,1) + \sum_{s'} p(s'|3,1)V_0(s') \right) = \operatorname{argmax} (-80 + (-80)) \\ &= \operatorname{argmax} (-160) = 1 \end{aligned}$$

$$V_1(3) = r(3,1) + \sum_{s'} p(s'|3,1)V_0(s') = -160$$

Then  $n = 2$ ,

$$\begin{aligned} f_0^*(1) &= a_1^{2-2} = \operatorname{argmax} \left( r(1,1) + \sum_{s'} p(s'|1,1)V_1(s'), r(1,2) + \sum_{s'} p(s'|1,2)V_1(s') \right) \\ &= \operatorname{argmax} (13 + (-39.1), 7 + 26) = \operatorname{argmax} (33) = 2 \end{aligned}$$

$$V_2(1) = r(1,2) + \sum_{s'} p(s'|1,2)V_1(s') = 7 + 26 = 33$$

$$\begin{aligned} f_0^*(2) &= a_2^{2-2} = \operatorname{argmax} \left( r(2,1) + \sum_{s'} p(s'|2,1)V_1(s') \right) \\ &= \operatorname{argmax} (13 + 26(0.3) - 160(0.7)) = \operatorname{argmax} (-91.2) = 1 \end{aligned}$$

$$V_2(2) = r(2,1) + \sum_{s'} p(s'|2,1)V_1(s') = -91.2$$

$$\begin{aligned} f_0^*(3) &= a_3^{2-2} = \operatorname{argmax} \left( r(3,1) + \sum_{s'} p(s'|3,1)V_1(s') \right) = \operatorname{argmax} (-80 + (-160)) \\ &= \operatorname{argmax} (-240) = 1 \end{aligned}$$

$$V_2(3) = r(3,1) + \sum_{s'} p(s'|3,1)V_1(s') = -80 - 160 = -240$$

Έχουμε πλέον δείξει ότι  $f_0^*(1) = 2$ , που σημαίνει ότι  $f_0 = f_0$  και

$$f_1^*(1) = f_2^*(1) = 1, \text{ που σημαίνει ότι } f_1 = f_2 = f_1$$

Δηλαδή έχουμε επιβεβαιώσει ότι η  $\pi = (f_0, f_1, f_1)$  του παραδείγματος 1.2.1 είναι πράγματι μια βέλτιστη στρατηγική.

### Σχόλια

- Παρά την απλότητα στην έκφραση του Αλγορίθμου 1.2.1 και του αντίστοιχου θεωρήματος, ίσως μπορεί να διαισθανεί κανείς πως η χρήση τους ενδέχεται να επιφέρει υπολογιστικές δυσκολίες. Το αναδρομικό βήμα του αλγορίθμου εμπεριέχει  $T \times N$  μεγιστοποιήσεις, αριθμός που προφανώς είναι ανάλογος τόσο του χρονικού ορίζοντα  $T$ , όσο και του πλήθους των καταστάσεων  $N$ . Ακόμη, εάν έχουμε έναν ‘μεγάλο’ χώρο ενεργειών  $A(s)$ , τότε καθε μία από αυτές τις μεγιστοποιήσεις ενδέχεται να είναι αρκετά χρονοβόρα. Λόγω αυτών, ο παραπάνω αλγόριθμος και το σχετικό θεώρημα, φαίνεται να μην είναι ο καταλληλότερος τρόπος αντιμετώπισης προβλημάτων με πολλές καταστάσεις, πολλές ενέργειες ή/και μεγάλο χρονικό ορίζοντα.
- Ένα πρακτικό πρόβλημα που συναντάται στη μοντελοποίηση με πεπερασμένο ορίζοντα, απορρέει από το γεγονός πως, εν γένει, οι βέλτιστες στρατηγικές είναι μη-στάσιμες. Το βήμα 3 του αλγορίθμου υποδεικνύει πως μία βέλτιστη Μαρκοβιανή στρατηγική προϋποθέτει από τον decision-maker να θυμάται και την τρέχουσα κατάσταση της διαδικασίας, αλλά και το τρέχον στάδιο. Παρόλο που αυτό φαίνεται απολύτως λογικό από μαθηματικής σκοπιάς, σε πολλές πρακτικές εφαρμογές, αυτό το επίπεδο πολυπλοκότητας της βέλτιστης στρατηγικής ενδέχεται να αποθαρρύνει τον d-m από τη χρήση της, ο οποίος θα προτιμούσε μία αρκετά ‘απλούστερη’ βέλτιστη πολιτική.

## 1.3 Γραμμικός Προγραμματισμός και Αθροίσιμες Μαρκοβιανές Διαδικασίες Απόφασης

Στην παράγραφο 1.1 θεωρήσαμε τη διαδικασία  $\Gamma_\beta$  με συντελεστή προεξόφλησης  $\beta$ . Το πρόβλημα βελτιστοποίησης που είναι συνδεδεμένο με την διαδικασία αυτή, αφορά στη μεγιστοποίηση της τιμής, ενώ βρισκόμαστε στο σύνολο των στάσιμων στρατηγικών, δηλαδή:

$$\max \mathbf{v}_\beta(\mathbf{f})$$

με περιορισμούς

(I)

$$\mathbf{f} \in F_s$$

Μία στρατηγική  $\mathbf{f}_0 \in F_s$  για την οποία επιτυγχάνεται το παραπάνω μέγιστο, θα καλείται βέλτιστη και ,τότε, το διάνυσμα αξίας για την  $\mathbf{f}_0$  θα είναι:

$$\mathbf{v}_\beta \doteq \mathbf{v}_\beta(\mathbf{f}_0) = \max_{\mathbf{f}} \mathbf{v}_\beta(\mathbf{f})$$

Φυσικά ακόμη δεν γνωρίζουμε καν εάν υπάρχει τέτοια βέλτιστη  $\mathbf{f}_0$  και το αντίστοιχο διάνυσμα τιμής. Θα δούμε στη συνέχεια αυτής της παραγράφου όχι μόνο ότι υπάρχουν τέτοιες βέλτιστες στρατηγικές, αλλά και ότι αντιστοιχούν σε λύσεις κατάλληλα κατασκευασμένων γραμμικών προγραμμαμάτων.

Έχουμε δείξει τη σχέση

$$\mathbf{v}_\beta(\mathbf{f}) = [\mathbf{I} - \beta \mathbf{P}(\mathbf{f})]^{-1} \mathbf{r}(\mathbf{f}) \quad (1.3.1)$$

Λόγω της μη γραμμικότητας της παραπάνω σχέσης ως προς  $\mathbf{f}$ , δεν μπορούμε να γνωρίζουμε με βεβαιότητα την ύπαρξη του  $\max$  στο (I), πόσο μάλλον την υπολογιστική δυσκολία αυτού του εγχειρήματος.

Πολλαπλασιάζοντας την (1.3.1) επί  $\mathbf{I} - \beta \mathbf{P}(\mathbf{f})$ , παίρνουμε:

$$[\mathbf{I} - \beta \mathbf{P}(\mathbf{f})] \mathbf{v}_\beta(\mathbf{f}) = \mathbf{r}(\mathbf{f}) \Rightarrow \mathbf{v}_\beta(\mathbf{f}) = \mathbf{r}(\mathbf{f}) + \beta \mathbf{P}(\mathbf{f}) \mathbf{v}_\beta(\mathbf{f}) \quad (1.3.2)$$

Συνεπώς, εάν υπάρχουν τα  $\mathbf{f}_0 \in F_s$  και  $\mathbf{v}_\beta(\mathbf{f}_0)$ , θα πρέπει να ικανοποιούν την (1.3.2).

Ακόμη, από την αρχή της βελτιστοποίησης, θα πρέπει το  $\mathbf{v}_\beta$  να ικανοποιεί και την:

$$v_\beta(s) = \max_{a \in A(s)} \left\{ r(s, a) + \beta \sum_{s'=1}^N p(s'|s, a) v_\beta(s') \right\} \quad (1.3.3)$$

όπου με  $v_\beta(s)$  συμβολίζουμε το  $s$ -οστό στοιχείο του  $\mathbf{v}_\beta$ .

Η τελευταία ισότητα γεννά ένα σύνολο γραμμικών ανισοτήτων. Έστω  $\mathbf{v}$  ένα αυθαίρετο διάνυσμα που τις ικανοποιεί, δηλαδή έστω:

$\mathbf{v} = (v(1), v(2), \dots, v(N))^T$ , τέτοιο ώστε

$$v(s) \geq r(s, \mathbf{f}) + \beta \sum_{s'=1}^N p(s'|s, \mathbf{f}) v(s'), \quad \text{για κάθε } s \in S$$

ή, σε μορφή πινάκων:

$$\mathbf{v} \geq \mathbf{r}(\mathbf{f}) + \beta \mathbf{P}(\mathbf{f}) \mathbf{v}$$

Αντικαθιστώντας την παραπάνω ανισότητα στον εαυτό της  $k$ -φορές, παίρνουμε ότι:



$$\mathbf{v} \geq \mathbf{r}(\mathbf{f}) + \beta \mathbf{P}(\mathbf{f})\mathbf{r}(\mathbf{f}) + \dots + \beta^{k-1} \mathbf{P}^{k-1}(\mathbf{f})\mathbf{r}(\mathbf{f}) + \beta^k \mathbf{P}^k(\mathbf{f})\mathbf{v}$$

Αφήνουμε τώρα  $k \rightarrow \infty$ , και προκύπτει:

$$\mathbf{v} \geq [\mathbf{I} - \beta \mathbf{P}(\mathbf{f})]^{-1} \mathbf{r}(\mathbf{f}) = \mathbf{v}_\beta(\mathbf{f})$$

Δηλαδή, το αυθαίρετο  $\mathbf{v}$  που επιλέξαμε να ικανοποιεί τις παραπάνω ανισότητες, αποτελεί ένα άνω φράγμα για την τιμή με προεξόφληση, στο σύνολο των στάσιμων στρατηγικών  $F_S$ .

Ακριβώς αυτή η τελευταία παρατήρηση μας επιτρέπει να αρχίσουμε να υποψιαζόμαστε πως το διάνυμα τιμής  $\mathbf{v}_\beta$  της διαδικασίας  $\Gamma_\beta$  θα μπορούσε να αποτελεί την βέλτιστη λύση για το γραμμικό πρόβλημα:

$$\min \sum_{s=1}^N \frac{1}{N} v(s)$$

με περιορισμούς

$(P_\beta)$

$$v(s) \geq r(s, a) + \beta \sum_{s'=1}^N p(s'|s, a) v(s'), \quad a \in A(s), s \in S$$

*Σχόλιο*

Οι συντελεστές  $\frac{1}{N}$  στην αντικειμενική συνάρτηση του  $(P_\beta)$  ερμηνεύονται από το γεγονός πως θεωρούμε ισοπίθανο η αρχική κατάσταση να είναι οποιαδήποτε από τις  $\{1, 2, \dots, N\}$ .

Πριν προχωρήσουμε στη διατύπωση του δυϊκού προβλήματος που αντιστοιχεί στο  $(P_\beta)$ , ας δούμε πρώτα τον τρόπο με τον οποίο μπορούμε να κατασκευάσουμε ένα δυϊκό πρόβλημα γραμμικού προγραμματισμού, στηριζόμενοι στο αντίστοιχο πρωτεύον αυτού.

Αρχικά, θέλουμε το πρωτεύον πρόγραμμα να είναι γραμμένο σε κανονική μορφή. Δηλαδή, να είναι της μορφής:

$$\begin{aligned} \max \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0} \end{aligned} \tag{P}$$

όπου  $\mathbf{c} \in \mathbb{R}_{n \times 1}$ ,  $\mathbf{x} \in \mathbb{R}_{n \times 1}$ ,  $\mathbf{b} \in \mathbb{R}_{m \times 1}$ ,  $\mathbf{A} \in \mathbb{R}_{m \times n}$

Τότε, το δυϊκό πρόβλημα  $(D)$  που αντιστοιχεί στο  $(P)$  θα έχει τη μορφή:

$$\begin{aligned}
& \min \mathbf{b}^T \mathbf{w} \\
& \mathbf{A}^T \mathbf{w} \geq \mathbf{c} \\
& \mathbf{w} \geq \mathbf{0}
\end{aligned} \tag{D}$$

όπου  $\mathbf{w} \in \mathbb{R}_{m \times 1}$ .

Η διαδικασία την οποία ακολουθούμε για να κατασκευάσουμε τους περιορισμούς και τις μεταβλητές του (D) είναι η εξής:

- (i) Ορίζουμε μια δυϊκή μεταβλητή  $w_i$  για καθέναν από τους  $m$  περιορισμούς του πρωτεύοντος.
- (ii) Για κάθε μία από τις  $n$  μεταβλητές του πρωτεύοντος, ορίζεται ένας δυϊκός περιορισμός.
- (iii) Οι συντελεστές των μεταβλητών στους περιορισμούς του δυϊκού είναι ίσοι με τους συντελεστές της συνδεόμενης μεταβλητής του πρωτεύοντος (ανεστραμμένοι). Το δεξιό μέλος των περιορισμών του δυϊκού θα είναι ίσο με τους συντελεστές της αντικειμενικής συνάρτησης του πρωτεύοντος.
- (iv) Οι συντελεστές της αντικειμενικής συνάρτησης του δυϊκού είναι ίσοι με τα δεξιά μέλη των περιορισμών του πρωτεύοντος.

### Παράδειγμα 1.3.1

Έστω το παρακάτω πρόγραμμα γραμμικού προγραμματισμού:

$$\max z = 5x_1 + 6x_2 + 2x_3$$

με περιορισμούς:

$$4x_1 + 3x_2 + 2x_3 \leq 20$$

$$2x_1 + 5x_2 + x_3 \geq 8$$

$$0.5x_1 + 3.5x_2 + 2x_3 \leq 12$$

Για να υπολογίσουμε το δυϊκό του παραπάνω προβλήματος, πρέπει πρώτα να μετατρέψουμε τον 2<sup>ο</sup> περιορισμό, έτσι ώστε να συμφωνεί με την κανονική μορφή. Οπότε ο 2<sup>ος</sup> περιορισμός θα γίνει:

$$-2x_1 - 5x_2 - x_3 \leq -8$$

Οπότε το δυϊκό πρόβλημα θα έχει τη μορφή:

$$\min u = 20w_1 - 8w_2 + 12w_3$$

με περιορισμούς:

$$4w_1 - 2w_2 + 0.5w_3 \geq 5$$

$$3w_1 - 5w_2 + 3.5w_3 \geq 6$$

$$2w_1 - w_2 + 2w_3 \geq 2$$

Ωστόσο, πολλές φορές η διαδικασία υπολογισμού της κανονικής μορφής ενός γραμμικού προβλήματος είναι ιδιαίτερα χρονοβόρα. Για αυτές τις περιπτώσεις, κατασκευάζουμε το δυϊκό πρόγραμμα σύμφωνα με τον παρακάτω πίνακα:

Πρωτεύον		Δυϊκό
maximize $z$		minimize $u$
$i$ -περιορισμός μορφής $=$		$i$ -μεταβλητή $w_i \geq 0$
$i$ -περιορισμός μορφής $\leq$		$i$ -μεταβλητή $w_i \in \mathbb{R}$
$i$ -περιορισμός μορφής $\geq$	$\Leftrightarrow$	$i$ -μεταβλητή $w_i \leq 0$
$i$ -μεταβλητή $x_i \geq 0$		$i$ -περιορισμός μορφής $\geq$
$i$ -μεταβλητή $x_i \in \mathbb{R}$		$i$ -περιορισμός μορφής $=$
$i$ -μεταβλητή $x_i \leq 0$		$i$ -περιορισμός μορφής $\leq$

Δηλαδή, ο παραπάνω πίνακας ερμηνεύεται ως εξής:

- Αν μια μεταβλητή του ενός δεν έχει περιορισμό στο πρόσημο, ο αντίστοιχος περιορισμός του άλλου είναι εξίσωση και αντίστροφα.
- Αν μια μεταβλητή του ενός είναι μη θετική, τότε ο αντίστοιχος περιορισμός του άλλου είναι ανίσωση με φορά αντίθετη της αναμενόμενης και αντίστροφα.

Τα σύνολα  $\{\mathbf{x} \in \mathbb{R}^n | \mathbf{Ax} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$  και  $\{\mathbf{w} \in \mathbb{R}^m | \mathbf{w}^T \mathbf{A} \geq \mathbf{c}^T, \mathbf{w} \geq \mathbf{0}\}$  ονομάζονται *εφικτά σύνολα* και, όταν αυτά είναι μη-κενά, το αντίστοιχο πρόβλημα θα καλείται *εφικτό*.

### Ορισμός 1.3.1.

Ένα στοιχείο  $\mathbf{x}$  του εφικτού συνόλου  $\{\mathbf{x} \in \mathbb{R}^n | \mathbf{Ax} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$  θα καλείται *εφικτή λύση*. Μία εφικτή λύση  $\mathbf{x}$  θα καλείται *βασική εφικτή λύση* οποτεδήποτε είναι δυνατό να επιλέξουμε έναν  $r \times r$  υποπίνακα  $A_r$  του  $A$  τέτοιο ώστε:

- (i) Ο  $A_r$  είναι μη-ιδιάζων
- (ii)  $A_r x_r = b_r$
- (iii)  $x_j = 0$  για όλα τα  $j$  πέρα των  $r$  επιλεγμένων στηλών

Παρατηρούμε ότι, όποτε υπάρχει μία βασική εφικτή λύση, τότε αυτή καθορίζεται πλήρως από τα  $A_r$  και  $b_r$ . Η επόμενη πρόταση φανερώνει την σημασία της ύπαρξης τέτοιων βασικών εφικτών λύσεων.

### Πρόταση 1.3.1.

- (i) Εάν το εφικτό σύνολο είναι μη-κενό, τότε υπάρχει βασική εφικτή λύση
- (ii) Εάν υπάρχει βέλτιστη εφικτή λύση, τότε υπάρχει βέλτιστη βασική εφικτή λύση

Η σημασία του (ii) της παραπάνω πρότασης έγκειται στο γεγονός πως αναζητώντας μια βέλτιστη λύση, θα μπορούσαμε να περιοριστούμε στο πεπερασμένο σύνολο των βασικών εφικτών λύσεων και έπειτα, μέσω της γνωστής μεθόδου *Simplex* να φτάσουμε στο ζητούμενο αποτέλεσμα.

Η γεωμετρική ερμηνεία του εφικτού συνόλου είναι η ακόλουθη. Κάθε γραμμικός περιορισμός της μορφής:

$$\sum_{j=1}^n a_{ij}x_j \leq b_i, \quad i = 1, 2, \dots, m$$

ορίζει ένα κλειστό ‘ημι-επίπεδο’ στον  $\mathbb{R}^n$ . Τότε, μπορούμε να δούμε το εφικτό σύνολο ως την τομή  $m$  στο πλήθος τέτοιων ‘ημι-επιπέδων’, το οποίο εάν είναι μη-κενό, θα καλείται *πολύεδρο*. Τα ακραία σημεία του πολυέδρου αυτού βρίσκονται σε 1-1 αντιστοιχία με τις βασικές εφικτές λύσεις.

### Πρόταση 1.3.2. (Δυϊκό Θεώρημα)

- (i) Εάν είτε το Πρωτεύον είτε το Δυϊκό πρόβλημα έχει πεπερασμένη βέλτιστη λύση, τότε το ίδιο ισχύει και για το άλλο πρόβλημα, και οι αντίστοιχες τιμές της αντικειμενικής συνάρτησης είναι ίσες.
- (ii) Εάν είτε το Πρωτεύον είτε το Δυϊκό πρόβλημα είναι μη φραγμένο, τότε το άλλο πρόβλημα δεν έχει εφικτές λύσεις.

*Πρόταση 1.3.3. (Συμπληρωματική Χαλαρότητα)*

Μια αναγκαία και ικανή συνθήκη ώστε ένα ζεύγος εφικτών λύσεων  $(\mathbf{w}, \mathbf{x})$  να είναι βέλτιστο είναι, για όλα τα  $i, j$  :

$$(a) \quad w_i > 0 \Rightarrow \sum_{i=1}^m a_{ij} w_i = c_j$$

$$(b) \quad \sum_{j=1}^n a_{ij} x_j < b_i \Rightarrow w_i = 0$$

$$(c) \quad x_j > 0 \Rightarrow \sum_{i=1}^m a_{ij} w_i = c_j$$

$$(d) \quad \sum_{i=1}^m a_{ij} w_i > c_j \Rightarrow x_j = 0$$

Είμαστε πλέον σε θέση να συνεχίσουμε με την ανάλυση του  $(P_\beta)$ . Εάν θεωρήσουμε το  $(P_\beta)$  ως ένα πρωτεύον (Primal) πρόβλημα γραμμικού προγραμματισμού, και συσχετίσουμε με κάθε περιορισμό τη δυϊκή μεταβλητή  $x_{sa}$ , τότε το δυϊκό (Dual) πρόγραμμα  $(D_\beta)$  του  $(P_\beta)$  θα έχει την ακόλουθη μορφή :

$$\max \sum_{s=1}^N \sum_{a=1}^{m(s)} r(s, a) x_{sa}$$

με περιορισμούς

$(D_\beta)$

$$\sum_{s=1}^N \sum_{a=1}^{m(s)} [\delta_{ss'} - \beta p(s'|s, a)] x_{sa} = \frac{1}{N}, s' \in S$$

$$x_{sa} \geq 0, a \in A(s), s \in S$$

*Θεώρημα 1.3.1*

(i) Τα γραμμικά προγράμματα  $(P_\beta)$  και  $(D_\beta)$  έχουν πεπερασμένες βέλτιστες λύσεις

(ii) Έστω  $\mathbf{v}^0 = (v^0(1), v^0(2), \dots, v^0(N))^T$  μία βέλτιστη λύση του  $(P_\beta)$ . Τότε  $\mathbf{v}^0 = \mathbf{v}_\beta$ , όπου  $\mathbf{v}_\beta$  είναι το διάνυσμα-τιμής της διαδικασίας  $\Gamma_\beta$ .

(iii) Έστω  $\mathbf{x}^0 = \{x_{sa}^0 | a \in A(s), s \in S\}$  μία βέλτιστη λύση του  $(D_\beta)$  και έστω

$$x_0^s \doteq \sum_{a=1}^{m(s)} x_{sa}^0, \text{ για κάθε } s \in S$$

Τότε  $x_0^s > 0$  και η  $f^0 \in F_s$  που ορίζεται ως:

$$f^0(s, a) \doteq \frac{x_{sa}^0}{x_s^0}, \text{ για όλα τα } a \in A(s), s \in S$$

είναι μία βέλτιστη στάσιμη στρατηγική για την διαδικασία  $\Gamma_\beta$ .

### Πόρισμα 1.3.1

(i) Το διάνυσμα τιμής  $v_\beta$  είναι η μοναδική λύση της εξίσωσης βελτιστοποίησης (1.3.3).

(ii) Για κάθε  $s \in S$ , έστω μια ενέργεια  $a_s \in A(s)$  για την οποία επιτυγχάνεται το  $\max$  στην (1.3.3), δηλαδή:

$$v(s) = r(s, a_s) + \beta \sum_{s'=1}^N p(s'|s, a_s) v(s') = \max_{a \in A(s)} \left\{ r(s, a) + \beta \sum_{s'=1}^N p(s'|s, a) v(s') \right\}$$

όπου  $v = (v(1), v(2), \dots, v(N))$  είναι η λύση της (1.3.3). Ορίζουμε την  $f^* \in F_s$  ως εξής:

$$f^*(s, a) = \begin{cases} 1, & \text{αν } a = a_s \\ 0, & \text{αλλιώς} \end{cases}, \text{ για όλα τα } s \in S$$

Τότε η  $f^*$  είναι μία βέλτιστη στάσιμη καθαρή (ντετερμινιστική) λύση.

### Παρατηρήσεις

- Η εγκυρότητα της εξίσωσης βελτιστοποίησης (1.3.3) υποδεικνύει πως το πρόβλημα εύρεσης βέλτιστης στρατηγικής είναι πλήρως αντιμετωπίσιμο από τη στιγμή που γίνει γνωστό το διάνυσμα  $v_\beta$ , εφόσον το μόνο που απαιτείται είναι ο υπολογισμός  $N$  στο πλήθος  $\max$ ima, όπως φαίνεται από το (ii) του Πορίσματος
- Η ανάλυση που προηγήθηκε για το μοντέλο της διαδικασίας  $\Gamma_\beta$ , ισχύει και στην περίπτωση μοντελοποίησης της τερματιζόμενης διαδικασίας  $\Gamma_T$ . Η ειδοποιός διαφορά είναι πως στην ανάλυση της  $\Gamma_T$  θα παραλείπεται ο παράγοντας προεξόφλησης  $\beta$  από τις σχετικές εξισώσεις.

Για παράδειγμα, το αντίστοιχο πρωτεύον γραμμικό πρόγραμμα του  $(P_\beta)$  θα είναι το:

$$\min \sum_{i=1}^N \frac{1}{N} v(s)$$

με περιορισμούς

$$v(s) = r(s, a) + \sum_{s'=1}^N p(s'|s, a) v(s'), \quad a \in A(s), s \in S \quad (P_T)$$

#### 1.4 Η Αδιαχώριστη διαδικασία Οριακού Μέσου

Στην παράγραφο αυτή θα χρησιμοποιήσουμε ένα νέο κριτήριο απόδοσης το οποίο είναι διαφορετικό τόσο από το κριτήριο με προεξόφληση όσο και από εκείνο της τερματιζόμενης διαδικασίας. Θα θέλαμε να εστιάσουμε στη μέση απόδοση μιας στρατηγικής μακροπρόθεσμα (in the long run).

Συγκεκριμένα, ορίζουμε την *τιμή οριακού μέσου* της στάσιμης στρατηγικής  $\mathbf{f}$ , ξεκινώντας από μια αρχική κατάσταση  $s$ , ως εξής:

$$v_a(s, \mathbf{f}) = \lim_{T \rightarrow \infty} \left( \left( \frac{1}{T+1} \right) \sum_{t=0}^T E_{sf} [R_t] \right) \quad (1.3.4)$$

Το μοντέλο που χρησιμοποιεί το παραπάνω ως κριτήριο απόδοσης θα το ονομάσουμε *Μαρκοβιανή Διαδικασία Απόφασης Οριακού Μέσου* και θα το συμβολίζουμε με  $\Gamma_a$ .

Κατ' αναλογία με τις προηγούμενες παραγράφους, ορίζουμε το διάνυσμα τιμής:

$$\mathbf{v}_a(\mathbf{f}) \doteq ((v_a(1, \mathbf{f}), v_a(2, \mathbf{f}), \dots, v_a(N, \mathbf{f})))^T$$

Μπορούμε τώρα να συνδέσουμε το  $\Gamma_a$  με το πρόβλημα βελτιστοποίησης:

$$\max \mathbf{v}_a(\mathbf{f})$$

με περιορισμούς:

$$\mathbf{f} \in F_s$$

Η σχέση (1.3.4) μπορεί να γραφεί σε διανυσματική μορφή:

$$\mathbf{v}_a(\mathbf{f}) = \lim_{T \rightarrow \infty} \left( \frac{1}{T+1} \sum_{t=0}^T \mathbf{P}^t(\mathbf{f}) \mathbf{r}(\mathbf{f}) \right) = \left( \lim_{T \rightarrow \infty} \frac{1}{T+1} \sum_{t=0}^T \mathbf{P}^t(\mathbf{f}) \right) \mathbf{r}(\mathbf{f}) = \mathbf{Q}(\mathbf{f}) \mathbf{r}(\mathbf{f})$$

Η τελευταία ισότητα προκύπτει από γνωστή ιδιότητα των Μαρκοβιανών Αλυσίδων, η οποία εξασφαλίζει την ύπαρξη ενός πίνακα  $\mathbf{Q}(\mathbf{f})$  τέτοιου ώστε:

$$\mathbf{Q}(\mathbf{f}) \doteq \lim_{T \rightarrow \infty} \frac{1}{T+1} \sum_{t=0}^T \mathbf{P}^t(\mathbf{f})$$

Ο πίνακας  $\mathbf{Q}(\mathbf{f})$  ονομάζεται *όριο κατά Cesaro* πίνακας του  $\mathbf{P}(\mathbf{f})$ . Αυτός ο πίνακας θα μας απασχολήσει ιδιαίτερα στο 3<sup>ο</sup> Κεφάλαιο.

Οπότε για  $\mathbf{f} \in F_s$  έχουμε ότι:

$$v_a(\mathbf{f}) = \mathbf{Q}(\mathbf{f})\mathbf{r}(\mathbf{f})$$

Ωστόσο, ένα πρόβλημα που θα μπορούσε να δημιουργηθεί σε αυτή την διαδικασία, αποτελεί η ύπαρξη απορροφητικών καταστάσεων. Είναι πιθανόν, ανάλογα με την επιλογή της στρατηγικής, η διαδικασία να εγκλωβιστεί σε μια τέτοια κατάσταση. Προς αποφυγή αυτού, αρχικά θα περιοριστούμε σε *αδιαχώριστες* Μαρκοβιανές Αλυσίδες, υποθέτοντας ότι για κάθε  $\mathbf{f} \in F_s$ , ο πίνακας πιθανοτήτων μετάβασης  $\mathbf{P}(\mathbf{f})$  ορίζει μια *αδιαχώριστη* Μαρκοβιανή Αλυσίδα.

Από μαθηματική σκοπιά, αυτό σημαίνει ότι για οποιεσδήποτε δύο καταστάσεις  $(s, s')$ , υπάρχει θετικός ακέραιος  $t$  τέτοιος ώστε  $[\mathbf{P}^t(\mathbf{f})]_{s,s'} > 0$ . Προφανώς το  $t$  εξαρτάται από την επιλογή των  $s, s'$ . Ουσιαστικά, μια αδιαχώριστη αλυσίδα συνεπάγεται για μία στρατηγική  $\mathbf{f}$ , ότι η διαδικασία θα επισκεφθεί κάθε κατάσταση του συστήματος άπειρες το πλήθος φορές. Στη συνέχεια παραθέτουμε ένα βασικό λήμμα που αφορά σε αυτή την κατηγορία μαρκοβιανών αλυσίδων.

#### Λήμμα 1.4.1.

Έστω  $\mathbf{P}$  ο πίνακας πιθανοτήτων μετάβασης μιας αδιαχώριστης ΜΑ και  $\mathbf{Q}$  το αντίστοιχο *όριο κατά Cesaro* του  $\mathbf{P}$ . Τότε

- (i) Οι γραμμές του  $\mathbf{Q}$  ταυτίζονται
- (ii) Έστω  $\mathbf{q} = (q_1, q_2, \dots, q_N)$  μία γραμμή του  $\mathbf{Q}$ . Τότε κάθε στοιχείο του  $\mathbf{q}$  θα είναι αυστηρά θετικό και  $\mathbf{q}$  είναι η μοναδική λύση του συστήματος γραμμικών εξισώσεων:

$$\mathbf{x} \mathbf{P} = \mathbf{x}$$

$$\mathbf{x} \mathbf{1} = 1$$

όπου με  $\mathbf{1}$  συμβολίζουμε το  $N$ - διάστατο διάνυσμα στήλης που έχει μονάδα σε κάθε στοιχείο του. Το διάνυσμα  $\mathbf{q}$  καλείται *στάσιμη κατανομή* της αδιαχώριστης ΜΑ.



Ας θεωρήσουμε μια αδιαχώριστη ΜΑ που προσδιορίζεται από κάποια  $\mathbf{f} \in F_s$  μέσω του πίνακα μεταβάσεων  $\mathbf{P}(\mathbf{f})$  και έστω  $\mathbf{q}(\mathbf{f})$  η στάσιμη κατανομή αυτής, όπως στο παραπάνω Λήμμα. Για κάθε  $a \in A(s), s \in S$  ορίζουμε:

$$x_{sa}(\mathbf{f}) \doteq q_s(\mathbf{f})f(s, a) \quad (1.4.1)$$

$$x_s(\mathbf{f}) \doteq \sum_{a \in A(s)} x_{sa}(\mathbf{f}) = q_s(\mathbf{f}) \quad \left( \text{αφού} \sum_{a \in A(s)} f(s, a) = 1 \right)$$

Ερμηνεύουμε το  $x_s(\mathbf{f})$  ως τη συχνότητα επίσκεψης της κατάστασης  $s$  μακροπρόθεσμα και το  $x_{sa}(\mathbf{f})$  ως τη συχνότητα εμφάνισης του ζεύγους  $(s, a)$  μακροπρόθεσμα. Ακόμη, θα συμβολίζουμε με  $\mathbf{x}(\mathbf{f})$  το διάνυσμα συχνότητας εμφάνισης του ζεύγους (κατάσταση-ενέργεια) μακροπρόθεσμα, του οποίου το  $s$ -οστό μπλοκ θα είναι το διάνυσμα στήλη:

$$\mathbf{x}_s(\mathbf{f}) = (x_{s1}(\mathbf{f}), x_{s2}(\mathbf{f}), \dots, x_{sm(s)}(\mathbf{f}))^T$$

Τέλος, ορίζουμε το διάνυσμα-γραμμή συχνότητας εμφάνισης κάθε κατάστασης μακροπρόθεσμα:

$$\bar{\mathbf{x}}(\mathbf{f}) = (x_1(\mathbf{f}), x_2(\mathbf{f}), \dots, x_N(\mathbf{f}))$$

Από τη σχέση (1.4.1) μπορούμε να ορίσουμε μια απεικόνιση

$$\begin{cases} M: F_s \rightarrow \mathbb{R}^m \\ M(\mathbf{f}) \rightarrow \mathbf{x}(\mathbf{f}) \end{cases}, \text{ όπου } m \doteq \sum_{s=1}^N m(s)$$

Από το Λήμμα 1.4.1 έχουμε ότι:

$$\mathbf{q}(\mathbf{f})\mathbf{P}(\mathbf{f}) = \mathbf{q}(\mathbf{f}) \Rightarrow \mathbf{q}(\mathbf{f})(\mathbf{I} - \mathbf{P}(\mathbf{f})) = \mathbf{0}$$

Η τελευταία σχέση μπορεί να γραφεί σε μη-διανυσματική μορφή ως

$$\sum_{s=1}^N (\delta_{ss'} - p(s'|s, \mathbf{f}))q_s(\mathbf{f}) = 0, \quad s, s' \in S$$

Αναλύοντας το  $p(s'|s, \mathbf{f})$  στις επιμέρους ενέργειες της  $\mathbf{f}$ , παίρνουμε:

$$\sum_{s=1}^N (\delta_{ss'} - p(s'|s, a))q_s(\mathbf{f})f(s, a) = 0 \Rightarrow \sum_{s=1}^N \sum_{a \in A(s)} (\delta_{ss'} - p(s'|s, a))x_{sa}(\mathbf{f}) = 0$$

Επιπλέον, έχουμε ότι

$$\sum_{s=1}^N \sum_{a \in A(s)} x_{sa}(\mathbf{f}) = \sum_{s=1}^N \sum_{a \in A(s)} q_s(\mathbf{f})f(s, a) = \sum_{s=1}^N q_s(\mathbf{f}) = 1$$

Οπότε τελικά καταλήγουμε στο πολύεδρο σύνολο  $X$  που ορίζεται από τους παρακάτω περιορισμούς:

$$\sum_{s=1}^N \sum_{a \in A(s)} (\delta_{ss'} - p(s'|s, a)) x_{sa} = 0, \quad s' \in S$$

$$\sum_{s=1}^N \sum_{a \in A(s)} x_{sa} = 1$$

$$x_{sa} \geq 0, \quad a \in A(s), \quad s \in S$$

Εναλλακτικά, σε μορφή πινάκων:

$$X = \{x | Wx = 0, 1^T x = 1, x \geq 0\}$$

Όπου τα  $x, 1$ , είναι  $m$ -διάστατα διανύσματα, και  $W$  είναι ένας  $N \times m$  πίνακας του οποίου το  $(s', (s, a))$ -οστό στοιχείο είναι:

$$w_{s'}(s, a) \doteq \delta_{ss'} - p(s'|s, a)$$

Στο σημείο αυτό θα θέλαμε να δείξουμε ότι το σύνολο  $X$  αντιστοιχεί στο ‘χώρο συχνοτήτων’ του  $\Gamma_\alpha$ , και στη συνέχεια ότι ο χώρος αυτός βρίσκεται σε 1-1 αντιστοιχία με το χώρο των στάσιμων στρατηγικών  $F_S$ . Παραθέτουμε, χωρίς απόδειξη, ορισμένα ‘τεχνικά’ αποτελέσματα.

#### Λήμμα 1.4.2

Έστω  $\Gamma_\alpha$  μια αδιαχώριστη διαδικασία οριακού μέσου και  $X$  το αντίστοιχο πολύεδρο όπως ορίζεται από τα τρία παραπάνω σύνολα περιορισμών. Έστω  $x \in X$  και ορίζουμε το  $\bar{x} = (x_1, x_2, \dots, x_N)$  ως εξής:

$$x_s = \sum_{a \in A(s)} x_{sa}, \quad s \in S$$

Τότε θα ισχύει ότι  $\bar{x} > 0$  (δηλαδή  $x_s > 0$  για κάθε  $s \in S$ )

#### Θεώρημα 1.4.1

Έστω  $\Gamma_\alpha$  αδιαχώριστο μοντέλο, το σύνολο  $X$  και η συνάρτηση  $M: F_S \rightarrow \mathbb{R}^m$ , όπως ακριβώς έχουν οριστεί. Τότε η  $M$  είναι αντιστρέψιμη απεικόνιση από το  $F_S$  στο  $X$  και η αντίστροφη της είναι η:

$$M^{-1}(x) = f_x, \quad \text{όπου}$$

$$f_x(s, a) \doteq \frac{x_{sa}}{x_s}, \quad a \in A(s), s \in S$$

Το τελευταίο θεώρημα υποδεικνύει πως μπορούμε να μετασχηματίσουμε το πρόβλημα

$$(I) \left\{ \max_{f \in F_s} v_a(f) \right\}, \text{ που ορίστηκε στην αρχή της παραγράφου, στο πρόβλημα}$$

$$\left\{ \max_{x \in X} v_a(M^{-1}(x)) \right\} \text{ πάνω στο χώρο συχνοτήτων } X.$$

Το επόμενο αποτέλεσμα υποδεικνύει τον τρόπο εύρεσης βέλτιστης λύσης για το (I) από μία βέλτιστη λύση ενός γραμμικού προγράμματος.

#### Θεώρημα 1.4.2

Έστω  $\Gamma_\alpha$  το αδιαχώριστο μοντέλο, το σύνολο  $X$  και η συνάρτηση  $M: F_s \rightarrow \mathbb{R}^m$  όπως στο Θεώρημα 1.4.1. Έστω  $x_0$  μια βέλτιστη λύση στο γραμμικό πρόγραμμα

$$\max \sum_{s=1}^N \sum_{a \in A(s)} r(s, a) x_{sa}$$

με περιορισμούς

$$Wx = 0$$

$$\mathbf{1}^T x = 0$$

$$x \geq 0$$

Τότε, η  $f^0 \doteq f_{x_0} = M^{-1}(x_0)$  αποτελεί μια βέλτιστη στρατηγική για το αρχικό πρόβλημα βελτιστοποίησης (I).

### 1.5 Συμπεριφορικές και Μαρκοβιανές στρατηγικές

Όπως μπορεί να αντιληφθεί κανείς, στα προβλήματα που έχουμε έως τώρα δει, δεν είναι αρκετό να περιοριστούμε μόνο στη χρήση καθαρών στρατηγικών, καθ' ότι κρίνεται αναγκαία η ύπαρξη στοχαστικότητας στην πολιτική που ακολουθείται, προκειμένου να φτάσουμε στη βέλτιστη λύση. Σε ακόμη γενικότερο πλαίσιο, θα μπορούσε κανείς να αναρωτηθεί εάν και το σύνολο  $F_s$  των στάσιμων στρατηγικών είναι επαρκές ή αρκετά μεγάλο, όσο θα θέλαμε.

Αντικείμενο αυτής της παραγράφου αποτελεί η εισαγωγή δύο νέων συνόλων από στρατηγικές, ευρύτερων από το  $F_S$ , καθώς και η σύγκριση της απόδοσης του μοντέλου για στρατηγικές που ανήκουν στα δύο αυτά νέα σύνολα.

Θα συμβολίζουμε με  $h_t = (s_0, a_0, s_1, a_1, \dots, a_{t-1}, s_t)$  την ιστορία της διαδικασίας μέχρι το χρόνο  $t$ , και με  $H_t$  το σύνολο όλων των πιθανών ιστοριών μέχρι το χρόνο  $t$ . Έστω  $A = \bigcup_{s \in S} A(s)$  ο χώρος όλων των ενεργειών της διαδικασίας και  $P(A)$  το σύνολο όλων των κατανομών πιθανότητας πάνω στο (πεπερασμένο) σύνολο  $A$ .

Ορίζουμε ως κανόνα απόφασης στο χρόνο  $t$ , τη συνάρτηση:

$f_t: H_t \rightarrow P(A)$  τέτοια ώστε

$$f_t(h_t, a) = \begin{cases} \mathbb{P}_{f_t}(A_t = a | h_t), & \text{αν } a \in A(s_t) \\ 0, & \text{αν } a \notin A(s_t) \end{cases}$$

Τότε, ορίζουμε ως συμπεριφορική στρατηγική  $\pi$  να είναι η ακολουθία από κανόνες απόφασης:

$$\pi \doteq (f_0, f_1, f_2, \dots)$$

Θα συμβολίζουμε με  $F_B$  την κλάση όλων των συμπεριφορικών στρατηγικών. Μια συμπεριφορική στρατηγική  $\pi$  θα καλείται *Μαρκοβιανή* εάν κάθε κανόνας απόφασης  $f_t$  εξαρτάται μόνο από την παρούσα κατάσταση, για κάθε  $t = 0, 1, 2, \dots$ , δηλαδή:

$$\begin{aligned} f_t(s_t, a) &= \mathbb{P}_{f_t}(A_t = a | S_t = s_t) = \mathbb{P}_{f_t}(A_t = a | S_0 = s_0, A_0 = a_0, \dots, A_{t-1} = a_{t-1}, S_t = s_t) \\ &= f_t(h_t, a) \end{aligned}$$

για όλες τις ιστορίες  $h_t = (s_0, a_0, s_1, a_1, \dots, a_{t-1}, s_t) \in H_t$

Η κλάση των *Μαρκοβιανών* στρατηγικών συμβολίζεται με  $F_M$ .

Η κλάση  $F_M$  είναι μεγαλύτερη από την κλάση  $F_S$  των στάσιμων στρατηγικών. Πράγματι, μια στάσιμη στρατηγική είναι μια Μαρκοβιανή στρατηγική για την οποία όλοι οι κανόνες απόφασης είναι ανεξάρτητοι από το χρόνο  $t$ , δηλαδή  $f_t = f$  για κάθε  $t$ , όπου:

$$f(s, a) = \mathbb{P}_f(A_t = a | S_t = s), \text{ για κάθε } t = 0, 1, 2, \dots$$

Υπενθυμίζουμε ότι μία καθαρή στάσιμη στρατηγική  $f$  είναι μία στάσιμη στρατηγική για την οποία, για κάθε  $s \in S$  υπάρχει  $a_s \in A(s)$  τέτοιο ώστε  $f(s, a) = 0$  για κάθε  $a \neq a_s$ . Το σύνολο των καθαρών στάσιμων στρατηγικών θα συμβολίζεται με  $F_D$ . Φυσικά αυτό το σύνολο είναι ένα πεπερασμένο σύνολο, με πληθικό αριθμό:

$$|F_D| = \prod_{s=1}^N m(s)$$

Από τον τρόπο που κατασκευάστηκαν όλες οι παραπάνω κλάσεις, γίνεται φανερό η σχέση εγκλεισμού που τις συνδέει. Ισχύει δηλαδή ότι:

$$F_B \supseteq F_M \supseteq F_S \supseteq F_D$$

### Σχόλιο

Παρατηρούμε ότι, ανεξάρτητα από την κλάση στρατηγικών, το γεγονός ότι ο χώρος καταστάσεων και ενεργειών είναι πεπερασμένα σύνολα, μας εξασφαλίζει ότι οι αναμενόμενες αμοιβές για κάθε  $t$  είναι καλά ορισμένες και ικανοποιούν την :

$$E_{s_0\pi}[R_t] = \sum_{s=1}^N \sum_{a \in A(s)} \mathbb{P}_{\pi}(S_t = s, A_t = a | S_0 = s_0) r(s, a) \quad , \pi \in F_B, s_0 \in S$$

Ο ορισμός για την τιμή μιας στρατηγικής με συντελεστή προεξόφλησης, καθώς και για την τιμή του τερματιζόμενου μοντέλου, επεκτείνονται με φυσικό τρόπο στο σύνολο στρατηγικών  $F_B$ .

Για την τιμή στο μοντέλο οριακού μέσου, ωστόσο, θα χρειαστεί να τροποποιήσουμε ελαφρώς τον ορισμό, κατά τον εξής τρόπο:

Για κάθε  $\pi \in F_B$  και αρχική κατάσταση  $s$ , ορίζουμε:

$$v_a(s, \pi) \doteq \liminf_{T \rightarrow \infty} \left( \frac{1}{T+1} \sum_{t=0}^T E_{s\pi}[R_t] \right)$$

Το επόμενο Θεώρημα μας λέει ότι μπορούμε να ‘προσαρμόσουμε’ σε κάθε συμπεριφορική στρατηγική μία Μαρκοβιανή στρατηγική

### Θεώρημα 1.5.1

Έστω  $\pi \in F_B$  μια αυθαίρετη συμπεριφορική στρατηγική. Τότε, για κάθε αρχική κατάσταση  $s_0 \in S$ , υπάρχει Μαρκοβιανή στρατηγική  $\bar{\pi} \in F_M$  τέτοια ώστε, για  $a \in A, s \in S, t = 0, 1, 2, \dots$ ,

$$\mathbb{P}_{\pi}(S_t = s, A_t = a | S_0 = s_0) = \mathbb{P}_{\bar{\pi}}(S_t = s, A_t = a | S_0 = s_0) \quad (1.5.1)$$

(Εν γένει, η  $\bar{\pi}$  εξαρτάται από την αρχική κατάσταση)

### Απόδειξη

Παρατηρούμε ότι η ιστορία μέχρι τον χρόνο  $t = 0$  είναι απλά  $h_0 = (s_0)$ .

Εάν  $\pi = (f_0, f_1, \dots, f_t, \dots)$  και ορίσουμε την  $\bar{\pi} \doteq (\bar{f}_0, \bar{f}_1, \dots, \bar{f}_t, \dots)$ , όπου  $\bar{f}_0 \doteq f_0$ , τότε προφανώς η σχέση (1.5.1) ισχύει για  $t = 0$ .

Θα δείξουμε ότι αν ορίσουμε την  $\bar{f}_t$  ως:

$$\bar{f}_t(s_t, a) = \mathbb{P}_{\pi}[A_t = a | S_0 = s_0, S_t = s_t] \quad (1.5.2)$$

για όλα τα  $a \in A(s_t), s_t, s_0 \in S, t = 1, 2, 3, \dots$

τότε η σχέση (1.5.1) ισχύει για όλους τους μη-αρνητικούς ακεραίους  $t$ , με επαγωγή. Πράγματι, είδαμε ότι ισχύει για  $t = 0$  και υποθέτουμε ότι ισχύει για όλα τα  $t = 0, 1, 2, \dots, T - 1$ .

Το αριστερό μέλος της (1.5.1) για  $t = T$  γράφεται ως εξής:

$$\begin{aligned}\mathbb{P}_{\pi}(S_T = s_T, A_T = a | S_0 = s_0) &= \sum_{s=1}^N \sum_{a \in A(s)} \mathbb{P}_{\pi}[(S_{T-1} = s, A_{T-1} = a | S_0 = s_0)] p(s_T | s, a) \\ &= \sum_{s=1}^N \sum_{a \in A(s)} \mathbb{P}_{\bar{\pi}}[(S_{T-1} = s, A_{T-1} = a | S_0 = s_0)] p(s_T | s, a) \\ &= \mathbb{P}_{\bar{\pi}}[S_T = s_T | S_0 = s_0]\end{aligned}$$

για όλα τα  $s_T, s_0$ . Συνεπώς, η (1.5.1) ισχύει και για την περίπτωση  $t = T$ , το οποίο ολοκληρώνει την απόδειξη. ■

### Πόρισμα 1.5.1.

Έστω  $s \in S$  και  $\pi \in F_B$  μια αυθαίρετη συμπεριφορική στρατηγική. Θεωρούμε την  $\bar{\pi}$  η οποία κατασκευάζεται από την  $\pi$  όπως στη σχέση (1.5.2). Έστω  $\Gamma_B, \Gamma_T, \Gamma_a$  το μοντέλο με προεξόφληση, το τερματιζόμενο και το μοντέλο οριακού μέσου, αντίστοιχα, όπως αυτά έχουν οριστεί.

Τότε, χωρίς περιορισμό της γενικότητας, μπορούμε να περιοριστούμε στο σύνολο στρατηγικών  $F_M$ , γιατί:

- (i)  $v_{\beta}(s, \pi) = v_{\beta}(s, \bar{\pi}), v_T(s, \pi) = v_T(s, \bar{\pi}), v_a(s, \pi) = v_a(s, \bar{\pi})$
- (ii)  $\sup_{F_B} v_{\beta}(s, \pi) = \sup_{F_M} v_{\beta}(s, \pi), \sup_{F_B} v_T(s, \pi) = \sup_{F_M} v_T(s, \pi), \sup_{F_B} v_a(s, \pi) = \sup_{F_M} v_a(s, \pi)$

### Απόδειξη

Αρχικά παρατηρούμε ότι οποιοδήποτε από τα τρία κριτήρια απόδοσης χρησιμοποιεί την ακολουθία αναμενόμενων αμοιβών  $\{E_{s,\pi}[R_t]\}_{t=0}^{\infty}$ ,  $\pi \in F_B, s \in S$

Από το προηγούμενο θεώρημα έχουμε:

$$\begin{aligned}E_{s_0, \pi}[R_t] &= \sum_{s=1}^N \sum_{a \in A(s)} \mathbb{P}_{\pi}[(S_t = s, A_t = a | S_0 = s_0)] r(s, a) \\ &= \sum_{s=1}^N \sum_{a \in A(s)} \mathbb{P}_{\bar{\pi}}[(S_t = s, A_t = a | S_0 = s_0)] r(s, a) \\ &= E_{s_0, \bar{\pi}}[R_t]\end{aligned}$$

για όλα τα  $s_0 \in S$ . Οπότε τα (i),(ii) έπονται άμεσα από τους ορισμούς των αντίστοιχων κριτηρίων απόδοσης και την τελευταία ισότητα.

■

### *Σχόλιο*

Συνοψίζοντας, το Θεώρημα 1.5.1 και το Πόρισμα 1.5.1 μας επιτρέπουν να περιοριστούμε στην κλάση των Μαρκοβιανών στρατηγικών, χωρίς αυτό να περιορίζει την γενικότητα.

## Κεφάλαιο 2

### Στοχαστικά παίγνια και γραμμικός προγραμματισμός

#### 2.0 Εισαγωγή

Τα στοχαστικά παιχνίδια αποτελούν γενίκευση των Μαρκοβιανών διαδικασιών απόφασης. Πλέον, η εξέλιξη της διαδικασίας ελέγχεται από δύο ή περισσότερους ελεγκτές, οι οποίοι θα καλούνται παίκτες. Στα πλαίσια της εργασίας αυτής, θα μας απασχολήσουν μόνον παίγνια δύο παικτών και μηδενικού αθροίσματος, το οποίο πολύ απλά ερμηνεύεται ως “το κέρδος του ενός εκ των δύο παικτών θα ισούται με τη ζημία του δεύτερου παίκτη”. Τα παιχνίδια που θα μελετήσουμε είναι μη-συνεργατικά (ανταγωνιστικά), δηλαδή ο κάθε παίκτης επιδιώκει να παίξει με τέτοιο τρόπο ώστε να πετύχει την μεγιστοποίηση του κέρδους του ή την ελαχιστοποίηση της ζημίας του. Φυσικά, όπως και στο προηγούμενο κεφάλαιο, θα μας απασχολήσουν περιπτώσεις όπου τα σύνολα καταστάσεων και ενεργειών είναι πεπερασμένα. Και πάλι το ενδιαφέρον μας θα στραφεί σε συγκεκριμένες κατηγορίες παιγνίων, τα οποία διαφέρουν ως προς τη δομή ή/και το κριτήριο πληρωμής των παικτών. Ο περιορισμός αυτός κρίνεται απαραίτητος, εφόσον δεν είναι ακόμα γνωστό εάν υπάρχει αναλυτική λύση για οποιοδήποτε γενικό Στοχαστικό παίγνιο, πόσο μάλλον όταν περιοριζόμαστε στα παίγνια εκείνα που επιδέχονται λύση με μεθόδους Γραμμικού Προγραμματισμού.

Αξιοσημείωτο είναι το γεγονός πως η θεωρία περί Στοχαστικών Παιγνίων αναπτύχθηκε ανεξάρτητα και, μάλιστα, προγενέστερα της θεωρίας των Μαρκοβιανών Διαδικασιών απόφασης ως ένα εντελώς ξεχωριστό αντικείμενο, προτού βρεθεί η σύνδεση μεταξύ των δύο κλάδων.

#### 2.1 Στοχαστικά παίγνια με συντελεστή προεξόφλησης

Στην παράγραφο αυτή θα γενικεύσουμε το Μαρκοβιανό μοντέλο απόφασης  $\Gamma_\beta$  με συντελεστή προεξόφλησης  $\beta$ , όπως αυτό ορίστηκε στην παράγραφο 1.1, για την περίπτωση που έχουμε 2 παίκτες-ελεγκτές.

Θεωρούμε πως, αν η διαδικασία βρίσκεται στην κατάσταση  $s \in S = \{1, 2, \dots, N\}$  στον χρόνο  $t$ , οι παίκτες επιλέγουν ανεξάρτητα ο ένας από τον άλλον ενέργειες  $a^1 \in A^1(s)$  και  $a^2 \in A^2(s)$  και λαμβάνουν αμοιβές  $r^1(s, a^1, a^2)$  και  $r^2(s, a^1, a^2)$ , αντίστοιχα.

Επιπλέον, οι πιθανότητες μετάβασης γενικεύονται έτσι, ώστε:

$$p(s'|s, a^1, a^2) \doteq \mathbb{P}\{S_{t+1} = s' | S_t = s, A_t^1 = a_1, A_t^2 = a_2\} \text{ , για κάθε } t = 0, 1, 2, \dots$$

Σε αναλογία με τον συμβολισμό του προηγούμενου κεφαλαίου, θα συμβολίζουμε με  $F_s$  το σύνολο των στάσιμων στρατηγικών του παίκτη 1 και με  $G_s$  το αντίστοιχο σύνολο για τον παίκτη



2. Δηλαδή, αν  $\mathbf{g} = (g(1), g(2), \dots, g(N)) \in G_s$ , τότε κάθε  $g(s)$  είναι ένα  $m^2(s)$ -διάστατο διάνυσμα πιθανοτήτων, όπου  $m^2(s) = |A^2(s)|$ .

Ακόμη, θα συμβολίζουμε με  $R_t^1$  ( $R_t^2$ ) τη συνολική αμοιβή του παίκτη 1 (παίκτη 2) στο χρόνο  $t$ , ενώ με  $r^1(\mathbf{f}, \mathbf{g})$  ( $r^2(\mathbf{f}, \mathbf{g})$ ) την αναμενόμενη αμοιβή ενός βήματος του παίκτη 1 (παίκτη 2) που αντιστοιχεί στο ζεύγος στρατηγικών  $(\mathbf{f}, \mathbf{g}) \in F_s \times G_s$ . Η αναμενόμενη αμοιβή του παίκτη  $k$  στο στάδιο  $t$  για τις στρατηγικές  $(\mathbf{f}, \mathbf{g})$  και αρχική κατάσταση  $s$ , θα συμβολίζεται με

$E_{s\mathbf{f}\mathbf{g}}[R_t^k]$ . Συνεπώς, η συνολική τιμή με προεξόφληση για ένα ζεύγος στρατηγικών  $(\mathbf{f}, \mathbf{g}) \in F_s \times G_s$  για τον παίκτη  $k$  ( $k = 1$  ή  $k = 2$ ) θα δίνεται από την:

$$v_\beta^k(s, \mathbf{f}, \mathbf{g}) \doteq \sum_{t=0}^{\infty} \beta^t E_{s\mathbf{f}\mathbf{g}}[R_t^k], \quad \text{όπου } \beta \in [0,1) \text{ και } k \in \{1,2\}$$

Θα θέλαμε στο σημείο αυτό να βρούμε, με κάποιο τρόπο, ένα ζευγάρι στάσιμων στρατηγικών  $(\mathbf{f}, \mathbf{g})$ , το οποίο να αποτελεί 'λύση' του παιχνιδιού. Στο αντίστοιχο σημείο του προηγούμενου κεφαλαίου είχαμε ορίσει ως λύση για τον παίκτη 1 την βελτιστοποίηση του προβλήματος:

$$\begin{cases} \max v_\beta^1(\mathbf{f}, \mathbf{g}) \\ \mathbf{f} \in F_s \end{cases}$$

Προφανώς μια τέτοια αντιμετώπιση θα κρίνονταν ανεπαρκής στην περίπτωση των δύο παικτών, εφόσον η λύση του παιχνιδιού, εν γένει, εξαρτάται και από την στρατηγική  $\mathbf{g} \in G_s$  του 2<sup>ου</sup> παίκτη.

Η κλασική υπόθεση που κάνουμε στα μη-συνεργατικά παίγνια τα οποία θα μας απασχολήσουν, έγκειται στο ότι οι παίκτες επιλέγουν ανεξάρτητα (και κρυφά) τις ενέργειές τους, καθώς και ότι μοναδικός τους στόχος είναι η μεγιστοποίηση της συνολικής προσωπικής τους αμοιβής. Εάν ακόμη υποθέσουμε ότι οι παίκτες γνωρίζουν επακριβώς το σύνολο στρατηγικών και τη συνάρτηση πληρωμής των υπόλοιπων παικτών, τότε η πλέον διαδεδομένη λύση του παιχνιδιού που ψάχνουμε είναι γνωστή στη βιβλιογραφία ως *σημείο Nash* ή *σημείο Ισορροπίας κατά Nash* (*Nash Equilibrium*).

### Ορισμός 2.1.1.

Θα λέμε ότι το ζεύγος στρατηγικών  $(\mathbf{f}^0, \mathbf{g}^0) \in F_s \times G_s$  αποτελεί ένα *σημείο Ισορροπίας κατά Nash* για το στοχαστικό παίγνιο  $\Gamma_\beta$  με συντελεστή προεξόφλησης, εάν:

$$v_\beta^1(\mathbf{f}, \mathbf{g}^0) \leq v_\beta^1(\mathbf{f}^0, \mathbf{g}^0) \quad , \text{για όλα τα } \mathbf{f} \in F_s$$

και

$$v_\beta^2(\mathbf{f}^0, \mathbf{g}) \leq v_\beta^2(\mathbf{f}^0, \mathbf{g}^0) \quad , \text{για όλα τα } \mathbf{g} \in G_s$$

Από τον παραπάνω ορισμό μπορούμε να δούμε πως η καταλληλότητα του σημείου Nash ως λύση του παιχνιδιού, έγκειται στο γεγονός πως κανένας παίκτης δεν ωφελείται εάν επιχειρήσει να παρεκκλίνει από την στρατηγική του. Με άλλα λόγια, εάν  $(f^0, g^0)$  είναι ένα σημείο Nash και ο παίκτης 1 γνωρίζει πως ο παίκτης 2 θα παίξει την στρατηγική  $g^0$ , τότε το καλύτερο που μπορεί να κάνει ο παίκτης 1 για να μεγιστοποιήσει το κέρδος του, είναι να επιλέξει την στρατηγική  $f^0$ .

Ωστόσο, ένα μεγάλο πρόβλημα που ενδέχεται να συναντήσουμε σε πολλά παίγνια, είναι η ύπαρξη πολλαπλών σημείων Nash με διαφορετικές πληρωμές για τους παίκτες.

Μία ευρεία κατηγορία παιχνιδιών στην οποία θα εστιάσουμε για το υπόλοιπο της εργασίας, αποτελούν τα παιχνίδια μηδενικού αθροίσματος. Θα δούμε, μάλιστα, πως για αυτή την κλάση παιγνίων, το πρόβλημα της ύπαρξης πολλαπλών σημείων Nash δεν υφίσταται.

Ένα στοχαστικό παίγνιο θα καλείται μηδενικού-αθροίσματος εάν:

$$r^1(s, a^1, a^2) + r^2(s, a^1, a^2) = 0, \text{ για κάθε } s \in S, a^1 \in A^1(s), a^2 \in A^2(s)$$

Οπότε μπορούμε να απλοποιήσουμε ελαφρώς τον συμβολισμό και να έχουμε:

$$\begin{cases} r(s, a^1, a^2) \doteq r^1(s, a^1, a^2) = -r^2(s, a^1, a^2), & s \in S, a^1 \in A^1(s), a^2 \in A^2(s) \\ v_\beta(f, g) \doteq v_\beta^1(f, g) = -v_\beta^2(f, g), & (f, g) \in F_s \times G_s \end{cases}$$

Έστω ότι  $(f^0, g^0)$  είναι ένα σημείο Nash. Τότε, σύμφωνα με τα παραπάνω, οι δύο ανισότητες του ορισμού του σημείου Nash ανάγονται στην :

$$v_\beta(f, g^0) \leq v_\beta(f^0, g^0) \leq v_\beta(f^0, g), \quad f \in F_s, g \in G_s \quad (2.1.1)$$

Θα καλούμε την στρατηγική  $f^0$  βέλτιστη στάσιμη στρατηγική για τον παίκτη 1 (Αντίστοιχα για τον παίκτη 2).

Στη συνέχεια παραθέτουμε ένα ιδιαίτερα χρήσιμο θεώρημα το οποίο θα απαντήσει στο πρόβλημα περί ύπαρξης πολλαπλών σημείων Nash σε ένα παίγνιο με 2 παίκτες.

### Θεώρημα 2.1.1.

Έστω ότι η ανισότητα (2.1.1) ικανοποιείται για τα ζεύγη στρατηγικών  $(f^0, g^0)$  και  $(f', g')$  που ανήκουν στο  $F_s \times G_s$ . Τότε:

$$(i) \quad v_\beta = v_\beta(f^0, g^0) = v_\beta(f', g') = v_\beta(f^0, g') = v_\beta(f', g^0)$$

$$(ii) \quad v_\beta(s, f^0, g^0) = \max_{F_s} \min_{G_s} v_\beta(s, f, g) = \min_{G_s} \max_{F_s} v_\beta(s, f, g), \quad \text{για κάθε } s \in S$$

- (iii) Αντίστροφα, εάν υπάρχουν τα  $\max_{F_s} \min_{G_s} v_\beta(s, \mathbf{f}, \mathbf{g})$ ,  $\min_{G_s} \max_{F_s} v_\beta(s, \mathbf{f}, \mathbf{g})$  και είναι ίσα για κάποιο  $s_0 \in S$ , τότε υπάρχουν στάσιμες στρατηγικές  $\mathbf{f}^0, \mathbf{g}^0$  οι οποίες ικανοποιούν την (2.1), για την κατάσταση  $s_0$ .

### Απόδειξη

- (i) Υποθέτουμε ότι τα ζεύγη  $(\mathbf{f}^0, \mathbf{g}^0)$  και  $(\mathbf{f}', \mathbf{g}')$  ικανοποιούν την (2.1.1). Αυτό σημαίνει πως:

$$v_\beta(\mathbf{f}, \mathbf{g}^0) \leq v_\beta(\mathbf{f}^0, \mathbf{g}^0) \leq v_\beta(\mathbf{f}^0, \mathbf{g}), \quad \text{για όλα τα } \mathbf{f} \in F_s, \mathbf{g} \in G_s$$

Οπότε θα πρέπει να ισχύει και ότι:

$$v_\beta(\mathbf{f}', \mathbf{g}^0) \leq v_\beta(\mathbf{f}^0, \mathbf{g}^0) \leq v_\beta(\mathbf{f}^0, \mathbf{g}') \quad (2.1.2)$$

Από την άλλη μεριά, για το ζεύγος  $(\mathbf{f}', \mathbf{g}')$  παίρνουμε ότι:

$$v_\beta(\mathbf{f}, \mathbf{g}') \leq v_\beta(\mathbf{f}', \mathbf{g}') \leq v_\beta(\mathbf{f}', \mathbf{g}), \quad \text{για όλα τα } \mathbf{f} \in F_s, \mathbf{g} \in G_s$$

Οπότε θα πρέπει να ισχύει και ότι:

$$v_\beta(\mathbf{f}^0, \mathbf{g}') \leq v_\beta(\mathbf{f}', \mathbf{g}') \leq v_\beta(\mathbf{f}', \mathbf{g}^0) \quad (2.1.3)$$

Συνδυάζοντας τις (2.1.2) και (2.1.3) λαμβάνουμε το ζητούμενο αποτέλεσμα.

- (ii) Από την (2.1.1) γνωρίζουμε ότι :

$$\max_{\mathbf{f} \in F_s} v_\beta(s, \mathbf{f}, \mathbf{g}^0) = v_\beta(s, \mathbf{f}^0, \mathbf{g}^0) = \min_{\mathbf{g} \in G_s} v_\beta(s, \mathbf{f}^0, \mathbf{g})$$

Επιπλέον:

$$\max_{\mathbf{f} \in F_s} \min_{\mathbf{g} \in G_s} v_\beta(s, \mathbf{f}, \mathbf{g}) \geq \min_{\mathbf{g} \in G_s} v_\beta(s, \mathbf{f}^*, \mathbf{g})$$

και

$$\min_{\mathbf{g} \in G_s} \max_{\mathbf{f} \in F_s} v_\beta(s, \mathbf{f}, \mathbf{g}) \leq \max_{\mathbf{f} \in F_s} v_\beta(s, \mathbf{f}, \mathbf{g}^*)$$

για όλα τα  $\mathbf{f}^* \in F_s$  και  $\mathbf{g}^* \in G_s$ .

Οπότε παίρνουμε ότι:

$$\max_{\mathbf{f} \in F_s} \min_{\mathbf{g} \in G_s} v_\beta(s, \mathbf{f}, \mathbf{g}) \geq \min_{\mathbf{g} \in G_s} v_\beta(s, \mathbf{f}^0, \mathbf{g}) = \max_{\mathbf{f} \in F_s} v_\beta(s, \mathbf{f}, \mathbf{g}^0) \geq \min_{\mathbf{g} \in G_s} \max_{\mathbf{f} \in F_s} v_\beta(s, \mathbf{f}, \mathbf{g})$$

και:

$$\min_{\mathbf{g} \in G_s} \max_{\mathbf{f} \in F_s} v_\beta(s, \mathbf{f}, \mathbf{g}) \geq \min_{\mathbf{g} \in G_s} v_\beta(s, \mathbf{f}^0, \mathbf{g}) = \max_{\mathbf{f} \in F_s} v_\beta(s, \mathbf{f}, \mathbf{g}^0) \geq \max_{\mathbf{f} \in F_s} \min_{\mathbf{g} \in G_s} v_\beta(s, \mathbf{f}, \mathbf{g})$$

Οι δύο παραπάνω ανισοτικές σχέσεις συνεπάγον την ζητούμενη ισότητα.

(iii) Ορίζουμε τις ακόλουθες ποσότητες:

$$v_\beta \doteq v_\beta(s^0, f^0, g^0), \quad F_\beta(f) \doteq \min_{g \in G_s} v_\beta(s^0, f, g), \quad G_\beta(g) \doteq \max_{f \in F_s} v_\beta(s^0, f, g)$$

όπου  $f \in F_s, g \in G_s$

Από την υπόθεσή μας γνωρίζουμε ότι θα υπάρξει ένα  $f^0 \in F_s$ , τέτοιο ώστε  $F_\beta(f^0) = v_\beta$

Αυτό σημαίνει ότι:

$$v_\beta = \min_{g \in G_s} v_\beta(s^0, f^0, g) \leq v_\beta(s^0, f^0, g), \quad \text{για όλα τα } g \in G_s$$

Με ανάλογο επιχείρημα, γνωρίζουμε ότι θα υπάρξει  $g^0 \in G_s$  ώστε:

$$v_\beta = G_\beta(g^0) = \max_{f \in F_s} v_\beta(s^0, f, g^0) \geq v_\beta(s^0, f, g^0), \quad \text{για όλα τα } f \in F_s$$

Οπότε το ζεύγος στρατηγικών  $(f^0, g^0)$  αποτελεί ένα ζεύγος στάσιμων βέλτιστων στρατηγικών.

■

## Σχόλιο

Από το προηγούμενο θεώρημα συμπεραίνουμε πως τα διανύσματα προεξοφλητικής τιμής ταυτίζονται για όλα τα ζεύγη βέλτιστων στρατηγικών. Αυτό το διάνυσμα θα καλείται *διάνυσμα τιμής του παιγνίου μηδενικού αθροίσματος*  $\Gamma_\beta$ , και θα συμβολίζεται με :

$$v_\beta = (v_\beta(1), v_\beta(2), \dots, v_\beta(N))^T$$

Έως τώρα έχουμε δει τα στοχαστικά παίγνια με προεξόφληση από την σκοπιά της θεωρίας αποφάσεων. Ορίστηκε, δηλαδή, το  $\Gamma_\beta$  ως μια γενίκευση της Μαρκοβιανής διαδικασίας απόφασης για περισσότερους από έναν ελεγκτές-παίκτες, και συγκεκριμένα για την περίπτωση των δύο παικτών. Μια εναλλακτική σκοπιά, από την οποία θα μπορούσαμε να κοιτάξουμε τα εν λόγω παίγνια, αποτελεί εκείνη των *στατικών πινακοπαιχνιδιών* (matrix games). Μπορούμε, με αυτό το σκεπτικό, να ορίσουμε το παίγνιο  $\Gamma_\beta$  ως γενίκευση αυτών των πινακοπαιχνιδιών σε πολλαπλά στάδια.

Έστω  $m^1(s)$  η πληθικότητα του  $A^1(s)$  (αντίστοιχα  $m^2(s)$  η πληθικότητα του  $A^2(s)$ ), με  $s \in S$ . Τότε, ορίζουμε  $N$  στο πλήθος πινακοπαιχνίδια ως εξής:

$$R(s) = [r(s, a^1, a^2)]_{a^1=1, a^2=1}^{m^1(s), m^2(s)}$$

τα οποία βρίσκονται σε 1-1 αντιστοιχία με τις καταστάσεις του  $\Gamma_\beta$ .

Ένα πινακοπαιχνίδι  $R(s)$  προσδιορίζεται πλήρως από έναν πίνακα  $\mathbf{M}$  διάστασης  $m^1(s) \times m^2(s)$ , όπου το  $m_{ij}$  στοιχείο του πίνακα είναι η αμοιβή που θα λάβει ο παίκτης 1 εάν αυτός επιλέξει την ενέργεια  $i$  και ο παίκτης 2 επιλέξει την ενέργεια  $j$ . Φυσικά, εφόσον αναφερόμαστε σε παίγνια δύο παικτών και μηδενικού αθροίσματος, από τον πίνακα  $\mathbf{M}$  είναι γνωστές και οι αμοιβές του δεύτερου παίκτη. Η  $i$ -γραμμή του πίνακα  $\mathbf{M}$  αντιπροσωπεύει την  $i$ -ενέργεια του παίκτη 1, ενώ η  $j$ -στήλη αντιπροσωπεύει την  $j$ -ενέργεια του παίκτη 2. Τα πινακοπαιχνίδια αποτελούν την πιο απλή κατηγορία παιχνιδιών και ιστορικά αποτέλεσαν τα πρώτα παιχνίδια που μελετήθηκαν. Για παράδειγμα, το «Πέτρα-Ψαλίδι-Χαρτί» αποτελεί ένα πινακοπαιχνίδι, όπου ο πίνακας  $\mathbf{M}$  είναι διάστασης  $3 \times 3$ .

Θα λέμε ότι ένα πινακοπαιχνίδι δύο παικτών και μηδενικού αθροίσματος έχει τιμή  $v$ , αν και μόνο αν

$$\sup_f \inf_g r(f, g) = v = \inf_g \sup_f r(f, g)$$

Για ένα τέτοιο παιχνίδι με τιμή  $v$  οι στρατηγικές  $f_\varepsilon, g_\varepsilon$  θα καλούνται  $\varepsilon$ -βέλτιστες,  $\varepsilon \geq 0$ , αν  $\inf_g r(f_\varepsilon, g) \geq v - \varepsilon$  και  $\sup_f r(f, g_\varepsilon) \leq v - \varepsilon$ . Οι  $0$ -βέλτιστες στρατηγικές θα καλούνται απλά βέλτιστες.

Η ακόλουθη πρόταση που αφορά στα πινακοπαιχνίδια και της βέλτιστες στρατηγικές οφείλεται στον J. von Neumann (1928).

### Πρόταση 2.1.1.

Για όλους τους πίνακες  $\mathbf{A}$ , το πινακοπαιχνίδι  $[\mathbf{A}]$  έχει τιμή και οι παίκτες διαθέτουν βέλτιστες στρατηγικές.

Φυσικά, οι στρατηγικές δεν θα είναι υποχρεωτικά ντετερμινιστικές, αλλά, εν γένει, μεικτές.

Μπορούμε να ερμηνεύσουμε το 'παίξιμο' ενός μόνο σταδίου, στην κατάσταση  $s$ , σαν οι παίκτες να παίζουν το πινακοπαιχνίδι  $R(s)$  μία φορά, όπου οι επιλογές των ενεργειών  $a^1 \in A^1(s)$  και  $a^2 \in A^2(s)$  προσδιορίζουν τόσο την αμοιβή ενός βήματος  $r(s, a^1, a^2)$  όσο και την πιθανότητα μετάβασης  $p(s'|s, a^1, a^2)$  στο επόμενο πινακοπαιχνίδι  $R(s')$  που θα παιχτεί στο επόμενο στάδιο.

Ακολουθώντας όμοια συλλογιστική πορεία με τη θεωρία των διαδικασιών απόφασης, τότε εάν υποθέσουμε πως η τιμή του παιχνιδιού,  $v_\beta$ , υπάρχει και ότι γνωρίζουμε ποιο είναι το βέλτιστο παίξιμο από το επόμενο στάδιο και έπειτα, τότε στο παρόν στάδιο έχουμε να αντιμετωπίσουμε το ακόλουθο βοηθητικό πινακοπαιχνίδι:

$$R(s, \mathbf{v}_\beta) = \left[ r(s, a^1, a^2) + \beta \sum_{s' \in S} p(s'|s, a^1, a^2) v_\beta(s') \right]_{a^1=1, a^2=1}^{m^1(s), m^2(s)} \quad (2.1.4)$$

Μπορούμε τώρα να υποθέσουμε ότι η σχέση (1.2.1) του 1<sup>ου</sup> Κεφαλαίου γενικεύεται και στην περίπτωση των στοχαστικών παιχνιδιών.

Εάν, λοιπόν, υπάρχει η τιμή  $\mathbf{v}_\beta$  του παιχνιδιού, τότε για κάθε  $s \in S$  θα πρέπει να ικανοποιείται ότι:

$$v_\beta(s) = \text{val}[R(s, \mathbf{v}_\beta)] \quad (2.1.5)$$

όπου με  $\text{val}[A]$  θα συμβολίζουμε την τιμή του πινακοπαιχνιδιού  $A$ , ενώ  $v_\beta(s)$  είναι το  $s$ -οστό στοιχείο του διανύσματος  $\mathbf{v}_\beta$ .

Όλα τα παραπάνω συνοψίζονται στο επόμενο βασικό θεώρημα, το οποίο παραθέτουμε χωρίς απόδειξη.

### Θεώρημα 2.1.2. (Shapley)

Το στοχαστικό παίγνιο μηδενικού αθροίσματος  $\Gamma_\beta$  με συντελεστή προεξόφλησης  $\beta$ , έχει ως διάνυσμα τιμής,  $\mathbf{v}_\beta$ , την μοναδική λύση των εξισώσεων:

$$v(s) = \text{val}[R(s, \mathbf{v})], \text{ για όλα τα } s \in S, \text{ όπου } \mathbf{v}^T = (v(1), v(2), \dots, v(N))^T$$

Ακόμη, εάν  $(f^0(s), g^0(s))$  είναι ένα ζευγάρι βέλτιστων στρατηγικών (εν γένει, όχι καθαρών) στο πινακοπαιχνίδι  $R(s, \mathbf{v}_\beta)$  για κάθε  $s \in S$ , τότε :

Η  $\mathbf{f}^0 = (f^0(1), f^0(2), \dots, f^0(N))$  είναι μια βέλτιστη στάσιμη στρατηγική για τον παίκτη 1 στο  $\Gamma_\beta$  και η  $\mathbf{g}^0 = (g^0(1), g^0(2), \dots, g^0(N))$  είναι μια βέλτιστη στάσιμη στρατηγική για τον παίκτη 2 στο  $\Gamma_\beta$ .

Κλείνοντας αυτή την παράγραφο, επεκτείνουμε τον συμβολισμό του προηγούμενου κεφαλαίου και παραθέτουμε τον ορισμό κάποιων ποσοτήτων που θα χρησιμοποιήσουμε εν συνεχεία.

Έστω ένα σταθεροποιημένο ζεύγος στάσιμων στρατηγικών  $\mathbf{f} = (\mathbf{f}(1), \mathbf{f}(2), \dots, \mathbf{f}(N))$  και  $\mathbf{g} = (\mathbf{g}(1), \mathbf{g}(2), \dots, \mathbf{g}(N))$  των παικτών 1 και 2 αντίστοιχα. Θα κάνουμε τη σύμβαση πως η  $\mathbf{f}$  είναι ένα μπλοκ διάνυσμα-γραμμή, ενώ η  $\mathbf{g}$  θα είναι ένα μπλοκ διάνυσμα-στήλη.

Συνεπώς, αν

$$m^2 \doteq \sum_{s=1}^N m^2(s)$$

τότε η  $\mathbf{g}$  θα είναι ένα  $m^2$ -διάστατο διάνυσμα-στήλη όπου το  $s$ -οστό μπλοκ  $\mathbf{g}(s)$  θα είναι  $m^2(s)$ -διάστατο.

Στη συνέχεια ορίζουμε :

(i) Τις πιθανότητες μετάβασης στην επόμενη κατάσταση:

$$\begin{aligned} p(s'|s, \mathbf{f}, a^2) &\doteq \sum_{a^1=1}^{m^1(s)} p(s'|s, a^1, a^2) f(s, a^1), \quad s, s' \in S, a^2 \in A^2(s) \\ p(s'|s, a^1, \mathbf{g}) &\doteq \sum_{a^2=1}^{m^2(s)} p(s'|s, a^1, a^2) g(s, a^2), \quad s, s' \in S, a^1 \in A^1(s) \\ p(s'|s, \mathbf{f}, \mathbf{g}) &\doteq \sum_{a^1=1}^{m^1(s)} \sum_{a^2=1}^{m^2(s)} p(s'|s, a^1, a^2) f(s, a^1) g(s, a^2), \quad s, s' \in S \end{aligned}$$

(ii) Τον πίνακα πιθανοτήτων μετάβασης που επάγεται από τις  $(\mathbf{f}, \mathbf{g})$  :

$$\mathbf{P}(\mathbf{f}, \mathbf{g}) \doteq [p(s'|s, \mathbf{f}, \mathbf{g})]_{s, s' \in S}$$

(iii) Τις αναμενόμενες αμοιβές ενός βήματος:

$$\begin{aligned} r(s, \mathbf{f}, a^2) &\doteq \sum_{a^1=1}^{m^1(s)} r(s, a^1, a^2) f(s, a^1) = [\mathbf{f}(s) \mathbf{R}(s)]_{a^2}, \quad s \in S, a^2 \in A^2(s) \\ r(s, a^1, \mathbf{g}) &\doteq \sum_{a^2=1}^{m^2(s)} r(s, a^1, a^2) g(s, a^2) = [\mathbf{R}(s) \mathbf{g}(s)]_{a^1}, \quad s \in S, a^1 \in A^1(s) \\ r(s, \mathbf{f}, \mathbf{g}) &\doteq \sum_{a^1=1}^{m^1(s)} \sum_{a^2=1}^{m^2(s)} r(s, a^1, a^2) f(s, a^1) g(s, a^2) = \mathbf{f}(s) \mathbf{R}(s) \mathbf{g}(s), \quad s \in S \end{aligned}$$

(iv) Το  $N$ -διάστατο διάνυσμα-στήλη αμοιβών ενός βήματος:

$$\mathbf{r}(\mathbf{f}, \mathbf{g}) \doteq (r(1, \mathbf{f}, \mathbf{g}), r(2, \mathbf{f}, \mathbf{g}), \dots, r(N, \mathbf{f}, \mathbf{g}))^T$$

(v) Το  $N$ -διάστατο διάνυσμα τιμής με προεξόφληση για το ζεύγος  $(\mathbf{f}, \mathbf{g})$  :

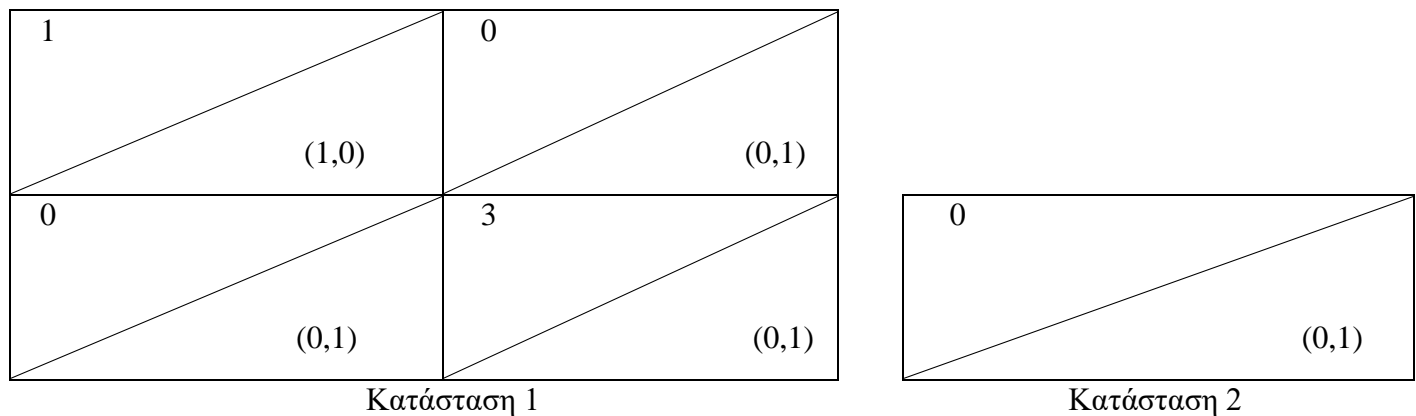
$$\mathbf{v}_\beta(\mathbf{f}, \mathbf{g}) \doteq (\mathbf{I} - \beta \mathbf{P}(\mathbf{f}, \mathbf{g}))^{-1} \mathbf{r}(\mathbf{f}, \mathbf{g})$$

## 2.2 Γραμμικός Προγραμματισμός και Στοχαστικά Παιγνία με συντελεστή προεξόφλησης

Όπως έχουμε ήδη διαπιστώσει, ο γραμμικός προγραμματισμός αποτελεί ένα ιδιαίτερος χρήσιμο εργαλείο για την επίλυση Μαρκοβιανών διαδικασιών Απόφασης. Όπως θα γίνει φανερό μέσω του επόμενου παραδείγματος, δεν θα πρέπει να αναμένουμε να φτάσουμε στη λύση οποιουδήποτε στοχαστικού παιγνίου με μεθόδους γραμμικού προγραμματισμού.

### Παράδειγμα 2.2.1

Έστω  $S = \{1, 2\}$ ,  $A^1(1) = A^2(1) = \{1, 2\}$ ,  $A^1(2) = A^2(2) = \{1\}$ ,  $\beta = \frac{1}{2}$ , ενώ οι αμοιβές και οι πιθανότητες μετάβασης φαίνονται στο επόμενο γράφημα:



Σκοπός του παραπάνω παραδείγματος είναι να δείξουμε πως, πράγματι, οι μέθοδοι γραμμικού προγραμματισμού κρίνονται ανεπαρκείς για την επίλυση ενός οποιουδήποτε στοχαστικού παιγνίου. Για το λόγο αυτό, χωρίς να δούμε καθόλου το θεωρητικό υπόβαθρο που απαιτείται για



τον υπολογισμό της τιμής του παιγνίου αυτού, παραθέτουμε απλώς το αποτέλεσμα. Αναφέρουμε, μόνο, πως το αποτέλεσμα αυτό είναι συνέπεια ενός θεωρήματος γνωστού στη βιβλιογραφία ως **Shapley-Snow Theorem**. Από το θεώρημα προκύπτει ότι τελικά το ζητούμενο  $v$  θα είναι λύση της εξίσωσης:

$$v = \frac{(1 + \frac{1}{2}v)(3 + \frac{1}{2}v)}{4 + v} \Rightarrow 3v^2 + 8v - 12 = 0 \Rightarrow v = \frac{1}{3}(-4 + 2\sqrt{13})$$

Οπότε το διάνυσμα τιμής θα είναι το :

$$v = \left( \frac{1}{3}(-4 + 2\sqrt{13}), 0 \right)^T$$

### Σχόλιο

Αξίζει να παρατηρήσουμε πως, εάν το παραπάνω παράδειγμα επιδέχονταν λύση με μεθόδους γραμμικού προγραμματισμού, με συντελεστές που είναι ρητές συναρτήσεις των (ρητών) δεδομένων του προβλήματος, τότε η βέλτιστη λύση θα έπρεπε να βρίσκεται σε κάποιο οριακό σημείο της εφικτής περιοχής και θα είχε υποχρεωτικά ρητή τιμή. Συνεπώς, το γεγονός ότι στο συγκεκριμένο παράδειγμα καταλήξαμε σε διάνυσμα τιμής με άρρητους αριθμούς, υποδηλώνει πως δεν είναι εφικτή η επίλυση μέσω γραμμικού προγραμματισμού. Το ερώτημα που θα μας απασχολήσει στη συνέχεια είναι το εξής: « Μήπως υπάρχουν ορισμένες ειδικές κλάσεις στοχαστικών παιγνίων, τα οποία μπορούν να αντιμετωπιστούν με μεθόδους γραμμικού προγραμματισμού;»

#### 2.2.1. Παίγνια με προεξόφληση με 1 ελεγκτή

Για αυτή την κατηγορία παιγνίων, η πιθανότητα μετάβασης μεταξύ των καταστάσεων εξαρτάται από τις ενέργειες που επιλέγει μόνον ο ένας παίκτης

Θα συμβολίζουμε με  $\Gamma_\beta(1)$  το παιχνίδι που ‘ελέγχεται’ από τον παίκτη 1 και που ορίζεται από την ιδιότητα:

$$p(s'|s, a^1, a^2) \equiv p(s'|s, a^1) \text{ , όπου } s, s' \in S, a^1 \in A^1(s), a^2 \in A^2(s)$$

Με εντελώς όμοιο τρόπο ορίζεται και το παιχνίδι  $\Gamma_\beta(2)$ .

Δοθέντος του παιχνιδιού  $\Gamma_\beta(1)$ , μπορούμε να ‘ελαφρύνουμε’ τον συμβολισμό που έχουμε δει ως εξής:

$$P(f, g) = P(f) \text{ και}$$

$$v_\beta(f, g) = (I - \beta P(f))^{-1} r(f, g) \text{ , για όλα τα } (f, g) \in F_s \times G_s$$

Λόγω της κατασκευής του  $\Gamma_\beta(1)$ , είναι λογικό να αναμένουμε αυτό να συμπεριφέρεται περισσότερο σαν Μαρκοβιανή διαδικασία Απόφασης του παίκτη 1, παρά σαν γενικό στοχαστικό παίγνιο με προεξόφληση.

Υποθέτουμε πως ο παίκτης 2 ακολουθεί μία στάσιμη στρατηγική  $\mathbf{g} \in G_s$ . Τότε, κατ'αναλογία με τα γραμμικά προγράμματα  $(P_\beta)$  και  $(D_\beta)$  του προηγούμενου κεφαλαίου, θεωρούμε το ακόλουθο ζεύγος (Πρωτεύοντος – Δυϊκού) προβλημάτων :

$$\min \sum_{s'=1}^N \frac{1}{N} v(s')$$

με περιορισμούς: ( $P_\beta(1)$ )

$$v(s) \geq [R(s)\mathbf{g}(s)]_{a^1} + \beta \sum_{s'=1}^N p(s'|s, a^1) v(s'), \quad s \in S, a^1 \in A^1(s)$$

$$\sum_{a^2 \in A^2(s)} g(s, a^2) = 1, \quad s \in S$$

$$g(s, a^2) \geq 0, \quad s \in S, a^2 \in A^2(s)$$

$$\max \sum_{s'=1}^N z(s)$$

με περιορισμούς: ( $D_\beta(1)$ )

$$\sum_{s=1}^N \sum_{a^1 \in A^1(s)} (\delta_{ss'} - \beta p(s'|s, a^1)) x_{sa^1} = \frac{1}{N}, \quad s' \in S$$

$$z(s) \leq [\mathbf{x}(s)R(s)]_{a^2}, \quad s \in S, a^1 \in A^1(s)$$

$$x(s, a^1) \geq 0, \quad s \in S, a^1 \in A^1(s)$$

όπου  $\mathbf{x}(s) = (x(s, 1), x(s, 2), \dots, x(s, m^1(s)))$  για κάθε  $s \in S$ .

### Σχόλιο

Όπως και στην παράγραφο 1.3, τα παραπάνω επιχειρήματα θα ήταν έγκυρα, ακόμη και αν αντικαθιστούσαμε τους συντελεστές  $\frac{1}{N}$  με κάποιες θετικές πιθανότητες εκκίνησης  $\gamma(s')$ , οι οποίες αθροίζουν στη μονάδα.

Το προηγούμενο ζεύγος γραμμικών προγραμμάτων μπορεί να χρησιμοποιηθεί για την επίλυση του στοχαστικού παιγνίου με προεξόφληση με έναν ελεγκτή.

Η επόμενη παρατήρηση θα φανεί ιδιαίτερα χρήσιμη στην συνέχεια του κεφαλαίου.

### Παρατήρηση

Παρατηρούμε πως ορισμένοι από τους περιορισμούς στο  $(P_\beta(1))$  και  $(D_\beta(1))$  είναι σε ομάδες που αντιστοιχούν ακριβώς στη μπλοκ δομή των στάσιμων στρατηγικών  $\mathbf{f}$  ή  $\mathbf{g}$  του παίκτη 1 ή 2. Επειδή αυτή ακριβώς η δομή πρόκειται να εμφανιστεί επανειλημμένα στη συνέχεια, είναι βολικό να χρησιμοποιούμε τη φράση «αναμιγνύοντας τους περιορισμούς ως προς την  $\mathbf{f}$  ή την  $\mathbf{g}$ ». Αυτή η φράση θα σημαίνει ότι κάθε περιορισμός μιας τέτοιας ομάδας που αντιστοιχεί στην  $(s, a^1)$  (αντίστοιχα  $(s, a^2)$ ) πολλαπλασιάζεται επί  $f(s, a^1)$  (αντίστοιχα επί  $g(s, a^2)$ ) και ότι όλοι οι περιορισμοί αυτής της ομάδας θα αθροίζονται για όλα τα  $a^1 \in A^1(s)$  (αντίστοιχα  $a^2 \in A^2(s)$ ).

Για παράδειγμα, έστω μια αυθαίρετη  $\mathbf{f} \in F_s$  και αναμιγνύουμε την πρώτη ομάδα περιορισμών στο  $(P_\beta(1))$  ως προς την  $\mathbf{f}$ . Τότε, για  $s \in S$ , θα πάρουμε:

$$\sum_{a^1 \in A^1(s)} v(s) f(s, a^1) \geq \mathbf{f}(s) R(s) \mathbf{g}(s) + \beta \sum_{s'=1}^N p(s'|s, \mathbf{f}) v(s')$$

ή, ισοδύναμα, αφού  $\sum_{a^1 \in A^1(s)} f(s, a^1) = 1$

$$v(s) \geq r(s, \mathbf{f}, \mathbf{g}) + \beta [\mathbf{P}(\mathbf{f})\mathbf{v}]_s$$

(όπου  $\mathbf{v} = (v(1), v(2), \dots, v(N))^T$ )

### Θεώρημα 2.2.1

Θεωρούμε το στοχαστικό παίγνιο με προεξόφληση με 1-ελεγκτή  $\Gamma_\beta(1)$  και το ζεύγος  $(P_\beta(1))$  και  $(D_\beta(1))$  γραμμικών προγραμμάτων. Ακόμη, έστω  $(\mathbf{v}^0, \mathbf{g}^0)$  μια βέλτιστη λύση του  $(P_\beta(1))$  και  $(\mathbf{z}^0, \mathbf{x}^0)$  μια βέλτιστη λύση του  $(D_\beta(1))$ . Τότε:

- (i) Το διάνυσμα αξίας του  $\Gamma_\beta(1)$  είναι  $\mathbf{v}^0$  και η  $\mathbf{g}^0$  είναι μία βέλτιστη στρατηγική για τον παίκτη 2
- (ii) Εάν θέσουμε:

$$x_s^o \doteq \sum_{a^1 \in A^1(s)} x_{sa^1}^o$$

και μία στάσιμη στρατηγική  $f^0$  για τον παίκτη 1 ορίζεται ως:

$$f^0(s, a^1) = \frac{x_{sa^1}^0}{x_s^0}, \quad s \in S, a^1 \in A^1(s)$$

τότε η  $f^0$  είναι μια βέλτιστη στρατηγική του παίκτη 1

### Παράδειγμα 2.2.2

Έστω  $S = \{1,2\}$ ,  $A^1(s) = A^2(s) = \{1,2\}$  για  $s = 1,2$ ,  $\beta = 0.6$  και οι αμοιβές και οι πιθανότητες μετάβασης είναι:

5	-2
(0.5,0.5)	(0.5,0.5)
-3	6
(0.7,0.3)	(0.7,0.3)

Κατάσταση 1

-2	7
(0.3,0.7)	(0.3,0.7)
3	-11
(0.8,0.2)	(0.8,0.2)

Κατάσταση 2

Από τη δομή των πιθανοτήτων μετάβασης μεταξύ των καταστάσεων, γίνεται φανερό πως το παραπάνω παίγνιο ελέγχεται από τον παίκτη 1. Το πρωτεύον γραμμικό πρόγραμμα  $P_{0.6}(1)$  θα έχει την ακόλουθη μορφή:

$$\min \left[ \frac{1}{2} v(1) + \frac{1}{2} v(2) \right]$$

με περιορισμούς:

$$v(1) \geq 5g(1,1) - 2g(1,2) + 0.3v(1) + 0.3v(2)$$

$$v(1) \geq -3g(1,1) + 6g(1,2) + 0.42v(1) + 0.18v(2)$$

$$v(2) \geq -2g(2,1) + 7g(2,2) + 0.18v(1) + 0.42v(2)$$

$$v(2) \geq 3g(2,1) - 11g(2,2) + 0.48v(1) + 0.12v(2)$$

$$g(1,1) + g(1,2) = 1$$

$$g(2,1) + g(2,2) = 1$$

$$g(1,1), g(1,2), g(2,1), g(2,2) \geq 0$$

Χρησιμοποιώντας το λογισμικό MATLAB, παίρνουμε τα ακόλουθα αποτελέσματα:

Αντικειμενική συνάρτηση	1.9225
$v(1)$	2.7393
$v(2)$	1.1058
$g(1,1)$	0.5123
$g(1,2)$	0.4877
$g(2,1)$	0.7613
$g(2,2)$	0.2387

Οπότε συμπεραίνουμε ότι :

$$v_{0.6}^0 = (2.7393, 1.1058) \text{ και } g^0 = ((0.5123, 0.4877), (0.7613, 0.2387))$$

Το δυϊκό πρόβλημα θα έχει τη μορφή:

$$\max[z(1) + z(2)]$$

με περιορισμούς:

$$0.7x_{11} + 0.58x_{12} - 0.18x_{21} - 0.48x_{22} = \frac{1}{2}$$

$$-0.3x_{11} - 0.18x_{12} + 0.58x_{21} + 0.88x_{22} = \frac{1}{2}$$

$$z(1) \leq 5x_{11} - 3x_{12}$$

$$z(1) \leq -2x_{11} + 6x_{12}$$

$$z(2) \leq -2x_{21} + 3x_{22}$$

$$z(2) \leq 7x_{21} - 11x_{22}$$

$$x_{11}, x_{12}, x_{21}, x_{22} \geq 0$$

Χρησιμοποιώντας και πάλι το MATLAB, λαμβάνουμε τα αποτελέσματα:

Αντικειμενική συνάρτηση	1.9225
$z(1)$	1.9740
$z(2)$	-0.0515
$x_{11}$	0.7403
$x_{12}$	0.5758
$x_{21}$	0.7207
$x_{22}$	0.4633

Σύμφωνα με το (ii) του Θεωρήματος 2.2.1, η βέλτιστη στρατηγική  $f^0$  θα είναι:

$$f^0 = \left( \left( \frac{x_{11}^0}{x_{11}^0 + x_{12}^0}, \frac{x_{12}^0}{x_{11}^0 + x_{12}^0} \right), \left( \frac{x_{21}^0}{x_{21}^0 + x_{22}^0}, \frac{x_{22}^0}{x_{21}^0 + x_{22}^0} \right) \right) \\ = ((0.5625, 0.4375), (0.6087, 0.3913))$$

### 2.2.2. Παίγνια με προεξόφληση διαχωρίσιμης αμοιβής (SEparable Reward) και ανεξάρτητα της κατάστασης μετάβασης (State Independent Transition) (SER-SIT)

Σ' αυτή την κλάση παιχνιδιών, θα κάνουμε ορισμένες υποθέσεις που αφορούν στη δομή του παιχνιδιού, ούτως ώστε να μπορέσουμε να αξιοποιήσουμε αποτελέσματα από τη θεωρία του γραμμικού προγραμματισμού και να φτάσουμε στον υπολογισμό του διανύσματος-τιμής καθώς και βέλτιστων στρατηγικών των παικτών.

Συγκεκριμένα, η δομή του παιχνιδιού θέλουμε να έχει την ακόλουθη μορφή:

$$r(s, a^1, a^2) = c(s) + \rho(a^1, a^2), \quad s \in S, a^1 \in A^1(s), a^2 \in A^2(s) \quad (SER)$$

$$p(s'|s, a^1, a^2) = p(s'|a^1, a^2), \quad s \in S, a^1 \in A^1(s), a^2 \in A^2(s) \quad (SIT)$$

Προφανώς, για να έχει νόημα η 2<sup>η</sup> υπόθεση (SIT), θα πρέπει υποχρεωτικά να ισχύει ότι  $m^1(s) = \mu$  και  $m^2(s) = \nu$  για όλα τα  $s \in S$ . Έτσι, ένα ζεύγος ενέργειων  $a^1, a^2$  θα καθορίζει την ίδια μετάβαση, ανεξαρτήτως από την παρούσα κατάσταση.

Η 1<sup>η</sup> υπόθεση (SER) καθιστά όλες τις αμοιβές των παικτών να προκύπτουν από το άθροισμα μιας ποσότητας που εξαρτάται μόνο από την παρούσα κατάσταση ( $c(s)$ ) και μιας άλλης που εξαρτάται μόνο από το ζεύγος των ενεργειών που επιλέχθηκαν ( $\rho(a^1, a^2)$ ).

Θα συμβολίζουμε :  $\mathbf{c}^T = (c(1), c(2), \dots, c(N))$

Ο συλλογισμός που θα ακολουθήσουμε για την επίλυση παιγνίων τύπου SER-SIT έχει ως εξής: Αντιστοιχίζουμε στο διάνυσμα  $\mathbf{c}$  το βοηθητικό πινακοπαιχνίδι:

$$R(\mathbf{c}) \doteq \left[ \rho(a^1, a^2) + \beta \sum_{s' \in S} p(s' | a^1, a^2) c(s') \right]_{a^1=1, a^2=1}^{\mu, \nu}$$

Έστω  $\rho \doteq \text{val}[R(\mathbf{c})]$  και  $\mathbf{x}^0 = (x_1^0, x_2^0, \dots, x_\mu^0)$ ,  $\mathbf{y}^0 = (y_1^0, y_2^0, \dots, y_\nu^0)$  ένα ζεύγος από βέλτιστες στρατηγικές για το παιχνίδι  $R(\mathbf{c})$ .

Θεωρούμε τώρα τη στρατηγική  $\mathbf{g}^0 \in G_s$  (αντίστοιχα  $\mathbf{f}^0 \in F_s$ ) η οποία ορίζεται ως:  $\mathbf{g}^0(s) = \mathbf{y}^0$  (αντίστοιχα  $\mathbf{f}^0(s) = \mathbf{x}^0$ ) για  $s \in S$ .

Οι στρατηγικές αυτές είναι βέλτιστες για το SER-SIT παίγνιο και μάλιστα θα δείξουμε ότι ισχύει:

$$\mathbf{v}_\beta = \mathbf{c} + \left( \frac{\rho}{1 - \beta} \right) \mathbf{1} \quad (2.2.1)$$

Αρχικά, θα αναδιατυπώσουμε το παιχνίδι χρησιμοποιώντας τη σχέση (1.1.4). Θα έχουμε, δηλαδή, ότι :

$$\mathbf{v}_\beta(\mathbf{f}, \mathbf{g}) = \mathbf{r}(\mathbf{f}, \mathbf{g}) + \beta \mathbf{P}(\mathbf{f}, \mathbf{g}) \mathbf{v}_\beta(\mathbf{f}, \mathbf{g})$$

Λόγω των περιορισμών SER και SIT , για κάθε  $s \in S$  θα έχουμε:

$$\begin{aligned} v_\beta(s, \mathbf{f}, \mathbf{g}) &= r(s, \mathbf{f}, \mathbf{g}) + \beta \sum_{s' \in S} p(s', \mathbf{f}, \mathbf{g}) v_\beta(s, \mathbf{f}, \mathbf{g}) \\ &= \sum_{a^1} \sum_{a^2} r(s, a^1, a^2) f(s, a^1) g(s, a^2) + \beta \sum_{s' \in S} p(s', \mathbf{f}, \mathbf{g}) v_\beta(s, \mathbf{f}, \mathbf{g}) \\ &= \sum_{a^1} \sum_{a^2} (c(s) + \rho(a^1, a^2)) f(s, a^1) g(s, a^2) + \beta \sum_{s' \in S} p(s', \mathbf{f}, \mathbf{g}) v_\beta(s, \mathbf{f}, \mathbf{g}) \\ &= c(s) \sum_{a^1} \sum_{a^2} f(s, a^1) g(s, a^2) + \sum_{a^1} \sum_{a^2} \rho(a^1, a^2) f(s, a^1) g(s, a^2) \\ &\quad + \beta \sum_{s' \in S} p(s', \mathbf{f}, \mathbf{g}) v_\beta(s, \mathbf{f}, \mathbf{g}) \end{aligned}$$

$$\begin{aligned}
&= c(s) + \sum_{a^1} \sum_{a^2} \left( [R(c)]_{a^1, a^2} - \beta \sum_{s' \in S} p(s' | a^1, a^2) c(s') \right) f(s, a^1) g(s, a^2) \\
&\quad + \beta \sum_{s' \in S} p(s', \mathbf{f}, \mathbf{g}) v_\beta(s, \mathbf{f}, \mathbf{g}) \\
&= c(s) + \mathbf{f}(s) R(c) \mathbf{g}(s) - \beta \sum_{a^1} f(s, a^1) \sum_{a^2} g(s, a^2) \sum_{s' \in S} p(s' | a^1, a^2) c(s') \\
&\quad + \beta \sum_{s' \in S} p(s', \mathbf{f}, \mathbf{g}) v_\beta(s, \mathbf{f}, \mathbf{g}) \\
&= c(s) + \mathbf{f}(s) R(c) \mathbf{g}(s) + \beta \sum_{s' \in S} p(s' | \mathbf{f}, \mathbf{g}) (v_\beta(s', \mathbf{f}, \mathbf{g}) - c(s'))
\end{aligned}$$

Εάν θέσουμε  $r(\mathbf{c}, \mathbf{f}, \mathbf{g}) \doteq \mathbf{f}(s) R(\mathbf{c}) \mathbf{g}(s)$ , για  $s \in S$ , τότε η παραπάνω σχέση σε διανυσματική μορφή γίνεται:

$$\mathbf{v}_\beta(\mathbf{f}, \mathbf{g}) = \mathbf{c} + r(\mathbf{c}, \mathbf{f}, \mathbf{g}) \mathbf{1} + \beta \mathbf{P}(\mathbf{f}, \mathbf{g}) (\mathbf{v}_\beta(\mathbf{f}, \mathbf{g}) - \mathbf{c})$$

και λύνοντας ως προς  $\mathbf{v}_\beta(\mathbf{f}, \mathbf{g})$ , λαμβάνουμε:

$$\begin{aligned}
&(\mathbf{I} - \beta \mathbf{P}(\mathbf{f}, \mathbf{g})) \mathbf{v}_\beta(\mathbf{f}, \mathbf{g}) = (\mathbf{I} - \beta \mathbf{P}(\mathbf{f}, \mathbf{g})) \mathbf{c} + r(\mathbf{c}, \mathbf{f}, \mathbf{g}) \mathbf{1} \\
&\Rightarrow \mathbf{v}_\beta(\mathbf{f}, \mathbf{g}) = \mathbf{c} + [\mathbf{I} - \beta \mathbf{P}(\mathbf{f}, \mathbf{g})]^{-1} [r(\mathbf{c}, \mathbf{f}, \mathbf{g}) \mathbf{1}] \quad (2.2.2)
\end{aligned}$$

Όπως είδαμε και στο 1<sup>ο</sup> κεφάλαιο, ισχύει ότι :

$$[\mathbf{I} - \beta \mathbf{P}(\mathbf{f}, \mathbf{g})]^{-1} = \sum_{t=0}^{\infty} \beta^t \mathbf{P}^t(\mathbf{f}, \mathbf{g})$$

Όμως, λόγω της δομής SER-SIT, θα έχουμε ότι για κάθε  $t$ ,  $\mathbf{P}^t(\mathbf{f}, \mathbf{g}) = \mathbf{P}(\mathbf{f}, \mathbf{g})$ , οπότε:

$$[\mathbf{I} - \beta \mathbf{P}(\mathbf{f}, \mathbf{g})]^{-1} = \mathbf{P}(\mathbf{f}, \mathbf{g}) \sum_{t=0}^{\infty} \beta^t = \frac{1}{1-\beta} \mathbf{P}(\mathbf{f}, \mathbf{g})$$

Συνεπώς, έχουμε καταλήξει στη σχέση:

$$\mathbf{v}_\beta(\mathbf{f}, \mathbf{g}) = \mathbf{c} + \frac{1}{1-\beta} \mathbf{P}(\mathbf{f}, \mathbf{g}) [r(\mathbf{c}, \mathbf{f}, \mathbf{g}) \mathbf{1}]$$

Οπότε αν θεωρήσουμε το ζεύγος βέλτιστων στρατηγικών  $(\mathbf{f}, \mathbf{g})$  τέτοιες ώστε  $\mathbf{f}(s) = \mathbf{x}^T$  και  $\mathbf{g}(s) = \mathbf{y}$  για όλα τα  $s \in S$ , και επειδή ο πίνακας  $\mathbf{P}$  είναι στοχαστικός, παίρνουμε :



$$v_\beta(f, g) = c + \frac{r(c, f, g)}{1 - \beta} \mathbf{1} = c + \frac{x^T R(c) y}{1 - \beta} \mathbf{1} = c + \left( \frac{\rho}{1 - \beta} \right) \mathbf{1}$$

και έτσι, έχουμε καταλήξει στο ζητούμενο.

### 2.2.3 Παίγνια με προεξόφληση εναλλασσόμενου ελεγκτή

Η κατηγορία αυτή αποτελεί την 3<sup>η</sup> και τελευταία κλάση παιγνίων με συντελεστή προεξόφλησης που θα μας απασχολήσει σε αυτή την εργασία. Πρόκειται για γενίκευση των παιχνιδιών με έναν ελεγκτή που αναλύθηκε προηγουμένως. Για αυτά τα παίγνια, οι πιθανότητες μετάβασης μεταξύ των καταστάσεων εξαρτώνται από τις ενέργειες του ενός παίκτη για ένα σύνολο καταστάσεων, και από τις ενέργειες του άλλου παίκτη, για τις υπόλοιπες καταστάσεις.

Για το λόγο αυτό, καλούμαστε να διαμερίσουμε το σύνολο καταστάσεων  $S$ . Θα έχουμε δηλαδή:

$S = S^1 \cup S^2$ , όπου τα  $S^1, S^2$  είναι μη-κενά και ξένα μεταξύ τους σύνολα, έτσι ώστε για κάθε  $a^1 \in A^1(s), a^2 \in A^2(s)$ , να ισχύει ότι :

$$p(s'|s, a^1, a^2) = \begin{cases} p(s'|s, a^1), & s \in S^1 \\ p(s'|s, a^2), & s \in S^2 \end{cases} \quad (SW)$$

Προκειμένου να αντιμετωπιστεί ένα τέτοιο παίγνιο εναλλασσόμενου ελεγκτή μέσω γραμμικού προγραμματισμού, θα χρειαστεί πρώτα να λυθεί ένα (πεπερασμένο) σύνολο γραμμικών προγραμμάτων, σε αντίθεση με την κατηγορία του ενός ελεγκτή όπου η λύση ενός γραμμικού προγράμματος ήταν αρκετή.

Προτού δοθεί ο σχετικός αλγόριθμος, αξίζει να παρατηρήσουμε μια ιδιότητα των παιγνίων τύπου  $SW$ . Έστω ότι ο παίκτης 1 σταθεροποιεί την στρατηγική  $\hat{f}(s)$  για κάθε  $s \in S^1$ . Τότε, μπορούμε να ορίσουμε ένα παίγνιο  $\Gamma_\beta(2, \hat{f})$  με ελεγκτή τον παίκτη 2, χώρο καταστάσεων και ενεργειών ίδιο με του  $\Gamma_\beta$ , αμοιβές και πιθανότητες μετάβασης ίδιες για τα  $s \in S^2$ , αλλά τροποποιημένες για τα  $s \in S^1$ , ως εξής :

$$\hat{r}(s, a^1, a^2) \doteq \sum_{a^1 \in A^1(s)} r(s, a^1, a^2) \hat{f}(s, a^1)$$

$$\hat{p}(s'|s, a^1, a^2) \doteq \sum_{a^1 \in A^1(s)} p(s'|s, a^1, a^2) \hat{f}(s, a^1)$$

για όλα τα  $a^1 \in A^1(s), a^2 \in A^2(s), s \in S^1, s' \in S^2$

Φυσικά, για να ισχύουν τα παραπάνω, θα πρέπει όλες οι ενέργειες του παίκτη 1 στο  $\Gamma_\beta(2, \hat{f})$ , για  $s \in S^1$ , να ταυτίζονται · όσον αφορά τόσο τις μεταβάσεις όσο και τις αμοιβές, ακριβώς όπως φαίνεται στο επόμενο παράδειγμα.

### Παράδειγμα 2.2.3.

Έστω  $S = \{1,2\}$ ,  $A^1(1) = A^1(2) = A^2(2) = \{1,2\}$ ,  $A^2(1) = \{1,2,3\}$

ενώ οι αμοιβές ενός βήματος και οι μεταβάσεις φαίνονται στο παρακάτω διάγραμμα.

1	2	3
(0.6,0.4)	(0.6,0.4)	(0.6,0.4)
2	3	4
(0,1)	(0,1)	(0,1)

Κατάσταση 1

1	2
(1,0)	(0,1)
2	4
(1,0)	(0,1)

Κατάσταση 2

Βλέπουμε πως ο παίκτης 1 ελέγχει τις μεταβάσεις στο σύνολο  $S^1 = \{1\}$ , ενώ ο παίκτης 2 στο σύνολο  $S^2 = \{2\}$ . Ακόμη, εάν ο παίκτης 1 σταθεροποιήσει μία στρατηγική για την κατάσταση 1, έστω την  $\hat{f}(1) = (p, 1 - p)$ , τότε σύμφωνα με τα προηγούμενα, μπορούμε να θεωρήσουμε το ακόλουθο παιχνίδι  $\Gamma_\beta(2, \hat{f}(1))$ , με ελεγκτή τον παίκτη 2 :

$2 - p$	$3 - p$	$4 - p$
$(0.6p, 1 - 0.6p)$	$(0.6p, 1 - 0.6p)$	$(0.6p, 1 - 0.6p)$
$2 - p$	$3 - p$	$4 - p$
$(0.6p, 1 - 0.6p)$	$(0.6p, 1 - 0.6p)$	$(0.6p, 1 - 0.6p)$

Κατάσταση 1

1	2
(1,0)	(0,1)
2	4
(1,0)	(0,1)

Κατάσταση 2

Ο αλγόριθμος που θα χρησιμοποιήσουμε προκειμένου να φτάσουμε στη λύση ενός παιχνιδιού τύπου SW, όπως κατασκευάστηκε παραπάνω, είναι ο ακόλουθος:

### Αλγόριθμος 2.2.1.

**Βήμα 1:** Θέσε  $k = 0$ , επίλεξε ένα αυθαίρετο  $\mathbf{v}^0 = (v^0(1), \dots, v^0(N))^T$  και βρες μια βέλτιστη στρατηγική  $\mathbf{f}^0(s)$  για τον παίκτη 1 στο πινακοπαιχνίδι  $R(s, \mathbf{v}^0)$ , για κάθε  $s \in S^1$ .

**Βήμα 2:** Θέσε  $k = k + 1$ . Λύσε το παίγνιο  $\Gamma_\beta(2, \hat{\mathbf{f}}^{k-1})$  με ελεγκτή τον παίκτη 2. Η τιμή του θα συμβολίζεται με  $\mathbf{v}_\beta$ .

Θέσε  $\mathbf{v}^k = \mathbf{v}_\beta$ .

**Βήμα 3:** Αν  $v^k(s) = \text{val}[R(s, \mathbf{v}^k)]$ , για κάθε  $s \in S$ , τότε σταματάμε.

Αλλιώς, βρες μια βέλτιστη στρατηγική  $\mathbf{f}^k(s)$  για τον παίκτη 1 στο πινακοπαιχνίδι  $R(s, \mathbf{v}^k)$ , για κάθε  $s \in S^1$  και επίστρεψε στο Βήμα 2.

## 2.3 Στοχαστικά παίγνια Οριακού Μέσου

Στην ανάλυση που έγινε στο προηγούμενο κεφάλαιο, πάνω στις Μαρκοβιανές διαδικασίες απόφασης, είδαμε πως το κριτήριο πληρωμής οριακού μέσου θεωρείται ‘πιο δύσκολο’ στην αντιμετώπιση, συγκριτικά με το κριτήριο προεξόφλησης. Παρ’ όλα αυτά, ξεπερνώντας ορισμένες τεχνικές δυσκολίες, καταλήξαμε στη διατύπωση σημαντικών αποτελεσμάτων και για τα δύο κριτήρια απόδοσης, όπως είναι η ύπαρξη βέλτιστων στάσιμων στρατηγικών, η εύρεση των οποίων προκύπτει να είναι η λύση σε κατάλληλα κατασκευασμένα γραμμικά προγράμματα.

Με αυτό το σκεπτικό, θα μπορούσε κανείς να εικάσει πως, με ίσως ορισμένες τεχνικές δυσκολίες, θα μπορούσαμε να γενικεύσουμε τα αποτελέσματα των Μαρκοβιανών διαδικασιών στα στοχαστικά παίγνια.

Δυστυχώς όμως, μια τέτοια γενίκευση δεν είναι εφικτή και αυτό θα γίνει φανερό μέσα από ένα πολύ γνωστό αντιπαράδειγμα, το οποίο απαντάται στη βιβλιογραφία ως ‘The Big Match’. Προτού περάσουμε στην ανάλυση του Big Match, ας ορίσουμε το στοχαστικό παίγνιο με κριτήριο πληρωμής τον οριακό μέσο,  $\Gamma_\alpha$ , περιορισμένοι στο σύνολο των στάσιμων στρατηγικών, με τρόπο ανάλογο της πληρωμής με προεξόφληση.

Η διαφοροποίηση των δύο παιγνίων έγκειται ουσιαστικά στον ορισμό της συνάρτησης απόδοσης των παικτών. Ορίζουμε, λοιπόν, την πληρωμή που θα λάβει ο παίκτης 1 από τον παίκτη 2, που αντιστοιχεί σε ζεύγος στρατηγικών  $(\mathbf{f}, \mathbf{g}) \in F^S \times G^S$  ως εξής:

$$v_a(s, \mathbf{f}, \mathbf{g}) \doteq \lim_{T \rightarrow \infty} \left[ \left( \frac{1}{T+1} \right) \sum_{t=0}^T E_{s, \mathbf{f}, \mathbf{g}}[R_t] \right] = [\mathbf{Q}(\mathbf{f}, \mathbf{g}) \mathbf{r}(\mathbf{f}, \mathbf{g})]_s$$

για κάθε  $s \in S$ , όπου  $\mathbf{Q}(\mathbf{f}, \mathbf{g})$  είναι το κατά Cesaro όριο του πίνακα  $\mathbf{P}(\mathbf{f}, \mathbf{g})$ .

Φυσικά, θα λέμε ότι  $\mathbf{f}_0 \in F_s$  και  $\mathbf{g}_0 \in G_s$  είναι βέλτιστες, στάσιμες στρατηγικές εάν , για κάθε  $s \in S, \mathbf{f} \in F_s, \mathbf{g} \in G_s$  , ισχύει ότι :

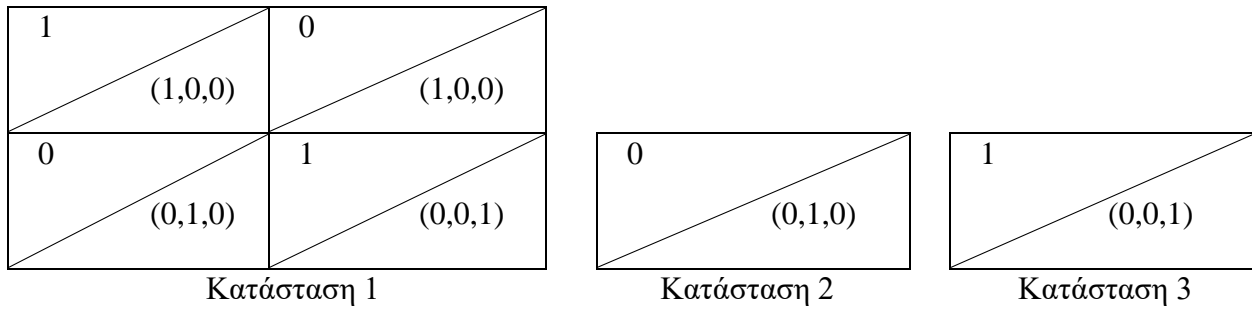
$$v_a(s, \mathbf{f}, \mathbf{g}_0) \leq v_a(s, \mathbf{f}_0, \mathbf{g}_0) \leq v_a(s, \mathbf{f}_0, \mathbf{g})$$

Εάν υπάρχει ένα τέτοιο ζεύγος στρατηγικών  $(\mathbf{f}_0, \mathbf{g}_0)$ , τότε υπάρχει και η μη-προεξοφλητική τιμή του παιχνιδιού  $\Gamma_\alpha$ , η οποία ορίζεται να είναι:

$$\mathbf{v}_\alpha \doteq \mathbf{v}_\alpha(\mathbf{f}_0, \mathbf{g}_0)$$

### Παράδειγμα 2.3.1. (The Big Match)

Έστω  $S = \{1,2,3\}$ ,  $A^1(1) = A^2(1) = \{1,2\}$ ,  $A^1(s) = A^2(s) = \{1\}$ , για  $s = 2,3$  και οι αμοιβές και πιθανότητες μετάβασης δίνονται στο ακόλουθο σχήμα:

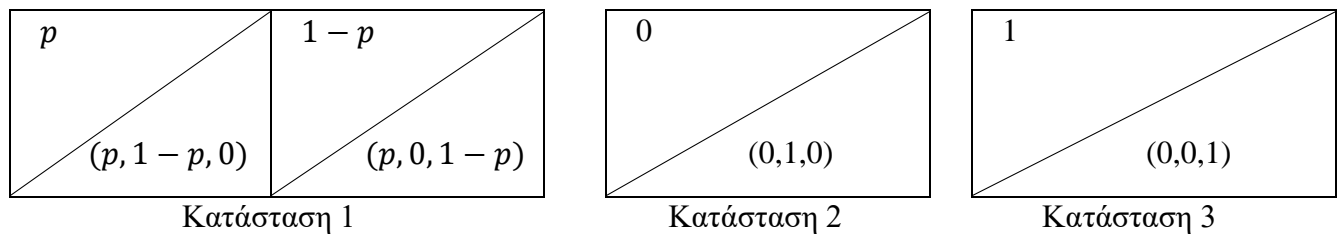


Βλέπουμε πως οι καταστάσεις 2 και 3 είναι απορροφητικές, συνεπώς  $v_a(2) = 0$  και  $v_a(3) = 1$

Έστω ότι ο παίκτης 1 διαθέτει μία βέλτιστη στάσιμη στρατηγική:

$$\mathbf{f}_p = ((p, 1 - p), (1), (1))$$

για κάποιο σταθερό  $p \in [0,1]$ . Έχοντας να αντιμετωπίσει αυτή τη στρατηγική, ο παίκτης 2 στην ουσία προσπαθεί να ελαχιστοποιήσει το κόστος για την Μαρκοβιανή διαδικασία απόφασης:



Διαχωρίζουμε 2 περιπτώσεις:

- A.  $p = 1$ , δηλαδή ο παίκτης 1 δεν ρισκάρει και επιλέγει να παραμείνει στην κατάσταση 1. Ωστόσο, αυτό σημαίνει ότι εναντίον της στρατηγικής  $\mathbf{g}_0 = ((0,1), (1), (1))^T$ , ο παίκτης 1 θα λαμβάνει σχεδόν πάντα 0, και έτσι,  $v_a(1, \mathbf{f}_1, \mathbf{g}_0) = 0$ .
- B.  $0 \leq p < 1$ , δηλαδή ο παίκτης 1 παίρνει το ρίσκο να επιλέξει την 2<sup>η</sup> ενέργεια στην κατάσταση 1 με πιθανότητα  $1 - p > 0$  κάθε φορά που η διαδικασία επισκέπτεται την 1<sup>η</sup> κατάσταση. Ωστόσο, σε αυτή την περίπτωση ο παίκτης 2 μπορεί να χρησιμοποιήσει την στρατηγική  $\mathbf{g}_1 = ((1,0), (1), (1))^T$  και έτσι η διαδικασία θα εγκλωβιστεί στην κατάσταση 2 με πιθανότητα 1. Οπότε και εδώ πάλι θα έχουμε ότι  $v_a(1, \mathbf{f}_p, \mathbf{g}_1) = 0$ .

Δηλαδή, για όλα τα  $p \in [0,1]$  μπορούμε να συμπεράνουμε πως

$$\min_{\mathbf{g} \in G_s} v_a(1, \mathbf{f}_p, \mathbf{g}) = 0$$

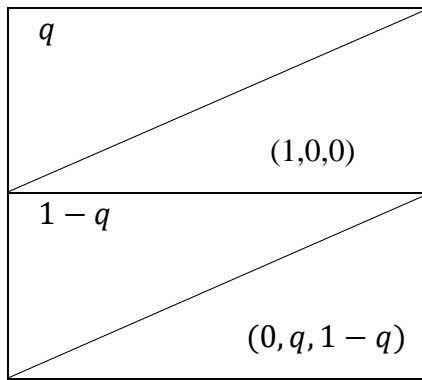
Απ' την άλλη μεριά, εάν ο παίκτης 2 χρησιμοποιήσει την στρατηγική  $\mathbf{g}^* = \left(\left(\frac{1}{2}, \frac{1}{2}\right), (1), (1)\right)^T$ , βλέπουμε πως, ανεξάρτητα από την επιλογή του παίκτη 1 στην κατάσταση 1, θα έχουμε ότι:

$$v_a(1, \mathbf{f}, \mathbf{g}^*) = \frac{1}{2}$$

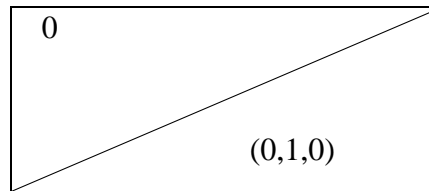
Ακόμη, παρατηρούμε ότι κάθε στάσιμη στρατηγική του παίκτη 2 μπορεί να γραφεί στη μορφή:

$$\mathbf{g}_q = ((q, 1 - q), (1), (1))^T, \quad q \in [0,1]$$

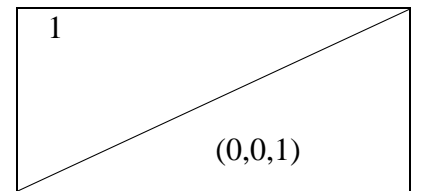
Φυσικά, όπως και πριν, για μια σταθεροποιημένη στρατηγική  $\mathbf{g}_q$ , ο παίκτης 1 καλείται να αντιμετωπίσει την διαδικασία απόφασης:



Κατάσταση 1



Κατάσταση 2

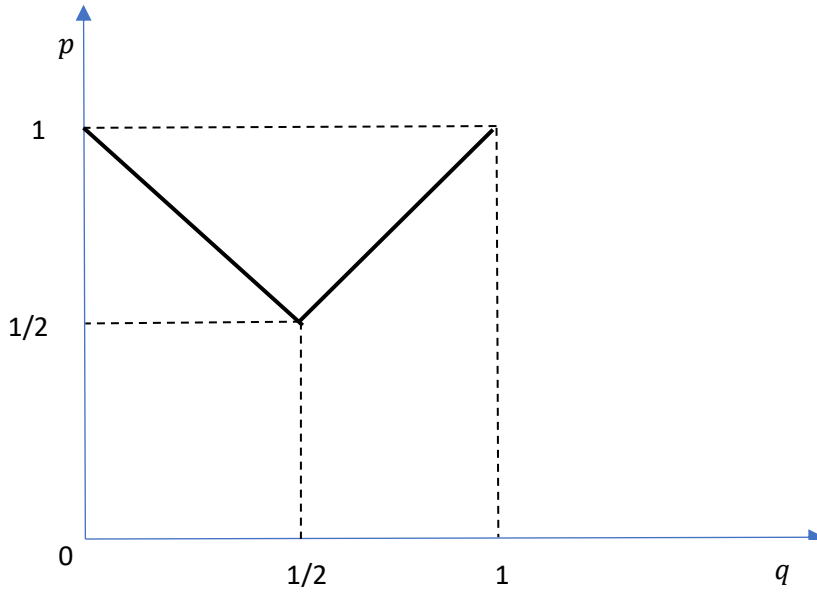


Κατάσταση 3

Παρατηρούμε πως, εάν ο παίκτης 1 χρησιμοποιεί την στρατηγική  $f_p$  με  $p < 1$ , τότε η απορρόφηση στις καταστάσεις 2 και 3 θα συμβεί με πιθανότητα  $q$  και  $1 - q$ , αντίστοιχα. Εάν  $p = 1$ , τότε η κατάσταση 1 θα επαναλαμβάνεται συνεχώς. Συνεπώς, λαμβάνουμε ότι:

$$v_a(1, f_p, g_q) = \begin{cases} q, & \text{εάν } p = 1 \\ 1 - q, & \text{εάν } p < 1 \end{cases}$$

Οπότε, το γράφημα της  $\max_{f \in F_s} v_a(1, f, g_q)$  φαίνεται παρακάτω:



Το ελάχιστο της παραπάνω συνάρτησης έχει τιμή  $\frac{1}{2}$ . Από την έως τώρα ανάλυση έχουμε στην ουσία δείξει ότι :

$$0 = \max_{f \in F_s} \min_{g \in G_s} v_a(1, f, g) < \frac{1}{2} = \min_{g \in G_s} \max_{f \in F_s} v_a(1, f, g)$$

Ακριβώς λόγω αυτής της αυστηρής ανισοτικής σχέσης, μπορούμε να συμπεράνουμε πως δεν υπάρχουν βέλτιστες στάσιμες στρατηγικές για το Big Match.

Ακόμη και αν θέλαμε να αναζητήσουμε ένα ζεύγος βέλτιστων στρατηγικών στο, εν γένει, ευρύτερο σύνολο των Μαρκοβιανών στρατηγικών, θα βλέπαμε ότι ούτε και αυτό είναι εφικτό.

Πράγματι, έστω  $\pi^1 = (f_0, f_1, \dots, f_t, \dots) \in F_M$  και  $\pi^2 = (g_0, g_1, \dots, g_t, \dots) \in G_M$  ένα ζεύγος Μαρκοβιανών στρατηγικών των παικτών, αντίστοιχα.

Εάν η πιθανότητα για τον παίκτη 1 να επιλέξει την 2<sup>η</sup> ενέργεια είναι ίση με 1, τότε παίρνοντας  $g_t = (1, 0)$  για όλα τα  $t$ , θα έχουμε ότι  $v_a(1, \pi^1, \pi^2) = 0$

Εάν όχι, τότε έστω  $m \geq 0$  ο μικρότερος αριθμός από τον οποίο και έπειτα ο παίκτης 1 επιλέγει την 2<sup>η</sup> ενέργεια με θετική πιθανότητα, έστω  $\varepsilon > 0$ . Τώρα, ορίζουμε  $\pi_\varepsilon^2$  να είναι η στρατηγική κατά την οποία επιλέγεται η  $\mathbf{g} = (1,0)$  για τα πρώτα  $m$  στάδια, έπειτα  $\mathbf{g}_{m+1} = (0,1)$  και  $\mathbf{g} = (\frac{1}{2}, \frac{1}{2})$  από εκεί και πέρα. Θα έχουμε, τότε, ότι  $v_a(1, \pi^1, \pi_\varepsilon^2) = \frac{(1-\varepsilon)}{2}$ .

Από την άλλη, παίρνοντας  $\mathbf{g}_t = (\frac{1}{2}, \frac{1}{2})$  για όλα τα  $t$ , μας δίνει  $v_a(1, \pi^1, \pi^2) = \frac{1}{2}$  για όλα τα  $\pi^1 \in F_M$ . Εάν θεωρήσουμε την  $\mathbf{g}_t = (q, 1 - q)$  όπου  $q \leq \frac{1}{2}$  για όλα τα  $t$ , τότε η  $\pi^1$  για την οποία  $\mathbf{f}_t = (1,0)$  για όλα τα  $t$ , μας δίνει ότι  $v_a(1, \pi^1, \pi^2) \geq \frac{1}{2}$ . Εάν τώρα σε κάποιο στάδιο  $n$  έχουμε ότι  $q > \frac{1}{2}$  για πρώτη φορά, επιλέγοντας την στρατηγική  $\pi^1$  όπου  $\mathbf{f}_t = (1,0)$  για  $t < n$  και  $\mathbf{f}_n = (0,1)$ , δίνει ότι  $v_a(1, \pi^1, \pi^2) > \frac{1}{2}$ .

Η παραπάνω κατασκευαστική απόδειξη δείχνει πως ακόμα και στο σύνολο των Μαρκοβιανών στρατηγικών, δεν μπορούμε να βρούμε ένα ζεύγος βέλτιστων στρατηγικών για τους 2 παίκτες.

### Σχόλιο

Στην πραγματικότητα, από την παραπάνω ανάλυση συμπεραίνουμε πως η τιμή του Big Match είναι ίση με  $\frac{1}{2}$  και ότι ο 2<sup>ος</sup> παίκτης διαθέτει βέλτιστη στρατηγική, την  $\mathbf{g}_t = (\frac{1}{2}, \frac{1}{2})$ . Ωστόσο, όπως δείξαμε, δεν είναι εφικτό να βρεθεί στρατηγική του παίκτη 1 που να του εξασφαλίζει  $\frac{1}{2}$ , αλλά μόνο ότι μπορεί να φτάσει «ε-κοντά» στην τιμή αυτή.

Στη συνέχεια, θα μας απασχολήσει να δούμε πώς μπορούν να βρεθούν βέλτιστες (στάσιμες) στρατηγικές, για εκείνα τα παιχνίδια με πληρωμή οριακού μέσου, για τα οποία υφίστανται τέτοιες στρατηγικές.

## 2.4 Στοχαστικά Παίγνια οριακού μέσου με έναν ελεγκτή

Το ερώτημα στο οποίο καλούμαστε να απαντήσουμε στην τελευταία παράγραφο αυτού του κεφαλαίου είναι το εξής: «Υπάρχει κάποια υποκατηγορία παιχνίγων οριακού μέσου τα οποία επιδέχονται λύση με μεθόδους γραμμικού προγραμματισμού;». Εάν η απάντηση στο παραπάνω είναι θετική, τότε θα είμαστε σε θέση να αξιοποιήσουμε ήδη υπάρχοντες αλγόριθμους προκειμένου να φτάσουμε στη λύση των παιχνίγων αυτών.

Η αναμενόμενη ιδέα τώρα είναι να κάνουμε την υπόθεση πως μόνον ο ένας εκ των δύο παικτών ‘ελέγχει’ τις πιθανότητες μετάβασης, ακριβώς όπως πράξαμε και για την περίπτωση με προεξόφληση.

Σκοπός μας είναι να διατυπώσουμε έναν αλγόριθμο επίλυσης στοχαστικών παιγνίων οριακού μέσου με ελεγκτή τον παίκτη 1, καθώς θα γίνει και μια προσπάθεια επιβεβαίωσης της εγκυρότητάς του.

Αρχικά, υπενθυμίζουμε πως για ένα στοχαστικό παίγνιο οριακού μέσου  $\Gamma_\alpha$  και ένα ζεύγος στάσιμων στρατηγικών  $(f, g)$ , η τιμή του παιχνιδιού, για αυτό το ζεύγος, είναι:

$$v_\alpha(f, g) = Q(f, g)r(f, g)$$

όπου με  $Q(f, g)$  συμβολίζουμε το *Cesaro-όριο* του  $P(f, g)$ .

Το παίγνιο  $\Gamma_\alpha(1)$  προκύπτει με φυσιολογικό τρόπο από το  $\Gamma_\alpha$ , θέτοντας τον επιπλέον περιορισμό:

$$p(s'|s, a^1, a^2) \equiv p(s'|s, a^1), \quad \text{όπου } a^1 \in A^1(s), a^2 \in A^2(s), s, s' \in S$$

Προφανώς, ανάλογα μπορεί να οριστεί και το παίγνιο  $\Gamma_\alpha(2)$ .

Προτού αναπτύξουμε την θεωρία για τον σχετικό αλγόριθμο, θα χρειαστεί να παραθέσουμε τον εξής συμβολισμό:

### Συμβολισμός

Για κάθε κατάσταση  $s \in S$ , ορίζουμε έναν  $N \times m(s)$  μπλοκ πίνακα  $W_s$ , του οποίου το  $(s', (s, a))$ -οστό στοιχείο θα είναι:

$$w_{s'(s,a)} \doteq \delta_{ss'} - p(s'|s, a), \quad s' \in S, \quad a = 1, 2, \dots, m(s) \in A(s)$$

Για αυτό το  $s$ -οστό μπλοκ, ορίζουμε τα ακόλουθα διανύσματα στήλης, διάστασης  $m(s) \times 1$ :

$$\mathbf{x}_s = (x_{s1}, x_{s2}, \dots, x_{sm(s)})^T$$

$$\mathbf{y}_s = (y_{s1}, y_{s2}, \dots, y_{sm(s)})^T$$

$$\mathbf{r}_s = (r(s, 1), r(s, 2), \dots, r(s, m(s)))^T$$

$$\mathbf{1}_s = (1, 1, \dots, 1), \quad \mathbf{0}_s = (0, 0, \dots, 0)$$

και συνθέτοντας αυτά τα μπλοκ, παίρνουμε:



$$\mathbf{x}^T = (\mathbf{x}_1^T, \mathbf{x}_2^T, \dots, \mathbf{x}_N^T)$$

$$\mathbf{y}^T = (\mathbf{y}_1^T, \mathbf{y}_2^T, \dots, \mathbf{y}_N^T)$$

$$\mathbf{r}^T = (\mathbf{r}_1^T, \mathbf{r}_2^T, \dots, \mathbf{r}_N^T)$$

Ορίζουμε τους μπλοκ-πίνακες:

$$\mathbf{J}_1^T = (\mathbf{1}_1^T, \mathbf{0}_2^T, \dots, \mathbf{0}_N^T)$$

$$\mathbf{J}_2^T = (\mathbf{0}_1^T, \mathbf{1}_2^T, \dots, \mathbf{0}_N^T)$$

...

$$\mathbf{J}_N^T = (\mathbf{0}_1^T, \mathbf{0}_2^T, \dots, \mathbf{1}_N^T)$$

Επιπλέον, θα χρησιμοποιήσουμε τους  $N \times m$  -διάστατους πίνακες:

$$\mathbf{W} \doteq (\mathbf{W}_1 \vdots \mathbf{W}_2 \vdots \dots \vdots \mathbf{W}_N)$$

$$\mathbf{J} \doteq (\mathbf{J}_1 \vdots \mathbf{J}_2 \vdots \dots \vdots \mathbf{J}_N)^T$$

καθώς και τα διανύσματα-γραμμή διάστασης  $1 \times N$ :

$$\mathbf{v}^T = (v(1), \dots, v(N))$$

$$\mathbf{u}^T = (u(1), \dots, u(N))$$

$$\boldsymbol{\gamma}^T = (\gamma(1), \dots, \gamma(N))$$

όπου  $\gamma(s) > 0$  και  $\sum_{s=1}^N \gamma(s) = 1$

Εφόσον οι αμοιβές εξαρτώνται από τις ενέργειες και των δύο παικτών, θα ορίσουμε έναν μπλοκ-διαγώνιο πίνακα

$$\mathbf{R} = \text{diag}[R(1), R(2), \dots, R(N)]$$

όπου  $R(s) = [r(s, a^1, a^2)]_{a^1=1, a^2=1}^{m^1(s), m^2(s)}$

Ο πίνακας  $R$  έχει διάσταση  $m^1 \times m^2$ , όπου με  $m^k$  συμβολίζουμε το άθροισμα του πλήθους των ενεργειών του παίκτη  $k$  ( $k = 1$  ή  $k = 2$ ) για όλες τις καταστάσεις.

Για μια στάσιμη στρατηγική του παίκτη 2 θα έχουμε:

$$\mathbf{R}\mathbf{g} = \left[ (R(1)\mathbf{g}(1))^T, (R(2)\mathbf{g}(2))^T, \dots, (R(N)\mathbf{g}(N))^T \right]^T$$

το οποίο είναι ένα  $m^1 \times 1$  μπλοκ διάνυσμα στήλης.

Όμοια, για μια στάσιμη στρατηγική  $\mathbf{f}$  του παίκτη 1 θα έχουμε:

$$\mathbf{fR} = [\mathbf{f}(1)R(1), \mathbf{f}(2)R(2), \dots, \mathbf{f}(N)R(N)]$$

το οποίο είναι ένα  $1 \times m^2$  μπλοκ διάνυσμα γραμμής.

Είμαστε πλέον σε θέση να ορίσουμε ένα ζεύγος γραμμικών προγραμμάτων (Primal-Dual),  $P_a(1)$  και  $D_a(1)$ , το οποίο αντιστοιχεί στο στοχαστικό παίγνιο οριακού μέσου με έναν ελεγκτή,  $\Gamma_a(1)$ .

$$\min [\gamma^T \mathbf{v}]$$

με περιορισμούς:

$$P_a(1)$$

$$a) \quad (\mathbf{u}^T, \mathbf{v}^T, \mathbf{g}^T) \begin{pmatrix} \mathbf{W} & \vdots & \mathbf{0} \\ \dots & \vdots & \dots \\ \mathbf{J} & \vdots & \mathbf{W} \\ \dots & \vdots & \dots \\ -\mathbf{R}^T & \vdots & \mathbf{0} \end{pmatrix} \geq (\mathbf{0}^T, \mathbf{0}^T)$$

$$b) \quad \mathbf{1}^T \mathbf{g}(s) = 1, \quad s \in S$$

$$c) \quad \mathbf{g}(s) \geq \mathbf{0}, \quad s \in S$$

Για το δυϊκό πρόγραμμα, τα διανύσματα δυϊκών μεταβλητών  $\mathbf{x}, \mathbf{y}$  αντιστοιχούν στα δύο μπλοκ περιορισμών στο  $a)$  και η δυϊκή μεταβλητή  $\mathbf{z}$  αντιστοιχεί στους περιορισμούς του  $\beta)$ .

$$\max[\mathbf{1}^T \mathbf{z}]$$

με περιορισμούς:

$$D_a(1)$$

$$d) \quad \begin{pmatrix} \mathbf{W} & \vdots & \mathbf{0} \\ \dots & \vdots & \dots \\ \mathbf{J} & \vdots & \mathbf{W} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{r} \end{pmatrix}$$

$$e) \quad [-\mathbf{R}^T \mathbf{x}]_s + z_s \mathbf{1}_{m^2(s)} \leq \mathbf{0}_{m^2(s)}, \quad s \in S$$

$$f) \quad \mathbf{x} \geq \mathbf{0}, \mathbf{y} \geq \mathbf{0}$$

Η ποσότητα  $[-\mathbf{R}^T \mathbf{x}]_s = -R(s)^T \mathbf{x}(s)$  είναι ένα  $(m^2(s) \times 1)$  – διάσταστο διάνυσμα για κάθε  $s \in S$ .

Τέλος, παρουσιάζουμε τον αλγόριθμο επίλυσης του  $\Gamma_a(1)$  μέσω αυτού του ζεύγους γραμμικών προγραμμάτων.

### Αλγόριθμος 2.4.1.

Βήμα 1: Βρες μια βέλτιστη λύση  $(\hat{\mathbf{u}}^T, \hat{\mathbf{v}}^T, \hat{\mathbf{g}}^T)$  του  $P_a(1)$  και μια βέλτιστη λύση  $(\hat{\mathbf{x}}^T, \hat{\mathbf{y}}^T, \hat{\mathbf{z}}^T)$  του  $D_a(1)$ .

Βήμα 2: Όρισε το σύνολο καταστάσεων:

$$S^* \doteq \{s \in S \mid \hat{x}_s \doteq \sum_{a^1 \in A^1(s)} \hat{x}_{sa^1} > 0\}$$

Βήμα 3: Κατασκεύασε μια στάσιμη στρατηγική:

$$\hat{\mathbf{f}} = (\hat{\mathbf{f}}(1), \hat{\mathbf{f}}(2), \dots, \hat{\mathbf{f}}(N))$$

σύμφωνα με:

$$\hat{f}(s, a^1) = \begin{cases} \frac{\hat{x}_{sa^1}}{\hat{x}_s}, & s \in S^*, a^1 \in A^1(s) \\ \frac{\hat{y}_{sa^1}}{\hat{y}_s}, & s \in S \setminus S^*, a^1 \in A^1(s) \end{cases}$$

όπου:

$$\hat{y}_s \doteq \sum_{a^1 \in A^1(s)} \hat{y}_{sa^1}$$

Παράδειγμα 2.4.1.

Έστω  $S = \{1,2\}$ ,  $A^1(1) = A^2(2) = \{1,2\}$ ,  $A^1(2) = A^2(1) = \{1,2,3\}$ , ενώ οι αμοιβές και οι πιθανότητες μετάβασης φαίνονται στο ακόλουθο σχήμα:

-2	-6	0
(1,0)	(1,0)	(1,0)
-3	0	-2
(0,1)	(0,1)	(0,1)

Κατάσταση 1

1	-2
(1,0)	(1,0)
-4	-3
(0,1)	(0,1)
-5	0
(1,0)	(1,0)

Κατάσταση 2

Παίρνουμε το διάνυσμα  $\gamma^T = (\frac{1}{2}, \frac{1}{2})$

Τότε το πρωτεύον πρόβλημα ανάγεται στο :

$$\min \left[ \frac{1}{2} v(1) + \frac{1}{2} v(2) \right]$$

με περιορισμούς:

$P_a(1)$

$$a) \begin{pmatrix} 0 & 1 & -1 & 0 & -1 & \vdots & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 1 & \vdots & 0 & 0 & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 1 & 0 & 0 & 0 & \vdots & 0 & 1 & -1 & 0 & -1 \\ 0 & 0 & 1 & 1 & 1 & \vdots & 0 & -1 & 1 & 0 & 1 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 2 & 3 & 0 & 0 & 0 & \vdots & 0 & 0 & 0 & 0 & 0 \\ 6 & 0 & 0 & 0 & 0 & \vdots & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & \vdots & 0 & 0 & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \vdots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & -1 & 4 & 5 & \vdots & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 3 & 0 & \vdots & 0 & 0 & 0 & 0 & 0 \end{pmatrix}^T \begin{pmatrix} u(1) \\ u(2) \\ v(1) \\ v(2) \\ g(1,1) \\ g(1,2) \\ g(1,3) \\ g(2,1) \\ g(2,2) \end{pmatrix} \geq 0$$

$$b) \quad g(1,1) + g(1,2) + g(1,3) = 1$$

$$g(2,1) + g(2,2) = 1$$

$$c) \quad g(s, a^2) \geq 0, s \in \{1,2\}, a^2 \in A^2(s)$$

Χρησιμοποιώντας το λογισμικό MATLAB, παίρνουμε τα ακόλουθα αποτελέσματα:

Αντικειμενική συνάρτηση	-2.0909
$u(1)$	0
$u(2)$	0.8409
$v(1)$	-2.0909
$v(2)$	-2.0909
$g_{11}$	0.9773
$g_{12}$	0.0227
$g_{13}$	0
$g_{21}$	0.25
$g_{22}$	0.75

Το αντίστοιχο δυϊκό πρόγραμμα θα είναι ως εξής:

$$\max[z(1) + z(2)]$$

με περιορισμούς:

$$D_a(1)$$

$$d) \quad \begin{pmatrix} 0 & 1 & -1 & 0 & -1 & \vdots & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 1 & \vdots & 0 & 0 & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 1 & 0 & 0 & 0 & \vdots & 0 & 1 & -1 & 0 & -1 \\ 0 & 0 & 1 & 1 & 1 & \vdots & 0 & -1 & 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_{11} \\ x_{12} \\ x_{21} \\ x_{22} \\ x_{23} \\ y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \\ y_{23} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1/2 \\ 1/2 \end{pmatrix}$$

$$e) \quad 2x_{11} + 3x_{12} + z(1) \leq 0$$

$$6x_{11} + z(1) \leq 0$$

$$2x_{12} + z(1) \leq 0$$

$$-x_{21} + 4x_{22} + 5x_{23} + z(2) \leq 0$$

$$2x_{21} + 3x_{22} + z(2) \leq 0$$

$$f) \quad x_{ij} \geq 0, y_{ij} \geq 0 \text{ για } i = 1, 2, j = 1, 2, 3$$

Χρησιμοποιώντας το λογισμικό MATLAB, παίρνουμε τα ακόλουθα αποτελέσματα:

Αντικειμενική συνάρτηση	-2.0909
$z(1)$	-1.6364
$z(2)$	-0.4545
$x_{11}$	0.2727
$x_{12}$	0.3636
$x_{21}$	0.2273
$x_{22}$	0
$x_{23}$	0.1364
$y_{11}$	0
$y_{12}$	0
$y_{21}$	0.1364
$y_{22}$	0
$y_{23}$	0

Από το Βήμα 2 του Αλγορίθμου:

$$S^* = \{1, 2\} = S$$

Χρησιμοποιώντας τα παραπάνω αποτελέσματα, μπορούμε τώρα να υπολογίσουμε την στρατηγική  $\hat{f}$ :

$$\hat{f}(1,1) = \frac{x_{11}}{x_{11} + x_{12}} = 0.4286$$

$$\hat{f}(1,2) = \frac{x_{12}}{x_{11} + x_{12}} = 0.5714$$

$$\hat{f}(2,1) = \frac{x_{21}}{x_{21} + x_{22} + x_{23}} = 0.625$$

$$\hat{f}(2,2) = \frac{x_{22}}{x_{21} + x_{22} + x_{23}} = 0$$

$$\hat{f}(2,3) = \frac{x_{23}}{x_{21} + x_{22} + x_{23}} = 0.375$$

Φυσικά, η βέλτιστη στρατηγική  $\hat{\mathbf{g}}$ , για τον παίκτη 2, δίνεται στα αποτελέσματα της λύσης του  $P_a(1)$ .

Συνεπώς, έχουμε καταλήξει σε ένα ζεύγος βέλτιστων στρατηγικών για το αρχικό παίγνιο

$$(\hat{\mathbf{f}}, \hat{\mathbf{g}}) = (((0.4286, 0.5714), (0.625, 0, 0.375)), ((0.9773, 0.0227, 0), (0.25, 0.75)))$$

Για το υπόλοιπο αυτής την ενότητας, θα ασχοληθούμε με τον σχολιασμό του Αλγορίθμου 2.4.1.. Αρχικά, θα δούμε γιατί ο αλγόριθμος αυτός είναι καλά ορισμένος.

Πράγματι, ας πάρουμε ένα αυθαίρετο  $\mathbf{g} \in G_s$ ,  $\mathbf{u} = \mathbf{0}$  και  $\mathbf{v} = M\mathbf{1}_N$ , όπου:

$$M \doteq \max_{s,a^1,a^2} \{|r(s, a^1, a^2)|\}$$

Ισχυριζόμαστε πως θα καταλήξουμε σε μια εφικτή λύση του  $P_a(1)$ .

Βλέπουμε πως οι περιορισμοί  $b)$  και  $c)$  ικανοποιούνται τετριμμένα

Το δεύτερο μπλοκ περιορισμών του  $a)$  ανάγεται σε:

$$\mathbf{v}^T \mathbf{W} = M\mathbf{1}^T \mathbf{W} \geq \mathbf{0}^T$$

(λόγω της κατασκευής του  $\mathbf{W}$ , πολλαπλασιάζονται μόνο τα διαγώνια στοιχεία, τα οποία έχουν μη αρνητικό πρόσημο)

ενώ το πρώτο μπλοκ γίνεται:

$$\mathbf{u}^T \mathbf{W} + \mathbf{v}^T \mathbf{J} - \mathbf{g}^T \mathbf{R}^T = \mathbf{0}^T + M[\mathbf{1}_N^T \mathbf{J}] - [\mathbf{R}\mathbf{g}]^T$$

η τελευταία έκφραση θα είναι πάντα μεγαλύτερη ή ίση με  $\mathbf{0}^T$ , αφού:

$$[M\mathbf{1}_{m^1(s)}]^T - [\mathbf{R}(s)\mathbf{g}(s)]^T \geq \mathbf{0}_{m^1(s)}^T, \quad \text{για κάθε } s \in S$$

Προκειμένου να εξασφαλιστεί η ύπαρξη πεπερασμένης βέλτιστης λύσης για το  $P_a(1)$ , και άρα λόγω δυϊκού θεωρήματος και για το  $D_a(1)$ , μένει να δείξουμε πως το  $\mathbf{v}^T$ -μπλοκ κάθε εφικτής λύσης του  $P_a(1)$  είναι κάτω φραγμένο. Το αποτέλεσμα αυτό θα προκύψει ως πόρισμα της ακόλουθης πρότασης.

*Πρόταση 2.4.1.*

Έστω  $(\bar{\mathbf{u}}^T, \bar{\mathbf{v}}^T, \bar{\mathbf{g}}^T)$  μια αυθαίρετη εφικτή λύση του  $P_a(1)$  και  $\mathbf{f}$  οποιαδήποτε στάσιμη στρατηγική του παίκτη 1. Τότε:

$$\bar{\mathbf{v}} \geq \mathbf{v}_a(\mathbf{f}, \bar{\mathbf{g}})$$

*Απόδειξη*

Αρχικά θεωρούμε το μπλοκ περιορισμών από το α) :

$$\bar{\mathbf{v}}^T \mathbf{W} \geq \mathbf{0}_{m^1}^T$$

Από τη στιγμή που το  $\mathbf{W}$  έχει την ίδια μπλοκ δομή με την  $\mathbf{f}$ , αναμιγνύοντας το  $s$ -οστό μπλοκ του  $\mathbf{W}$  ως προς την  $\mathbf{f}$ , παίρνουμε ότι για όλα τα  $s \in S$  :

$$v(s) \geq [\mathbf{P}(\mathbf{f})\bar{\mathbf{v}}]_s$$

ή, ισοδύναμα, σε διανυσματική μορφή:

$$\bar{\mathbf{v}} \geq \mathbf{P}(\mathbf{f})\bar{\mathbf{v}}, \text{ όπου } \mathbf{f} \in F_s$$

Ισχυριζόμαστε ότι:

$$\bar{\mathbf{v}} \geq \mathbf{P}(\mathbf{f})\bar{\mathbf{v}} \Rightarrow \bar{\mathbf{v}} \geq \mathbf{Q}(\mathbf{f})\bar{\mathbf{v}}$$

Πράγματι, έχουμε ότι:

$$\bar{\mathbf{v}} \geq \mathbf{P}(\mathbf{f})\bar{\mathbf{v}} \geq \mathbf{P}^2(\mathbf{f})\bar{\mathbf{v}} \geq \dots \geq \mathbf{P}^k(\mathbf{f})\bar{\mathbf{v}}$$

Ακόμη, για κάθε  $T \in \mathbb{N}$ , παίρνουμε ότι:

$$T\bar{\mathbf{v}} \geq \sum_{t=1}^T \mathbf{P}^t(\mathbf{f})\bar{\mathbf{v}} \Rightarrow T\bar{\mathbf{v}} + \bar{\mathbf{v}} \geq \sum_{t=1}^T \mathbf{P}^t(\mathbf{f})\bar{\mathbf{v}} + \bar{\mathbf{v}} \Rightarrow \bar{\mathbf{v}} \geq \frac{1}{T+1} \sum_{t=0}^T \mathbf{P}^t(\mathbf{f})\bar{\mathbf{v}}$$

Παίρνοντας  $T \rightarrow \infty$ , προκύπτει το ζητούμενο αποτέλεσμα:

$$\bar{\mathbf{v}} \geq \lim_{T \rightarrow \infty} \left( \frac{1}{T+1} \sum_{t=0}^T \mathbf{P}^t(\mathbf{f})\bar{\mathbf{v}} \right) = \mathbf{Q}(\mathbf{f})\bar{\mathbf{v}} \quad (2.4.1)$$

Όμοια, αναμιγνύοντας ως προς την  $\mathbf{f}$  το  $s$ -οστό μπλοκ περιορισμών της:

$$\bar{\mathbf{u}}^T \mathbf{W} + \bar{\mathbf{v}}^T \mathbf{J} - \bar{\mathbf{g}}^T \mathbf{R}^T \geq \mathbf{0}_{m^1}^T$$

λαμβάνουμε ότι:

$$\bar{v}(s) + \bar{u}(s) \geq \mathbf{f}(s)R(s)\bar{\mathbf{g}}(s) + [\mathbf{P}(\mathbf{f})\bar{\mathbf{u}}]_s$$



για κάθε  $s \in S$ . Ισοδύναμα, σε διανυσματική γραφή:

$$\bar{v} + \bar{u} \geq r(f, \bar{g}) + P(f)\bar{u}$$

Πολλαπλασιάζοντας και τα 2 μέλη επί  $Q(f)$ , η τελευταία σχέση δίνει:

$$Q(f)\bar{v} + Q(f)\bar{u} \geq Q(f)r(f, \bar{g}) + Q(f)P(f)\bar{u}$$

Όμως γνωρίζουμε ότι  $QP = Q$ , οπότε:

$$Q(f)\bar{v} \geq Q(f)r(f, \bar{g})$$

και χρησιμοποιώντας τη σχέση (2.4.1) παίρνουμε το ζητούμενο, δηλαδή ότι:

$$\bar{v} \geq v_a(f, \bar{g}) \quad (2.4.2)$$

■

### Πόρισμα 2.4.1

- (i) Η αντικειμενική συνάρτηση  $[\gamma^T v]$  στο  $P_a(1)$  είναι κάτω φραγμένη.
- (ii) Τα προβλήματα  $P_a(1)$  και  $D_a(1)$  έχουν πεπερασμένες βέλτιστες λύσεις.

### Απόδειξη

- (i) Έστω  $M_L \doteq \min_{s, a^1, a^2} \{r(s, a^1, a^2)\}$ . Από τον ορισμό του κριτηρίου οριακού μέσου, μπορούμε να δούμε ότι για όλα τα  $(f, g) \in F_s \times G_s$  ισχύει:

$$v_a(f, g) \geq M_L \mathbf{1}_N$$

Τότε, από την (2.5.2), για οποιαδήποτε εφικτά  $(\bar{u}^T, \bar{v}^T, \bar{g}^T)$  για το  $P_a(1)$ , έπεται:

$$[\gamma^T \bar{v}^T] = \sum_{s \in S} \gamma(s) \bar{v}(s) \geq M_L \left( \sum_{s \in S} \gamma(s) \right) = M_L$$

- (ii) Εφόσον το  $P_a(1)$  είναι ένα εφικτό και φραγμένο γραμμικό πρόγραμμα, θα επιδέχεται πεπερασμένη βέλτιστη λύση. Από το Δυϊκό θεώρημα του γραμμικού προγραμματισμού, συμπεραίνουμε πως το ίδιο ισχύει και για το  $D_a(1)$ .

■

Έχοντας πλέον δείξει το ζητούμενο αποτέλεσμα, συνεχίζοντας την ανάλυση του Αλγορίθμου 2.5.1., παρατηρούμε ότι το δεύτερο μπλοκ περιορισμών του  $d$ ), δηλαδή  $J\mathbf{x} + W\mathbf{y} = \gamma$ , έχει ως  $s$ -οστή είσοδο (μετά από αναδιάταξη όρων) :

$$x_{s'} + y_{s'} = \left[ \sum_{s \in S} \sum_{a^1 \in A^1(s)} p(s'|s, a^1) y_{sa^1} \right] + \gamma(s') \geq \gamma(s') > 0, \quad s' \in S$$

Οπότε, θα έχουμε ότι  $y_{s'} > 0$  κάθε φορά που  $x_{s'} = 0$ , και έτσι, η  $\hat{f}$  του Βήματος 3 του Αλγορίθμου είναι καλώς ορισμένη.

#### Πρόταση 2.4.2.

Έστω  $(\hat{\mathbf{u}}^T, \hat{\mathbf{v}}^T, \hat{\mathbf{g}}^T)$  και  $(\hat{\mathbf{x}}^T, \hat{\mathbf{y}}^T, \hat{\mathbf{z}}^T)$  να είναι ένα ζεύγος δυϊκών βέλτιστων λύσεων στα  $P_a(1)$  και  $D_a(1)$ , αντίστοιχα. Έστω, ακόμα  $\hat{f} \in F_S$  μία στάσιμη στρατηγική του παίκτη 1 και  $S^* \subset S$  κατασκευασμένα όπως στο Βήμα 3 και Βήμα 2 του Αλγορίθμου, αντίστοιχα. Τότε:

(i) Το  $S^*$  είναι το σύνολο των επαναλαμβανόμενων καταστάσεων της Μαρκοβιανής αλυσίδας που επάγεται από τον  $P(\hat{f})$ .

(ii)  $\hat{\mathbf{v}} = \mathbf{P}(\hat{f})\hat{\mathbf{v}} = \mathbf{Q}(\hat{f})\hat{\mathbf{v}}$

(iii)  $\hat{v}(s) + \left[ (I - \mathbf{P}(\hat{f})) \hat{\mathbf{u}} \right]_s = \hat{f}(s)R(s)\hat{g}(s), \quad s \in S^*$

Έστω τώρα  $\mathbf{Q}(\hat{f}) = [q(s'|s, \hat{f})]_{s, s' \in S}$ . Από το (i) της προηγούμενης πρότασης παίρνουμε ότι :

$$q(s'|s, \hat{f}) = 0, \quad \text{αν } s' \notin S^*$$

Ακόμη, για όλα τα  $s \in S$  έχουμε ότι:

$$v_a(s, \hat{f}, \hat{g}) = [\mathbf{Q}(\hat{f})\mathbf{r}(\hat{f}, \hat{g})]_s = \sum_{s' \in S^*} q(s'|s, \hat{f})r(s, \hat{f}, \hat{g})$$

Από το (iii) της πρότασης, η παραπάνω σχέση γίνεται:

$$\begin{aligned} v_a(s, \hat{f}, \hat{g}) &= \sum_{s' \in S^*} q(s'|s, \hat{f})\hat{v}(s) + \sum_{s' \in S^*} q(s'|s, \hat{f}) \left[ (I - \mathbf{P}(\hat{f})) \hat{\mathbf{u}} \right]_s \\ &= [\mathbf{Q}(\hat{f})\hat{\mathbf{v}}]_s + [\mathbf{Q}(\hat{f})(I - \mathbf{P}(\hat{f}))\hat{\mathbf{u}}]_s = \hat{v}(s) \end{aligned}$$

όπου η τελευταία ισότητα προκύπτει από το (ii) και το γεγονός ότι  $\mathbf{Q}(\hat{f}) = \mathbf{Q}(\hat{f})\mathbf{P}(\hat{f})$

Έτσι, έχουμε δείξει πως μια βέλτιστη λύση  $(\hat{\mathbf{u}}^T, \hat{\mathbf{v}}^T, \hat{\mathbf{g}}^T)$  του  $P_a(1)$  ικανοποιεί την:

$$\hat{\mathbf{v}} = \mathbf{v}_a(\hat{\mathbf{f}}, \hat{\mathbf{g}})$$

Σε συνδυασμό με την Πρόταση 2.4.1. , έχουμε μέχρι τώρα δείξει το ένα σκέλος της συνθήκης βελτιστοποίησης, δηλαδή ότι:

$$\mathbf{v}_a(\mathbf{f}, \hat{\mathbf{g}}) \leq \mathbf{v}_a(\hat{\mathbf{f}}, \hat{\mathbf{g}}), \text{ με } \mathbf{f} \in F_s \quad (2.4.3)$$

Για το δεύτερο σκέλος της συνθήκης βελτιστοποίησης θα χρειαστούμε μία τελευταία πρόταση, την οποία παραθέτουμε χωρίς απόδειξη.

### Πρόταση 2.4.3.

Έστω  $(\hat{\mathbf{u}}^T, \hat{\mathbf{v}}^T, \hat{\mathbf{g}}^T)$ ,  $(\hat{\mathbf{x}}^T, \hat{\mathbf{y}}^T, \hat{\mathbf{z}}^T)$  και  $\hat{\mathbf{f}} \in F_s$  όπως είναι στη διατύπωση της Πρότασης 2.4.2. Τότε:

(i) Για όλα τα  $\mathbf{g} \in G_s$  :

$$\sum_{s \in S} \hat{z}_s = \gamma^T \mathbf{v}_a(\hat{\mathbf{f}}, \hat{\mathbf{g}}) \leq \gamma^T \mathbf{v}_a(\hat{\mathbf{f}}, \mathbf{g})$$

(ii) Για όλα τα  $\mathbf{g} \in G_s$  :

$$\mathbf{v}_a(\hat{\mathbf{f}}, \hat{\mathbf{g}}) \leq \mathbf{v}_a(\hat{\mathbf{f}}, \mathbf{g})$$

Φυσικά, το (ii) της Πρότασης μαζί με την (2.4.3) συμπληρώνουν τα δύο μέλη της ζητούμενης ανισότητας. Έχουμε καταλήξει δηλαδή στην :

$$\mathbf{v}_a(\mathbf{f}, \hat{\mathbf{g}}) \leq \mathbf{v}_a(\hat{\mathbf{f}}, \hat{\mathbf{g}}) \leq \mathbf{v}_a(\hat{\mathbf{f}}, \mathbf{g})$$

για όλα τα  $\mathbf{f} \in F_s$  και  $\mathbf{g} \in G_s$ .

## Κεφάλαιο 3

### Ημιμαρκοβιανά παίγνια τέλειας πληροφόρησης

#### 3.0 Εισαγωγή

Το τελευταίο κεφάλαιο της εργασίας αυτής πραγματεύεται ημιμαρκοβιανά παίγνια τέλειας πληροφόρησης. Πρόκειται για μια γενίκευση των “κλασικών” στοχαστικών παιγνίων, όπου πλέον θεωρούμε ότι η επιλογή ενεργειών από τους παίκτες συνεπάγεται και έναν χρόνο παραμονής στην παρούσα κατάσταση, προτού πραγματοποιηθεί μετάβαση στην επόμενη κατάσταση και ούτω καθεξής. Ο χρόνος παραμονής θα προέρχεται από γνωστή κατανομή, χωρίς να επεκταθούμε παραπάνω στο συγκεκριμένο ζήτημα, εφόσον το ενδιαφέρον μας στρέφεται γύρω από τη μέση τιμή της παραμονής αυτής. Πιο συγκεκριμένα, θα γίνει ανάλυση των παιγνίων που διαθέτουν το χαρακτηριστικό της “τέλειας πληροφόρησης”. Αφού δοθεί το κατάλληλο θεωρητικό υπόβαθρο σχετικά με αυτήν την κατηγορία παιγνίων, παρατίθεται ο αλγόριθμος επίλυσης ενός τέτοιου παιγνίου μέσω γραμμικού προγραμματισμού, το οποίο πρώτα έχει υποβιβαστεί σε ημιμαρκοβιανή διαδικασία απόφασης. Τέλος, γίνεται εφαρμογή του συγκεκριμένου αλγορίθμου σε αριθμητικό παράδειγμα, το οποίο επιλύουμε με δύο διαφορετικούς τρόπους.

#### 3.1 Πεπερασμένα ημιμαρκοβιανά παίγνια δύο παικτών και μηδενικού αθροίσματος

Ας ξεκινήσουμε από τον ορισμό ενός ημιμαρκοβιανού παιγνίου 2 παικτών και μηδενικού αθροίσματος, όπου τόσο ο χώρος καταστάσεων όσο και ο χώρος των διαθέσιμων ενεργειών για κάθε παίκτη, είναι πεπερασμένα σύνολα. Ορίζουμε, λοιπόν, ένα τέτοιο παίγνιο  $\Gamma$ , ως μία συλλογή από 7 αντικείμενα :

$$\Gamma = \langle S, \{A^1(s) \mid s \in S\}, \{A^2(s) \mid s \in S\}, p, Q, r \rangle, \text{ όπου:}$$

$S = \{1, 2, \dots, N\}$  είναι το μη-κενό πεπερασμένο σύνολο καταστάσεων

$A^1(s) = \{1, 2, \dots, m^1(s)\}$  είναι το μη-κενό πεπερασμένο σύνολο των επιτρεπτών ενεργειών για τον παίκτη 1 στην κατάσταση  $s$

$A^2(s) = \{1, 2, \dots, m^2(s)\}$  είναι το μη-κενό πεπερασμένο σύνολο των επιτρεπτών ενεργειών για τον παίκτη 2 στην κατάσταση  $s$

$K = \{(s, a^1, a^2) \mid s \in S, a^1 \in A^1(s), a^2 \in A^2(s)\}$  είναι το σύνολο που περιλαμβάνει όλες τις επιτρεπτές 3-αδες καταστάσεων-ενεργειών των 2 παικτών

$p(\cdot \mid s, a^1, a^2)$  είναι η πιθανότητα μετάβασης στην επόμενη κατάσταση (νόμος της κίνησης)

$Q_{ss'}(. | a^1, a^2)$  είναι μια κατανομή πιθανότητας στο  $[0, \infty)$ , δεδομένου του  $K \times S$ , την οποία θα καλούμε κατανομή των ενδιάμεσων χρόνων παραμονής

$r$  είναι η πραγματική συνάρτηση πληρωμής ενός βήματος για τον παίκτη 1, που ορίζεται πάνω στο  $K$  (η συνάρτηση πληρωμής ενός βήματος για τον παίκτη 2 θα είναι η  $r_2(s, a^1, a^2) = -r(s, a^1, a^2)$ )

Το ημιμαρκοβιανό παίγνιο εξελίσσεται επ' άπειρον ως εξής:

Στην μηδενική εποχή απόφασης, το παιχνίδι ξεκινάει από την αρχική κατάσταση  $s_0 \in S$  και οι παίκτες 1 και 2 ταυτόχρονα και ανεξάρτητα ο ένας απ' τον άλλο, επιλέγουν ενέργειες.

$a_0^1 \in A^1(s_0)$  και  $a_0^2 \in A^2(s_0)$ . Συνεπώς, οι παίκτες λαμβάνουν αμοιβές ενός βήματος  $r(s_0, a_0^1, a_0^2)$  και  $r_2(s_0, a_0^1, a_0^2)$ , αντίστοιχα και το παίγνιο μεταβαίνει στην επόμενη κατάσταση  $s_1$  με πιθανότητα μετάβασης  $p(s_1 | s_0, a_0^1, a_0^2)$ . Ο χρόνος παραμονής στην κατάσταση  $s_0$  έως ότου γίνει η μετάβαση στην κατάσταση  $s_1$  ορίζεται από τη συνάρτηση κατανομής  $Q_{s_0 s_1}(. | a_0^1, a_0^2)$ . Εφόσον φτάσουμε στην κατάσταση  $s_1$  στην επόμενη εποχή απόφασης, επαναλαμβάνεται η ίδια διαδικασία με την κατάσταση  $s_0$  να έχει πλέον αντικατασταθεί από την  $s_1$ .

### Στρατηγικές

Στο 1<sup>ο</sup> κεφάλαιο ορίσαμε μία συμπεριφορική στρατηγική για τον παίκτη 1 να είναι μια ακολουθία

$$\mathbf{f} \doteq (f_0, f_1, f_2, \dots, f_t, \dots)$$

όπου τα  $f_t$  είναι κανόνες απόφασης από το σύνολο των ιστοριών της διαδικασίας,  $H_t$ , στο  $P(A^1)$ , το οποίο είναι το σύνολο των κατανομών πιθανότητας πάνω στο πεπερασμένο σύνολο ενεργειών  $A^1$ .

Με όμοιο τρόπο μπορεί να οριστεί η συμπεριφορική στρατηγική  $\mathbf{g}$  για τον παίκτη 2.

Οι χώροι όλων των συμπεριφορικών στρατηγικών των παικτών θα συμβολίζονται με  $F_B^1$  και  $F_B^2$ , αντίστοιχα. Οι συμπεριφορικές στρατηγικές αποτελούν την πιο ευρεία κλάση στρατηγικών που θα μας απασχολήσουν.

Εν συνεχεία, θα καλούμε μία συμπεριφορική στρατηγική  $\mathbf{f} = \{f_t\}_{t=0}^\infty$  του παίκτη 1 *ημιμαρκοβιανή στρατηγική* εάν για κάθε  $t$ , το  $f_t$  εξαρτάται από τα  $s_0, s_t$  και την εποχή απόφασης  $t$ . Ανάλογα ορίζεται και η *ημιμαρκοβιανή στρατηγική*  $\mathbf{g} = \{g_t\}_{t=0}^\infty$  για τον παίκτη 2.

Όπως έχουμε ήδη δει, μία *στάσιμη στρατηγική* θα εξαρτάται μόνο από την παρούσα κατάσταση. Δηλαδή, μία στάσιμη στρατηγική του παίκτη 1 (ανάλογα και για τον παίκτη 2) ορίζεται ως μία  $N$ -άδα,

$f = (f(1), f(2), \dots, f(N))$ , όπου  $f(s) = (f(s, 1), f(s, 2), \dots, f(s, m^1(s)))$  και το  $f(s, i)$  δηλώνει την πιθανότητα να επιλεγεί η ενέργεια  $i$  στην κατάσταση  $s$ . Θα συμβολίζουμε με  $F_s^1$  και  $F_s^2$  τα σύνολα των στάσιμων στρατηγικών των παικτών, αντίστοιχα.

Υπενθυμίζουμε πως μία στάσιμη στρατηγική θα καλείται *καθαρή* εάν σε κάθε κατάσταση επιλέγεται μία ενέργεια με πιθανότητα 1 (και όλες οι υπόλοιπες με πιθανότητα 0). Θα συμβολίζουμε με  $F_{sp}^1$  και  $F_{sp}^2$  τα σύνολα των στάσιμων καθαρών στρατηγικών των παικτών, αντίστοιχα.

Μία *ημιστάσιμη στρατηγική* είναι μία ημιμαρκοβιανή στρατηγική, η οποία είναι ανεξάρτητη της εποχής απόφασης  $t$ . Θα συμβολίζουμε με  $\xi_1$  και  $\xi_2$  τα σύνολα των ημιστάσιμων στρατηγικών των παικτών.

### Ορισμός 3.1.1.

Ένα ημιμαρκοβιανό παίγνιο μηδενικού αθροίσματος και δύο παικτών  $\Gamma = \langle S, \{A^1(s) \mid s \in S\}, \{A^2(s) \mid s \in S\}, p, Q, r \rangle$  θα καλείται παίγνιο *τέλειας πληροφόρησης*, εάν ισχύουν τα ακόλουθα :

- (i)  $S = S_1 \cup S_2, S_1 \cap S_2 = \emptyset$
- (ii)  $|A^2(s)| = 1$ , για όλα τα  $s \in S_1$ , δηλαδή στο  $S_1$  ο παίκτης 2 έχει μόνο μία διαθέσιμη ενέργεια
- (iii)  $|A^1(s)| = 1$ , για όλα τα  $s \in S_2$ , δηλαδή στο  $S_2$  ο παίκτης 1 έχει μόνο μία διαθέσιμη ενέργεια

### 3.2. Ημιμαρκοβιανά παίγνια με κριτήριο πληρωμής τον οριακό λόγο των μέσων

Έστω  $(X_1, A_1^1, A_1^2, X_2, A_2^1, A_2^2, \dots)$  μία διατεταγμένη ακολουθία στο  $S \times (A^1 \times A^2 \times S)^\infty$ . Δοθέντος ζεύγους συμπεριφορικών στρατηγικών  $(f, g) \in F_B^1 \times F_B^2$  και αρχικής κατάστασης  $s$ , ορίζουμε ως αναμενόμενη πληρωμή τον οριακό λόγο των μέσων, για τον παίκτη 1, να είναι :

$$v_a(s, f, g) = \liminf_{t \rightarrow \infty} \frac{E_{f,g}[\sum_{m=1}^t (r(X_m, A_m^1, A_m^2) \mid X_1 = s)]}{E_{f,g}[\sum_{m=1}^t (\bar{r}(X_m, A_m^1, A_m^2) \mid X_1 = s)]} \quad (3.2.1)$$

Όπου το  $\bar{r}(s, a^1, a^2) = \sum_{s' \in S} \left( p(s' \mid a^1, a^2) \int_0^\infty t dQ_{ss'}(t \mid a^1, a^2) \right)$  δηλώνει τον αναμενόμενο

χρόνο παραμονής στην κατάσταση  $s$  για το ζευγάρι ενεργειών  $(a^1, a^2) \in A^1(s) \times A^2(s)$ .

*Ορισμός 3.2.1.*

Για ένα ζεύγος στάσιμων στρατηγικών  $(\mathbf{f}, \mathbf{g}) \in F_S^1 \times F_S^2$  ορίζουμε τον πίνακα πιθανοτήτων μετάβασης ως:

$$\mathbf{P}(\mathbf{f}, \mathbf{g}) = [p(s' | s, \mathbf{f}, \mathbf{g})]_{N \times N}$$

όπου

$$p(s' | s, \mathbf{f}, \mathbf{g}) = \sum_{a^1 \in A^1(s)} \sum_{a^2 \in A^2(s)} p(s' | s, a^1, a^2) \text{ δηλώνει την πιθανότητα}$$

ότι, ξεκινώντας από την κατάσταση  $s$ , η επόμενη κατάσταση θα είναι η  $s'$  για το ζεύγος στρατηγικών  $(\mathbf{f}, \mathbf{g})$  των παικτών, αντίστοιχα.

Για τις στάσιμες στρατηγικές  $(\mathbf{f}, \mathbf{g})$ , θα γράφουμε την συνολική πληρωμή των παικτών ως :

$$v_a(s, \mathbf{f}, \mathbf{g}) = \liminf_{t \rightarrow \infty} \frac{\sum_{m=1}^t r^m(s, \mathbf{f}, \mathbf{g})}{\sum_{m=1}^t \bar{r}^m(s, \mathbf{f}, \mathbf{g})}, \text{ για κάθε } s \in S.$$

όπου τα  $r^m(s, \mathbf{f}, \mathbf{g})$  και  $\bar{r}^m(s, \mathbf{f}, \mathbf{g})$  είναι αντίστοιχα η αναμενόμενη αμοιβή για τον παίκτη 1 και ο αναμενόμενος χρόνος παραμονής στην  $m$ -οστή εποχή απόφασης, όταν βρισκόμαστε στην κατάσταση  $s$  και έχουν επιλεγεί οι στρατηγικές  $\mathbf{f}, \mathbf{g}$  από τους παίκτες. Ορίζουμε :

$$\mathbf{r}(\mathbf{f}, \mathbf{g}) = [r(s, \mathbf{f}, \mathbf{g})]_{N \times 1}, \quad \bar{\mathbf{r}}(\mathbf{f}, \mathbf{g}) = [\bar{r}(s, \mathbf{f}, \mathbf{g})]_{N \times 1}, \quad \mathbf{v}_a(\mathbf{f}, \mathbf{g}) = [v_a(s, \mathbf{f}, \mathbf{g})]_{N \times 1}$$

ως τα διανύσματα-στήλες της αναμενόμενης αμοιβής ενός βήματος, του αναμενόμενου χρόνου παραμονής και της αναμενόμενης πληρωμής του οριακού λόγου των μέσων αντίστοιχα, για το ζεύγος στρατηγικών  $(\mathbf{f}, \mathbf{g}) \in F_S^1 \times F_S^2$ . Έχουμε, ακόμη, ότι:

$$\begin{aligned} r^m(s, \mathbf{f}, \mathbf{g}) &= \sum_{s' \in S} \mathbb{P}_{\mathbf{f}, \mathbf{g}}(X_m = s' | X_1 = s) r(s', \mathbf{f}, \mathbf{g}) = \sum_{s' \in S} r(s', \mathbf{f}, \mathbf{g}) p^{m-1}(s' | s, \mathbf{f}, \mathbf{g}) \\ &= [\mathbf{P}^{m-1}(\mathbf{f}, \mathbf{g}) \mathbf{r}(\mathbf{f}, \mathbf{g})](s) \\ \bar{r}^m(s, \mathbf{f}, \mathbf{g}) &= \sum_{s' \in S} \mathbb{P}_{\mathbf{f}, \mathbf{g}}((X_m = s' | X_1 = s) \bar{r}(s', \mathbf{f}, \mathbf{g})) = \sum_{s' \in S} \bar{r}(s', \mathbf{f}, \mathbf{g}) p^{m-1}(s' | s, \mathbf{f}, \mathbf{g}) \\ &= [\mathbf{P}^{m-1}(\mathbf{f}, \mathbf{g}) \bar{\mathbf{r}}(\mathbf{f}, \mathbf{g})](s) \end{aligned}$$

Στο σημείο αυτό παραθέτουμε το ακόλουθο Λήμμα, όπως αυτό είναι διατυπωμένο στο [7] (σελ.175)

*Λήμμα 3.2.1.*

Έστω  $\mathbf{P} = [p(s'|s)]_{n \times n}$  να είναι ένας πίνακας πιθανοτήτων μετάβασης . Τότε υπάρχει στοχαστικός πίνακας  $\mathbf{P}^* = [p^*(s'|s)]_{n \times n}$  , τον οποίο θα καλούμε *όριο κατά Cesaro* του πίνακα  $\mathbf{P}$  και ο οποίος ορίζεται να είναι :

$$p^*(s'|s) = \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{k=1}^m p^k(s'|s) \quad , \quad s, s' = 1, 2, \dots, n$$

Εφόσον ο πίνακας  $\mathbf{P}(\mathbf{f}, \mathbf{g})$  είναι ένας τέτοιος πίνακας πιθανοτήτων μετάβασης, έχουμε ότι το όριο

$$\lim_{t \rightarrow \infty} \frac{1}{t+1} \sum_{m=0}^t \mathbf{P}^m(\mathbf{f}, \mathbf{g})$$

υπάρχει και είναι ίσο με  $\mathbf{P}^*(\mathbf{f}, \mathbf{g})$ , όπου με  $\mathbf{P}^*(\mathbf{f}, \mathbf{g})$  συμβολίζουμε τον πίνακα που προκύπτει ως όριο κατά Cesaro από τον  $\mathbf{P}(\mathbf{f}, \mathbf{g})$

Επιπλέον έχουμε ότι:

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{m=1}^t r^m(\mathbf{f}, \mathbf{g}) = [\mathbf{P}^*(\mathbf{f}, \mathbf{g})\mathbf{r}(\mathbf{f}, \mathbf{g})](s)$$

και

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{m=1}^t \bar{r}^m(\mathbf{f}, \mathbf{g}) = [\mathbf{P}^*(\mathbf{f}, \mathbf{g})\bar{\mathbf{r}}(\mathbf{f}, \mathbf{g})](s)$$

Οπότε, τελικά, για ένα ζεύγος στάσιμων στρατηγικών  $(\mathbf{f}, \mathbf{g}) \in F_s^1 \times F_s^2$  , θα έχουμε ότι :

$$v_a(s, \mathbf{f}, \mathbf{g}) = \frac{[\mathbf{P}^*(\mathbf{f}, \mathbf{g})\mathbf{r}(\mathbf{f}, \mathbf{g})](s)}{[\mathbf{P}^*(\mathbf{f}, \mathbf{g})\bar{\mathbf{r}}(\mathbf{f}, \mathbf{g})](s)} \quad , \text{για όλα τα } s \in S \quad (3.2.2)$$

*Ορισμός 3.2.2.*

Θα λέμε ότι το ημιμαρκοβιανό παίγνιο με πληρωμές τον οριακό λόγο των μέσων , δύο-παικτών και μηδενικού αθροίσματος, θα έχει διάνυσμα τιμής  $\mathbf{v}_a = [v_a(s)]_{N \times 1}$  αν

$$\sup_{\mathbf{f} \in F_b^1} \inf_{\mathbf{g} \in F_b^2} v_a(s, \mathbf{f}, \mathbf{g}) = v_a(s) = \inf_{\mathbf{g} \in F_b^2} \sup_{\mathbf{f} \in F_b^1} v_a(s, \mathbf{f}, \mathbf{g}) \quad , \text{για όλα τα } s \in S$$



Ένα ζεύγος στρατηγικών  $(f^*, g^*) \in F_b^1 \times F_b^2$  θα λέμε ότι είναι *βέλτιστο* για τους παίκτες, εάν

$$v_a(s, f^*, g) \geq v_a(s) \geq v_a(s, f, g^*) , \text{ για όλα τα } s \in S \text{ και όλα τα } (f, g) \in F_b^1 \times F_b^2$$

### 3.3 Πεπερασμένες ημιμαρκοβιανές διαδικασίες απόφασης

Κατ' αναλογία με την προηγούμενη παράγραφο καθώς και με την σχετική παράγραφο του κεφαλαίου 1 για μαρκοβιανές διαδικασίες, θα περιγράψουμε εν συντομία την διαδικασία εξέλιξης μιας τέτοιας ημιμαρκοβιανής διαδικασίας.

Ορίζουμε λοιπόν μια τέτοια πεπερασμένη διαδικασία (πεπερασμένος χώρος καταστάσεων και ενεργειών) ως μία συλλογή 5 αντικειμένων ως εξής :

$$\hat{F} = \langle S, \hat{A} = \{A(s) | s \in S\}, \hat{p}, \hat{Q}, \hat{r} \rangle$$

Ο συμβολισμός που θα χρησιμοποιήσουμε ταυτίζεται με τον αντίστοιχο για ημιμαρκοβιανά παίγνια. Θα χρησιμοποιούμε τον συμβολισμό  $(\cdot)$  για να δηλώσουμε την διαφοροποίηση από τον αντίστοιχο συμβολισμό των παιγνίων.

Η διαδικασία και πάλι εξελίσσεται στο διηγεκές, ακριβώς όπως και στο ανάλογο παίγνιο , με τη μόνη διαφορά να είναι πως τώρα οι δύο παίκτες αντικαθίστανται από έναν που λαμβάνει αποφάσεις.

Φυσικά οι ορισμοί που αφορούν τους χώρους στρατηγικών είναι ακριβώς οι ίδιοι με εκείνους των ημιμαρκοβιανών παιχνιδιών. Θα συμβολίζουμε με  $F_B, F_S, F_{SP}, \xi$  τα σύνολα των συμπεριφορικών, στάσιμων, καθαρών στάσιμων και ημιστάσιμων στρατηγικών, αντίστοιχα, του decision maker.

Για μια συμπεριφορική στρατηγική  $f \in F_B$ , η αναμενόμενη πληρωμή του οριακού λόγου των μέσων ορίζεται από την:

$$\widehat{v}_a(s, f) = \liminf_{t \rightarrow \infty} \frac{E_f[\sum_{m=1}^t (\hat{r}(X_m, A_m) | X_1 = s)]}{E_f[\sum_{m=1}^t (\bar{\tau}(X_m, A_m) | X_1 = s)]} , \text{ για όλα τα } s \in S \quad (3.3.1)$$

όπου

$$\bar{\tau}(s, a) = \sum_{s' \in S} \left( \hat{p}(s' | a) \int_0^\infty t d\hat{Q}_{ss'}(t | a) \right)$$

είναι ο αναμενόμενος χρόνος παραμονής στην κατάσταση  $s$  όταν ο decision-maker έχει επιλέξει την ενέργεια  $a \in \hat{A}(s)$ .

### Παρατήρηση

Ο τρόπος που ορίσαμε την αναμενόμενη πληρωμή των παικτών, μέσω της (3.2.1) για τα παίγνια και μέσω της (3.3.1) για διαδικασίες απόφασης, προφανώς δεν είναι μοναδικός. Ωστόσο, είναι αυτός που απαντάται συχνότερα στη βιβλιογραφία και θα παραθέσουμε ορισμένους λόγους για τους οποίους αυτό συμβαίνει. Ως εναλλακτική συνάρτηση πληρωμής των παικτών θα μπορούσαμε να θεωρήσουμε την:

$$\varphi(s, f) = E_f \left[ \liminf_{n \rightarrow \infty} \frac{\sum_{m=0}^n r(X_m, A_m)}{\sum_{m=0}^n \tau(X_m, A_m)} \mid X_0 = s \right]$$

ή και την :

$$\psi(s, f) = \liminf_{n \rightarrow \infty} E_f \left[ \frac{\sum_{m=0}^n r(X_m, A_m)}{\sum_{m=0}^n \tau(X_m, A_m)} \mid X_0 = s \right]$$

Ο κυριότερος λόγος που επιλέγουμε η πληρωμή να ορίζεται από την σχέση (3.3.1), έναντι των  $\varphi$  και  $\psi$ , είναι πως γνωρίζουμε την ύπαρξη βέλτιστης ημι-στάσιμης στρατηγικής για αυτό το κριτήριο, από την εργασία [14], ενώ κάτι τέτοιο δεν είναι ακόμα γνωστό για τα δύο εναλλακτικά κριτήρια πληρωμής. Επιπλέον, η σχέση (3.3.1) φαίνεται να είναι πιο ‘εύκολη’ στον χειρισμό, εφόσον τόσο ο αριθμητής όσο και ο παρονομαστής μπορούν να υπολογιστούν ακολουθώντας τη λογική μιας Μαρκοβιανής διαδικασίας πεπερασμένου ορίζοντα, όπως αυτή ορίστηκε στο 1<sup>ο</sup> Κεφάλαιο, και έπειτα η συνολική αναμενόμενη πληρωμή θα προκύπτει ως το  $\liminf$  του πηλίκου που υπολογίστηκε. Απ’την άλλη, κάτι τέτοιο δεν είναι εφικτό για την  $\varphi$ , αφού θα πρέπει να υπολογιστεί το  $\liminf$  για κάθε άπειρου μήκους ιστορία της διαδικασίας και έπειτα να πάρουμε την μέση τιμή αυτής, ενώ για το κριτήριο  $\psi$  θα πρέπει να υπολογιστεί η μέση τιμή για μια άπειρη ακολουθία τυχαίων μεταβλητών, ώστε να πάρουμε το  $\liminf$  αυτής.

Τέλος, ένας εναλλακτικός ορισμός θα ήταν να αντικαταστήσουμε στη σχέση (3.3.1) (και στην (3.2.1)) το  $\liminf$  με  $\limsup$ . Αυτό δεν θα είχε καμία διαφορά στην κλάση των στάσιμων στρατηγικών (επειδή το συγκεκριμένο όριο υπάρχει), αλλά αποτελεί, εν γένει, ένα διαφορετικό κριτήριο πληρωμής για την ευρύτερη κλάση των συμπεριφορικών στρατηγικών.

### Ορισμός 3.3.1

Μία στρατηγική  $f^*$  θα καλείται *βέλτιστη*, με κριτήριο πληρωμής τον οριακό λόγο των μέσων, εάν

$$\widehat{v}_a(s, f^*) \geq \widehat{v}_a(s, f)$$

για όλα τα  $f \in F_B$  και όλα τα  $s \in S$ . Έστω

$$\widehat{v}_a(s) = \sup_{f \in F_B} \{\widehat{v}_a(s, f)\}$$

Τότε, η τιμή  $\widehat{v}_a(s)$  θα καλείται η *τιμή* της ημιμαρκοβιανής διαδικασίας απόφασης για την αρχική κατάσταση  $s$ , ενώ το διάνυσμα

$$\widehat{v}_a = [\widehat{v}_a(s)]_{s \in S}$$

θα είναι το *διάνυσμα τιμής* για την διαδικασία.

### Θεώρημα 3.3.1

Σε μια πεπερασμένη ημιμαρκοβιανή διαδικασία απόφασης με κριτήριο πληρωμής τον οριακό λόγο των μέσων, υπάρχει καθαρή ημιστάσιμη στρατηγική  $f_p^* \in \xi$ , η οποία είναι βέλτιστη για τον decision-maker.

Το κυριότερο αποτέλεσμα αυτού του κεφαλαίου συνοψίζεται στο επόμενο σημαντικό θεώρημα. Αρχικά διατυπώνουμε το θεώρημα και έπειτα παραθέτουμε την απόδειξή του. Η συλλογιστική πορεία της απόδειξης θα φανεί χρήσιμη στη συνέχεια, όταν θα είμαστε σε θέση να εφαρμόσουμε τον αλγόριθμο επίλυσης τέτοιου τύπου παιγνίων.

### Θεώρημα 3.3.2

Οποιοδήποτε ημιμαρκοβιανό παίγνιο δύο παικτών μηδενικού αθροίσματος και τέλει πληροφόρησης, με κριτήριο πληρωμής τον οριακό λόγο των μέσων, έχει λύση στο σύνολο των καθαρών ημιστάσιμων στρατηγικών.

#### Απόδειξη

Θα δείξουμε το παραπάνω θεώρημα δείχνοντας την στρατηγική ισοδυναμία μεταξύ του ημιμαρκοβιανού παιγνίου και της αντίστοιχης ημιμαρκοβιανής διαδικασίας απόφασης. Έστω λοιπόν το ημιμαρκοβιανό παίγνιο

$$\Gamma = \langle S = S_1 \cup S_2, A^1 = \{A^1(s) \mid s \in S_1\}, A^2 = \{A^2(s) \mid s \in S_2\}, p, Q, r \rangle$$

δύο παικτών, μηδενικού αθροίσματος και τέλει πληροφόρησης με πληρωμή που προκύπτει από τον οριακό λόγο των μέσων, όπου σε  $|S_1|$  στο πλήθος καταστάσεις ο παίκτης 2 διαθέτει μόνο μία ενέργεια για να επιλέξει, ενώ το ίδιο ισχύει για τον παίκτη 1 για τις καταστάσεις  $\{|S_1| + 1, |S_1| + 2, \dots, |S_1| + |S_2|\}$ . Από αυτό το ημιμαρκοβιανό παίγνιο θα κατασκευάσουμε μία διαδικασία απόφασης  $\hat{\Gamma}$  με τον εξής τρόπο. Έστω :

$$\hat{\Gamma} = \langle S = S_1 \cup S_2, \hat{A} = \{\hat{A}(s) = A^1(s) \mid s \in S_1\} \cup \{\hat{A}(s) = A^2(s) \mid s \in S_2\}, \hat{p} = p, \hat{Q} = Q, \hat{r} \rangle$$

όπου το  $\hat{r}(s, a)$  ορίζεται ως:

$$\hat{r}(s, a) = \begin{cases} r(s, a), & \text{εάν } s \in S_1, a \in \hat{A}(s) = A^1(s) \\ -r(s, a), & \text{εάν } s \in S_2, a \in \hat{A}(s) = A^2(s) \end{cases}$$

Από το Θεώρημα 3.3.1 γνωρίζουμε πως υπάρχει βέλτιστη καθαρή ημιστάσιμη στρατηγική στο  $\hat{\Gamma}$ . Έστω  $\widehat{f}^*$  μια καθαρή ημιστάσιμη στρατηγική, η οποία είναι βέλτιστη για τη διαδικασία  $\hat{\Gamma}$ . Έστω  $s_0$  μία αρχική κατάσταση, οπότε :

$$\widehat{f}^* = \{\widehat{f}^*(s_0, s) \mid \widehat{f}^*(s_0, s) \in P(\hat{A}(s)) \text{ για κάθε } s_0, s \in S\}$$

όπου με  $P(\hat{A}(s))$  συμβολίζουμε το σύνολο όλων των κατανομών πιθανότητας πάνω στο χώρο  $\hat{A}(s)$ . Στη συνέχεια, κατασκευάζουμε ένα ζεύγος καθαρών ημιστάσιμων στρατηγικών  $(f^*, g^*)$  στο παίγνιο  $\Gamma$ , που προκύπτουν από την  $\widehat{f}^*$ , ως εξής:

$$f^*(s_0, s) = \begin{cases} \widehat{f}^*(s_0, s) & s \in S_1 \\ 1 & s \in S_2 \end{cases}$$

$$g^*(s_0, s) = \begin{cases} 1 & s \in S_1 \\ \widehat{f}^*(s_0, s) & s \in S_2 \end{cases}$$

Συμβολίζουμε με  $v_a$  και  $\widehat{v}_a$  τις συναρτήσεις πληρωμών για το παίγνιο  $\Gamma$  (για τον παίκτη 1) και για την διαδικασία  $\hat{\Gamma}$ , αντίστοιχα.

Από το θεώρημα 9 του [11] γνωρίζουμε ότι

Εάν υπάρχει ζεύγος καθαρών ημιστάσιμων στρατηγικών  $(f^*, g^*)$  τέτοιο ώστε για όλα τα  $s \in S$  να ισχύει:

$$v_a(s, f, g^*) \leq v_a(s, f^*, g^*) \leq v_a(s, f^*, g) \text{ για όλα τα } (f, g) \in \xi_1 \times \xi_2$$

τότε το  $(f^*, g^*)$  είναι ένα ζεύγος βέλτιστων στρατηγικών του παιχνιδιού. Μπορούμε, λοιπόν, να εστιάσουμε στο σύνολο των ημιστάσιμων στρατηγικών αντί του μεγαλύτερου συνόλου των συμπεριφορικών στρατηγικών.

Έστω τώρα ότι σταθεροποιούμε την στρατηγική  $g^*$  για τις καταστάσεις  $\{|S_1| + 1, |S_1| + 2, \dots, |S_1| + |S_2|\}$  στην διαδικασία απόφασης  $\hat{\Gamma}$ . Τότε αναγώμαστε σε ένα νέο μοντέλο ημιμαρκοβιανής διαδικασίας απόφασης  $\hat{\Gamma}_1$ , όπου για τις καταστάσεις  $\{|S_1| + 1, |S_1| + 2, \dots, |S_1| + |S_2|\}$ , ο decision-maker ακολουθεί την στρατηγική  $g^*(s_0, |S_1| + j)$ , όπου  $j \in \{1, 2, \dots, |S_2|\}$ . Ομοίως, μπορούμε να θεωρήσουμε την διαδικασία  $\hat{\Gamma}_2$ , σταθεροποιώντας την στρατηγική  $f^*$  για τις καταστάσεις  $\{1, 2, \dots, |S_1|\}$  της διαδικασίας  $\hat{\Gamma}$ . Εάν μπορέσουμε να αποδείξουμε ότι η  $f^*$  είναι μία βέλτιστη καθαρή ημιστάσιμη στρατηγική για την  $\hat{\Gamma}_1$ , με σταθεροποιημένη την  $g^*$ , και ότι η  $g^*$  είναι μία βέλτιστη καθαρή ημιστάσιμη στρατηγική για την  $\hat{\Gamma}_2$ , με σταθεροποιημένη την  $f^*$ , τότε θα έχουμε καταλήξει στο γεγονός ότι το  $(f^*, g^*)$  αποτελεί ένα ζεύγος βέλτιστων καθαρών ημιστάσιμων στρατηγικών για το παίγνιο  $\Gamma$ .

### Λήμμα 3.3.1

Η  $\mathbf{f}^*$  είναι βέλτιστη καθαρή ημιστάσιμη στρατηγική για τον παίκτη 1 στην διαδικασία  $\hat{\Gamma}_1$ .

#### Απόδειξη (Λήμμα 3.3.1)

Με απαγωγή σε άτοπο. Ας υποθέσουμε ότι η  $\mathbf{f}^*$  δεν είναι βέλτιστη καθαρή ημιστάσιμη στρατηγική στην  $\hat{\Gamma}_1$ . Έστω  $\widehat{v}_1$  η συνάρτηση συνολικής πληρωμής για την διαδικασία  $\hat{\Gamma}_1$ . Έστω, ακόμη, ότι η  $\mathbf{f}_1^*$  είναι μία βέλτιστη καθαρή ημιστάσιμη στρατηγική στην  $\hat{\Gamma}_1$ . Τότε, θα έχουμε ότι :

$$\widehat{v}_1(s, \mathbf{f}_1^*) \geq \widehat{v}_1(s, \mathbf{f}^*) \quad \forall s \in (S_1 \cup S_2)$$

και θα πρέπει να υπάρχει τουλάχιστον μία κατάσταση για την οποία η ανισότητα είναι αυστηρή. Ας υποθέσουμε ότι για την  $s_1 \in S$  ισχύει ότι  $\widehat{v}_1(s_1, \mathbf{f}_1^*) > \widehat{v}_1(s_1, \mathbf{f}^*)$

Αυτό όμως σημαίνει ότι  $v_a(s_1, \mathbf{f}_1^*, \mathbf{g}^*) > v_a(s_1, \mathbf{f}^*, \mathbf{g}^*)$

Μπορούμε τώρα να κατασκευάσουμε την καθαρή ημιστάσιμη στρατηγική  $\widehat{\mathbf{f}}_1^*$ , η οποία θα ταυτίζεται με την  $\mathbf{f}_1^*$  στο  $S_1$  και με την  $\mathbf{g}^*$  στο  $S_2$ , δηλαδή :

$$\widehat{\mathbf{f}}_1^*(s_0, s) = \begin{cases} \mathbf{f}_1^*(s_0, s) & , s \in S_1 \\ \mathbf{g}^*(s_0, s) & , s \in S_2 \end{cases}$$

Όμως, τότε έχουμε ότι  $\widehat{v}_a(s_1, \widehat{\mathbf{f}}_1^*) > \widehat{v}_a(s_1, \widehat{\mathbf{f}}^*)$ , το οποίο είναι αντίφαση, εφόσον γνωρίζουμε ότι η  $\widehat{\mathbf{f}}^*$  είναι βέλτιστη για την  $\hat{\Gamma}$ .

Οπότε συμπεραίνουμε πως, πράγματι, η  $\mathbf{f}^*$  είναι βέλτιστη καθαρή ημιστάσιμη στρατηγική στην  $\hat{\Gamma}_1$ . ■

#### Σχόλιο

Με όμοια επιχειρήματα μπορεί να δείξει κανείς ότι η στρατηγική  $\mathbf{g}^*$  είναι μια βέλτιστη καθαρή ημιστάσιμη στρατηγική για την ημιμαρκοβιανή διαδικασία  $\hat{\Gamma}_2$ .

Τέλος, χρησιμοποιώντας το Λήμμα (3.3.1) και το παραπάνω σχόλιο, μπορούμε πλέον να συμπεράνουμε πως το  $(\mathbf{f}^*, \mathbf{g}^*)$  αποτελεί ζεύγος βέλτιστων καθαρών ημιστάσιμων στρατηγικών και για τους δύο παίκτες στο ημιμαρκοβιανό παίγνιο τέλειας πληροφόρησης  $\hat{\Gamma}$ . ■

### 3.4 Αλγόριθμος επίλυσης ημιμαρκοβιανών παιγνίων τέλειας πληροφόρησης

Στη βιβλιογραφία έχει αναπτυχθεί ένας αλγόριθμος επίλυσης ημιμαρκοβιανών διαδικασιών απόφασης με κριτήριο πληρωμής τον οριακό λόγο των μέσων. Για την επίλυση ενός παιγνίου με τα χαρακτηριστικά που έχουμε μελετήσει, αρχικά θα το υποβιβάσουμε σε ημιμαρκοβιανή διαδικασία και έπειτα θα χρησιμοποιήσουμε τον αλγόριθμο από το [12]. Προκειμένου να υπολογιστεί η τιμή του παιχνιδιού καθώς και βέλτιστες στρατηγικές των παικτών, θα ακολουθήσουμε τη λογική της απόδειξης του προηγούμενου Θεωρήματος.

#### Αλγόριθμος 3.4.1.

Έστω  $s_0$  μια αυθαίρετη αρχική κατάσταση. Θεωρούμε το ακόλουθο πρόβλημα γραμμικού προγραμματισμού για τις μεταβλητές  $v(s_0)$ ,  $g = (g_s \mid s \in S)$ ,  $h = (h_s \mid s \in S)$  ως εξής :

$$LP : \min v(s_0)$$

υπό τους περιορισμούς

$$g_s \geq \sum_{s' \in S} \hat{p}(s'|s, a) g_{s'} \quad \forall s \in S, \forall a \in \hat{A}(s)$$

$$g_s + h_s \geq r(s, a) - v(s_0) \bar{r}(s, a) + \sum_{s' \in S} \hat{p}(s'|s, a) h_{s'} \quad \forall s \in S, \forall a \in \hat{A}(s)$$

$$g_{s_0} \leq 0$$

Οι μεταβλητές  $v(s_0)$ ,  $(g_s \mid s \in S \setminus \{s_0\})$ ,  $(h_s \mid s \in S)$  δεν έχουν περιορισμό στο πρόσημο.

Το αντίστοιχο δυϊκό γραμμικό πρόγραμμα για τις μεταβλητές

$$x = (x_{sa} \mid s \in S, a \in \hat{A}(s)), y = (y_{sa} \mid s \in S, a \in \hat{A}(s)) \quad \text{και } t, \text{ θα είναι :}$$

$$DLP : \max R_s$$

$$\text{όπου } R_s = \sum_{s \in S} \sum_{a \in \hat{A}(s)} \hat{r}(s, a) x_{sa} \quad (3a)$$

υπό τους περιορισμούς :

$$\sum_{s \in S} \sum_{a \in \hat{A}(s)} \{ \delta_{ss'} - \hat{p}(s'|s, a) \} x_{sa} = 0, \quad \forall s' \in S \quad (3b)$$

$$\sum_{a \in \hat{A}(s)} x_{sa} + \sum_{s \in S} \sum_{a \in \hat{A}(s)} \{ \delta_{ss'} - \hat{p}(s'|s, a) \} y_{sa} = 0, \quad \forall s' \in S \setminus \{s_0\} \quad (3c)$$

$$\sum_{a \in \hat{A}(s_0)} x_{s_0 a} + \sum_{s \in S} \sum_{a \in \hat{A}(s)} \{ \delta_{ss_0} - \hat{p}(s_0|s, a) \} y_{sa} - t = 0 \quad (3d)$$

$$\sum_{s \in S} \sum_{a \in \hat{A}(s)} \bar{\tau}(s, a) x_{sa} = 1 \quad (3e)$$

$$x_{sa} \geq 0, y_{sa} \geq 0 \quad \forall s \in S$$

$$a \in \hat{A}(s)$$

$$t \geq 0$$

(όπου το  $\delta_{ss'}$  είναι το Δέλτα του Kronecker)

Για την εύρεση μιας εφικτής λύσης  $(x, y, t)$  του *DLP*, ορίζουμε τα ακόλουθα σύνολα :

$$S_x = \left\{ s \in S \mid \sum_{a \in \hat{A}(s)} x_{sa} > 0 \right\}$$

$$S_y = \left\{ s \in S \mid \sum_{a \in \hat{A}(s)} x_{sa} = 0 \text{ και } \sum_{a \in \hat{A}(s)} y_{sa} > 0 \right\}$$

$$S_{xy} = \left\{ s \in S \mid \sum_{a \in \hat{A}(s)} x_{sa} = 0 \text{ και } \sum_{a \in \hat{A}(s)} y_{sa} = 0 \right\}$$

Οπότε θα ισχύει ότι  $S = S_x \cup S_y \cup S_{xy}$ , και τα σύνολα  $S_x, S_y$  και  $S_{xy}$  είναι ανά δύο ξένα μεταξύ τους. Μία καθαρή στάσιμη στρατηγική που αντιστοιχεί στην εφικτή λύση  $(x, y, t)$  του *DLP* ορίζεται ως  $f_{xyt}^{ps_0}$ , όπου  $s_0$  είναι η αυθαίρετη αρχική κατάσταση που σταθεροποιήσαμε εξ αρχής.

$f_{xyt}^{ps_0}(s) = a_s, s \in S$ , τέτοια ώστε:

$$a_s = \begin{cases} a, & \text{εάν } s \in S_x \text{ και } x_{sa} > 0 \\ a', & \text{εάν } s \in S_y \text{ και } y_{sa'} > 0 \\ \text{αυθαίρετο,} & \text{εάν } s \in S_{xy} \end{cases}$$

Στο σημείο αυτό, θα χρειαστούμε το παρακάτω σημαντικό αποτέλεσμα, το οποίο παρατίθεται στο [12].

#### Θεώρημα 3.4.1

Έστω  $(x^*, y^*, t^*)$  μία βέλτιστη λύση του  $DLP$ . Τότε η  $f_{x^*y^*t^*}^{ps_0}$  είναι μία βέλτιστη καθαρή στάσιμη στρατηγική της ημιμαρκοβιανής διαδικασίας, για την αρχική κατάσταση  $s_0$ .

Από το παραπάνω θεώρημα εξάγουμε το συμπέρασμα πως μία βέλτιστη καθαρή ημιστάσιμη στρατηγική μιας ημιμαρκοβιανής διαδικασίας απόφασης μπορεί να βρεθεί μέσω της βέλτιστης λύσης του δυϊκού γραμμικού προγράμματος. Για το λόγο αυτό, ο συγκεκριμένος αλγόριθμος είναι ιδιαίτερα χρήσιμος, καθώς θα μας επιτρέψει να εξάγουμε βέλτιστες στρατηγικές και για τους δύο παίκτες στο ημιμαρκοβιανό παίγνιο τέλειας πληροφόρησης, ακολουθώντας την συλλογιστική πορεία της απόδειξης του Θεωρήματος 3.3.2.

Για μία σταθεροποιημένη αρχική κατάσταση  $s_0$ , έχουμε μία βέλτιστη καθαρή στάσιμη στρατηγική, την  $f_{x^*y^*t^*}^{ps_0}$ , όπου  $f_{x^*y^*t^*}^{ps_0}(s) = a_s^*$  και το  $a_s^*$  είναι όπως ορίστηκε προηγουμένως. Έτσι, για διαφορετικές αρχικές καταστάσεις, λαμβάνουμε μία βέλτιστη στρατηγική για τον decision-maker, στην διαδικασία απόφασης. Τέλος, θα χρειαστούμε ακόμη ένα πόρισμα, το Πόρισμα 1 από το [12].

#### Πόρισμα 3.4.1

Έστω  $f_{s_0}^p$  να είναι μία καθαρή στάσιμη στρατηγική που εξήχθη από μία βέλτιστη λύση του  $DLP$  για αρχική κατάσταση  $s_0 \in S$ . Τότε η καθαρή ημιστάσιμη στρατηγική

$$\xi^{p^*} = (f_1^{p^*}, f_2^{p^*}, \dots, f_N^{p^*})$$

θα είναι βέλτιστη για την ημιμαρκοβιανή διαδικασία με κριτήριο τον οριακό λόγο των μέσων.

Οπότε από το τελευταίο πόρισμα, λαμβάνουμε μία βέλτιστη στρατηγική  $\hat{f}^*$  του decision-maker στην διαδικασία  $\hat{\Gamma}$ , όπου :

$$\hat{f}^* = (f_1^*, f_2^*, \dots, f_s^*, \dots, f_N^*)$$

και  $f_s^* = f_{x^*y^*t^*}^{ps}$  αποτελεί μία βέλτιστη καθαρή στάσιμη στρατηγική της ημιμαρκοβιανής διαδικασίας με αρχική κατάσταση  $s$ .

Το ζευγάρι  $(f^*, g^*)$  βέλτιστων καθαρών στρατηγικών για τους δύο παίκτες στο παίγνιο  $\Gamma$ , δίνεται όπως είδαμε και παραπάνω από :



$$f^*(s_0, s) = \begin{cases} \widehat{f^*}(s_0, s) & s \in S_1 \\ 1 & s \in S_2 \end{cases}$$

$$g^*(s_0, s) = \begin{cases} 1 & s \in S_1 \\ \widehat{f^*}(s_0, s) & s \in S_2 \end{cases}$$

Εάν  $R_s$  είναι η αντικειμενική συνάρτηση του  $DLP$  για αρχική κατάσταση  $s$ , τότε το

$$\hat{v} = (\hat{v}(1), \hat{v}(2), \dots, \hat{v}(s), \dots, \hat{v}(N))$$

θα είναι το διάνυσμα τιμής της διαδικασίας  $\hat{F}$ , όπου  $\hat{v}(s) = \max R_s$ .

Το διάνυσμα τιμής για το παίγνιο  $\Gamma$  θα δίνεται από :

$$v = (\hat{v}(1), \hat{v}(2), \dots, \hat{v}(|S_1|), -\hat{v}(|S_1| + 1), \dots, -\hat{v}(|S_1| + |S_2|))$$

Το αρνητικό πρόσημο, που αφορά στις καταστάσεις  $\{|S_1| + 1, \dots, |S_1| + |S_2|\}$ , δικαιολογείται από το γεγονός ότι ο παίκτης 2 θέλει να ελαχιστοποιήσει το κόστος για αυτές τις καταστάσεις.

### *Υπολογισμός βέλτιστης στρατηγικής μέσω πλήρους απαρίθμησης*

Σύμφωνα με την μέθοδο της πλήρους απαρίθμησης, θα υπολογίσουμε την τιμή της ημιμαρκοβιανής διαδικασίας για μία καθαρή στάσιμη στρατηγική  $f$ , κατ' αναλογία με την σχέση (3.2.1) που αφορά τα παίγνια, ως εξής:

$$\hat{v}(s, f) = \frac{[P^*(f)\hat{r}(f)](s)}{[P^*(f)\bar{\tau}(f)](s)}, \quad \text{για όλα τα } s \in S$$

όπου  $\hat{r}(f)$  είναι το διάνυσμα αμοιβής και  $\bar{\tau}(f)$  είναι το διάνυσμα του αναμενόμενου χρόνου παραμονής.

Εν συνεχεία, παραθέτουμε σε βήματα τον αλγόριθμο υπολογισμού του πίνακα ορίου κατά Cesaro για μια ημιμαρκοβιανή διαδικασία με  $N$  καταστάσεις,

**Είσοδος:** Ο πίνακας πιθανοτήτων μετάβασης  $P \in M_N(\mathbb{R})$

**Έξοδος:** Ο πίνακας όριο κατά Cesaro  $P^* \in M_N(\mathbb{R})$  του  $P$

**Βήμα 1:** Ορίζουμε το χαρακτηριστικό πολυώνυμο  $Ch_p(z) = |P - zI_n|$

**Βήμα 2:** Διαιρούμε το  $Ch_p(z)$  με το πολυώνυμο  $(z - 1)^{m(1)}$ , όπου  $m(1)$  είναι η αλγεβρική πολλαπλότητα της ιδιοτιμής  $z_0 = 1$ . Θα συμβολίζουμε το πηλίκο της διαίρεσης με  $D(z)$ .

**Βήμα 3:** Υπολογίζουμε τον πίνακα  $R = D(Q)$

**Βήμα 4:** Ορίζουμε τον πίνακα όριο  $P^*$  διαιρώντας τον πίνακα  $R$  με το άθροισμα των στοιχείων μιας αυθαίρετης γραμμής του.

Κλείνοντας αυτό το κεφάλαιο, παραθέτουμε ένα παράδειγμα, αναλυτικά λυμένο στο προγραμματιστικό περιβάλλον MATLAB. Αρχικά, θα γίνει επίλυση μέσω γραμμικού προγραμματισμού και έπειτα θα συγκρίνουμε τη λύση με εκείνη που θα λάβουμε από τη μέθοδο πλήρους απαρίθμησης. Φυσικά, αναμένεται τα αποτελέσματα όσον αφορά τόσο την τιμή του παιγνίου όσο και τη βέλτιστη στρατηγική, να ταυτίζονται.

### 3.5 Επίλυση μέσω γραμμικού προγραμματισμού και μέσω πλήρους απαρίθμησης

Έστω το ακόλουθο παίγνιο  $\Gamma$ , όπου οι παίκτες 1 και 2 θα παίζουν σύμφωνα με τους παρακάτω πίνακες:

1.5	$(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$	1
1	$(\frac{1}{3}, 0, \frac{2}{3}, 0)$	0.8

Κατάσταση 1

3	$(\frac{1}{2}, \frac{1}{2}, 0, 0)$	1.5
2	$(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, 0)$	1

Κατάσταση 2

3	$(\frac{1}{2}, 0, \frac{1}{2}, 0)$	1
7	$(\frac{1}{2}, 0, 0, \frac{1}{2})$	2

Κατάσταση 3

2	$(\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$	2
1	$(\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$	1.2

Κατάσταση 4

Η δομή του παραπάνω παιγνίου ερμηνεύεται ως εξής:

$r$
$(p_1, p_2, p_3, 1 - p_1 - p_2 - p_3)$
$\bar{r}$

όπου με  $r$  συμβολίζεται η αμοιβή που θα λάβει ο παίκτης 1 (και ο παίκτης 2 θα λαμβάνει αμοιβή  $-r$ ), με  $p_i$  συμβολίζεται η πιθανότητα η επόμενη μετάβαση να γίνει στην κατάσταση  $i$  (υποχρεωτικά, για κάθε κατάσταση, τα  $p_i$  θα έχουν άθροισμα ίσο με μονάδα) και με  $\bar{r}$  συμβολίζεται η αναμενόμενη τιμή του χρόνου παραμονής στην κατάσταση που βρίσκεται η διαδικασία.

Αρχικά, προκειμένου να γίνει η ανάλυση του παραπάνω παιγνίου, θα το μετατρέψουμε πρώτα σε διαδικασία απόφασης για τον παίκτη 1. Με ανάλογο τρόπο μπορεί να γίνει η μετατροπή σε διαδικασία απόφασης για τον δεύτερο παίκτη. Ορίζεται, λοιπόν η διαδικασία  $\hat{I}_1$ , με τον τρόπο που περιγράφεται στην απόδειξη του Θεωρήματος 3.3.2., ως εξής:

1.5 $(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$ 1	3 $(\frac{1}{2}, \frac{1}{2}, 0, 0)$ 1.5	-3 $(\frac{1}{2}, 0, \frac{1}{2}, 0)$ 1	-2 $(\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$ 2
1 $(\frac{1}{3}, 0, \frac{2}{3}, 0)$ 0.8	2 $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, 0)$ 1	-7 $(\frac{1}{2}, 0, 0, \frac{1}{2})$ 2	-1 $(\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$ 1.2
Κατάσταση 1	Κατάσταση 2	Κατάσταση 3	Κατάσταση 4

α) Επίλυση μέσω γραμμικού προγραμματισμού

Ακολουθώντας τα βήματα του αλγορίθμου που περιγράφηκαν παραπάνω, θα αναζητήσουμε λύση της διαδικασίας  $\hat{I}_1$  επιλύοντας το αντίστοιχο δυϊκό πρόβλημα γραμμικού προγραμματισμού.

Το ζητούμενο είναι η μεγιστοποίηση της  $R_s$ , η οποία από την σχέση (3α) προκύπτει να είναι:

$$R_s = 1.5x_{11} + x_{12} + 3x_{21} + 2x_{22} - 3x_{31} - 7x_{32} - 2x_{41} - x_{42}$$

Το σύνολο περιορισμών που προκύπτουν από την (3b) είναι:

$$9x_{11} + 8x_{12} - 6x_{12} - 4x_{22} - 6x_{31} - 6x_{32} - 3x_{41} - 6x_{42} = 0$$

$$-3x_{11} + 6x_{21} + 8x_{22} - 6x_{41} - 3x_{42} = 0$$

$$-3x_{11} - 8x_{12} - 4x_{22} + 6x_{31} + 12x_{32} - 3x_{41} - 3x_{42} = 0$$

$$-x_{11} - 2x_{32} + 4x_{41} + 4x_{42} = 0$$

Το σύνολο περιορισμών που προκύπτει από τις (3c) και (3d) είναι:

$$12x_{11} + 12x_{22} + 9y_{11} + 8y_{12} - 6y_{21} - 4y_{22} - 6y_{31} - 6y_{32} - 3y_{41} - 6y_{42} - 12\delta_{s_01}t = 0$$

$$12x_{21} + 12x_{22} - 3y_{11} + 6y_{21} + 8y_{22} - 6y_{41} - 3y_{42} - 12\delta_{s_02}t = 0$$

$$12x_{31} + 12x_{32} - 3y_{11} - 8y_{12} - 4y_{22} + 6y_{31} + 12y_{32} - 3y_{41} - 3y_{42} - 12\delta_{s_03}t = 0$$

$$4x_{41} + 4x_{42} - y_{11} - 2y_{32} + 4y_{41} + 4y_{42} - 4$$

Τέλος, από την (3e) λαμβάνουμε τον περιορισμό:

$$x_{11} + 0.8x_{12} + 1.5x_{21} + x_{22} + x_{31} + 2x_{32} + 2x_{41} + 1.2x_{42} = 1$$

και

$$t \geq 0, x_{sa} \geq 0, y_{sa} \geq 0 \text{ για κάθε } s \in S, a \in \hat{A}(s)$$

Χρησιμοποιώντας την συνάρτηση linprog, λαμβάνουμε από το MATLAB, για κάθε πιθανή αρχική κατάσταση, τα ακόλουθα αποτελέσματα:

(i)  $s_0 = 1$

$x_{11}$	0	$y_{11}$	0	$t$	0.7724
$x_{12}$	0.3239	$y_{12}$	1.0590		
$x_{21}$	0	$y_{21}$	0		
$x_{22}$	0.0498	$y_{22}$	0		
$x_{31}$	0	$y_{31}$	0		
$x_{32}$	0.2658	$y_{32}$	0.4651		
$x_{41}$	0	$y_{41}$	0.0997		
$x_{42}$	0.1329	$y_{42}$	0		

$$\max R_1 = 1.5698$$

(ii)  $s_0 = 2$

$x_{11}$	0	$y_{11}$	0	$t$	0.7724
$x_{12}$	0.3239	$y_{12}$	0.2554		
$x_{21}$	0	$y_{21}$	0		
$x_{22}$	0.0498	$y_{22}$	1.0839		
$x_{31}$	0	$y_{31}$	0		
$x_{32}$	0.2658	$y_{32}$	0.2658		
$x_{41}$	0	$y_{41}$	0		
$x_{42}$	0.1329	$y_{42}$	0		

$\max R_2 = 1.5698$

(iii)  $s_0 = 3$

$x_{11}$	0	$y_{11}$	0	$t$	0.7724
$x_{12}$	0.3239	$y_{12}$	0		
$x_{21}$	0	$y_{21}$	0		
$x_{22}$	0.0498	$y_{22}$	0		
$x_{31}$	0	$y_{31}$	0.1329		
$x_{32}$	0.2658	$y_{32}$	0.4651		
$x_{41}$	0	$y_{41}$	0.0997		
$x_{42}$	0.1329	$y_{42}$	0		

$\max R_3 = 1.5698$

(iv)  $s_0 = 4$

$x_{11}$	0	$y_{11}$	0	$t$	0.7724
$x_{12}$	0.3239	$y_{12}$	0		
$x_{21}$	0	$y_{21}$	0		
$x_{22}$	0.0498	$y_{22}$	0.4049		
$x_{31}$	0	$y_{31}$	0.0581		
$x_{32}$	0.2658	$y_{32}$	0		
$x_{41}$	0	$y_{41}$	0.6395		
$x_{42}$	0.1329	$y_{42}$	0		

$\max R_4 = 1.5698$

Από τα παραπάνω αποτελέσματα, γίνεται φανερό ότι το διάνυσμα τιμής της διαδικασίας  $\widehat{f}_1$  θα είναι ίσο με :

$$\widehat{\varphi} = (1.5698, 1.5698, 1.5698, 1.5698)$$

Έστω τώρα  $\widehat{f}_1^* = (f_1^*, f_2^*, f_3^*, f_4^*)$  μια βέλτιστη καθαρή ημι-στάσιμη στρατηγική της διαδικασίας  $\widehat{f}_1$ , όπου με  $f_i^*, i = 1, 2, 3, 4$ , θα συμβολίζουμε μια βέλτιστη στρατηγική της διαδικασίας, δεδομένου ότι η αρχική κατάσταση είναι η  $i$ .

Τότε, σύμφωνα με την σχέση  $f_{x^*y^*t^*}^{ps_0}(s) = a_s^*$ , υπολογίζεται ότι :

$$f_1^* = f_{x^*y^*t^*}^{p1}(s) = a_{s1}^* = (2, 2, 2, 2)$$

$$f_2^* = f_{x^*y^*t^*}^{p2}(s) = a_{s2}^* = (2, 2, 2, 2)$$

$$f_3^* = f_{x^*y^*t^*}^{p3}(s) = a_{s3}^* = (2, 2, 2, 2)$$

$$f_4^* = f_{x^*y^*t^*}^{p4}(s) = a_{s4}^* = (2, 2, 2, 2)$$

Υπολογίστηκε, λοιπόν, η βέλτιστη στρατηγική για την διαδικασία να είναι η

$$\widehat{f}_1^* = (f_1^*, f_2^*, f_3^*, f_4^*) = (f_1^*, f_1^*, f_1^*, f_1^*)$$

και η  $f_1^*$  είναι η καθαρή στρατηγική που υποδεικνύει την επιλογή της 2<sup>ης</sup> διαθέσιμης ενέργειας, για κάθε κατάσταση  $s \in \{1, 2, 3, 4\}$ , για κάθε αρχική κατάσταση  $s_0$ .

### β) Επίλυση μέσω της μεθόδου πλήρους απαρίθμησης

Σε αυτή την μέθοδο και πάλι ζητούμενο είναι ο υπολογισμός της αξίας του παιχνιδιού και η εύρεση μιας βέλτιστης στρατηγικής για κάθε πιθανής αρχική κατάσταση  $s_0$ , για την διαδικασία  $\widehat{f}_1$ .

Αρχικά, παρατηρούμε ότι υπάρχουν συνολικά  $2^4$  διαθέσιμες καθαρές στρατηγικές για την διαδικασία  $\widehat{f}_1$ , ακριβώς όσοι είναι οι τρόποι με τους οποίους μπορούμε να κατανείμουμε τις 2 διαθέσιμες ενέργειες κάθε κατάστασης σε ένα 4-διάστατο διάνυσμα, κάθε συντεταγμένη του οποίου αντιστοιχεί σε μία εκ των 4 καταστάσεων της διαδικασίας.

Υλοποιώντας τον αλγόριθμο που αφορά στον υπολογισμό του ορίου κατά Cesaro ενός πίνακα, και χρησιμοποιώντας την σχέση (3.2) για τον υπολογισμό της αναμενόμενης τιμής, λαμβάνονται τα ακόλουθα αποτελέσματα:

➤ Στρατηγική  $f_1 = (1,1,1,1)$

$$P^*(f_1) = \begin{pmatrix} 0.3810 & 0.2857 & 0.2381 & 0.0952 \\ 0.3810 & 0.2857 & 0.2381 & 0.0952 \\ 0.3810 & 0.2857 & 0.2381 & 0.0952 \\ 0.3810 & 0.2857 & 0.2381 & 0.0952 \end{pmatrix}$$

Τότε, για την  $f_1$ , παίρνουμε:

$$\widehat{v}_a(f_1) = (0.4231, 0.4231, 0.4231, 0.4231)$$

➤ Στρατηγική  $f_2 = (2,1,1,1)$

$$P^*(f_2) = \begin{pmatrix} 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \end{pmatrix}$$

Τότε, για την  $f_2$ , παίρνουμε:

$$\widehat{v}_a(f_2) = (1.4062, 1.4063, 1.4063, 1.4062)$$

➤ Στρατηγική  $f_3 = (1,2,1,1)$

$$P^*(f_3) = \begin{pmatrix} 0.3556 & 0.2 & 0.3556 & 0.0889 \\ 0.3556 & 0.2 & 0.3556 & 0.0889 \\ 0.3556 & 0.2 & 0.3556 & 0.0889 \\ 0.3556 & 0.2 & 0.3556 & 0.0889 \end{pmatrix}$$

Τότε, για την  $f_3$ , παίρνουμε:

$$\widehat{v}_a(f_3) = (0.2857, 0.2857, 0.2857, 0.2857)$$



➤ Στρατηγική  $f_4 = (2,2,1,1)$

$$P^*(f_4) = \begin{pmatrix} 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \end{pmatrix}$$

Τότε, για την  $f_4$ , παίρνουμε:

$$\widehat{v}_a(f_4) = (1.4062, 1.4062, 1.4063, 1.4062)$$

➤ Στρατηγική  $f_5 = (1,1,2,1)$

$$P^*(f_5) = \begin{pmatrix} 0.3684 & 0.3421 & 0.1316 & 0.1579 \\ 0.3684 & 0.3421 & 0.1316 & 0.1579 \\ 0.3684 & 0.3421 & 0.1316 & 0.1579 \\ 0.3684 & 0.3421 & 0.1316 & 0.1579 \end{pmatrix}$$

Τότε, για την  $f_5$ , παίρνουμε:

$$\widehat{v}_a(f_5) = (0.2342, 0.2342, 0.2342, 0.2342)$$

➤ Στρατηγική  $f_6 = (2,1,2,1)$

$$P^*(f_6) = \begin{pmatrix} 0.3962 & 0.1509 & 0.3019 & 0.1509 \\ 0.3962 & 0.1509 & 0.3019 & 0.1509 \\ 0.3962 & 0.1509 & 0.3019 & 0.1509 \\ 0.3962 & 0.1509 & 0.3019 & 0.1509 \end{pmatrix}$$

Τότε, για την  $f_6$ , παίρνουμε:

$$\widehat{v}_a(f_6) = (1.0807, 1.0807, 1.0807, 1.0807)$$

➤ Στρατηγική  $f_7 = (1,2,2,1)$

$$P^*(f_7) = \begin{pmatrix} 0.3265 & 0.2653 & 0.2177 & 0.1905 \\ 0.3265 & 0.2653 & 0.2177 & 0.1905 \\ 0.3265 & 0.2653 & 0.2177 & 0.1905 \\ 0.3265 & 0.2653 & 0.2177 & 0.1905 \end{pmatrix}$$

Τότε, για την  $f_7$ , παίρνουμε:

$$\widehat{v}_a(f_7) = (0.6280, 0.6280, 0.6280, 0.6280)$$

➤ Στρατηγική  $f_8 = (2,2,2,1)$

$$P^*(f_8) = \begin{pmatrix} 0.3750 & 0.1250 & 0.3333 & 0.1667 \\ 0.3750 & 0.1250 & 0.3333 & 0.1667 \\ 0.3750 & 0.1250 & 0.3333 & 0.1667 \\ 0.3750 & 0.1250 & 0.3333 & 0.1667 \end{pmatrix}$$

Τότε, για την  $f_8$ , παίρνουμε:

$$\widehat{v}_a(f_8) = (1.4327, 1.4327, 1.4327, 1.4327)$$

➤ Στρατηγική  $f_9 = (1,1,1,2)$

$$P^*(f_9) = \begin{pmatrix} 0.4000 & 0.2500 & 0.2500 & 0.1000 \\ 0.4000 & 0.2500 & 0.2500 & 0.1000 \\ 0.4000 & 0.2500 & 0.2500 & 0.1000 \\ 0.4000 & 0.2500 & 0.2500 & 0.1000 \end{pmatrix}$$

Τότε, για την  $f_9$ , παίρνουμε:

$$\widehat{v}_a(f_9) = (0.4367, 0.4367, 0.4367, 0.4367)$$

➤ Στρατηγική  $f_{10} = (2,1,1,2)$

$$P^*(f_{10}) = \begin{pmatrix} 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \end{pmatrix}$$

Τότε, για την  $f_{10}$ , παίρνουμε:

$$\widehat{v}_a(f_{10}) = (1.4062, 1.4063, 1.4063, 1.4063)$$

➤ Στρατηγική  $f_{11} = (1, 2, 1, 2)$

$$P^*(f_{11}) = \begin{pmatrix} 0.3765 & 0.1765 & 0.3529 & 0.0941 \\ 0.3765 & 0.1765 & 0.3529 & 0.0941 \\ 0.3765 & 0.1765 & 0.3529 & 0.0941 \\ 0.3765 & 0.1765 & 0.3529 & 0.0941 \end{pmatrix}$$

Τότε, για την  $f_{11}$ , παίρνουμε:

$$\widehat{v}_a(f_{11}) = (0.2309, 0.2309, 0.2309, 0.2309)$$

➤ Στρατηγική  $f_{12} = (2, 2, 1, 2)$

$$P^*(f_{12}) = \begin{pmatrix} 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \\ 0.4286 & 0 & 0.5714 & 0 \end{pmatrix}$$

Τότε, για την  $f_{12}$ , παίρνουμε:

$$\widehat{v}_a(f_{12}) = (1.4062, 1.4062, 1.4063, 1.4062)$$

➤ Στρατηγική  $f_{13} = (1, 1, 2, 2)$

$$P^*(f_{13}) = \begin{pmatrix} 0.4000 & 0.2857 & 0.1429 & 0.1714 \\ 0.4000 & 0.2857 & 0.1429 & 0.1714 \\ 0.4000 & 0.2857 & 0.1429 & 0.1714 \\ 0.4000 & 0.2857 & 0.1429 & 0.1714 \end{pmatrix}$$

Τότε, για την  $f_{13}$ , παίρνουμε:

$$\widehat{v}_a(f_{13}) = (0.2165, 0.2165, 0.2165, 0.2165)$$

➤ Στρατηγική  $f_{14} = (2,1,2,2)$

$$P^*(f_{14}) = \begin{pmatrix} 0.4286 & 0.0816 & 0.3265 & 0.1633 \\ 0.4286 & 0.0816 & 0.3265 & 0.1633 \\ 0.4286 & 0.0816 & 0.3265 & 0.1633 \\ 0.4286 & 0.0816 & 0.3265 & 0.1633 \end{pmatrix}$$

Τότε, για την  $f_{14}$ , παίρνουμε:

$$\widehat{v}_a(f_{14}) = (1.3509, 1.3509, 1.3509, 1.3509)$$

➤ Στρατηγική  $f_{15} = (1,2,2,2)$

$$P^*(f_{15}) = \begin{pmatrix} 0.3714 & 0.2143 & 0.2143 & 0.2000 \\ 0.3714 & 0.2143 & 0.2143 & 0.2000 \\ 0.3714 & 0.2143 & 0.2143 & 0.2000 \\ 0.3714 & 0.2143 & 0.2143 & 0.2000 \end{pmatrix}$$

Τότε, για την  $f_{15}$ , παίρνουμε:

$$\widehat{v}_a(f_{15}) = (0.5695, 0.5695, 0.5695, 0.5695)$$

➤ Στρατηγική  $f_{16} = (2,2,2,2)$

$$P^*(f_{16}) = \begin{pmatrix} 0.4194 & 0.0645 & 0.3441 & 0.1720 \\ 0.4194 & 0.0645 & 0.3441 & 0.1720 \\ 0.4194 & 0.0645 & 0.3441 & 0.1720 \\ 0.4194 & 0.0645 & 0.3441 & 0.1720 \end{pmatrix}$$

Τότε, για την  $f_{16}$ , παίρνουμε:

$$\widehat{v}_a(f_{16}) = (1.5698, 1.5698, 1.5698, 1.5698)$$

Όπως ήταν αναμενόμενο, η βέλτιστη στρατηγική για τη διαδικασία  $\widehat{F}_1$  είναι η  $f_{16} = (2,2,2,2)$ , δηλαδή το βέλτιστο ‘παίξιμο’ είναι να επιλέγεται η 2<sup>η</sup> ενέργεια για κάθε κατάσταση  $s \in \{1,2,3,4\}$ , και το διάνυσμα τιμής της διαδικασίας, ταυτίζεται με εκείνο που λάβαμε από τον αλγόριθμο γραμμικού προγραμματισμού και είναι ίσο με:

$$\widehat{v}_a(f_{16}) = (1.5698, 1.5698, 1.5698, 1.5698)$$

## Σχόλια-Παρατηρήσεις

- Οπότε έχουμε καταλήξει και με τις 2 μεθόδους στο ίδιο διάνυσμα τιμής για την διαδικασία  $\bar{I}_1$  και επιστρέφοντας στο αρχικό παίγνιο  $\Gamma$ , εάν  $(\mathbf{f}^*, \mathbf{g}^*)$  είναι ένα ζεύγος βέλτιστων στρατηγικών, τότε η  $\mathbf{f}^*$  υποδεικνύει στον 1<sup>ο</sup> παίκτη να επιλέγει την 2<sup>η</sup> ενέργεια για κάθε αρχική κατάσταση για τις δύο πρώτες καταστάσεις του παιγνίου και η  $\mathbf{g}^*$  υποδεικνύει στον 2<sup>ο</sup> παίκτη να επιλέγει την 2<sup>η</sup> ενέργεια για κάθε αρχική κατάσταση για τις δύο τελευταίες καταστάσεις του παιγνίου.
- Κατά την υλοποίηση των παραπάνω μέσω του λογισμικού MATLAB, έγινε μια προσπάθεια εκτίμησης του υπολογιστικού χρόνου που απαιτείται για την επίλυση του συγκεκριμένου παραδείγματος με τις δύο διαφορετικές μεθόδους, κάνοντας χρήση της εντολής “tic-tac”. Οι μετρήσεις έδειξαν πως για τη μέθοδο πλήρους απαρίθμησης, οι υπολογισμοί ολοκληρώνονταν κατά μέσο όρο περίπου 10 φορές ταχύτερα από εκείνους της μεθόδου γραμμικού προγραμματισμού. Ενώ αυτό φαίνεται ως ένα μη αναμενόμενο αποτέλεσμα, μπορεί να αποδοθεί στις “μικρές” διαστάσεις του προβλήματος. Σε ανάλογες μετρήσεις που έγιναν σε αριθμητικό παράδειγμα με 8 καταστάσεις, παρατηρήσαμε πως ο χρόνος εκτέλεσης της μεθόδου πλήρους απαρίθμησης, αυτή τη φορά ήταν κατά μέσο όρο περίπου 3 φορές πιο γρήγορος από εκείνον της μεθόδου γραμμικού προγραμματισμού. Για το λόγο αυτό και γνωρίζοντας πως η μέθοδος Simplex ενδείκνυται υπολογιστικά για προβλήματα μεγάλης διάστασης, εικάζουμε, πώς όσο μεγαλώνει η διάσταση του προβλήματος, τόσο περισσότερο υπολογιστικά συμφέρουσα θα είναι η μέθοδος αυτή. Έτσι, θεωρούμε πως θα υπάρχει μία διάσταση προβλήματος για την οποία οι δύο τρόποι υπολογισμού που συγκρίνονται παραπάνω θα έχουν το ίδιο υπολογιστικό κόστος και πως από εκείνο το σημείο και έπειτα, η μέθοδος πλήρους απαρίθμησης θα αρχίσει να είναι ιδιαίτερα χρονοβόρα, ώστε να αποθαρρύνει κάποιον από τη χρήση της. Το γεγονός αυτό αποτελεί αντικείμενο μελλοντικής έρευνας.

## ***Παράρτημα***

Για την υλοποίηση της εφαρμογής της παραγράφου 3.5, έγινε χρήση των παρακάτω, μέσω του λογισμικού MATLAB:

i) Συνάρτηση linprog για επίλυση μέσω γραμμικού προγραμματισμού

```
S = 4; % Number of States
d = zeros(S,1);
e = eye(S);
A = [];
b = [];
Aeq = @(d)[9 8 -6 -4 -6 -6 -3 -6 0 0 0 0 0 0 0 0
           -3 0 6 8 0 0 -6 -3 0 0 0 0 0 0 0 0
           -3 -8 0 -4 6 12 -3 -3 0 0 0 0 0 0 0 0
           -1 0 0 0 0 -2 4 4 0 0 0 0 0 0 0 0
           12 12 0 0 0 0 0 0 9 8 -6 -4 -6 -6 -3 -6 -12*d(1)
           0 0 12 12 0 0 0 0 -3 0 6 8 0 0 -6 -3 -12*d(2)
           0 0 0 0 12 12 0 0 -3 -8 0 -4 6 12 -3 -3 -12*d(3)
           0 0 0 0 0 0 4 4 -1 0 0 0 0 -2 4 4 -4*d(4)
           1 0.8 1.5 1 1 2 2 1.2 0 0 0 0 0 0 0 0 0];
beq = [0 0 0 0 0 0 0 0 1];
f = [1.5 1 3 2 -3 -7 -2 -1 0 0 0 0 0 0 0 0 0];

s = size(f,2); %number of variables
lb = zeros([1,s]); % lower bound of variables
ub = inf*ones([1,s]); % upper bound of variables
```

```

for i = 1:S
d = e(i,:);
[x,fval] = linprog(f,A,b,Aeq(d),beq,lb,ub);
disp(x)
disp(fval)
end

```

ii) Μέθοδος πλήρους απαρίθμησης

```
state_actions = [2,2,2,2];
```

```

rewards = { [1.5,1]
             [3,2]
             [-3,-7]
             [-2,-1]
             };

```

```

sojourn = {[1,0.8]
           [1.5,1]
           [1,2]
           [2,1.2]
           };

```

```

probs = {[1/4,1/4,1/4,1/4; 1/3, 0, 2/3,0]
         [1/2, 1/2, 0, 0 ; 1/3,1/3,1/3,0]
         [1/2,0,1/2,0; 1/2,0,0,1/2]
         [1/4,1/2,1/4,0 ; 1/2,1/4,1/4,0]

```

```

    };

states_count = size(state_actions,2);
policies = [ones(1,states_count)];
new_p = [];
for i = 1:states_count
    temp = [];
    for j = 2:state_actions(i)
        new_p = policies;
        new_p(:,i) = new_p(:,i) + j-1;
        temp = cat(1,temp,new_p);
    end
    policies = cat(1,policies,temp);
end
policies = int32(policies);
%%

```

```

phi = zeros(size(policies));

for i = 1:size(policies,1)
    pol = policies(i,:);
    P = zeros(states_count);
    re = zeros(states_count,1);
    soj = zeros(states_count,1);

    for j = 1:states_count
        P(j,:) = probs{j}(pol(j,:));
        re(j) = rewards{j}(pol(j));
    end
end

```



```

    soj(j) = sojourn{j}(pol(j));
end

Pstar = cesaro(P);
disp(Pstar);
phi(i,:) = ((Pstar*re)./(Pstar*soj))';

end

[val, argmax] = max(phi,[],1);
disp('-----phi-----')
disp(phi);

```

### iii) Συνάρτηση eigval

```

function [evalues, repeats] = eigval(A)

% eigval Eigenvalues and their algebraic multiplicity.
%
% evalues = eigval(A) returns the distinct eigenvalues of A,
% with duplicates removed and sorted in decreasing order.
%
% [evalues, repeats] = eigval(A) also returns the row vector
% repeats that gives the multiplicity of each eigenvalue.
% The sum of the multiplicities is n.
%
% Examples: Let A = eye(n) and B = diag([3 4]).
% For A, evalues is 1 and repeats is n.

```

```

% For B, evals is [4; 3] and repeats is [1 1].

tol = sqrt(eps);
lambda = sort(eig(A));
lambda = round(lambda/tol) * tol;
%
% lambda gives all n eigenvalues (repetitions included).
%
evals = unique(lambda);
evals = flipud(evals);
n = length(lambda);
d = length(evals);
A = ones(n, 1) * evals';
B = lambda * ones(1, d);
MATCH = abs(A-B) <= tol;
%
% MATCH is an n by d zero matrix except
% MATCH(i,j) = 1 when lambda(i) = evals(j).
% Summing the columns gives the row vector repeats.
%
repeats = sum(MATCH);

```

#### iv) Συνάρτηση Cesaro

```

function ces = cesaro(t)

% get unique eigenvalues (ev) and their multiplicity (mu)
[ev,mu] = eigval(t);
eig1 = find(ev==1);

```

```

%remove eigenvalue equal to 1 and its multiplicity from ev and rep
ev(eig1) = [];
mu(eig1) = [];

I = eye(size(t));
R = 1;
for i = 1:size(ev,1)
    R = R*(t-ev(i)*I)^mu(i);
end
% divide R by the sum of its first row
ces = R/sum(R(1,:));
end

```

## Βιβλιογραφία

- [1] Alexandru Lazari, Dmitrii Lozovanu. (2020). New algorithms for finding the limiting. *Buletinul Academiei de Stiinte a Moldovei*, pp. 92(1):75–88.
- [2] D. Blackwell, T. Ferguson. (1968). The big match. *The Annals of Mathematical Statistics*, pp. 39: 159-163.
- [3] *Function "eig"*. (Introduced before R2006a). Retrieved from MathWorks:  
<https://www.mathworks.com/help/matlab/ref/eig.html>
- [4] *Function "linprog"*. (Introduced before R2006a). Retrieved from MathWorks:  
<https://www.mathworks.com/help/optim/ug/linprog.html>
- [5] Howard, Ronald A. (1971). *Dynamic Probabilistic Systems Volume II: Semi-Markov and decision Processes*. John Wiley & Sons.
- [6] Inc., T. M. (2022). MATLAB. version 9.13.0.2114483 (R2022b). Natick, Massachusetts.
- [7] J.L.Doob. (1953). *Stochastic Processes*. John Willey & Sons.
- [8] Jerzy Filar, Koos Vrieze (1997). *Competitive Markov decision processes*. New York: Springer.
- [9] *MATLAB Teaching Codes*. (n.d.). Retrieved from web.mit.edu:  
<http://web.mit.edu/18.06/www/Course-Info/Tcodes.html?fbclid=IwAR2C8oO2pmQeTm6edQwNc-yFpCkYSAeiabfwNrkKzwrZVTPUAcMqMKyIAfo>
- [10] Mondal, P. (2015). Linear Programming and Zero-Sum Two-Person Undiscounted Semi-Markov Games. *Asia-Pacific Journal of Operational Research*, p. 150043(20).
- [11] Mondal, P. (2017). On zero-sum two-person undiscounted semi-markov games. *Advances in Applied Probability*, pp. 49: 826-849.
- [12] Mondal, P. (2017). Computing semi-stationary optimal policies for multichain semi-Markov decision processes. *Annals of Operations Research*, pp. 843-865.
- [13] S. Sinha, K. G. Bakshi (2022). *Zero-Sum Two Person Perfect Information Semi-Markov Games: A Reduction*. Retrieved from arXiv:2201.00179
- [14] Sagnik Sinha, P. Mondal (2017). Semi-Markov decision processes with limiting ratio average. *Journal of Mathematical Analysis and Applications*, pp. 864-871.
- [15] Strang, G. (2006). *Linear Algebra and Its Applications, 4th Edition*. Cengage Learning.
- [16] Κ.Α.Μηλολιδάκης (2009). *Θεωρία Παιγνίων*. Αθήνα: Εκδόσεις Σοφία.
- [17] Ν.Δ. Τσάντας, Π. -Χ.Γ.Βασιλείου (2000). *Εισαγωγή στην επιχειρησιακή έρευνα*. Θεσσαλονίκη: Εκδόσεις Ζήτη.