

Μελέτη και Υλοποίηση Ελεγκτών Ρομποτικών Συστημάτων με χρήση Αλγορίθμων Ενισχυτικής Μάθησης

Κόντες Γεώργιος

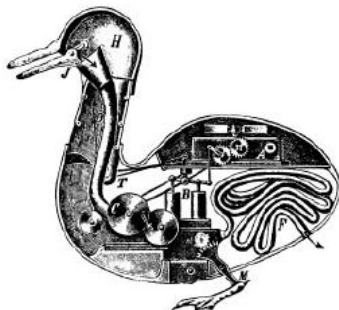
Πολυτεχνείο Κρήτης

22 Ιουλίου, 2009

- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- Ενισχυτική Μάθηση
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- Συμπεράσματα

- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- Ενισχυτική Μάθηση
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- Συμπεράσματα

- Ρομπότ στην αρχαιότητα \Rightarrow εργασία, ψυχαγωγία



- Σήμερα υπάρχει πληθώρα ρομπότ \Rightarrow οχήματα, ιπτάμενα, υποβρύχια, βιομηχανικά...



- Τελικός στόχος \implies Πραγματικά αυτόνομα ρομπότ
- Αυτονομία \implies Προσαρμοστικότητα \implies Μάθηση
- Ενισχυτική Μάθηση

- Εισαγωγή
- **Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)**
- Ενισχυτική Μάθηση
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- Συμπεράσματα

- Υπόθεση Μάρκοβ (π.χ. αυτοκίνητο)
 - αν $x_t \longrightarrow$ είναι η τρέχουσα κατάσταση
 - αν $u_t \longrightarrow$ τρέχουσα ενέργεια που εφαρμόζεται
 - $Pr\{x_{t+1} = x' | x_t, u_t, x_{t-1}, u_{t-1}, \dots, x_0, u_0\} = Pr\{x_{t+1} = x' | x_t, u_t\}$

- Στοχαστική διεργασία διακριτού χρόνου
- Περιγράφεται από το διάνυσμα (X, U, T, R, γ, D)
 - $X \longrightarrow$ Ο χώρος καταστάσεων του προβλήματος
 - $U \longrightarrow$ Ο χώρος των ενεργειών ή αποφάσεων της διεργασίας
 - $T \longrightarrow$ Το Μοντέλο Μετάβασης (Transition Model) της διεργασίας
 $\longrightarrow T(\mathbf{x}, u, \mathbf{x}') = T(\mathbf{x}'|\mathbf{x}, u)$

- Στοχαστική διεργασία διακριτού χρόνου
- Περιγράφεται από το διάνυσμα (X, U, T, R, γ, D)
 - R \longrightarrow Η συνάρτηση Ανταμοιβής (Reward Function) ή Συνάρτηση Κόστους (Cost Function) $\longrightarrow R : X \times U \times X \rightarrow \mathbb{R}$.
Συνηθέστερες μορφές: $r_t = R(\mathbf{x}, u, \mathbf{x}')$, $r_t = R(\mathbf{x}, u)$, $r_t = R(\mathbf{x})$
 - $\gamma \in (0, 1]$ \longrightarrow Ο παράγοντας έκπτωσης (discount factor) της διεργασίας
 - D \longrightarrow Καθορίζει την αρχική κατάσταση της διεργασίας.

- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- **Ενισχυτική Μάθηση**
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- Συμπεράσματα

- Πολιτική (Policy) ή ελεγκτής (controller) = $\pi(\mathbf{x}) : X \rightarrow U$ από καταστάσεις σε ενέργειες
- Μαρκοβιανό Σύστημα

$$\mathbf{x}_0 \sim p(\mathbf{x}_0), \quad (1)$$

$$\mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t, t), \quad (2)$$

$$r_{t+1} = R(\mathbf{x}_t, \mathbf{u}_t, \mathbf{x}_{t+1}, t + 1), \quad (3)$$

$$\mathbf{u}_t \sim \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t, t) = p(\mathbf{u}_t | \mathbf{x}_t, t, \theta), \quad (4)$$

- Τροχιές

$$\mathbf{x}_0 \xrightarrow[r_0]{u_0} \mathbf{x}_1 \xrightarrow[r_1]{u_1} \mathbf{x}_2 \xrightarrow[r_2]{u_2} \mathbf{x}_3 \xrightarrow[r_3]{u_3} \mathbf{x}_4 \xrightarrow[r_4]{u_4} \dots \mathbf{x}_{H-1} \xrightarrow[r_{H-1}]{u_{H-1}} \mathbf{x}_H \quad (5)$$

- Στόχος \implies Εύρεση πολιτικής π_{θ}^* , βέλτιστη προς:

$$J(\theta) = \mathbb{E} \left\{ \sum_{t=0}^H \gamma^t r_t; \theta \right\} \quad (6)$$

- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- Ενισχυτική Μάθηση
- **Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης**
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- Συμπεράσματα

- Συναρτήσεις Αξιολόγησης:

$$V^\pi(x) = \mathbb{E}_\pi \left\{ \sum_{t=0}^H \gamma^t r_t | x_0 = x \right\} \quad (7)$$

και

$$Q^\pi(x, u) = \mathbb{E}_\pi \left\{ \sum_{t=0}^H \gamma^t r_t | x_0 = x, u_0 = u \right\} \quad (8)$$

- Άπληστη Πολιτική

$$\pi'(x) = \arg \max_{u \in U} Q^\pi(x, u) \quad (9)$$

με

$$\forall x \in X, \quad Q^{\pi'}(x, \pi'(x)) \geq Q^\pi(x, \pi(x)) \quad (10)$$

- Γραμμική Εξίσωση Bellman

$$Q^\pi(x, u) = r(x, u) + \gamma \sum_{x' \in X} T(x, u, x') Q^\pi(x', \pi(x')) \quad (11)$$

- Bellman Optimality Equation

$$Q^*(x, u) = r(x, u) + \gamma \sum_{x' \in X} T(x, u, x') \max_{u' \in U} Q^*(x', u') \quad (12)$$

- Value Iteration \implies Bellman Optimality Equation
- Policy Iteration \implies Γραμμική Εξίσωση Bellman

- Προβλήματα:
 - Γνώση του Μοντέλου
 - Διακριτές Ενέργειες
- Λύσεις:
 - Συναρτήσεις Βάσης

$$Q^\pi(\mathbf{x}, \mathbf{u}) = \phi(\mathbf{x}, \mathbf{u})^T \mathbf{w} \quad (13)$$

- Neural Networks
- Βέλτιστες - Δυναμικές Διακριτοποιήσεις

- Αστάθεια Μάθησης
- Διακριτές Ενέργειες

- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- Ενισχυτική Μάθηση
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- **Αλγόριθμοι Πολιτικής Κλίσης**
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- Συμπεράσματα

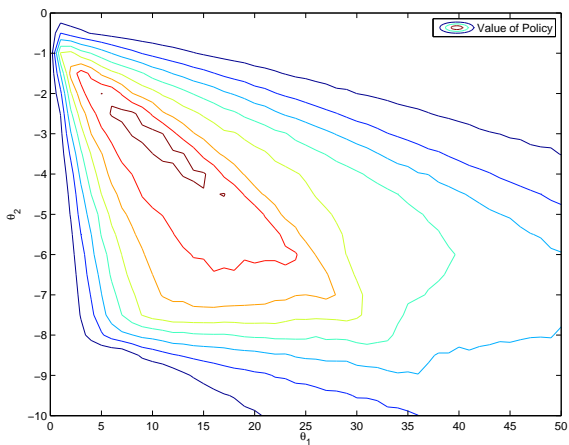
- Διδιάστατο Συνθετικό Πρόβλημα
- x_1 =θέση, x_2 =ταχύτητα, u =επιτάχυνση
- Δυναμικό Μοντέλο

$$x_2(t+1) = x_2(t) - u(t) + \kappa, \quad (14)$$

$$x_1(t+1) = x_1(t) - 0.1x_2(t+1) + \kappa, \quad (15)$$

- $\sigma_\kappa = 0.003$
- $u_t = 1/[1 + \exp(-(\boldsymbol{\theta} + \boldsymbol{\varepsilon}_t)^T \mathbf{x})] - 0.5$
- $\mathbf{x}_0 = [1, 0] + noise, \mathbf{x}_{goal} = [0, 0]$
- $r(t) = \exp(-0.5\|\mathbf{x}(t)\|_2/(0.03)^2)$

- Η συνάρτηση προς βελτιστοποίηση



- Απ' ευθείας στο χώρο των πολιτικών

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha_k \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_k} \quad (16)$$

- Μέθοδος Διακριτών Διαφορών

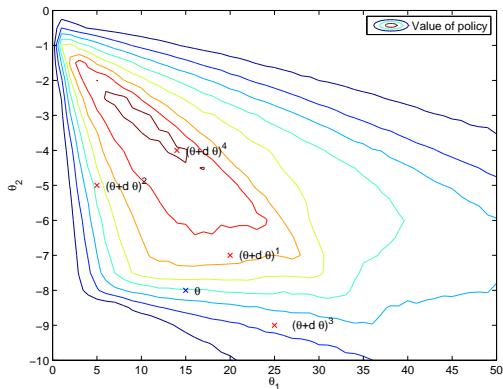
$$\Delta \hat{J} \approx J(\boldsymbol{\theta}_k + \Delta \boldsymbol{\theta}) - J_{ref}^1 \quad (17)$$

και στην συνέχεια

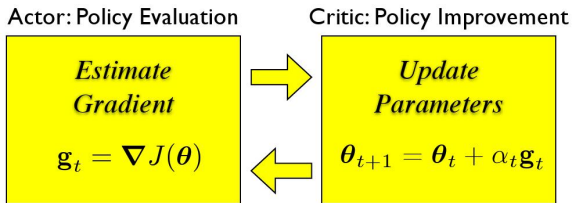
$$\mathbf{g}_{FD} = (\Delta \boldsymbol{\theta}^T \Delta \boldsymbol{\theta})^{-1} \Delta \boldsymbol{\theta}^T \Delta \hat{J} \quad (18)$$

¹Ως J_{ref} χρησιμοποιείται συνήθως το $J(\boldsymbol{\theta}_k)$ ή το $J(\boldsymbol{\theta}_k - \Delta \boldsymbol{\theta})$

- Η μέθοδος διακριτών διαφορών



- Actor - Critic



- Θετικά:
 - Συνεχείς καταστάσεις - ενέργειες
 - Ελεγκτές → Φυσική σημασία για το σύστημα
- Αρνητικά:
 - Ρύθμιση των ρυθμών μάθησης
 - Τοπικά ελάχιστα

- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- Ενισχυτική Μάθηση
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- Συμπεράσματα

- Έστω παρατηρήσιμες μεταβλητές \mathbf{X} , κρυφές μεταβλητές \mathbf{Z} , παράμετροι θ
- Σκοπός:

$$\max p(\mathbf{X}|\theta) = \max_{\mathbf{z}} \sum_{\mathbf{z}} p(\mathbf{X}, \mathbf{z}|\theta) \quad (19)$$

- Ορίζουμε $q(\mathbf{Z})$

- $\forall q(\mathbf{Z})$ ισχύει:

$$\ln p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + KL(q||p), \quad (20)$$

όπου

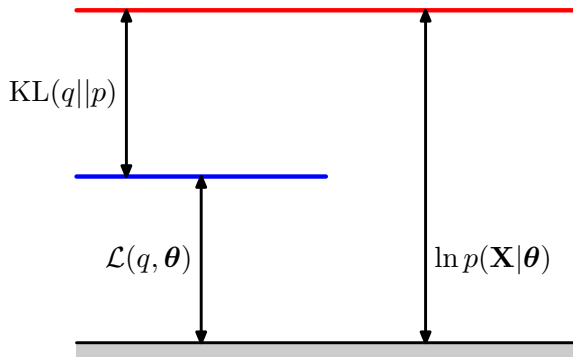
$$\mathcal{L}(q, \theta) = \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{z}|\theta)}{q(\mathbf{z})} \right\} \quad (21)$$

και

$$KL(q||p) = - \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left\{ \frac{p(\mathbf{z}|\mathbf{X}, \theta)}{q(\mathbf{z})} \right\} \quad (22)$$

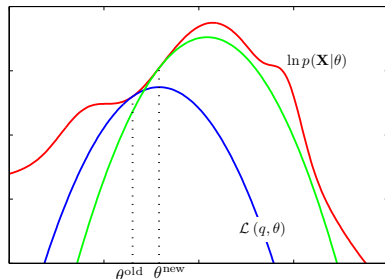
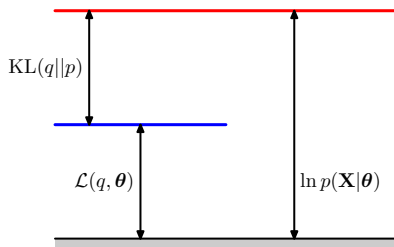
- $KL(q||p) = 0 \iff q(\mathbf{z}) = p(\mathbf{z}|\mathbf{X}, \theta)$

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + KL(q||p), \quad (23)$$



Ο αλγόριθμος EM

- E - step:
 - Παράμετροι θ_{old}
 - Μεγιστοποίηση $\mathcal{L}(q, \theta_{old})$ ως προς $q(\mathbf{Z})$, θ_{old} σταθερό
 - $\max \mathcal{L}(q, \theta_{old}) \rightarrow KL(q||p) = 0 \rightarrow q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta_{old})$
- M - step:
 - $q(\mathbf{Z})$ σταθερή
 - $\max \mathcal{L}(q, \theta)$ ως προς $\theta \rightarrow \theta_{new}$



- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- Ενισχυτική Μάθηση
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- **Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού**
- Πειράματα
- Συμπεράσματα

- Monte Carlo Expectation - Maximization για Ενισχυτική Μάθηση
- Ο ορίζοντας της ΜΔΑ είναι τυχαία μεταβλητή T
- Οι ανταμοιβές θεωρούνται πιθανότητες φανταστικών γεγονότων R
- Χρησιμοποιούμε όλες τις υποτροχιές μήκους t της κάθε τροχιάς
- Οι λανθάνουσες μεταβλητές είναι T, t, ξ

- Για το E-step αποδεικνύεται πως:

$$q^*(\xi, T, t) = p(\xi, T, t|R; \theta_{old}) \quad (24)$$

$$\propto \alpha(T)b(t)p(\xi|t; \theta_{old})p(R|\xi, T) \quad (25)$$

- Μοντέλο μη διαθέσιμο \rightarrow δειγματοληψία

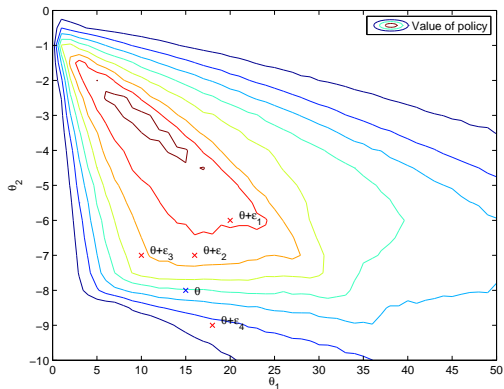
- $\alpha(t) = (1 - \delta)\delta^t, \gamma < \delta < 1$

- Στο M-step μεγιστοποιούμε και προκύπτει:

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \frac{\sum_{i=1}^m \frac{1}{|\xi_i|+1} \sum_{t=0}^{|\xi_i|} Q_{it} \boldsymbol{\varepsilon}_{it}}{\sum_{i=1}^m \frac{1}{|\xi_i|+1} \sum_{t=0}^{|\xi_i|} Q_{it}} \quad (26)$$

- $Q_{it} = \sum_{\tau=t}^{|\xi_i|} b(\tau) r_{i\tau}$
- $b(t) = (1 - \gamma/\delta)(\gamma/\delta)^t$
- Ελεγκτής: $u_t = (\boldsymbol{\theta} + \boldsymbol{\varepsilon}_t) \phi(\mathbf{x}_t)$, $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\boldsymbol{\varepsilon}_t; \mathbf{0}, \sigma^2 \mathbf{I}_d)$

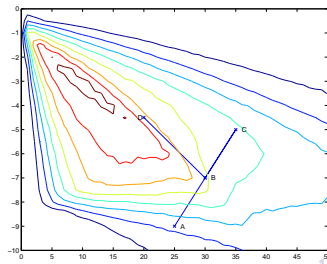
- MCEM



- Ο αλγόριθμος PoWER
- Σταθερό μήκος H
- Θέτουμε $b(t) = 1/(1 + H)$

Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού

- Χρήση των $p\%$ καλύτερων τροχιών
- Αναζήτηση επί γραμμής
 - Κλασική θεώρηση: αναζήτηση στην κατεύθυνση του διανύσματος κλίσης $\theta_{k+1} = \theta_k + \mu\alpha_k \mathbf{g}$
 - Θεωρούμε: $\mu\alpha_k \mathbf{g} = \theta_{k+1}^{MCEM} - \theta_k^{MCEM}$
 - $n \ll m$ τροχιές σταθερού μήκους H



- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- Ενισχυτική Μάθηση
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- Συμπεράσματα

- Διδιάστατο Συνθετικό Πρόβλημα
- x_1 =θέση, x_2 =ταχύτητα, u =επιτάχυνση
- Δυναμικό Μοντέλο

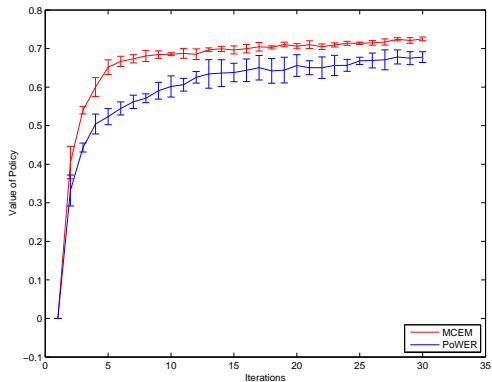
$$x_2(t+1) = x_2(t) - u(t) + \kappa, \quad (27)$$

$$x_1(t+1) = x_1(t) - 0.1x_2(t+1) + \kappa, \quad (28)$$

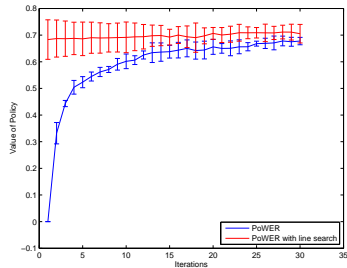
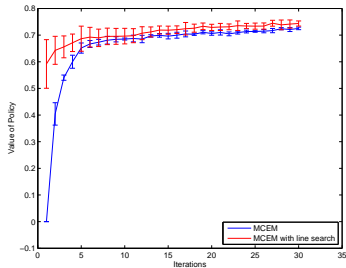
- $\sigma_\kappa = 0.003$
- $u_t = 1/[1 + \exp(-(\boldsymbol{\theta} + \boldsymbol{\varepsilon}_t)^T \mathbf{x})] - 0.5$

- $\mathbf{x}_0 = [1, 0] + noise, x_{goal} = [0, 0]$
- $r(t) = \exp(-0.5\|\mathbf{x}(t)\|_2/(0.03)^2)$
- MCEM
 - $\gamma = 0.95, \delta = 0.99, T_{max} = 100$
- PoWER
 - $H = 100$
- Αναζήτηση επί γραμμής
 - 10 τροχιές
- Χρήση καλύτερων τροχιών
 - 2 τροχιές

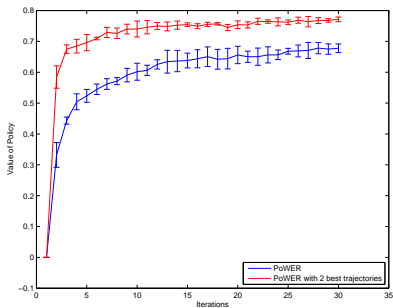
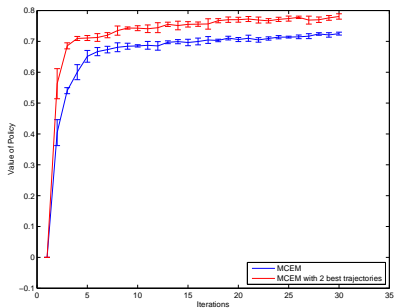
- MCEM vs PoWER



- Με και χωρίς αναζήτηση επί γραμμής

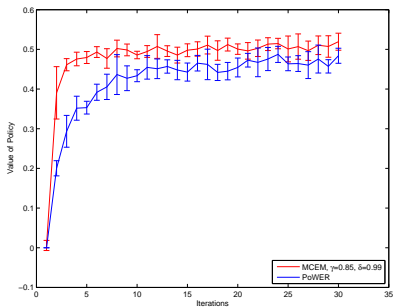
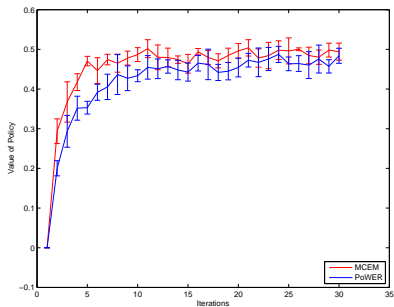


- Με και χωρίς τη χρήση των καλύτερων τροχιών

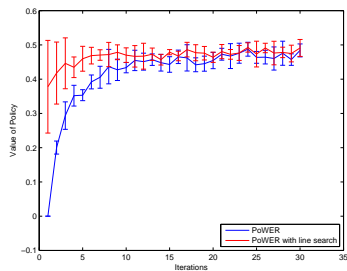
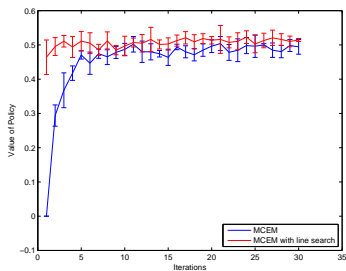


- $\sigma_{\kappa} = 0.01$

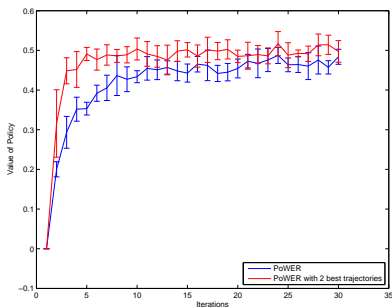
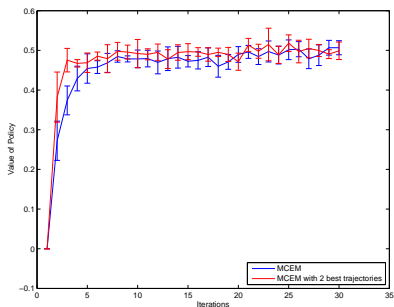
- MCEM vs PoWER



- $\sigma_{\kappa} = 0.01$
- Με και χωρίς αναζήτηση επί γραμμής

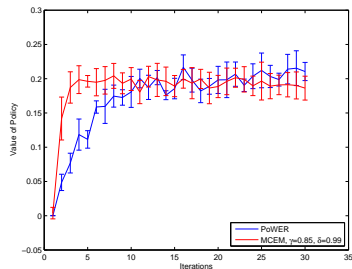
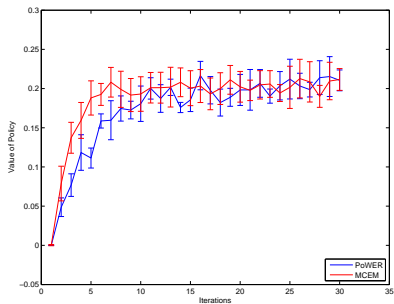


- $\sigma_{\kappa} = 0.01$
- Με και χωρίς τη χρήση των καλύτερων τροχιών

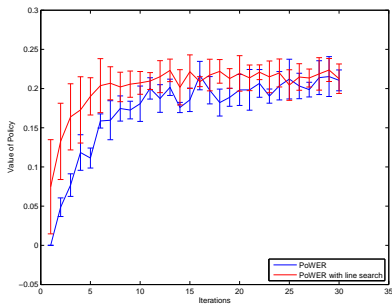
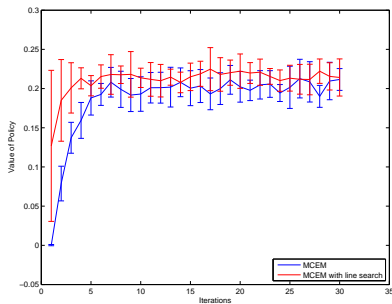


- $\sigma_{\kappa} = 0.03$

- MCEM vs PoWER

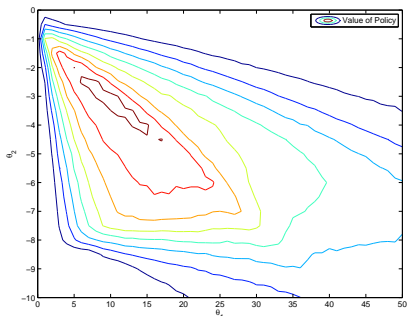


- $\sigma_{\kappa} = 0.03$
- Με και χωρίς αναζήτηση επί γραμμής

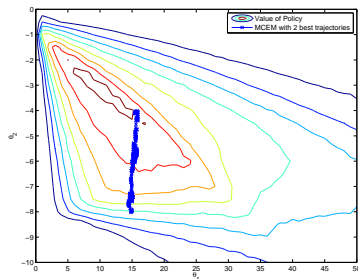
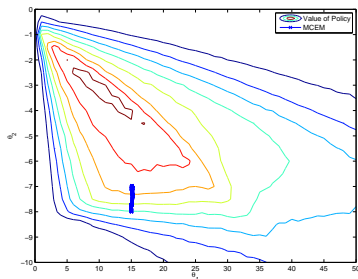


- $\sigma_{\kappa} = 0.02$
- Συνάρτηση Ανταμοιβής

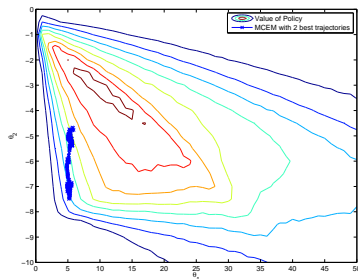
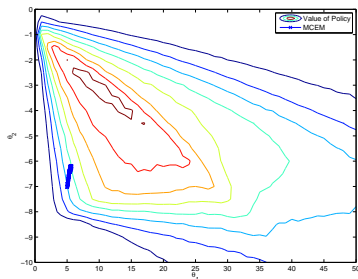
$$r(t) = \begin{cases} 1, & \text{αν } \|\mathbf{x}\|_2 < 0.1 \\ 0, & \text{αλλού} \end{cases} \quad (29)$$



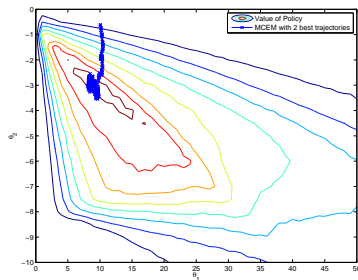
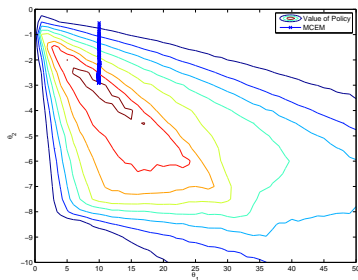
- 300 επαναλήψεις με και χωρίς τη χρήση των 2 καλύτερων τροχιών
- $\mathbf{x}_0 = [15, -8]$



- 300 επαναλήψεις με και χωρίς τη χρήση των 2 καλύτερων τροχιών
- $\mathbf{x}_0 = [5, -7]$



- 300 επαναλήψεις με και χωρίς τη χρήση των 2 καλύτερων τροχιών
- $\mathbf{x}_0 = [10, -0, 5]$



- Εξομοίωση αιώρησης ελικοπτέρου
- Μη γραμμικό, υψηλής διάστασης, στοχαστικό
- Χώρος καταστάσεων
 - Η θέση του ελικοπτέρου (x, y, z)
 - Η διεύθυνσή του (roll ϕ , pitch θ , yaw ω)
 - Η ταχύτητά του σε κάθε διεύθυνση $(\dot{x}, \dot{y}, \dot{z})$
- Χώρος ενεργειών
 - u_1 : Έλεγχος μπρος - πίσω κυκλικής κλίσης
 - u_2 : Έλεγχος αριστερά - δεξιά κυκλικής κλίσης
 - u_3 : Έλεγχος κλίσης ρότορα της κύριας έλικας
 - u_4 : Έλεγχος κλίσης ρότορα της ουραίας έλικας

- Διάνυσμα λάθους 9 διαστάσεων
- Ο ελεγκτής

$$u_1 = -w_1 x_{error} - w_2 \dot{x}_{error} - w_3 \theta_{error} + w_4 \quad (30)$$

$$u_2 = -w_5 y_{error} - w_6 \dot{y}_{error} - w_7 \phi_{error} + w_8 \quad (31)$$

$$u_3 = w_9 z_{error} - w_{10} \dot{z}_{error} + w_{11} \quad (32)$$

$$u_4 = -w_{12} \omega_{error} \quad (33)$$

- Η ανταμοιβή

$$r_1(t) = \exp(-x_{error}^2) + \exp(-\dot{x}_{error}^2) + \exp(-\theta_{error}^2) \quad (34)$$

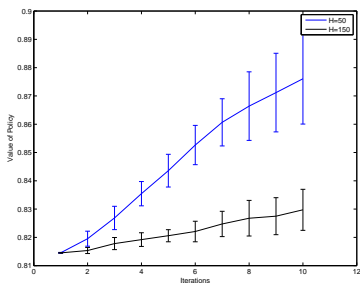
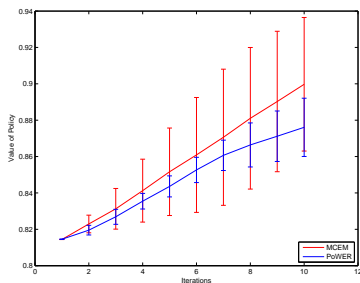
$$r_2(t) = \exp(-y_{error}^2) + \exp(-\dot{y}_{error}^2) + \exp(-\phi_{error}^2) \quad (35)$$

$$r_3(t) = \exp(-z_{error}^2) + \exp(-\dot{z}_{error}^2) \quad (36)$$

$$r_4(t) = \exp(-\omega_{error}^2) \quad (37)$$

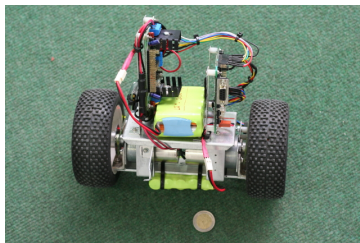
- 6000 βήματα \rightarrow 10 λεπτά πραγματικού χρόνου πτήσης
- 100 επεισόδια σε κάθε επανάληψη
- Αναζήτηση επί γραμμής \rightarrow 10 επεισόδια
- Χρήση καλύτερων τροχιών \rightarrow 5 καλύτερες
- PoWER
 - $H = 50$
- MCEM
 - $\gamma = 0.95$
 - $\delta = 0.99$
 - $T_{max} = 600$

- Ευαισθησία του PoWER στο μήκος επεισοδίου



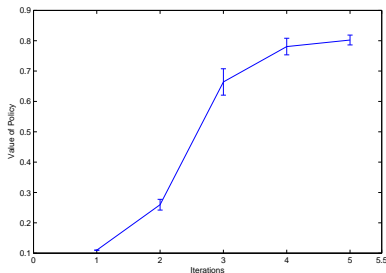
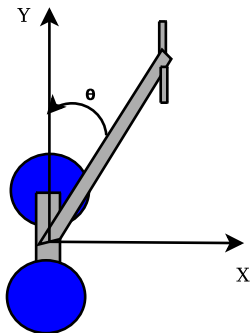
Πειράματα

- Δύο 12 Vdc, 152 RpM spur geared κινητήρες
- Έναν ooPIC microcontroller
- Δύο 64 παλμών ανά περιστροφή οδόμετρα
- Ένα επιταχυνσιόμετρο 2 αξόνων και ένα γυροσκόπιο ενός άξονα
- Δύο πηγές ρεύματος



Πειράματα

- Χώρος καταστάσεων: $\mathbf{x} = [x_1, x_2, x_3, x_4]$
- Ελεγκτής: $u_t = \theta^T \mathbf{x}$
- Ανταμοιβή: $r(t) = \exp(-x_1^2(t)) + \exp(-x_3^2(t))$



- Εισαγωγή
- Μαρκοβιανές Διεργασίες Απόφασης (ΜΔΑ)
- Ενισχυτική Μάθηση
- Ενισχυτική Μάθηση με Συναρτήσεις Αξιολόγησης
- Αλγόριθμοι Πολιτικής Κλίσης
- Ο αλγόριθμος EM
- Ενισχυτική Μάθηση μέσω Πιθανοτικού Συμπερασμού
- Πειράματα
- **Συμπεράσματα**

- MCEM
 - (+) Απλή υλοποίηση
 - (+) Ελάχιστη χειροκίνητη παραμετροποίηση
 - (+) Εύρωστη συμπεριφορά σε στοχαστικά περιβάλλοντα
 - (+) Χαμηλή υπολογιστική πολυπλοκότητα $O(mT)$
 - (-) Τοπικά βέλτιστα
 - (-) "Τεθλασμένα" πορεία στο χώρο παραμέτρων
 - (-) Μεγάλες απαιτήσεις σε μνήμη
- Εμπειρικοί κανόνες χρήσης των μεθόδων επιτάχυνσης
 - Χρήση πάντα των $p\%$ καλύτερων τροχιών
 - Μικρή στοχαστικότητα \longrightarrow Αναζήτηση επί γραμμής

- 1 Nikos Vlassis, Kontes Georgios and Savas Piperidis, "Reinforcement Learning of Robot Control via Probabilistic Inference", 1st Hellenic Robotics Conference (HEROC), February 23-24, Athens, 2009
- 2 Nikos Vlassis, Marc Toussaint, Georgios Kontes and Savas Piperidis, "Learning Model - free Robot Control by a Monte Carlo EM Algorithm", Autonomous Robots, Special issue on robot learning (accepted)

Ερωτήσεις???