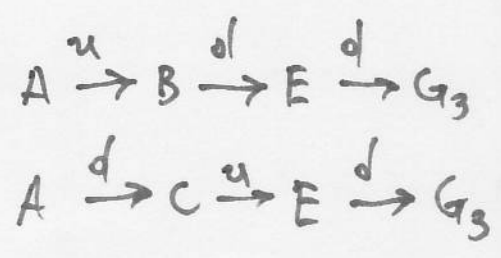
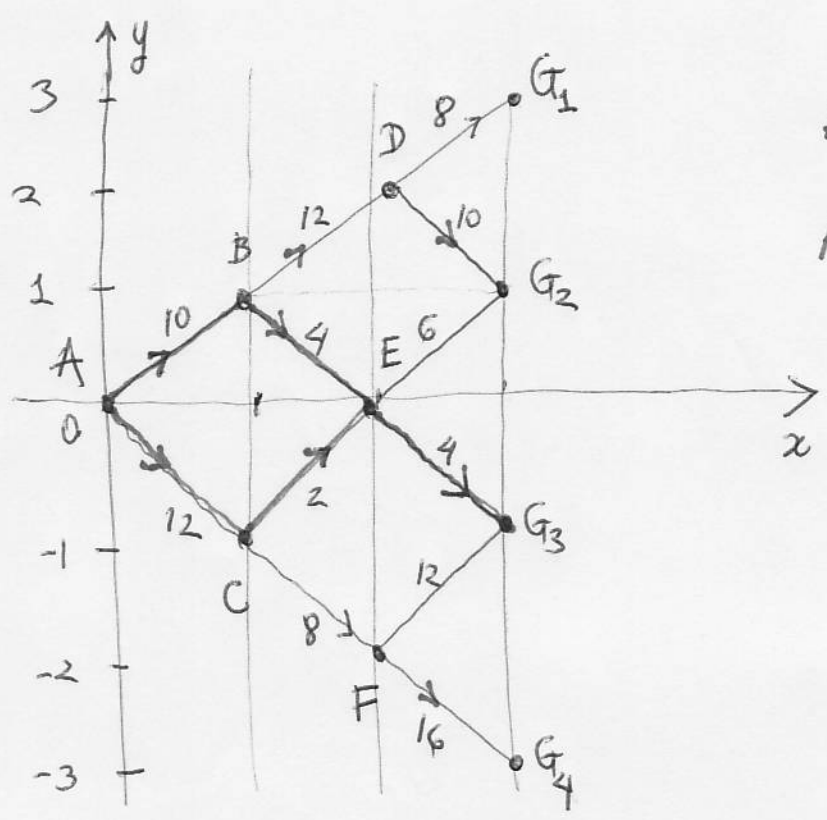


ΕΝΑ ΣΤΟΧΑΣΤΙΚΟ ΠΡΟΒΛΗΜΑ

ΔΙΑΔΡΟΜΗΣ ΕΛΑΧΙΣΤΟΥ ΚΟΣΤΟΥΣ.



↑ οι μη στοχαστικές
 άριστες πολιτικές

Θέλουμε να μεταβούμε από την κορυφή A στην "ευθεία" $\{G_i\}$ με ελάχιστο κόστος. Υπάρχει όμως αβεβαιότητα στο γεγονός εάν ακολουθείται η βέλπτη διαδρομή. Δηλαδή ενώ σε κάθε κορυφή $(x_t, y_t) = S_t$ εκδίδουμε την οδηγία $\beta_t^* \in \{u, d\}$ η οδηγία ακολουθείται με πιθανότητα: $p_t(x_t, y_t) = P\{\alpha_t = u \mid \beta_t^* = u\} = P\{\alpha_t = d \mid \beta_t^* = d\}$

και δεν ακολουθείται με πιθανότητα:

$$1 - p_t(x_t, y_t) = q_t(x_t, y_t) = P\{\alpha_t = d \mid \beta_t^* = u\} = P\{\alpha_t = u \mid \beta_t^* = d\}$$

$$V(x, y) = V_t(s_t) = \min_{\beta_t \in \{u, d\}} \left\{ \mathbb{E} \left[c(\alpha_t, s_t) + V_{t+1}(s_{t+1}(\alpha_t, s_t)) \mid \beta_t \right] \right\}$$

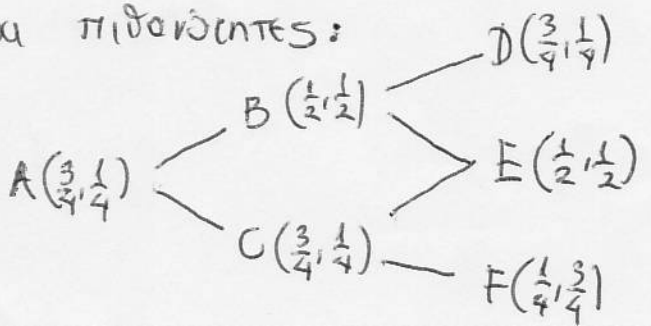
$$= \min \left\{ \mathbb{E} \left[c(\alpha_t, s_t) + V_{t+1}(s_{t+1}(\alpha_t, s_t)) \mid \beta_t = u \right], \right. \\ \left. \mathbb{E} \left[c(\alpha_t, s_t) + V_{t+1}(s_{t+1}(\alpha_t, s_t)) \mid \beta_t = d \right] \right\}$$

$$= \min \left\{ P\{\alpha_t = u \mid \beta_t = u\} \cdot (c(u, s_t) + V_{t+1}(s_{t+1}(u, s_t))) + \right. \\ \left. + P\{\alpha_t = d \mid \beta_t = u\} \cdot (c(d, s_t) + V_{t+1}(s_{t+1}(d, s_t))), \right. \\ \left. P\{\alpha_t = u \mid \beta_t = d\} \cdot (c(u, s_t) + V_{t+1}(s_{t+1}(u, s_t))) + \right. \\ \left. + P\{\alpha_t = d \mid \beta_t = d\} \cdot (c(d, s_t) + V_{t+1}(s_{t+1}(d, s_t))) \right\}$$

$$= \min \left\{ p_t(x, y) \cdot (u(x, y) + V(x_{t+1}, y_{t+1})) + \right. \\ \left. + q_t(x, y) \cdot (d(x, y) + V(x_{t+1}, y_{t-1})), \right. \\ \left. q_t(x, y) \cdot (u(x, y) + V(x_{t+1}, y_{t+1})) + \right. \\ \left. + p_t(x, y) \cdot (d(x, y) + V(x_{t+1}, y_{t-1})) \right\}$$

με Ευρωπαϊκές συνθήκες: $V(3, 3-2i) = V(G_{i+1}) = 0, i=0, \dots, 3$

και πιθανότητες:



$x=2$:
$$V(D) = \min \left\{ p_2(D)(u(D) + \hat{V}(G_1)) + q_2(D)(d(D) + \hat{V}(G_2)), \right. \\ \left. p_2(D)(d(D) + \hat{V}(G_2)) + q_2(D)(u(D) + \hat{V}(G_1)) \right\}$$

$$= \min \left\{ \left(\frac{17}{2}, \frac{19}{2}\right) \right\} = \frac{17}{2}, \quad \beta_2^* = u$$

$$V(E) = \min \left\{ p_1(E)(u(E) + V(G_2)) + q_1(E)(d(E) + V(G_3)), \right. \\ \left. p_1(E)(d(E) + V(G_3)) + q_1(E)(u(E) + V(G_2)) \right\}$$

$$= \min \left\{ (5, 5) \right\} = 5, \quad \beta_2^* \in \{u, d\}$$

$$V(F) = \min \left\{ p_1(F)(u(F) + V(G_3)) + q_1(F)(d(F) + V(G_4)), \right. \\ \left. p_1(F)(d(F) + V(G_4)) + q_1(F)(u(F) + V(G_3)) \right\}$$

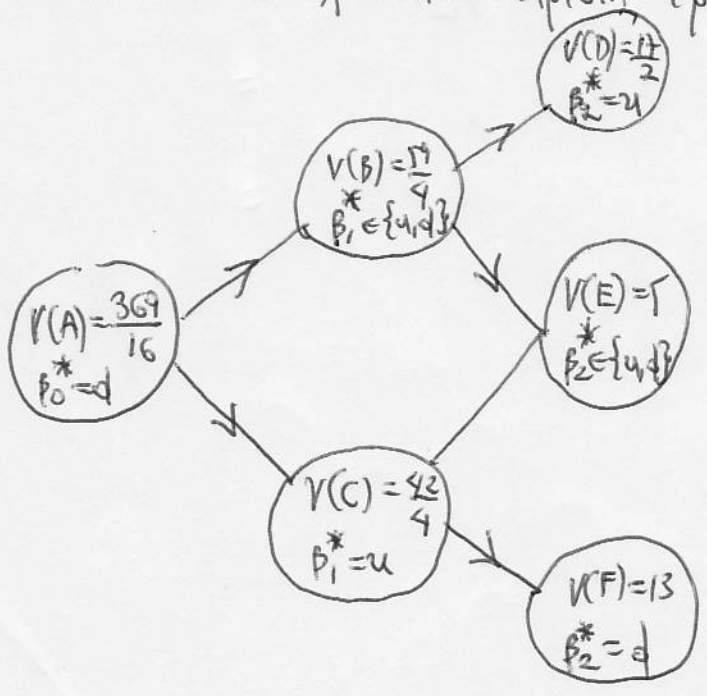
$$= \min \left\{ 15, (13) \right\} = 13, \quad \beta_2^* = d$$

$x=1$:
$$V(B) = \min \left\{ \left(\frac{59}{4}, \frac{59}{4}\right) \right\} = \frac{59}{4}, \quad \beta_1^* \in \{u, d\}$$

$$V(C) = \min \left\{ \left(\frac{42}{4}, \frac{70}{4}\right) \right\} = \frac{42}{4}, \quad \beta_1^* = u$$

$x=0$:
$$V(A) = \min \left\{ \frac{387}{16}, \left(\frac{369}{16}\right) \right\} = \frac{369}{16}, \quad \beta_0^* = d$$

FTG1 η στοχαστική άριστη τροχιά είναι:



ΕΝΑ ΑΙΤΙΟΚΡΑΤΙΚΟ ΠΡΟΒΛΗΜΑ ΦΟΡΤΩΣΗΣ

(Η' ΠΩΣ ΝΑ ΛΥΝΟΥΜΕ ΠΡΟΒΛΗΜΑΤΑ ΑΚΕΡΑΙΟΥ ΠΡΟΓΡΑΜΜΑΤΙΣΜΟΥ ΧΡΗΣΙΜΟΠΟΙΩΝΤΑΣ ΔΥΝΑΜΙΚΟ ΠΡΟΓΡΑΜΜΑ.)

Έστω $W = \text{σταθ}$ το μέγιστο δυνατό φορτίο που μπορεί να μεταφερθεί (από ένα ας πούμε αεροκράνος), και N είδη αντικειμένων τα οποία θέλουμε να μεταφέρουμε. Θετούμε

$$u_t = \text{αξία του } t\text{-αντικειμένου} \quad 1 \leq t \leq N$$

$$w_t = \text{το βάρος του } t\text{-αντικειμένου}$$

Θέλουμε η αξία των αντικειμένων που τελικώς θα μεταφερθεί να είναι μέγιστη. Δηλαδή θέλουμε να λύσουμε το εφεής πρόβλημα ακεραίου προγρ.

$$\max_{(\alpha_1, \dots, \alpha_N)} \sum_{t=1}^N u_t \alpha_t = \sum_{t=1}^N u_t \alpha_t^* = V(\alpha_1^*, \dots, \alpha_N^*)$$

έτσι ώστε: $\sum_{t=1}^N w_t \alpha_t \leq W$ και $\alpha_t \in \mathbb{Z}^{\geq 0}, 1 \leq t \leq N$

↑
υπό τον
περιορισμό.

Για την εξίσωση Bellman έχουμε:

$$V_t(s_t) = \max_{\alpha_t \in A_t(s_t)} \left\{ b_t(\alpha_t, s_t) + V_{t+1}(s_{t+1}(\alpha_t, s_t)) \right\}$$

$$s_t = \text{κοινότητα} \in \{0, 1, 2, \dots, W\}$$

$$A_t(s_t) = \{0, 1, \dots, \lfloor s_t/w_t \rfloor\} \ni \alpha_t = \# \text{ κομμάτιων } t\text{-είδος από το}$$

$$b_t(\alpha_t, s_t) = \text{benefit} = \alpha_t \cdot u_t$$

$$s_{t+1}(\alpha_t, s_t) = s_t - \alpha_t \cdot w_t, \quad 1 \leq t < N$$

$$t=N \Rightarrow \begin{cases} V_N(s_N) = \alpha_N \cdot u_N; & \alpha_N \in A_N(s_N) = \{ \lfloor s_N / w_N \rfloor \} \\ s_N \in \{0, 1, \dots, W\} & \text{(η ευρωπαϊκή συνθήκη)} \end{cases}$$

Αριθμητικό παράδειγμα: Για $W=5$, $U(\alpha_1, \alpha_2, \alpha_3)$

Δίνεται:

t	w _t	u _t
3	1	30
2	3	80
1	2	65

 $\Leftrightarrow \begin{cases} \max (\underbrace{65\alpha_1 + 80\alpha_2 + 30\alpha_3}_{U(\alpha_1, \alpha_2, \alpha_3)}) \\ \text{O.T.Π: } 2\alpha_1 + 3\alpha_2 + \alpha_3 \leq 5 \end{cases}$

\uparrow
W

t=3:

$$V_3(s_3) = u_3 \cdot \lfloor s_3 / w_3 \rfloor$$

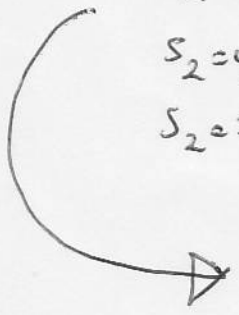
$$s_3 = \{0, \dots, 5\}$$

s ₃	0	1	2	3	4	5
V ₃ (s ₃)	0	30	60	90	120	150
α ₃ *	0	1	2	3	4	5

$$t=2: V_2(s_2) = \max_{\alpha_2 \in A_2(s_2)} \{ \alpha_2 u_2 + V_3(s_2 - w_2 \alpha_2) \} = \max_{\alpha_2 \in A_2(s_2)} \{ 80\alpha_2 + V_3(s_2 - 3\alpha_2) \}$$

$s_2 \in \{0, \dots, 5\} \Rightarrow$

- $s_2=0 \Rightarrow A_2(0) = \{0\} \Rightarrow V_2(0) = 80 \cdot 0 + V_3(0) = 0 \quad \alpha_2^* = 0$
- $s_2=1 \Rightarrow A_2(1) = \{0\} \Rightarrow V_2(1) = 80 \cdot 0 + V_3(1) = 30 \quad \alpha_2^* = 0$
- $s_2=2 \Rightarrow A_2(2) = \{0\} \Rightarrow V_2(2) = 80 \cdot 0 + V_3(2) = 60 \quad \alpha_2^* = 0$
- $s_2=3 \Rightarrow A_2(3) = \{0, 1\} \Rightarrow V_2(3) = \max\{90, 80\} = 90 \quad \alpha_2^* = 0$
- $s_2=4 \Rightarrow A_2(4) = \{0, 1\} \Rightarrow V_2(4) = \max\{120, 110\} = 120 \quad \alpha_2^* = 0$
- $s_2=5 \Rightarrow A_2(5) = \{0, 1\} \Rightarrow V_2(5) = \max\{150, 140\} = 150 \quad \alpha_2^* = 0$



s ₂	0	1	2	3	4	5
V ₂ (s ₂)	0	30	60	90	120	150
α ₂ *	0	0	0	0	0	0

t=1 Επειδή πάντα έχουμε ολοκληρω το βάρος

44

$$\begin{aligned} S_1 \in \{W=5\} &\Rightarrow V_1(S_1) = \max_{\alpha_1 \in A_1(S_1)} \{u_1 \alpha_1 + V_2(S_1 - w_1 \alpha_1)\} \\ &= \max_{\alpha_1 \in A_1(S_1)} \{65\alpha_1 + V_2(5 - 2\alpha_1)\} \end{aligned}$$

$$A_1(S_1) = A_1(5) = \{0, \dots, \lfloor 5/w_1 \rfloor\} = \{0, 1, 2\}$$

2

$$V_1(S_1) = \max \{0 + V_2(5), 65 + V_2(3), 130 + V_2(1)\} = 160$$

κ' $\alpha_1^* = 2$

Βέλτιστη τροχιά

$$\begin{array}{ccc} V_1(5) = 160 & \begin{array}{c} \downarrow \\ V_2(1) = 30 \end{array} & \begin{array}{c} \downarrow \\ V_3(1) = 30 \end{array} \\ \alpha_1^* = 2 \Rightarrow S_2^* = \underbrace{S_1^*}_{5} - w_1 \alpha_1^* = 1 & \alpha_2^* = 0 \Rightarrow S_3^* = \underbrace{S_2^*}_{1} - w_2 \alpha_2^* = 1 & \alpha_3^* = 1 \end{array}$$



$$\begin{cases} \max (65\alpha_1 + 80\alpha_2 + 30\alpha_3) = U(2, 0, 1) = 160 \\ (\alpha_1, \alpha_2, \alpha_3) \in \text{Ακέραια} \\ 2\alpha_1 + 3\alpha_2 + \alpha_3 \leq 5 \end{cases}$$

Απλοποιούσε την προηγούμενη εξίσωση Ιπώκτος
 οι συν/σεις $b_t(\alpha_t, s_t)$ να είναι μη' στοχαστικές κ'
 επειδή $x_t \in S = \mathcal{F}_t(\Omega)$, $A_t(s_t) = A(s_t)$.

$$\begin{aligned} \mathbb{E}_{\mathcal{F}_t} (V_{t+1}(s_{t+1}(\alpha_t, x_t))) &= \mathbb{E} [V_{t+1}(\mathcal{F}_{t+1}) \mid \mathcal{F}_t = \alpha(x)] = \\ &= \sum_{y \in S} P\{\mathcal{F}_{t+1} = y \mid \mathcal{F}_t = \alpha(x)\} V_{t+1}(y) = \sum_{y \in S} P_{xy}(\alpha) V_{t+1}(y) : (2.1) \end{aligned}$$

Έτσι η εξίσωση Bellman γίνεται:

$$V_t(x) = \sup_{\alpha \in A(x)} \left\{ b_t(\alpha(x), x) + \beta \sum_{y \in S} P_{xy}(\alpha) V_{t+1}(y) \right\} : (2.2)$$

Σε τέτοιου είδους προβλήματα μας ενδιαφέρει η συμπεριφορά της $V_t(x)$ όταν $t \rightarrow \infty$ (infinite horizon problems). Γνωρίζουμε άλλωστε ότι μια ερгодική (aperiodic κ' positive recurrent) κ' αδόγητη αλυσίδα Markov, συγχλίνει όταν $t \rightarrow \infty$, σε μια αδόγητη κατανομή ανεξάρτητη της αρχικής συνθήκης του συστήματος. Αυτή η ιδιότητα είναι που μας επιτρέπει να πούμε με μια πιθανότητα το σύστημα θα βρεθεί σε μια συγκεκριμένη κατάσταση εάν ακολουθήσουμε μια συγκεκριμένη πολιτική.

3

Δεχόμαστε την ύπαρξη του ορίου $V_t(x) \xrightarrow[t \rightarrow \infty]{} V(x)$, $\forall x \in S$
 κ' η σχέση (2.2) μας δίνει:

$$V(x) = \sup_{\alpha \in A} \left\{ b(\alpha(x), x) + \beta \sum_{y \in S} P_{xy}(\alpha) V(y) \right\}, \quad \forall x \in S$$

που είναι μια βαναλίστικη εξίσωση ως προς V .

ΛΥΝΟΝΤΑΣ ΩΣ ΠΡΟΣ V

Α' τρόπος (Policy iteration)

Βήμα 1: Διαλέγουμε αυθαίρετα ένα κανόνα που αντιστοιχεί μια ενέργεια $\alpha(x)$ σε κάθε κατάσταση $x \in S$ (the policy rule) δηλ. $\alpha(x) = x'$, $x, x' \in S = \{1, 2, \dots, r\}$

Βήμα 2: Χρησιμοποιούμε τις r εξισώσεις

$$V(x) = b(\alpha(x), x) + \beta \sum_{y \in S} P_{xy}(\alpha) V(y)$$

για να λύσουμε ως προς τους αγνώστους $V(x)$, $x \in S = \{1, \dots, r\}$

Βήμα 3: Για κάθε μια από τις r καταστάσεις ελέγχουμε εάν υπάρχει άλλο policy rule $\tilde{\alpha}(x) = x''$ που να μας οδηγεί σε μεγαλύτερο $V(x)$ δηλ. για την κατάσταση x ελέγχουμε:

$$b(\tilde{\alpha}(x), x) + \beta \sum_{y \in S} P_{xy}(\tilde{\alpha}) V(y) > b(\alpha(x), x) + \beta \sum_{y \in S} P_{xy}(\alpha) V(y)$$

εάν υπάρχει τέτοιο $\tilde{\alpha}$ θέτουμε $\alpha(x) \equiv \tilde{\alpha}(x)$ κ' συνεχίζουμε

$P = (P(x_i, x_j)) = (P\{\mathcal{J}_{t+1} = x_j | \mathcal{J}_t = x_i\}) =$

present						
	E	0.7	0.3	0	0	$P\{\mathcal{J}_{t+1} = G \mathcal{J}_t = G\} = 0.7$
	G	0	0.7	0.3	0	$P\{\mathcal{J}_{t+1} = B \mathcal{J}_t = A\} = 0.4$
	A	0	0	0.6	0.4	$P\{\mathcal{J}_{t+1} = B \mathcal{J}_t = B\} = 1.0$
	B	0	0	0	1	
	x_i/x_j	E	G	A	B	← future

Ενώ δίνεται ο παράγοντας προεξόφλησης $\beta = 0.9$

Αυθαίρετα λοιπόν διαλέγουμε το policy πλε :

$\alpha = \begin{pmatrix} E & G & A & B \\ E & G & E & E \end{pmatrix} \iff$ Αλλάζουμε με νέα μηχανή την παλιά μηχανή, μόνο εάν βρισκόμαστε σε κατάσταση A ή B.

Το σύστημα των γραμμικών εξισώσεων στο βήμα 2 γίνεται:

$$V(E) = b(E, E) + \beta \{ P\{\mathcal{J}_{t+1} = E | \mathcal{J}_t = E\} V(E) + P\{\mathcal{J}_{t+1} = G | \mathcal{J}_t = E\} V(G) \}$$

$$V(G) = b(G, G) + \beta \{ P\{\mathcal{J}_{t+1} = G | \mathcal{J}_t = G\} V(G) + P\{\mathcal{J}_{t+1} = A | \mathcal{J}_t = G\} V(A) \}$$

$$V(A) = b(E, E) + \beta \{ P\{\mathcal{J}_{t+1} = E | \mathcal{J}_t = E\} V(E) + P\{\mathcal{J}_{t+1} = G | \mathcal{J}_t = E\} V(G) \}$$

$$V(B) = b(E, E) + \beta \{ P\{\mathcal{J}_{t+1} = E | \mathcal{J}_t = E\} V(E) + P\{\mathcal{J}_{t+1} = G | \mathcal{J}_t = E\} V(G) \}$$

$$\Leftrightarrow \begin{bmatrix} V(E) \\ V(G) \\ V(A) \\ V(B) \end{bmatrix} = \begin{bmatrix} 100 \\ 80 \\ -100 \\ -100 \end{bmatrix} + (0.9) \begin{bmatrix} 0.7 & 0.3 & 0 & 0 \\ 0 & 0.7 & 0.3 & 0 \\ 0.7 & 0.3 & 0 & 0 \\ 0.7 & 0.3 & 0 & 0 \end{bmatrix} \begin{bmatrix} V(E) \\ V(G) \\ V(A) \\ V(B) \end{bmatrix} \Rightarrow$$

$$\Rightarrow \begin{bmatrix} V(E) \\ V(G) \\ V(A) \\ V(B) \end{bmatrix} = \begin{bmatrix} 687.81 \\ 572.19 \\ 487.81 \\ 487.81 \end{bmatrix} = (\phi, 1)$$

Στην γενική περίπτωση το γραμμικό σύστημα για τους αγνώστους $\underline{V} = [V(1), \dots, V(n)]^T$ θα είναι

$$V(x) = b(\alpha(x), x) + \beta \sum_{y=1}^n P_{xy}(\alpha) V(y) \quad ; \quad x=1, 2, \dots, n$$

με λύση: $\underline{V} = (\delta_{ij} - \beta P_{ij}(\alpha))^{-1} \underline{b}$

όπου $\underline{b} = (b(\alpha(1), 1), \dots, b(\alpha(n), n))^T$

Ελέγχουμε τώρα εάν η προηγούμεως αυθαίρετα επιλεγμένη policy rule είναι optimal:

(i) Εάν η μηχανή βρίσκεται σε κατάσταση E, δεν χρειάζεται να την αντικαταστήσουμε άρα $\alpha(E) = E$ είναι optimal

(ii) Εάν η μηχανή βρίσκεται σε G κ' εμείς την αλλάζουμε σε E (δηλ δοκιμάζουμε το εναλλακτικό policy rule

$\tilde{\alpha}(G) = E$) θα έχουμε:

$$\tilde{V}(G) = -100 + (0.9)((0.7) \cdot (687.81) + (0.3) \cdot (572.19))$$

$$= 487.81 < 572.19 = V(G) \Rightarrow$$

\Rightarrow Μένουμε στο παλιό policy rule. $\Leftrightarrow \alpha(G) = G$ optimal

(iii) Έαν η μηχ. βρίσκεται σε A κατάσταση κ' επίσης εναλλακτικό δεν την αλλάζουμε ($\alpha(A) = A$) τότε

$$\tilde{V}(A) = 50 + (0.9) \cdot [(0.6) \cdot (487.81) + (0.4) \cdot (487.81)]$$

$$= 489.03 > 487.81 = V(A) \Rightarrow$$

\Rightarrow Αλλάζουμε σε νέο policy rule \Leftrightarrow $\alpha(A) = A$ optimal.

(iv) Έαν η μηχ. βρίσκεται σε B κατάσταση κ' επίσης εναλλακτικό δεν την αλλάζουμε ($\alpha(B) = B$) τότε

$$\tilde{V}(B) = 10 + (0.9) \cdot [(487.81)(1)] = 449.03 < 487.81 = V(B)$$

\Rightarrow Μένουμε στο παλιό policy rule \Leftrightarrow $\alpha(B) = E$ optimal

Έτσι η νέα policy θα είναι: $\alpha = \begin{pmatrix} E & G & A & B \\ E & G & A & E \end{pmatrix} \Rightarrow$

$$\Rightarrow \underline{V} = \begin{bmatrix} 100 \\ 80 \\ 50 \\ -100 \end{bmatrix} + (0.9) \begin{bmatrix} 0.7 & 0.3 & 0 & 0 \\ 0 & 0.7 & 0.3 & 0 \\ 0 & 0 & 0.6 & 0.4 \\ 0.7 & 0.3 & 0 & 0 \end{bmatrix} \underline{V} \Rightarrow$$

$$\Rightarrow \underline{V} = \begin{bmatrix} 690.23 \\ 575.50 \\ 492.36 \\ 490.23 \end{bmatrix}$$

Θα πρέπει να κάνουμε έλεγχο για optimality στις

νέες τιμές:

- (i) Σαφώς $\alpha(E) = E$.
- (ii) Δέτορας $\tilde{\alpha}(G) = E \Rightarrow \tilde{V}(G) = -100 + (0.9)[(0.7)(690.23) + (0.3)(575.50)] = 490.23 < V(G) = 575.50 \Rightarrow$
 \Rightarrow Εμφέρουμε σε $\alpha(G) = G$.
- (iii) Κατά την ίδια έρροια εμφέρουμε σε $\alpha(A) = A$.
- (iv) Κ' εμφέρουμε σε $\alpha(B) = E$.

Τώρα μπορούμε να πούμε ότι :

βέλτιστη πολιτική : $\alpha = \begin{pmatrix} E & G & A & B \\ E & G & A & E \end{pmatrix}$

και $\underline{V}^T = (690.23, 575.50, 492.36, 490.23)$.

β' Τρόπος Προς τα εμπρός επαναλήψεις της \underline{V}
 (Value function iterations).

Βήμα 1 : Επιλέγουμε αυθαίρετα μια $\underline{V}^{(0)}$ συν/ση

$$\underline{V}^{(0)} = (V^{(0)}(1), \dots, V^{(0)}(n))^T \quad (\text{συνήθως την μηδενική})$$

$$\text{Βήμα 2 : } \begin{cases} V^{(n)}(x) = \max_{\alpha \in A(x)} \left\{ b(\alpha(x), x) + \beta \sum_{y=1}^n P_{xy}(\alpha) V^{(n-1)}(y) \right\} \\ \forall x \in \{1, \dots, n\} = S \end{cases}$$

Βήμα 3

Εάν $\sup_{x \in S} |V^{(n)}(x) - V^{(n-1)}(x)| \leq \epsilon$

όπου ϵ προκαθορισμένο βελάκι, σταματάμε κ' $V^{(n)}$ είναι η άγνωστη συν/ση αλλιώς αυξάνουμε n σε $n+1$ κ' πηγαίνουμε στο Βήμα 2.

Παράδειγμα (Το προηγούμενο παράδειγμα με Value-Function Iteration).

αρχικέ δέτω $V^{(0)}(E) = V^{(0)}(G) = V^{(0)}(A) = V^{(0)}(B) = 0$

Decision \ State	replace		$V^{(1)}(\text{state})$	$\alpha^{(1)} *$
	r	r'		
E	-100	100	100	r'
G	-100	80	80	r'
A	-100	50	50	r'
B	-100	10	10	r'

$$V^{(1)}(E) = \max \left\{ -100 + \beta \cdot \sum_{y \in S} P(\overset{E}{r}(E), y) \overset{0}{V^{(0)}(y)}, \right.$$

$$\left. 100 + \beta \sum_{y \in S} P(\underset{E}{r'}(E), y) \overset{0}{V^{(0)}(y)} \right\} = 100$$

$$V^{(1)}(G) = \max \left\{ -100 + \beta \sum_{y \in S'} P(\overbrace{r(G), y}^E) V^{(0)}(y), \right.$$

$$\left. 80 + \beta \sum_{y \in S} P(\overbrace{r'(G), y}^G) V^{(0)}(y) \right\} = 80$$

$$\underline{\text{optimal}}: \begin{cases} V^{(1)}(A) = 80 \\ \alpha^{(1)*}(A) = A \end{cases} \quad \text{vs} \quad \begin{cases} V^{(1)}(B) = 10 \\ \alpha^{(1)*}(B) = B \end{cases}$$

Decision \ State	r	r'	$V^{(2)}(\text{State})$	$\alpha^{(2)*}$
E	-15.4	184.6	184.6	r'
G	-15.4	143.9	143.9	r'
A	-15.4	80.6	80.6	r'
B	-15.4	19.0	19.0	r'

$$V^{(2)}(E) = \left\{ -100 + \beta \sum_{y \in S'} P(\overbrace{r(E), y}^E) V^{(1)}(y), 100 + \beta \sum_{y \in S'} P(\overbrace{r'(E), y}^E) V^{(1)}(y) \right\}$$

$$= \left\{ -100 + \beta \left[\underbrace{P(E, E)}_{0.7} \underbrace{V^{(1)}(E)}_{100} + \underbrace{P(E, G)}_{0.3} \underbrace{V^{(1)}(G)}_{80} + \underbrace{P(E, A)}_0 V^{(1)}(A) + \underbrace{P(E, B)}_0 V^{(1)}(B) \right], \right.$$

$$\left. 100 + \beta \left[P(E, E) V^{(1)}(E) + P(E, G) V^{(1)}(G) + P(E, A) V^{(1)}(A) + P(E, B) V^{(1)}(B) \right] \right\}$$

$$= \{-15.4, 184.6\} = 184.6, \quad \alpha^{(2)*} = r'$$

State \ Decision	r	r'	$V^{(3)}(\text{State})$	$\alpha^{(3)*}$
E	55.15	255.15	255.15	r/
G	55.15	192.42	192.42	r/
A	55.15	100.36	100.36	r/
B	55.15	27.10	55.15	r

⋮

State \ Decision	r	r'	$V^{(100)}(\text{State})$	$\alpha^{(100)*}$
E	490.22	690.22	690.22	r/
G	490.22	575.49	575.49	r/
A	490.22	492.34	492.34	r/
B	490.22	451.19	490.22	r

Παρατήρηση : Εάν ο αριθμός των states κ' των δυνατών actions είναι μεγάλος οι προς τα εμπρός επαναλήψεις της V δε είναι η καλύτερη μέθοδος