

ΚΕΦΑΛΑΙΟ 4

**Αλγόριθμοι για το πρόβλημα των μερικά παρατηρήσιμων
Μαρκοβιανών διαδικασιών απόφασης σε άπειρο χρονικό
ορίζοντα.**

Περίληψη

Οι Παπαδημητρίου και Τσιτσικλής [92] απέδειξαν ότι το πρόβλημα POMDP σε άπειρο χρονικά ορίζοντα δεν έχει στη γενική του έκφραση ακριβή λύση με πεπερασμένους αλγόριθμους. Αυτό αποτελεί για μας εφιαλτήριο για αναζήτηση προσεγγιστικών λύσεων. Σκοπός του κεφαλαίου είναι η εύρεση προσεγγίσεων της άριστης συνάρτησης τιμών και της άριστης πολιτικής (προσδιορισμός σχεδόν άριστων πολιτικών) σε άπειρο χρονικά ορίζοντα, επεκτείνοντας αποτελέσματα των Bertsekas [11], Puterman [99] και Hauskrecht [48]. Το κεφάλαιο αυτό οργανώνεται ως εξής:

Στην ενότητα 4.1. οι προσεγγίσεις υλοποιούνται με την επαναληπτική εφαρμογή του αλγορίθμου των ακροτάτων σημείων που περιγράψαμε στο κεφάλαιο 3.

Στην ενότητα 4.2. αναφέρονται διάφορα φράγματα για την άριστη συνάρτηση τιμών που απαντώνται στη βιβλιογραφία και τα οποία μπορούν να χρησιμεύσουν ως αρχικές προσεγγίσεις.

Στην ενότητα 4.3. εφαρμόζεται επαναληπτικά ο αλγόριθμος των ακροτάτων σημείων στα αρχικά φράγματα της ενότητας 4.2 για την δημιουργία νέων προσεγγίσεων της άριστης συνάρτησης τιμών και σχεδόν άριστων πολιτικών. Σε κάθε περίπτωση υπολογίζεται ο απαιτούμενος αριθμός των βημάτων (επαναλήψεων), καθώς

και το προκαθορισμένο σφάλμα του αλγορίθμου, έτσι ώστε να επιτύχουμε προσέγγιση με οποιαδήποτε επιθυμητή ακρίβεια..

4.1. Προσέγγιση της άριστης συνάρτησης τιμών για άπειρο χρονικό ορίζοντα και εύρεση σχεδόν άριστων πολιτικών στα πλαίσια της επαναληπτικής μεθόδου τιμών (Value-Iteration).

Στην ενότητα αυτή θα εφαρμόσουμε τη μέθοδο των διαδοχικών προσεγγίσεων ή επαναληπτική μέθοδο τιμών (successive approximations or value iteration method) για την προσέγγιση της άριστης συνάρτησης τιμών V^* και της άριστης πολιτικής για άπειρο χρονικό ορίζοντα (βλ. ενότητα 1.5). Όπως θα διαπιστώσουμε στη συνέχεια, οι προσεγγίσεις αυτές είναι δυνατόν να υλοποιηθούν με οποιαδήποτε επιθυμητή ακρίβεια εφαρμόζοντας επαναληπτικά τον αλγόριθμο των ακροτάτων σημείων που περιγράψαμε στο κεφάλαιο 3. Σημειώνουμε ακόμη ότι στο πλαίσιο του κριτηρίου βελτιστοποίησης σε άπειρο χρονικό ορίζοντα, για τον συντελεστή εκπτώσεως (discount factor) υποθέτουμε $\beta \in (0,1)$.

Το επόμενο λήμμα είναι πολύ βασικό για την εφαρμογή της μεθόδου των διαδοχικών προσεγγίσεων και οφείλεται στον Denardo [23] (βλ. επίσης Blackwell [15]).

Λήμμα 4.1.1: Εστω V^ η βέλτιστη συνάρτηση τιμών για άπειρο χρονικό ορίζοντα σε πρόβλημα εσόδων ή κόστους, H ο τελεστής μεγιστοποίησης ή ελαχιστοποίησης και $w \in B(P)$ (βλ. ενότητα 1.4,1.5). Τότε*

i) Η συνάρτηση V^ είναι το μοναδικό σταθερό σημείο του τελεστή H , δηλ.*

$$w \in B(P), Hw=w \text{ ή } w=V^*.$$

$$ii) \|H_n u - V^*\| \leq \beta^n \|u - V^*\|, n \geq 1.$$

$$iii) \|H_n u - V^*\| \leq \frac{\beta^n}{1-\beta} \|Hu - u\|, n \geq 1$$

όπου με H_n συμβολίζουμε την επαναληπτική χρήση του τελεστού H , n φορές και H_0 είναι ο ταυτοτικός τελεστής: $H_0 u = u$. □

Θα ασχοληθούμε με προσεγγίσεις σε προβλήματα εσόδων, οπότε ο H θεωρείται τελεστής μεγιστοποίησης. Οι προσεγγίσεις για προβλήματα κόστους γίνονται με ανάλογο τρόπο.

Θεωρούμε $V_0 \in B(\pi)$ και

$$V_n = HV_{n-1}, n = 1, 2, 3, \dots \quad \mathbf{4.1.1}$$

Η παραπάνω σχέση γράφεται:

$$V_n = H_n V_0, n \geq 1$$

Το λήμμα 4.1.1 (ii) δείχνει ότι η ακολουθία $\{V_n\}$ συγκλίνει ομαλά, όταν το $n \in \mathbb{N}$ στη βέλτιστη συνάρτηση τιμών V^* , ανεξάρτητα από την επιλογή της αρχικής συνάρτησης V_0 .

Αν $V_0 = 0$ (μηδενική συνάρτηση), τότε:

$$\|V_n - V^*\| \leq \frac{b^n L}{1-b}, n \geq 1, \quad \mathbf{4.1.2}$$

όπου

$$L = \max_{i,a} |q(i,a)|.$$

Πράγματι,

$$V_1(p) = HV_0(p) = \max_a \{p \cdot q^a\} \leq \max_{i,a} |q(i,a)| = L, \quad p \in P$$

Επομένως $\|V_1\| \leq L$.

Εφαρμόζοντας το λήμμα 4.1.1 (iii) παίρνουμε:

$$\|V_n - V^*\| = \|H_n V_0 - V^*\| \leq \frac{\beta^n}{1-\beta} \|HV_0 - V_0\| = \frac{\beta^n}{1-\beta} \|V_1\| \leq \frac{b^n}{1-b} L, n \geq 1$$

Η σχέση (4.1.2) είναι χρήσιμη στην περίπτωση που χρησιμοποιούμε τον αλγόριθμο των ακρότατων σημείων (Κεφ.3), για τον ακριβή υπολογισμό των συναρτήσεων V_n (προκαθορισμένο σφάλμα $\varepsilon=0$). Αν $h > 0$ είναι ένα επιθυμητό φράγμα για το σφάλμα προσέγγισης, δηλ.

$$\|V_n - V^*\| \leq h$$

τότε, μπορούμε να επιλέξουμε χρονικό ορίζοντα n , έτσι ώστε :

$$\frac{b^n}{1-b} L \leq h$$

δηλ.
$$n \geq \frac{\ln\left(\frac{(1-b)h}{L}\right)}{\ln b}.$$

Εργαζόμενοι με ανάλογο τρόπο μπορούμε να προσεγγίσουμε τη βέλτιστη συνάρτηση τιμών V^* εφαρμόζοντας τον προσεγγιστικό αλγόριθμο των ακρότατων σημείων με προκαθορισμένο σφάλμα $\varepsilon > 0$, σύμφωνα με την ακόλουθη πρόταση.

Πρόταση 4.1.2: Εστω \tilde{V}_n η προσέγγιση της V_n μέσω του αλγορίθμου των ακρότατων σημείων, με προκαθορισμένο σφάλμα $\varepsilon > 0$ (βλέπε ενότητα 3.4). Τότε,

$$\|\tilde{V}_n - V^*\| \leq \frac{1-b^n}{1-b} \varepsilon + \frac{b^n L}{1-b} \quad \mathbf{4.1.3}$$

Απόδειξη

Από την τριγωνική ιδιότητα της νόρμας του supremum παίρνουμε:

$$\|\tilde{V}_n - V^*\| \leq \|\tilde{V}_n - V_n\| + \|V_n - V^*\|.$$

Το συσσωρευμένο σφάλμα σ_n της προσέγγισης \tilde{V}_n για την συνάρτηση V_n δίνεται από την σχέση (3.4.7)

Επομένως

$$\|\tilde{V}_n - V_n\| \leq \sigma_n = \frac{1 - b^n}{1 - b} \cdot e$$

Η (4.1.3) συνάγεται άμεσα από την παραπάνω σχέση και την σχέση (4.1.2).

W

Αν $h > 0$ είναι ένα επιθυμητό άνω φράγμα για το σφάλμα προσέγγισης,

$$\|V_n - V^*\| \leq \eta$$

τότε μπορούμε να επιλέξουμε προκαθορισμένο σφάλμα ε και χρονικό ορίζοντα n έτσι ώστε:

$$\frac{1 - b^n}{1 - b} \cdot e + \frac{b^n \cdot L}{1 - b} \leq h$$

Μπορούμε π.χ να επιλέξουμε τα ε, η έτσι ώστε: $\frac{e}{1 - b} \leq \frac{h}{2}, \frac{b^n \cdot L}{1 - b} \leq \frac{h}{2},$

δηλαδή: $e \leq \frac{(1 - b) \cdot h}{2}, n \geq \frac{\ln\left(\frac{(1 - b) \cdot h}{2L}\right)}{\ln b}$

Έστω $d^\# = (\delta, \delta, \dots)$ μία στάσιμη πολιτική και $V(\pi/\delta)$, $\pi \in \Pi$, η συνάρτηση του αναμενόμενου ολικού εκπίπτοντος κέρδους για άπειρο χρονικό ορίζοντα εφαρμόζοντας την πολιτική $d^\#$.

Η συνάρτηση $V(./\delta)$ αναφέρεται επίσης ως συνάρτηση τιμών για την πολιτική δ^∞ .

Σημειώνουμε ακόμα ότι η $V(./\delta)$ είναι το μοναδικό σταθερό σημείο του τελεστή H_δ , ο οποίος είναι συστολή modulus β (βλέπε ενότητα 1.4):

$$H_\delta V(./\delta) = V(./\delta).$$

Με άλλα λόγια η συνάρτηση $V(./\delta)$ ικανοποιεί την εξίσωση:

$$V(\pi/\delta) = \pi \cdot q^{\delta(\pi)} + \beta \cdot \sum_{\theta} \{\theta/\pi, \delta\} \cdot V(T(\pi, \theta, \delta)/\delta), \pi \in \Pi.$$

Θα στρέψουμε τώρα το ενδιαφέρον μας στον προσδιορισμό στάσιμων πολιτικών που προσεγγίζουν την άριστη πολιτική $(\delta^*)^\infty$ (βλ. ενότητα 1.5).

Ορισμός 4.1.1: Μια πολιτική δ^∞ , λέγεται σχεδόν άριστη πολιτική με σφάλμα $h > 0$ (ή, h -άριστη πολιτική) αν η συνάρτηση τιμών για την πολιτική δ^∞ , $V(./\delta)$, αποτελεί προσέγγιση της άριστης συνάρτησης τιμών V^* με μέγιστο σφάλμα προσέγγισης μικρότερο ή ίσο του αριθμού h , δηλαδή:

$$\|V(./\delta) - V^*\| \leq h$$

Μπορούμε να προσεγγίσουμε την άριστη πολιτική $(d^*)^\#$ με τη στάσιμη πολιτική $(\delta^n)^\infty$, όπου δ^n , είναι η άριστη συνάρτηση ελέγχου στον χρονικό ορίζοντα n :

$$V_n = H V_{n-1} = H_{\delta^n} V_{n-1}. \quad \underline{\underline{4.1.4}}$$

(όπου $V_0 = 0$, μηδενική συνάρτηση)

δηλαδή : $\delta^n(\pi) = \arg \max_a \{p \cdot q^a + b \cdot \mathbf{E}_q \{q/p, a\} \cdot V_{n-1}(T(p, q, a))\}, p \in P$.

Η συνάρτηση $V(\cdot/\delta^n)$ είναι μια καλή προσέγγιση της V^* αν επιλέξουμε το n αρκετά μεγάλο, σύμφωνα με την ακόλουθη πρόταση.

Πρόταση 4.1.3: Έστω δ^n άριστη συνάρτηση ελέγχου για τον χρονικό ορίζοντα n σύμφωνα με την σχέση (4.1.4). Τότε η στάσιμη πολιτική $(\delta^n)^\infty$ είναι σχεδόν άριστη

με σφάλμα $\frac{2 \cdot \beta^n}{1 - \beta} \cdot A$ ή $\frac{2 \cdot b^n \cdot L}{1 - b}$ - άριστη δηλαδή:

$$\|V(\cdot/\delta^n) - V^*\| \leq \frac{2 \cdot \beta^n}{1 - \beta} \cdot A$$

όπου

$$L = \max_{i,a} |q(i, a)|.$$

Bertsekas [10]

□

Αν επιθυμούμε η πολιτική $(\delta^n)^\infty$ να είναι h -άριστη για δοσμένο $h > 0$, τότε μπορούμε να επιλέξουμε τον χρονικό ορίζοντα n έτσι ώστε:

$$\frac{2 \cdot \beta^n}{1 - \beta} \cdot A \leq h$$

δηλαδή :

$$n^3 \geq \frac{\ln\left(\frac{(1 - b) \cdot h}{2L}\right)}{\ln b}$$

Σημειώνουμε, ότι η εφαρμογή της πρότασης 4.1.3, είναι δυνατή μόνον στην περίπτωση όπου εφαρμόζουμε τον αλγόριθμο των ακρότατων σημείων με προκαθορισμένο σφάλμα $\epsilon=0$, οπότε επιτυγχάνεται ακριβής υπολογισμός της συνάρτησης V_n και ο καθορισμός της άριστης συνάρτησης ελέγχου δ^n , μέσω της

σχέσης (4.1.4) είναι εφικτός. Αν το προκαθορισμένο σφάλμα στον αλγόριθμο των ακρότατων σημείων είναι $e > 0$, τότε αντί για την δ^n προσδιορίζουμε μια συνάρτηση ελέγχου \tilde{d}^n σύμφωνα με τη σχέση

$$\tilde{V}_n = \tilde{H} \tilde{V}_{n-1} = H_{\tilde{d}^n} \tilde{V}_{n-1} \quad \underline{\underline{4.1.5}}$$

(βλέπε επίσης παρατήρηση 2 στην ενότητα 3.2).

Μπορούμε να προσεγγίσουμε την άριστη πολιτική $(\delta^*)^\infty$ με την στάσιμη πολιτική $(\tilde{d}^n)^\infty$. Η πρόταση που ακολουθεί παρέχει το σφάλμα αυτής της προσέγγισης και αποτελεί επέκταση της πρότασης 4.1.3.

Πρόταση 4.1.4: Έστω \tilde{V}_n η προσέγγιση της V_n μέσω του αλγορίθμου των ακρότατων σημείων με προκαθορισμένο σφάλμα $e > 0$ και \tilde{d}^n η συνάρτηση ελέγχου στον χρονικό ορίζοντα n σύμφωνα με τη σχέση (4.1.5). Τότε η στάσιμη πολιτική $(\tilde{d}^n)^\infty$ είναι $f(e, n)$ - *áristh*, όπου

$$f(e, n) := \frac{1 + b - 2b^n}{(1 - b)^2} \cdot e + \frac{2b^n \cdot L}{1 - b}, \quad \underline{\underline{4.1.6}}$$

δηλαδή:

$$\|V(\cdot/\tilde{d}^n) - V^*\| \leq f(e, n).$$

Απόδειξη

Από την τριγωνική ιδιότητα της νόρμας του supremum προκύπτει ότι:

$$\|V(. / \tilde{d}^n) - V^*\| \leq \|V(. / \tilde{d}^n) - \tilde{V}_n\| + \|\tilde{V}_n - V^*\|. \quad \underline{\mathbf{4.1.7}}$$

Εφαρμόζοντας πάλι την τριγωνική ιδιότητα παίρνουμε

$$\|V(. / \tilde{d}^n) - \tilde{V}_n\| \leq \|V(. / \tilde{d}^n) - H_{\tilde{d}^n} \tilde{V}_n\| + \|H_{\tilde{d}^n} \tilde{V}_n - \tilde{V}_n\|. \quad \underline{\mathbf{4.1.8}}$$

Λαμβάνοντας υπόψη ότι η συνάρτηση $V(. / \tilde{d}^n)$ είναι το σταθερό σημείο του τελεστή $H_{\tilde{d}^n}$ και ότι ο τελεστής είναι συστολή modulus β , έχουμε:

$$\|V(. / \tilde{d}^n) - H_{\tilde{d}^n} \tilde{V}_n\| = \|H_{\tilde{d}^n} V(. / \tilde{d}^n) - H_{\tilde{d}^n} \tilde{V}_n\| \leq \beta \cdot \|V(. / \tilde{d}^n) - \tilde{V}_n\|. \quad \underline{\mathbf{4.1.9}}$$

Από τις (4.1.8) και (4.1.9) παίρνουμε

$$: \quad \|V(. / \tilde{d}^n) - \tilde{V}_n\| \leq \frac{1}{1 - \beta} \cdot \|H_{\tilde{d}^n} \tilde{V}_n - \tilde{V}_n\|. \quad \underline{\mathbf{4.1.10}}$$

Λαμβάνοντας υπόψη τη σχέση (4.1.9) και την ιδιότητα της συστολής του τελεστή $H_{\tilde{d}^n}$ προκύπτει ότι:

$$\|H_{\tilde{d}^n} \tilde{V}_n - \tilde{V}_n\| = \|H_{\tilde{d}^n} \tilde{V}_n - H_{\tilde{d}^n} \tilde{V}_{n-1}\| \leq \beta \cdot \|\tilde{V}_n - \tilde{V}_{n-1}\|. \quad \underline{\mathbf{4.1.11}}$$

Από την τριγωνική ιδιότητα της νόρμας supremum έχουμε:

$$\|\tilde{V}_n - \tilde{V}_{n-1}\| \leq \|\tilde{V}_n - V_n\| + \|V_n - V_{n-1}\| + \|V_{n-1} - \tilde{V}_{n-1}\| \leq s_n + s_{n-1} + \|V_n - V_{n-1}\|, \quad \underline{\mathbf{4.1.12}}$$

όπου s_n, s_{n-1} είναι τα συσσωρευμένα σφάλματα των προσεγγίσεων $\tilde{V}_n, \tilde{V}_{n-1}$ για τις V_n, V_{n-1} αντίστοιχα (βλέπε ενότητα 3.4). Εφαρμόζοντας διαδοχικά την ιδιότητα συστολής του τελεστή μεγιστοποίησης H παίρνουμε:

$$\|V_n - V_{n-1}\| = \|H_{n-1} V_1 - H_{n-1} V_0\| \leq \beta^{n-1} \cdot \|V_1 - V_0\| \leq \beta^{n-1} \cdot \Lambda. \quad \underline{\mathbf{4.1.13}}$$

Η τελευταία ανισότητα προκύπτει από το γεγονός ότι έχουμε επιλέξει ως συνάρτηση οφέλους στον τερματισμό τη μηδενική συνάρτηση, δηλαδή $V_0=0$ και επομένως

$$\|V_1 - V_0\| = \|V_1\| = \|HV_0\| \leq \max_{i,a} |q(i,a)| = L$$

Από τις σχέσεις (4.1.10),(4.1.11),(4.1.12),(4.1.13) λαμβάνοντας υπόψη και την (3.4.7) παίρνουμε:

$$\|V(\tilde{d}^n) - \tilde{V}_n\| \leq \frac{b \cdot (s_n + s_{n-1})}{1-b} + \frac{b^n}{1-b} \cdot \Lambda = \frac{b \cdot (2 - b^n - b^{n-1})}{(1-b)^2} \cdot e + \frac{b^n}{1-b} \cdot \Lambda. \quad \mathbf{4.1.14}$$

Από τις σχέσεις (4.1.7),(4.1.14) και (4.1.3) έχουμε

$$\begin{aligned} \|V(\tilde{d}^n) - V^*\| &\leq \frac{b \cdot (2 - b^n - b^{n-1})}{(1-b)^2} \cdot e + \frac{1-b^n}{1-b} \cdot e + \frac{2b^n}{1-b} \cdot \Lambda = \\ &= \frac{2b^n}{1-b} \cdot \Lambda + \frac{1+b-2b^n}{(1-b)^2} \cdot e = f(e,n) \end{aligned}$$

Επομένως η στάσιμη πολιτική $(\tilde{d}^n)^\ddagger$ είναι $f(e,n)$ - *áristh*. W

Σημειώνουμε ότι για $\varepsilon=0$ η πρόταση 4.1.4 δίνει το ίδιο φράγμα όπως και η

$$\text{πρόταση 4.1.3: } f(0,n) = \frac{2b^n \cdot L}{1-b}.$$

Αν επιθυμούμε η πολιτική $(\tilde{d}^n)^\ddagger$ να είναι h - *áristh* για δοσμένο $h > 0$, τότε μπορούμε να επιλέξουμε προκαθορισμένο σφάλμα e και χρονικό ορίζοντα n έτσι ώστε:

$$f(e,n) \leq h.$$

Μπορούμε π.χ. να επιλέξουμε τα ε, n έτσι ώστε:

$$\frac{1+b}{(1-b)^2} \cdot e \leq \frac{h}{2}, \quad \frac{2b^n}{1-b} \cdot \Lambda \leq \frac{h}{2}$$

δηλαδή:

$$e \leq \frac{(1-b)^2 \cdot h}{2 \cdot (1+b)}, \quad n \geq \frac{\ln\left(\frac{(1-b) \cdot h}{4L}\right)}{\ln b}.$$

Ένας εναλλακτικός τρόπος εύρεσης σχεδόν άριστων πολιτικών σχετίζεται με το κατάλοιπο Bellman (Bellman residual) σε κάποιο χρονικό ορίζοντα n , που ορίζεται ως η μέγιστη απόλυτη διαφορά των συναρτήσεων V_n και V_{n-1} , δηλαδή η ποσότητα $\|V_n - V_{n-1}\|$. Αν αυτή η ποσότητα είναι "αρκούντως μικρή", τότε στην περίπτωση αυτή μπορούμε να προσεγγίσουμε την άριστη πολιτική $(\delta^*)^\infty$ με την πολιτική $(\delta^n)^\infty$, όπου δ^n είναι η άριστη συνάρτηση ελέγχου για τον χρονικό ορίζοντα σύμφωνα με τη σχέση (4.1.4). Με άλλα λόγια η συνάρτηση του αναμενόμενου ολικού εκπίπτοντος κέρδους για άπειρο χρονικό ορίζοντα εφαρμόζοντας την πολιτική $(\delta^n)^\infty, V(\cdot/d^n)$, προσεγγίζει ικανοποιητικά τη βέλτιστη συνάρτηση τιμών V^* .

Πρόταση 4.1.5: Έστω ότι για το κατάλοιπο Bellman σε κάποιο χρονικό ορίζοντα n ισχύει:

$$\|V_n - V_{n-1}\| \leq h \quad (\text{όπου } h > 0)$$

και δ^n είναι η άριστη συνάρτηση ελέγχου στον χρονικό ορίζοντα n , σύμφωνα με τη σχέση (4.1.4). Τότε η στάσιμη πολιτική $(\delta^n)^\infty$ είναι $\frac{2b \cdot h}{1-b}$ -άριστη, δηλαδή

$$\|V(\cdot/d^n) - V^*\| \leq \frac{2b \cdot h}{1-b} \tag{4.1.15}$$

Puterman [98]

W

Είναι φανερό ότι η παραπάνω πρόταση είναι εφαρμόσιμη μόνο στην περίπτωση όπου εφαρμόζουμε τον αλγόριθμο των ακρότατων σημείων με προκαθορισμένο

σφάλμα $e = 0$, οπότε επιτυγχάνεται ακριβής υπολογισμός των συναρτήσεων V_n , των καταλοίπων Bellman και των συναρτήσεων ελέγχου δ^n για $n=1,2,\dots$

Στην περίπτωση όπου εφαρμόζουμε τον αλγόριθμο των ακρότατων σημείων με προκαθορισμένο σφάλμα $e > 0$, χρησιμοποιούμε το "τροποποιημένο κατάλοιπο Bellman", που ορίζεται ως η μέγιστη απόλυτη διαφορά των συναρτήσεων \tilde{V}_n και \tilde{V}_{n-1} , δηλαδή η ποσότητα $\|\tilde{V}_n - \tilde{V}_{n-1}\|$.

Αν για κάποιο n η ποσότητα αυτή καθώς και το προκαθορισμένο σφάλμα e είναι "αρκούντως μικρά", τότε στην περίπτωση αυτή μπορούμε να προσεγγίσουμε την άριστη πολιτική $(\delta^*)^\infty$ με την πολιτική $(\tilde{d}^n)^\sharp$, όπου η συνάρτηση ελέγχου (\tilde{d}^n) καθορίζεται σύμφωνα με τη σχέση (4.1.5). Η επόμενη πρόταση καλύπτει αυτή την περίπτωση και αποτελεί επέκταση της πρότασης 4.1.5.

Πρόταση 4.1.6: Έστω ότι για το τροποποιημένο κατάλοιπο Bellman, σε κάποια χρονικό ορίζοντα n ισχύει

$$\|\tilde{V}_n - \tilde{V}_{n-1}\| \leq h \quad (\text{όπου } h > 0) \quad \underline{\underline{4.1.16}}$$

και \tilde{d}^n είναι η συνάρτηση ελέγχου στον χρονικό ορίζοντα n σύμφωνα με τη σχέση (4.1.5). Τότε η στάσιμη πολιτική $(\tilde{d}^n)^\sharp$ είναι $f(e,h)$ -άριστη, όπου

$$f(e,h) := \frac{1+b-2b^n}{(1-b)^2} \cdot e + \frac{2b \cdot h}{1-b}, \quad \underline{\underline{4.1.17}}$$

δηλαδή:

$$\|V(\tilde{d}^n) - V^*\| \leq f(e,h)$$

Απόδειξη

Από την τριγωνική ιδιότητα της νόρμας του supremum παίρνουμε

$$\|V(\cdot/d^n) - V^*\| \leq \|V(\cdot/d^n) - \tilde{V}_n\| + \|\tilde{V}_n - V_n\| + \|V_n - V^*\|$$

$$\leq \|V(\cdot/d^n) - \tilde{V}_n\| + s_n + \|V_n - V^*\| \quad \underline{\underline{4.1.18}}$$

Εφαρμόζοντας πάλι την τριγωνική ιδιότητα έχουμε

$$\|V(\cdot/d^n) - \tilde{V}_n\| \leq \|V(\cdot/d^n) - H_{d^n} \tilde{V}_n\| + \|H_{d^n} \tilde{V}_n - \tilde{V}_n\| \quad \underline{\underline{4.1.19}}$$

Λαμβάνοντας υπόψη ότι η συνάρτηση $V(\cdot/d^n)$ είναι το σταθερό σημείο του τελεστή H_{d^n} , ότι ο τελεστής είναι συστολή modulus β , και τις σχέσεις (4.1.5), (4.1.16),

από την (4.1.19) παίρνουμε:

$$\|V(\cdot/d^n) - \tilde{V}_n\| \leq \|H_{d^n} V(\cdot/d^n) - H_{d^n} \tilde{V}_n\| + \|H_{d^n} \tilde{V}_n - H_{d^n} \tilde{V}_{n-1}\| \leq$$

$$\leq \beta \cdot \|V(\cdot/d^n) - \tilde{V}_n\| + \beta \cdot \|\tilde{V}_n - \tilde{V}_{n-1}\| \leq$$

$$\leq \beta \cdot \|V(\cdot/d^n) - \tilde{V}_n\| + b \cdot h$$

Επομένως

$$\|V(\cdot/d^n) - \tilde{V}_n\| \leq \frac{b \cdot h}{1 - \beta} \quad \underline{\underline{4.1.20}}$$

Από το λήμμα 4.1.1 (iii) παίρνουμε:

$$\|V_n - V^*\| = \|H_{d^n} V_{n-1} - V^*\| \leq \frac{b}{1 - \beta} \|H_{d^n} V_{n-1} - V_{n-1}\|$$

Άρα

$$\|V_n - V^*\| \leq \frac{b}{1-b} \|V_n - V_{n-1}\| \quad \underline{4.1.21}$$

Από την τριγωνική ιδιότητα και τη σχέση (4.1.16) έχουμε

$$\|V_n - V_{n-1}\| \leq \|V_n - \tilde{V}_n\| + \|\tilde{V}_n - \tilde{V}_{n-1}\| + \|\tilde{V}_{n-1} - V_{n-1}\| \leq s_n + h + s_{n-1} \quad \underline{4.1.22}$$

Από τις σχέσεις (4.1.21),(4.1.22) έχουμε:

$$\|V_n - V^*\| \leq \frac{b}{1-b} \|V_n - V_{n-1}\| \leq \frac{b}{1-b} \cdot (s_n + h + s_{n-1}) \quad \underline{4.1.23}$$

Από τις σχέσεις (4.1.18),(4.1.20),(4.1.23) λαμβάνοντας υπόψη και την (3.4.7) παίρνουμε:

$$\begin{aligned} \|V(\tilde{d}^n) - V^*\| &\leq \frac{b \cdot h}{1-b} + s_n + \frac{b \cdot (s_n + s_{n-1} + h)}{1-b} = \\ &= \frac{s_n + b \cdot s_{n-1}}{1-b} + \frac{2bh}{1-b} = \frac{1+b-2b^n}{(1-b)^2} \cdot e + \frac{2b \cdot h}{1-b} = f(e, h) \end{aligned}$$

Επομένως η στάσιμη πολιτική $(\tilde{d}^n)^*$ είναι $f(e, h)$ - *áristh*. W

Σημειώνουμε ότι για $\varepsilon=0$ η πρόταση 4.1.6 δίνει το ίδιο φράγμα όπως και η

πρόταση 4.1.5 : $f(0, h) = \frac{2b \cdot h}{1-b}$. Αν επιθυμούμε η πολιτική $(\tilde{d}^n)^*$ να είναι l - *áristh*

για δοσμένο $l > 0$, τότε μπορούμε να επιλέξουμε προκαθορισμένο σφάλμα e και φράγμα h για το τροποποιημένο κατάλοιπο Bellman σε κάποιο χρονικό ορίζοντα έτσι ώστε:

$$f(e, h) \leq l$$

Μπορούμε π.χ. να επιλέξουμε τα ε, h έτσι ώστε:

$$\frac{1+b}{(1-b)^2} \cdot e \leq \frac{l}{2}, \quad \frac{2b \cdot h}{1-b} \leq \frac{l}{2}$$

δηλαδή:
$$e \leq \frac{(1-b)^2 l}{2 \cdot (1+b)}, \quad h \leq \frac{(1-b)l}{4b}.$$

Οι προτάσεις αυτής της ενότητας ισχύουν επίσης για προσεγγίσεις της άριστης συνάρτησης τιμών V^* και της άριστης πολιτικής $(\delta^*)^\infty$ σε προβλήματα κόστους (με την προφανή αλλαγή $L = \max_{i,a} |q(i,a)|$).

4.2. Κατασκευή φραγμάτων για την βέλτιστη συνάρτηση τιμών.

Στην ενότητα αυτή αναφέρουμε ορισμένα φράγματα για τη βέλτιστη συνάρτηση τιμών V^* που απαντώνται στη βιβλιογραφία. Η διαδικασία δημιουργίας προσεγγίσεων της V^* και της άριστης πολιτικής σε άπειρο χρονικό ορίζοντα από φράγματα θα μας απασχολήσει στην ενότητα 4.3. Θα περιοριστούμε στην κατασκευή φραγμάτων μόνο για προβλήματα POMDP μεγιστοποίησης εσόδων.

A) Κατασκευή άνω φραγμάτων για την V^*

1) Μέσω της βέλτιστης συνάρτησης τιμών της αντίστοιχης MDP. (Boutilier-Poole [17]).

Τα απλούστερα μη τετριμμένα άνω φράγματα για την συνάρτηση V^* μιας POMDP επιτυγχάνονται μέσω της βέλτιστης συνάρτησης τιμών V_{MDP}^* της αντίστοιχης (πλήρως παρατηρήσιμης Μαρκοβιανής διαδικασίας αποφάσεων (MDP) που ικανοποιεί την ακόλουθη εξίσωση αριστοποίησης.

$$V_{MDP}^*(i) = \max_a \{ q(i,a) + \beta \cdot \sum_{j=1}^{j=N} p_{ij}^\alpha \cdot V_{MDP}^*(j) \}, i=1,2,\dots,N$$

Η συνάρτηση V_{MDP}^* υπολογίζεται απλά και επακριβώς με την επαναληπτική μέθοδο πολιτικής (policy-iteration) σε πεπερασμένο αριθμό βημάτων (βλέπε Howard [51]). Σύμφωνα με τους Boutilier-Poole [17] έχουμε τα ακόλουθα άνω φράγματα για την συνάρτηση V^* : Για κάθε $\pi=(\pi_1,\pi_2,\dots,\pi_N) \in \Pi$,

$$V^*(\pi) \leq \max_a \sum_{i=1}^N \pi_i (q(i,a) + \beta \cdot \sum_{j=1}^N p_{ij}^\alpha \cdot V_{MDP}^*(j)) \quad \underline{4.2.1}$$

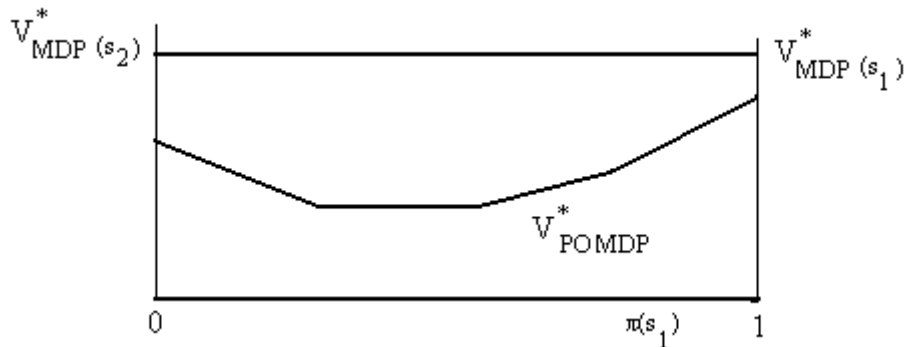
$$V^*(\pi) \leq \sum_{i=1}^N \pi_i \cdot V_{MDP}^*(i). \quad \underline{4.2.2}$$

Επειδή

$$\begin{aligned} \max_a \sum_{i=1}^N \pi_i (q(i,a) + \beta \cdot \sum_{j=1}^N p_{ij}^\alpha \cdot V_{MDP}^*(j)) &\leq \sum_{i=1}^N \pi_i \cdot \max_a (q(i,a) + \beta \cdot \sum_{j=1}^N p_{ij}^\alpha \cdot V_{MDP}^*(j)) = \\ &= \sum_{i=1}^N \pi_i \cdot V_{MDP}^*(i) \quad \forall \pi = (\pi_1, \pi_2, \dots, \pi_N) \in \Pi \end{aligned}$$

συνάγεται ότι το άνω φράγμα της (4.2.2) ως προσέγγιση της V^* υπολείπεται του αντίστοιχου άνω φράγματος της (4.2.1), όμως ο υπολογισμός του είναι απλούστερος.

Η ιδέα της προσέγγισης διευκρινίζεται στο σχήμα 4.1, όπου έχουμε μια POMDP δύο καταστάσεων s_1, s_2 και την MDP προσέγγισή της.



Σχήμα 4.1: Προσέγγιση βασιζόμενη σε μια πλήρως παρατηρήσιμη MDP για μια POMDP δύο καταστάσεων s_1, s_2 .

2) Μέσω των συναρτήσεων V_n (Lovejoy [75])

Μια ακολουθία άνω φραγμάτων για την V^* σχηματίζεται μέσω της ακολουθίας των συναρτήσεων $\{V_n\}$:

$$V_n = HV_{n-1}, n = 1, 2, \dots \quad \underline{4.2.3}$$

$$V_0 = 0 \text{ (μηδενική συνάρτηση),}$$

όπου H είναι ο τελεστής μεγιστοποίησης. Για $n=0,1,2,\dots$ έχουμε:

$$V^*(\pi) \leq V_n(\pi) + \frac{\beta^n}{1-\beta} \cdot q_{\max} \quad \forall \pi \in \Pi. \quad \underline{4.2.4}$$

όπου
$$q_{\max} = \max_{i,a} q(i,a).$$

Για $n=0$ έχουμε την ειδική περίπτωση

$$V^*(\pi) \leq \frac{1}{1-\beta} \cdot q_{\max} \quad \forall \pi \in \Pi.$$

B) Κατασκευή κάτω φραγμάτων για την V^*

1) Μέσω της βέλτιστης συνάρτησης τιμών της αντίστοιχης UMDP (unobservable MDP), μη παρατηρήσιμη Μαρκοβιανή διαδικασία αποφάσεων, χωρίς μηνύματα (Madani [78]).

Ένα κάτω φράγμα για την βέλτιστη συνάρτηση τιμών V^* μιας POMDP είναι η βέλτιστη συνάρτηση τιμών V^*_{UMDP} της αντίστοιχης μη παρατηρήσιμης διαδικασίας αποφάσεων (UMDP), που ικανοποιεί την ακόλουθη εξίσωση αριστοποίησης:

$$V^*_{UMDP}(\pi) = \max_a \{ \pi \cdot q(i,a) + \beta \cdot V^*_{UMDP}(T(\pi, \alpha)) \}, \pi \in \Pi$$

Όπου
$$T(\pi, \alpha) = \pi \cdot P^a, \pi \in \Pi, a \in A.$$

Αποδεικνύεται ότι
$$V^*_{UMDP} \leq V^*(\pi) \quad \forall \pi \in \Pi.$$

Παρόλο που το πρόβλημα UMDP είναι απλούστερο από το αντίστοιχο POMDP, ο υπολογισμός της V^*_{UMDP} δεν είναι απλός και σε γενικές γραμμές ακολουθούνται οι ίδιες μέθοδοι προσέγγισης με την V^* . Επομένως το κάτω φράγμα V^*_{UMDP} για την V^* έχει μόνο θεωρητική αξία.

2) Μέσω των συναρτήσεων V_n (Lovejoy [76])

Μια ακολουθία κάτω φραγμάτων για την V^* σχηματίζεται μέσω της ακολουθίας των συναρτήσεων $\{V_n\}$ που ορίζονται από την (9.2.3). Για $n=0,1,2,\dots$ έχουμε:

$$V_n(\pi) + \frac{\beta^n}{1-\beta} \cdot q_{\min} \leq V^*(\pi) \quad \forall \pi \in \Pi. \quad \underline{4.2.5}$$

$$q_{\min} = \min_{i,a} q(i,a).$$

Για $n=0$ έχουμε την ειδική περίπτωση

$$\frac{1}{1-\beta} \cdot q_{\min} \leq V^*(\pi) \quad \forall \pi \in \Pi.$$

3) Μέσω κάτω φραγμάτων των συναρτήσεων V_n (Lovejoy [76]).

Έστω $v(\pi), \pi \in \Pi$ μια κατά τμήματα γραμμική και κυρτή συνάρτηση. Όπως είναι γνωστό, ή συνάρτηση Hv είναι επίσης κατά τμήματα γραμμική και κυρτή συνάρτηση, πού καθορίζεται μέσω ενός πεπερασμένου συνόλου Γ_H από gradients (διανύσματα) του χώρου \mathbb{R}^N :

$$Hv(\pi) = \max_{\gamma \in \Gamma_H} \pi \cdot \gamma, \quad \pi \in \Pi.$$

Έστω $\gamma(\pi) \in \Gamma_H$ το gradient της Hv στο π , δηλαδή:

$$Hv(\pi) = \pi \cdot \gamma(\pi).$$

Το gradient $\gamma(\pi)$ προσδιορίζεται μέσω του αλγορίθμου του ενός βήματος (βλέπε ενότητα 2.2).

Ένα κάτω φράγμα για τη συνάρτηση Hv κατασκευάζεται ως εξής: Θεωρούμε $G \subseteq \Pi$ ένα πεπερασμένο σύνολο από δ.π (grid points) και Γ_L το σύνολο των gradients της συνάρτησης Hv στα δ.π του συνόλου G , δηλαδή

$$\Gamma_L = \{\gamma(\pi) : \pi \in G\}.$$

Ορίζουμε τη συνάρτηση

$$H_L v(\pi) := \max_{\gamma \in \Gamma_L} \pi \cdot \gamma, \quad \pi \in \Pi.$$

Επειδή $\Gamma_L \subseteq \Gamma_H$ συνάγεται ότι

$$H_L v(\pi) \leq Hv(\pi) \quad \forall \pi \in \Pi.$$

Ακολουθώς κατασκευάζεται η ακολουθία των συναρτήσεων $\{V_{L_n}\}$

$$V_{L_n} = H_L V_{L_{n-1}}, \quad n=1,2,3,\dots$$

Όπου $V_{L_0} = 0$ (μηδενική συνάρτηση).

Σημειώνουμε ότι το σύνολο G παραμένει το ίδιο για όλα τα n . Με άλλα λόγια, το σύνολο των gradients της συνάρτησης V_{L_n} είναι το σύνολο των gradients της $H_L V_{L_{n-1}}$ στα δ.π του G . Αποδεικνύεται επαγωγικά ότι:

$$V_{L_n}(\pi) \leq V_n(\pi) \quad \forall \pi \in \Pi.$$

Μια ακολουθία κάτω φραγμάτων για την V^* σχηματίζεται μέσω της ακολουθίας των συναρτήσεων $\{V_{L_n}\}$. Για $n=0,1,2,\dots$ έχουμε:

$$V_{L_n}(\pi) + \frac{\beta^n}{1-\beta} \cdot q_{\min} \leq V^*(\pi) \quad \forall \pi \in \Pi. \quad \underline{\underline{4.2.6}}$$

Προφανώς το κάτω φράγμα της (4.2.6) ως προσέγγιση της V^* υπολείπεται του αντίστοιχου κάτω φράγματος της (4.2.5), όμως ο υπολογισμός του είναι γενικά πολύ απλούστερος.

Απαλοιφή αποφάσεων πού δεν είναι άριστες για κάποιο δ.π

Η εύρεση άνω και κάτω φραγμάτων για την βέλτιστη συνάρτηση τιμών V^* είναι δυνατόν να συμβάλει στην απαλοιφή μη άριστων αποφάσεων για κάποιο δ.π στο πρόβλημα POMDP για άπειρο χρονικό ορίζοντα.

Θεωρούμε τη συνάρτηση $h: \Pi \times A \times B(\Pi) \rightarrow \mathbb{R}$

πού ορίζεται ως:

$$h(\pi, a, u) := \pi \cdot q^a + \beta \sum_{\theta} \{\theta / \pi, \alpha\} \cdot u(T(\pi, \theta, \alpha)), \quad \forall \pi \in \Pi, \alpha \in A, u \in B(\pi).$$

(βλέπε ενότητα 1.4).

Έστω V_L, V_U κάτω και άνω φράγμα αντίστοιχα για την V^* :

$$V_L(\pi) \leq V^*(\pi) \leq V_U(\pi) \quad \forall \pi \in \Pi.$$

Η επόμενη πρόταση παρέχει ένα κριτήριο απαλοιφής μη άριστων αποφάσεων για κάποιο δ.π στο πρόβλημα του άπειρου χρονικού ορίζοντα.

Πρόταση 4.2.1 :(MC Queen)[100]

Αν για κάποιο $\pi \in \Pi$ και κάποια απόφαση a ισχύει

$$h(\pi, a, V_U) < HV_L(\pi),$$

τότε η απόφαση a δεν είναι άριστη για το π στο πρόβλημα του άπειρου χρονικού ορίζοντα. □

4.3. Προσεγγίσεις της άριστης συνάρτησης τιμών για άπειρο χρονικό ορίζοντα και προσδιορισμός σχεδόν άριστων πολιτικών μέσω φραγμάτων.

Στην παρούσα ενότητα θα ασχοληθούμε με προσεγγίσεις της άριστης συνάρτησης τιμών V^* καθώς και με τον προσδιορισμό σχεδόν άριστων πολιτικών, που βασίζονται σε συναρτήσεις φραγμάτων της V^* .

Ο Hauskrecht [48], λαμβάνοντας δύο οποιαδήποτε άνω και κάτω φράγματα ως αρχικές προσεγγίσεις της V^* και εφαρμόζοντας την επαναληπτική μέθοδο τιμών (value-iteration) κατασκευάζει νέα φράγματα, που αποτελούν αυθαίρετα καλές προσεγγίσεις της V^* και οι προβαλλόμενες από αυτά πολιτικές (lookahead-controllers) είναι σχεδόν άριστες πολιτικές. Το πλήθος των βημάτων (επαναλήψεων) που απαιτούνται εξαρτάται από την «απόσταση» των αρχικών φραγμάτων και την επιθυμητή ακρίβεια της προσέγγισης. Εμείς επεκτείνουμε αυτά τα αποτελέσματα εφαρμόζοντας προσεγγιστικό αλγόριθμο των ακρότατων σημείων. Επιπλέον υπολογίζουμε τον αριθμό των βημάτων καθώς και το προκαθορισμένο σφάλμα του αλγορίθμου, ώστε να επιτύχουμε οποιαδήποτε επιθυμητή ακρίβεια προσέγγισης. Θα περιορισθούμε μόνο σε προσεγγίσεις προβλημάτων POMDP στα πλαίσια του κριτηρίου μεγιστοποίησης των ολικών εσόδων σε άπειρο χρονικό ορίζοντα. Οι αντίστοιχες προσεγγίσεις στα πλαίσια του κριτηρίου ελαχιστοποίησης του ολικού κόστους για άπειρο χρονικό ορίζοντα είναι ανάλογες.

Έστω $V_L(\pi)$, $V_U(\pi)$, $\pi \in \Pi$ συνάρτηση κάτω και άνω φράγματος αντίστοιχα για την συνάρτηση V^* :

$$V_L(\pi) \leq V^*(\pi) \leq V_U(\pi), \pi \in \Pi.$$

Πρόταση 4.3.1: Milos Hauskrecht [48]

Ας είναι $\eta = \|V_U - V_L\| = \sup_{\pi} |V_U(\pi) - V_L(\pi)|$ και V οποιοδήποτε από τα δύο φράγματα,

δηλαδή $V = V_L$ ή $V = V_U$. Τότε

i) $\|V - V^*\| \leq \eta$

ii) Αν δ είναι η συνάρτηση ελέγχου για την οποία

$$H_\delta V = HV$$

δηλαδή

$$\delta(\pi) := \arg \max_a \left\{ \pi q^a + \beta \cdot \sum_{\theta=1}^M \{\theta/\pi, \alpha\} \cdot V(T(\pi, \theta, \alpha)) \right\}, \pi \in \Pi$$

τότε η στάσιμη πολιτική δ^∞ είναι $\frac{2-\beta}{1-\beta} \cdot \eta$ -άριστη, δηλαδή:

$$\|V(./\delta) - V^*\| \leq \frac{2-b}{1-b} \cdot h \quad W$$

Ο αριθμός η εκφράζει την απόσταση (μέγιστη απόλυτη διαφορά) ανάμεσα στα φράγματα V_L και V_U . Η πρόταση 4.3.1 δηλώνει ότι οποιοδήποτε από τα φράγματα V_L , V_U μπορεί να θεωρηθεί προσέγγιση της βέλτιστης συνάρτησης τιμών V^* με μέγιστο σφάλμα το πολύ η . Επίσης οποιοδήποτε από τα δύο φράγματα V_L, V_U μπορεί να χρησιμοποιηθεί ως εφαλτήριο για την εύρεση $\frac{2-b}{1-b} \cdot h$ -άριστης πολιτικής.

Η δ αντιπροσωπεύει την βέλτιστη συνάρτηση ελέγχου «για ένα βήμα μπροστά» (one-step lookahead controller) αναφορικά με το φράγμα V που επιλέγουμε.

Η συνάρτηση τιμών που αντιστοιχεί στην πολιτική $\delta^\infty, V(./\delta)$, προσεγγίζει τη βέλτιστη συνάρτηση τιμών V^* με μέγιστο σφάλμα το πολύ $\frac{2-b}{1-b} \cdot h$. Η πολιτική

δ^∞ αναφέρεται ως προβαλλόμενη πολιτική από το φράγμα V (lookahead-policy).

Θεωρούμε τις αναγωγικές σχέσεις:

$$V_L^k = HV_L^{k-1}, V_U^k = HV_U^{k-1}, k = 1, 2, 3, \dots$$

$$V_L^0 = V_L, V_U^0 = V_U$$

Εισάγοντας τον τελεστή H_k που συμβολίζει την επαναληπτική χρήση του τελεστή H k φορές, οι παραπάνω σχέσεις γράφονται:

$$V_L^k = H_k V_L \leq H_k V^* = V^* \leq H_k V_U = V_U^k$$

Επομένως οι συναρτήσεις V_L^k, V_U^k αποτελούν αντίστοιχα κάτω και άνω φράγμα για την V^* . Θα αναφερόμαστε σε αυτές ως φράγματα τάξεως k . Επειδή ο τελεστής H_k είναι συστολή modulus β^k έχουμε:

$$\|V_U^k - V_L^k\| \leq \beta^k \|V_L - V_U\| = \beta^k \cdot \eta.$$

Συνοψίζοντας, από τα φράγματα V_L, V_U για την συνάρτηση V^* , εφαρμόζοντας επαναληπτικά τον τελεστή H_k φορές, σχηματίζονται τα φράγματα τάξεως k V_L^k, V_U^k για την V^* των οποίων η απόσταση (μέγιστη απόλυτη διαφορά) είναι το πολύ $\beta^k \cdot \eta$.

Από την πρόταση 4.3.1 (i) συνάγεται ότι τα φράγματα k τάξεως μπορούν να θεωρηθούν προσεγγίσεις της συνάρτησης V^* με μέγιστο σφάλμα το πολύ $\beta^k \cdot \eta$. Πιο συγκεκριμένα έχουμε:

$$\|V^k - V^*\| \leq \beta^k \cdot \eta, \quad \underline{\mathbf{4.3.1}}$$

όπου V^k είναι οποιοδήποτε από τα δύο φράγματα k τάξεως,

δηλαδή $V^k = V_L^k$ ή $V^k = V_U^k$.

Αν $\lambda > 0$ είναι ένα επιθυμητό άνω φράγμα για το σφάλμα προσέγγισης

$$\|V^k - V^*\| \leq \lambda$$

τότε μπορούμε να επιλέξουμε την τάξη k έτσι ώστε:

$$\beta^k \cdot \eta \leq \lambda$$

δηλαδή $k \geq \frac{\ln(\lambda / \eta)}{\ln \beta}$

Εστω δ^k η βέλτιστη συνάρτηση ελέγχου για ένα βήμα μπροστά (one-step lookahead controller) αναφορικά με το φράγμα V^k που επιλέγουμε, δηλαδή:

$$H_{\delta^k} V^k = H V^k \quad \underline{\mathbf{4.3.2}}$$

Από την πρόταση 4.3.1 (ii) συνάγεται ότι η προβαλλόμενη από το φράγμα V^k

πολιτική $(\delta^k)^\infty$ (lookahead policy) είναι $\frac{2-\beta}{1-\beta} \cdot \beta^k \cdot \eta$ - άριστη.

Αν επιθυμούμε η πολιτική $(\delta^k)^\infty$ να είναι λ -άριστη για δοσμένο $\lambda > 0$, τότε μπορούμε να επιλέξουμε την τάξη k έτσι ώστε:

$$\frac{2-\beta}{1-\beta} \cdot \beta^k \cdot \eta \leq \lambda$$

δηλαδή
$$k \geq \frac{\ln(\lambda \cdot (1-\beta) / \eta \cdot (2-\beta))}{\ln \beta}$$

Σημειώνουμε ότι οι προσεγγίσεις πού περιγράψαμε είναι εφικτές μόνο στην περίπτωση όπου εφαρμόζουμε διαδοχικά τον ακριβή αλγόριθμο των ακρότατων σημείων του κεφ. 3 (προκαθορισμένο σφάλμα $\varepsilon=0$), οπότε επιτυγχάνεται ακριβής υπολογισμός των φραγμάτων k τάξεως V_L^k, V_U^k καθώς επίσης ο καθορισμός της συνάρτησης ελέγχου δ^k μέσω της σχέσης (4.3.2) είναι εφικτός.

Ανάλογες προσεγγίσεις είναι δυνατές εφαρμόζοντας διαδοχικά τον προσεγγιστικό αλγόριθμο των ακρότατων σημείων. Επιλέγοντας ως αρχικές συναρτήσεις τις συναρτήσεις φραγμάτων V_L, V_U και εφαρμόζοντας k φορές τον αλγόριθμο των ακρότατων σημείων με προκαθορισμένο σφάλμα $\varepsilon > 0$ υπολογίζονται οι προσεγγίσεις $\tilde{V}_L^k, \tilde{V}_U^k$ των φραγμάτων k τάξης V_L^k, V_U^k αντίστοιχα, με συσσωρευμένο σφάλμα προσέγγισης σ_k σε κάθε περίπτωση (πρβλ. ενότητα 3.4). Συγκεκριμένα,

$$\tilde{V}_L^k = \tilde{H} \tilde{V}_L^{k-1}, \quad \tilde{V}_U^k = \tilde{H} \tilde{V}_U^{k-1}, k=1,2,3\dots$$

$$\tilde{V}_L^0 = V_L, \quad \tilde{V}_U^0 = V_U$$

Για $k=1,2,3,\dots$ έχουμε:

$$\tilde{V}_L^k(\pi) \leq V_L^k(\pi) \leq \tilde{V}_L^k(\pi) + \sigma_k \quad \forall \pi \in \Pi$$

$$\tilde{V}_U^k(\pi) \leq V_U^k(\pi) \leq \tilde{V}_U^k(\pi) + \sigma_k \quad \forall \pi \in \Pi \quad \mathbf{4.3.3}$$

Το συσσωρευμένο σφάλμα προσέγγισης σ_k ικανοποιεί την αναγωγική σχέση (3.4.6) και υπολογίζεται από την σχέση (3.4.7).

Σημειώνουμε ότι η προσέγγιση \tilde{V}_L^k του κάτω φράγματος k τάξεως V_L^k είναι επίσης κάτω φράγμα για την συνάρτηση V^* . Δεν μπορούμε όμως να ισχυριστούμε το ίδιο για την προσέγγιση \tilde{V}_U^k του άνω φράγματος k τάξεως V_U^k πού ενδέχεται να μην είναι άνω φράγμα της V^* . Ωστόσο αμφότερες οι συναρτήσεις $\tilde{V}_L^k, \tilde{V}_U^k$

μπορούν να ληφθούν ως προσεγγίσεις της συνάρτησης V^* . Το σφάλμα αυτών των προσεγγίσεων παρέχεται στην επόμενη πρόταση.

Παρατήρηση

Στα προβλήματα κόστους ισχύει το αντίστροφο. Πράγματι (βλέπε ενότητα 3.5) για $k=1,2,3\dots$ έχουμε:

$$\tilde{V}_L^k(\pi) - \sigma_k \leq V_L^k(\pi) \leq \tilde{V}_L^k(\pi)$$

$$\tilde{V}_U^k(\pi) - \sigma_k \leq V_U^k(\pi) \leq \tilde{V}_U^k(\pi)$$

Επομένως η προσέγγιση $\tilde{V}_U^k(\pi)$ του άνω φράγματος $V_U^k(\pi)$ εξακολουθεί να είναι άνω φράγμα για την συνάρτηση V^* , ενώ η προσέγγιση $\tilde{V}_L^k(\pi)$ του κάτω φράγματος $V_L^k(\pi)$, ενδέχεται να μην είναι κάτω φράγμα για την V^* .

Ωστόσο και στην περίπτωση αυτή αμφότερες οι συναρτήσεις $\tilde{V}_L^k, \tilde{V}_U^k$, μπορούν να ληφθούν ως προσεγγίσεις της V^* .

Πρόταση 4.3.2: Αν $\eta = \|V_U - V_L\| = \sup_{\pi} |V_U(\pi) - V_L(\pi)|$, τότε για $k=1,2,3\dots$

$$\text{i) } \|\tilde{V}_L^k - V^*\| \leq \frac{1-\beta^k}{1-\beta} \cdot \varepsilon + \beta^k \cdot \eta,$$

$$\text{ii) } \|\tilde{V}_U^k - V^*\| \leq \frac{1-\beta^k}{1-\beta} \cdot \varepsilon + \beta^k \cdot \eta.$$

Απόδειξη

i) Από την τριγωνική ιδιότητα της νόρμας supremum και τις σχέσεις (4.3.1), (4.3.3) και (3.4.7) παίρνουμε:

$$\|\tilde{V}_L^k - V^*\| \leq \|\tilde{V}_L^k - V_L^k\| + \|V_L^k - V^*\| \leq \sigma_k + \beta^k \cdot \eta =$$

$$= \frac{1-\beta^k}{1-\beta} \cdot \varepsilon + \beta^k \cdot \eta$$

ii) Παρόμοια □

Σημειώνουμε ότι είναι δυνατόν να επιτευχθούν προσεγγίσεις της συνάρτησης V^* με οποιαδήποτε επιθυμητή ακρίβεια. Συγκεκριμένα αν $\lambda > 0$ είναι ένα επιθυμητό άνω φράγμα για το μέγιστο σφάλμα προσέγγισης, δηλαδή

$$\|\tilde{V}^k - V^*\| \leq \lambda$$

όπου \tilde{V}^k είναι οποιαδήποτε από τις δύο προσεγγίσεις $\tilde{V}_L^k, \tilde{V}_U^k$, τότε μπορούμε να επιλέξουμε προκαθορισμένο σφάλμα ε και τάξη k έτσι ώστε

$$\frac{1-\beta^k}{1-\beta} \cdot \varepsilon + \beta^k \cdot \eta \leq \lambda$$

Μπορούμε π.χ να επιλέξουμε τα ε, k έτσι ώστε:

$$\frac{\varepsilon}{1-\beta} \leq \frac{\lambda}{2}, \quad \beta^k \cdot \eta \leq \frac{\lambda}{2}$$

δηλαδή
$$\varepsilon \leq \frac{\lambda \cdot (1-\beta)}{2}, \quad k \geq \frac{\ln(\frac{\lambda}{2\eta})}{\ln(\beta)}$$

Εστω $\tilde{\delta}^k$ η συνάρτηση ελέγχου που προκύπτει εφαρμόζοντας τον αλγόριθμο των ακρότατων σημείων με προκαθορισμένο σφάλμα $\varepsilon > 0$ στην προσέγγιση \tilde{V}^k , δηλαδή

$$H_{\tilde{\delta}^k} \tilde{V}^k = \tilde{H} \tilde{V}^k \quad \underline{\mathbf{4.3.4}}$$

(βλέπε επίσης παρατήρηση 2 στην ενότητα 3.2)

Η στάσιμη πολιτική $(\tilde{\delta}^k)^\infty$ μπορεί να ληφθεί ως προσέγγιση της άριστης πολιτικής $(\delta^*)^\infty$. Το σφάλμα αυτής της προσέγγισης παρέχεται στην πρόταση 4.3.3

Λήμμα 4.3.1: Αν $\eta = \|V_U - V_L\| = \sup_{\pi} |V_U(\pi) - V_L(\pi)|$, τότε για $k=1,2,3\dots$

$$\|HV^k - V^k\| \leq \beta^k \cdot \eta$$

όπου V^k είναι οποιοδήποτε από τα φράγματα k τάξεως V_L^k, V_U^k .

Απόδειξη

Επειδή

$$V_L^k \leq V^* \leq V_U^k \quad \text{και}$$

$$HV_L^k \leq HV^* = V^* \leq HV_U^k$$

συνάγεται ότι για κάθε $\pi \in \Pi$,

$$\begin{aligned} |HV^k(\pi) - V^k(\pi)| &\leq \max\{|HV^k(\pi) - V^*(\pi)|, |V^k(\pi) - V^*(\pi)|\} \leq \\ &\leq \max\{\|HV^k - V^*\|, \|V^k - V^*\|\} \end{aligned}$$

από την οποία προκύπτει ότι:

$$\|HV^k - V^k\| \leq \max\{\|HV^k - V^*\|, \|V^k - V^*\|\} \quad \underline{4.3.5}$$

Από τις σχέσεις (4.3.1),(4.3.5) και την ακόλουθη σχέση

$$\|HV^k - V^*\| = \|HV^k - HV^*\| \leq \beta \cdot \|V^k - V^*\| \leq \beta^{k+1} \cdot \eta$$

συνάγεται ότι

$$\|HV^k - V^k\| \leq \max\{\beta^{k+1} \cdot \eta, \beta^k \cdot \eta\} = \beta^k \cdot \eta \quad \square$$

Πρόταση 4.3.3: Αν $\eta = \|V_U - V_L\| = \sup_{\pi} |V_U(\pi) - V_L(\pi)|$, \tilde{V}^k είναι οποιαδήποτε από τις προσεγγίσεις $\tilde{V}_L^k, \tilde{V}_U^k$ των φραγμάτων k-τάξεως V_L^k, V_U^k και $\tilde{\delta}^k$ η συνάρτηση ελέγχου που ορίζεται από τη σχέση (4.3.4) τότε η στάσιμη πολιτική $(\tilde{\delta}^k)^\infty$ είναι $f(\varepsilon, \kappa)$ -άριστη, όπου

$$f(\varepsilon, \kappa) := \frac{3 - \beta - 2\beta^k}{(1 - \beta)^2} \cdot \varepsilon + \frac{2 - \beta}{1 - \beta} \cdot \beta^k \cdot \eta \quad \underline{4.3.6}$$

δηλαδή:

$$\|V(. / \tilde{\delta}^k) - V^*\| \leq f(\varepsilon, \kappa)$$

Απόδειξη

Από την τριγωνική ιδιότητα της νόρμας supremum παίρνουμε

$$\|V(. / \tilde{\delta}^k) - V^*\| \leq \|V(. / \tilde{\delta}^k) - \tilde{V}^k\| + \|\tilde{V}^k - V^*\| \quad \underline{4.3.7}$$

Εφαρμόζοντας πάλι την τριγωνική ιδιότητα έχουμε:

$$\|V(. / \tilde{\delta}^k) - \tilde{V}^k\| \leq \|V(. / \tilde{\delta}^k) - H_{\tilde{\delta}^k} \tilde{V}^k\| + \|H_{\tilde{\delta}^k} \tilde{V}^k - \tilde{V}^k\| \quad \underline{4.3.8}$$

Λαμβάνοντας υπόψη ότι η συνάρτηση τιμών για την πολιτική $(\tilde{\delta}^k)^\infty, V(\cdot/\tilde{\delta}^k)$, είναι το σταθερό σημείο του τελεστή $H_{\tilde{\delta}^k}$ και ότι ο τελεστής είναι συστολή modulus β , έχουμε:

$$\|V(\cdot/\tilde{\delta}^k) - H_{\tilde{\delta}^k} \tilde{V}^k\| = \|H_{\tilde{\delta}^k} V(\cdot/\tilde{\delta}^k) - H_{\tilde{\delta}^k} \tilde{V}^k\| \leq \beta \cdot \|V(\cdot/\tilde{\delta}^k) - \tilde{V}^k\| \quad \underline{4.3.9}$$

Από τις σχέσεις (4.3.8),(4.3.9) παίρνουμε

$$\|V(\cdot/\tilde{\delta}^k) - \tilde{V}^k\| \leq \|H_{\tilde{\delta}^k} \tilde{V}^k - \tilde{V}^k\| / (1 - \beta) \quad \underline{4.3.10}$$

Θεωρούμε $V^k = V_L^k$ αν $\tilde{V}^k = \tilde{V}_L^k$

$V^k = V_U^k$ αν $\tilde{V}^k = \tilde{V}_U^k$

Από την τριγωνική ιδιότητα, τις σχέσεις (4.3.3),(4.3.4) και το λήμμα 4.3.1 παίρνουμε:

$$\begin{aligned} \|H_{\tilde{\delta}^k} \tilde{V}^k - \tilde{V}^k\| &= \|\tilde{H} \tilde{V}^k - \tilde{V}^k\| \leq \|\tilde{H} \tilde{V}^k - H \tilde{V}^k\| + \|H \tilde{V}^k - V^k\| + \|V^k - \tilde{V}^k\| \leq \\ &\leq \sigma_{k+1} + \beta^k \cdot \eta + \sigma_k \end{aligned} \quad \underline{4.3.11}$$

Από τις σχέσεις (4.3.10), (4.3.11) παίρνουμε

$$\|V(\cdot/\tilde{\delta}^k) - \tilde{V}^k\| \leq \frac{\sigma_{k+1} + \sigma_k + \beta^k \cdot \eta}{1 - \beta} \quad \underline{4.3.12}$$

Από την πρόταση 4.3.1 έχουμε:

$$\|\tilde{V}^k - V^*\| \leq \sigma_k + \beta^k \cdot \eta \quad \underline{4.3.13}$$

Από τις σχέσεις (4.3.7),(4.3.13) και λαμβάνοντας υπόψη τη σχέση (3.4.7) συνάγεται ότι:

$$\begin{aligned} \|V(\cdot/\tilde{\delta}^k) - V^*\| &\leq \frac{\sigma_{k+1} + \sigma_k + \beta^k \cdot \eta}{1 - \beta} + \sigma_k + \beta^k \cdot \eta = \frac{(2 - \beta) \cdot \sigma_k + \sigma_{k+1}}{1 - \beta} + \frac{2 - \beta}{1 - \beta} \cdot \beta^k \cdot \eta \\ &= \frac{3 - \beta - 2\beta^k}{(1 - \beta)^2} \cdot \varepsilon + \frac{2 - \beta}{1 - \beta} \cdot \beta^k \cdot \eta = f(\varepsilon, k) \end{aligned}$$

Επομένως η στάσιμη πολιτική $(\tilde{\delta}^k)^\infty$ είναι $f(\varepsilon, k)$ -άριστη. \square

Αν επιθυμούμε η πολιτική $(\tilde{\delta}^k)^\infty$ να είναι λ -άριστη για δοσμένο $\lambda > 0$, τότε μπορούμε να επιλέξουμε το προκαθορισμένο σφάλμα ε και την τάξη k έτσι ώστε

$$f(\varepsilon, k) \leq \lambda$$

Μπορούμε π.χ να επιλέξουμε τα ε, k έτσι ώστε:

$$\frac{3-\beta}{(1-\beta)^2} \cdot \varepsilon \leq \frac{\lambda}{2} \quad , \quad \frac{2-\beta}{1-\beta} \cdot \beta^k \cdot \eta \leq \frac{\lambda}{2} ,$$

δηλαδή

$$\varepsilon \leq \frac{\lambda'(1-\beta)^2}{2(3-\beta)} \quad , \quad k \geq \frac{\ln\left(\frac{\lambda(1-\beta)}{2(2-\beta)\eta}\right)}{\ln \beta} . \quad \square$$

Παρατηρήσεις

1) Σε όλες τις περιπτώσεις, είτε χρησιμοποιούμε τον ακριβή είτε τον προσεγγιστικό αλγόριθμο των ακρότατων σημείων, ο ελάχιστος αριθμός βημάτων (επαναλήψεων) k πού απαιτείται ώστε να επιτύχουμε οποιαδήποτε επιθυμητή ακρίβεια στην προσέγγιση της V^* ή της άριστης πολιτικής εξαρτάται από την απόσταση η των αρχικών άνω και κάτω φραγμάτων. Επομένως η διαδικασία προσέγγισης μπορεί να επιταχυνθεί σημαντικά αν το η είναι μικρό, δεδομένου ότι το η αποτελεί μέτρο του σφάλματος της αρχικής προσέγγισης της V^* (πρόταση 4.3.1(i)).

2) Επειδή έχουμε τη δυνατότητα ως αρχική προσέγγιση της V^* να επιλέξουμε ένα από τα δύο αρχικά φράγματα (είτε το κάτω είτε το άνω φράγμα), προφανώς είναι λογικό να επιλέξουμε το απλούστερο από αυτά. Σημαντικό πλεονέκτημα ως αρχική προσέγγιση της V^* παρουσιάζουν τα φράγματα πού υπολογίζονται μέσω της V_{MDP}^* (βλέπε παράγραφο 4.2) λόγω της απλότητάς τους.

ΣΥΜΠΕΡΑΣΜΑΤΑ

Προσεγγίσεις της άριστης συνάρτησης τιμών και της άριστης πολιτικής σε άπειρο χρονικό ορίζοντα επιτυγχάνονται με οποιαδήποτε επιθυμητή ακρίβεια μέσω διαδοχικών επαναλήψεων του αλγόριθμου των ακροτάτων σημείων που περιγράψαμε στο κεφάλαιο 3. Καλές προσεγγίσεις επιτυγχάνονται επίσης αν το “κατάλοιπο Bellman” σε κάποιο βήμα (επανάληψη) είναι αρκούντως μικρό.

Η μέθοδος των φραγμάτων αναφέρεται στην επιλογή άνω και κάτω φραγμάτων ως αρχικών προσεγγίσεων της άριστης συνάρτησης τιμών η οποία συνοδεύεται από την επαναληπτική εφαρμογή του αλγόριθμου των ακροτάτων σημείων για την δημιουργία νέων προσεγγίσεων. Από τις προβαλλόμενες συναρτήσεις ελέγχου των νέων προσεγγίσεων (lookahead controllers) κατασκευάζονται σχεδόν άριστες πολιτικές. Επιτάχυνση της διαδικασίας είναι δυνατή αν η απόσταση των αρχικών φραγμάτων είναι μικρή.

Σε κάθε περίπτωση υπολογίζεται ο απαιτούμενος αριθμός επαναλήψεων καθώς και το προκαθορισμένο σφάλμα του αλγόριθμου των ακροτάτων σημείων, έτσι ώστε να επιτυγχάνεται προσέγγιση με οποιαδήποτε επιθυμητή ακρίβεια.