

Covariance Matrix Adaptation for Direct Reinforcement Learning

Freek Stulp

Cognitive Robotics Group – ENSTA-ParisTech

FLOWERS Team – INRIA/Bordeaux

Olivier Sigaud

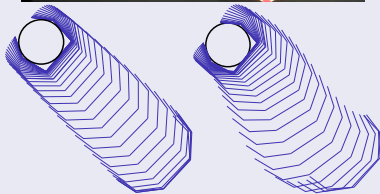
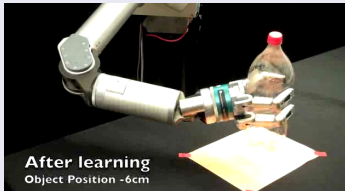
Institut des Systmes Intelligents et de Robotique – UPMC CNRS UMR 7222

JFPDA 23.05.2012

Motivation

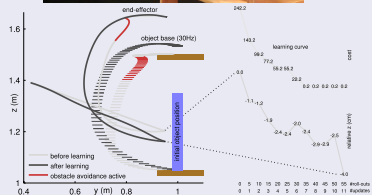
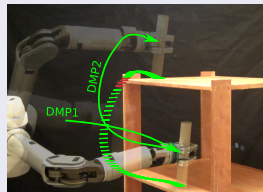
- PI² powerful algorithm for direct RL on robots

Grasping under Uncertainty



IROS'11, ICRA'11

Pick-and-Place Tasks



Humanoids'11

Motivation

- PI^2 powerful algorithm for direct RL on robots
 - Tuning exploration parameter in PI^2

Motivation

- PI^2 powerful algorithm for direct RL on robots
 - Tuning exploration parameter in PI^2
- Update rules for PI^2 and CEM/CMA-ES almost identical
 - But CEM and CMA-ES additionally tune exploration parameter automatically

Motivation

- PI^2 powerful algorithm for direct RL on robots
 - Tuning exploration parameter in PI^2
- Update rules for PI^2 and CEM/CMA-ES almost identical
 - But CEM and CMA-ES additionally tune exploration parameter automatically
- Goals
 - Analysis and comparison of PI^2 /CEM/CMA-ES
 - Novel algorithm PI^2 -CMA
 - Essentially PI^2 with automatic exploration parameter tuning

Outline

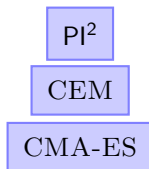
π^2

Outline

PI^2

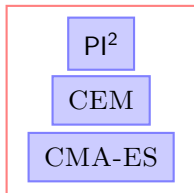
CEM

Outline



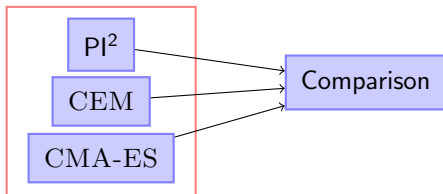
Outline

Reward-Weighted Averaging



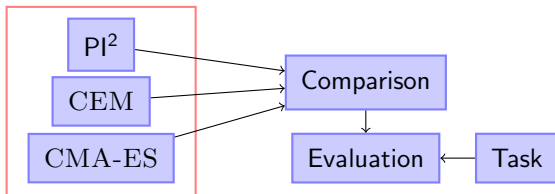
Outline

Reward-Weighted Averaging



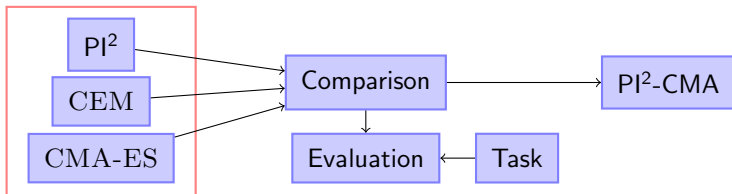
Outline

Reward-Weighted Averaging



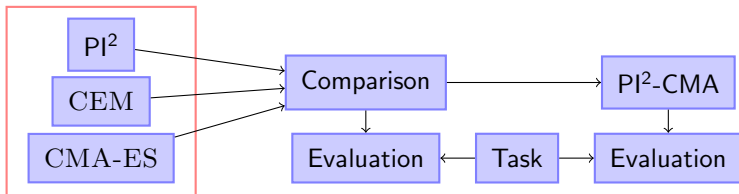
Outline

Reward-Weighted Averaging



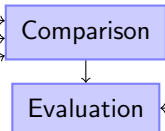
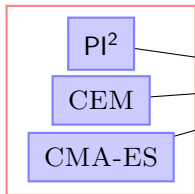
Outline

Reward-Weighted Averaging

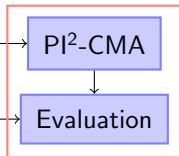


Outline

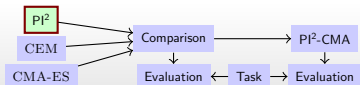
Reward-Weighted Averaging



Adaptive Exploration

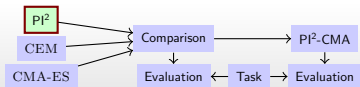


PI²- Derivation Outline



- From first principles of Stochastic Optimal Control
 - Start with Hamilton Jacobi Bellman equations
 - ① Log transformation + benign assumption
 - ⇒ Linear!
 - ② Transform PDE to path integral with Feynman-Kac theorem
 - ⇒ Roll-outs!
 - ③ Apply to parameterized policies
 - ⇒ Model free!
- ⇒ Iterative update rule for θ

PI²- Algorithm

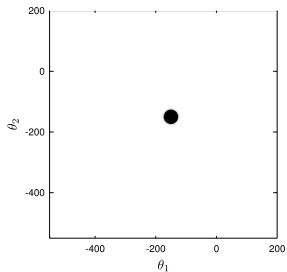
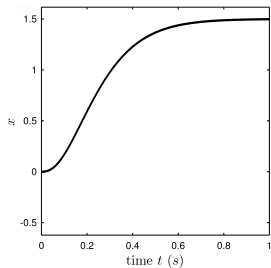


Words

Images

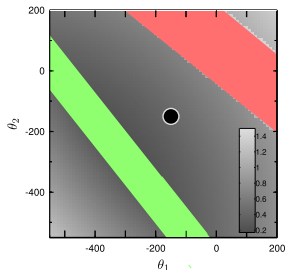
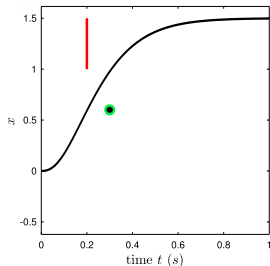
Formulae

- **Input:** DMP with initial parameters θ



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T \boldsymbol{\theta}$$

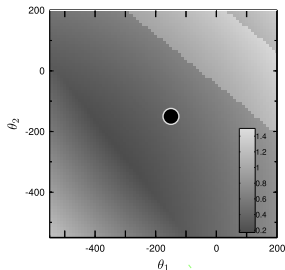
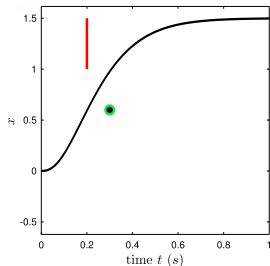
- **Input:** DMP with initial parameters θ , cost function J



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T \boldsymbol{\theta}$$

$$J(\boldsymbol{\tau}_i) = \phi_{t_N} + \int_{t_i}^{t_N} (q_t + \frac{1}{2} \boldsymbol{\theta}_t^T \mathbf{R} \boldsymbol{\theta}_t) dt$$

- **Input:** DMP with initial parameters θ , cost function J



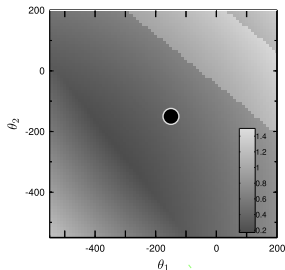
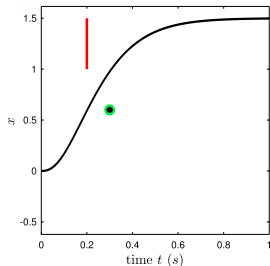
$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T \boldsymbol{\theta}$$

$$J(\boldsymbol{\tau}_i) = \phi_{t_N} + \int_{t_i}^{t_N} (q_t + \frac{1}{2} \boldsymbol{\theta}_t^T \mathbf{R} \boldsymbol{\theta}_t) dt$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

Explore

Update



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T \boldsymbol{\theta}$$

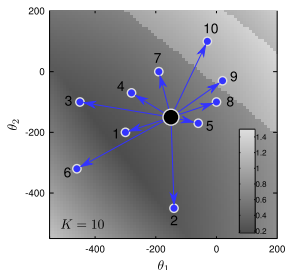
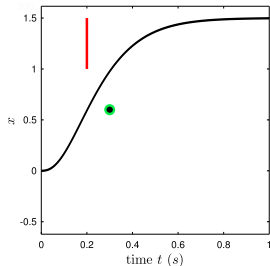
$$J(\boldsymbol{\tau}_i) = \phi_{t_N} + \int_{t_i}^{t_N} (q_t + \frac{1}{2} \boldsymbol{\theta}_t^T \mathbf{R} \boldsymbol{\theta}_t) dt$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

Explore

sample exploration vectors

Update



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(\mathbf{g} - x_t) - \dot{x}_t) + \mathbf{g}_t^T(\boldsymbol{\theta} + \boldsymbol{\epsilon}_{t,k})$$

$$J(\boldsymbol{\tau}_i) = \phi_{\tau_i N} + \int_{t_i}^{t_i + N} (q_t + \frac{1}{2} \boldsymbol{\theta}_t^T \mathbf{R} \boldsymbol{\theta}_t) dt$$

$$\boldsymbol{\epsilon}_{t,k} \sim \mathcal{N}(0, \boldsymbol{\Sigma}^\epsilon)$$

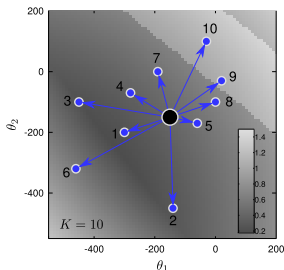
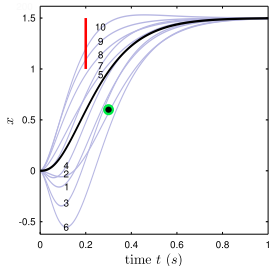
$$\boldsymbol{\theta}_k = \boldsymbol{\theta} + \boldsymbol{\epsilon}_k$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

Explore

sample exploration vectors
execute DMP

Update



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T(\theta + \epsilon_{t,k})$$

$$J(\tau_i) = \phi_{t_N} + \int_{t_i}^{t_N} (q_t + \frac{1}{2} \theta_t^T \mathbf{R} \theta_t) dt$$

$$\epsilon_{t,k} \sim \mathcal{N}(0, \Sigma^\epsilon)$$

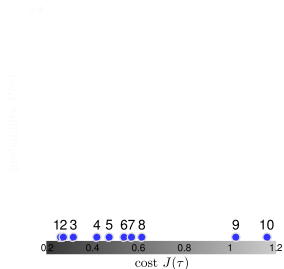
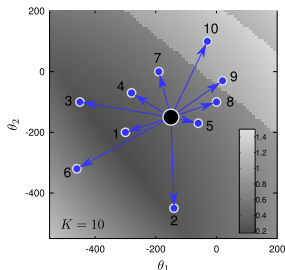
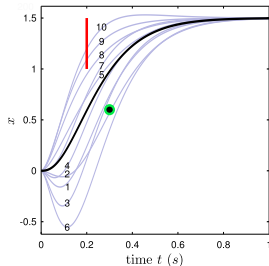
$$\theta_k = \theta + \epsilon_k$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

Explore

sample exploration vectors
 execute DMP
 determine cost

Update



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + g_t^T(\theta + \epsilon_{t,k})$$

$$J(\tau_i) = \phi_{\tau_i N} + \int_{t_i}^{\tau_i N} (q_t + \frac{1}{2} \theta_t^T R \theta_t) dt$$

$$\epsilon_{t,k} \sim \mathcal{N}(0, \Sigma^\epsilon)$$

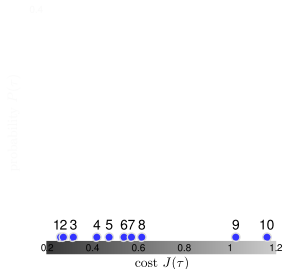
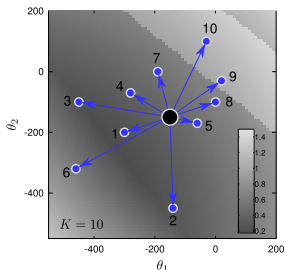
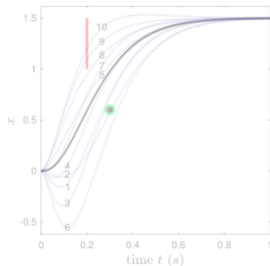
$$\theta_k = \theta + \epsilon_k$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

Explore

sample exploration vectors
 execute DMP
 determine cost

Update



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + g_t^T(\theta + \epsilon_{t,k})$$

$$J(\tau_i) = \phi_{t_N} + \int_{t_i}^{t_N} (q_t + \frac{1}{2} \theta_t^T R \theta_t) dt$$

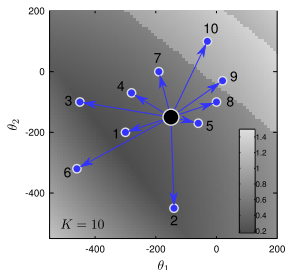
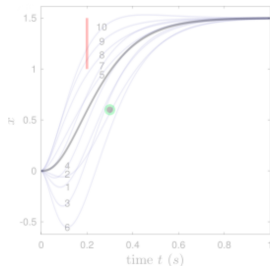
$$\epsilon_{t,k} \sim \mathcal{N}(0, \Sigma^\epsilon)$$

$$\theta_k = \theta + \epsilon_k$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

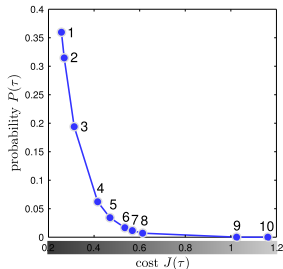
Explore

sample exploration vectors
execute DMP
determine cost



Update

compute prob. from cost



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T(\theta + \epsilon_{t,k})$$

$$J(\tau_i) = \phi_{t_N} + \int_{t_j}^{t_N} (q_t + \frac{1}{2} \theta_t^T R \theta_t) dt$$

$$\epsilon_{t,k} \sim \mathcal{N}(0, \Sigma^\epsilon)$$

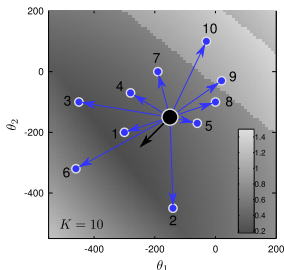
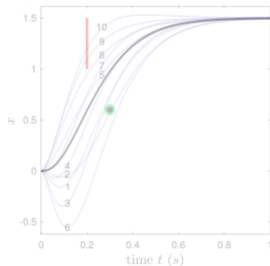
$$\theta_k = \theta + \epsilon_k$$

$$P(\tau_{i,k}) = \frac{e^{-\frac{1}{\lambda} J(\tau_{i,k})}}{\sum_{k=1}^K [e^{-\frac{1}{\lambda} J(\tau_{i,k})}]}$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

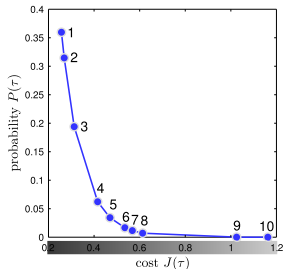
Explore

sample exploration vectors
execute DMP
determine cost



Update

compute prob. from cost
prob.-weighted averaging



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T(\theta + \epsilon_{t,k})$$

$$J(\tau_i) = \phi_{t_N} + \int_{t_i}^{t_N} (q_t + \frac{1}{2} \theta_t^T R \theta_t) dt$$

$$\epsilon_{t,k} \sim \mathcal{N}(0, \Sigma^\epsilon)$$

$$\theta_k = \theta + \epsilon_k$$

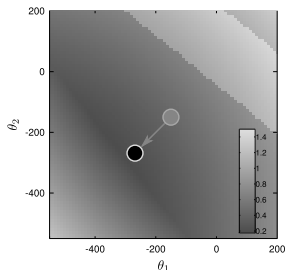
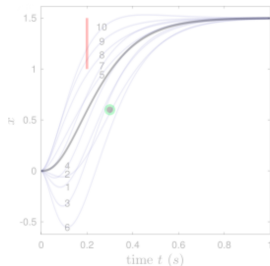
$$P(\tau_{i,k}) = \frac{e^{-\frac{1}{\lambda} J(\tau_{i,k})}}{\sum_{k=1}^K [e^{-\frac{1}{\lambda} J(\tau_{i,k})}]}$$

$$\delta \theta_{t_i} = \sum_{k=1}^K [P(\tau_{i,k}) \mathbf{M}_{t_i,k} \epsilon_{t_i,k}]$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

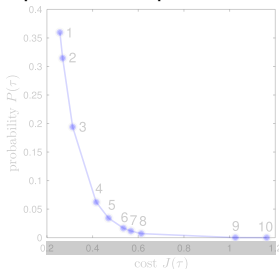
Explore

sample exploration vectors
execute DMP
determine cost



Update

compute prob. from cost
prob.-weighted averaging
parameter update



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T(\theta + \epsilon_{t,k})$$

$$J(\tau_i) = \phi_{t_N} + \int_{t_i}^{t_N} (q_t + \frac{1}{2} \theta_t^T \mathbf{R} \theta_t) dt$$

$$\epsilon_{t,k} \sim \mathcal{N}(0, \Sigma^\epsilon)$$

$$\theta_k = \theta + \epsilon_k$$

$$\theta \leftarrow \theta + \delta \theta$$

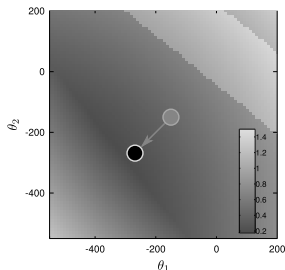
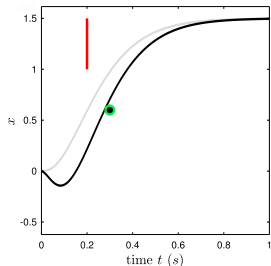
$$P(\tau_{i,k}) = \frac{e^{-\frac{1}{\lambda} J(\tau_{i,k})}}{\sum_{k=1}^K [e^{-\frac{1}{\lambda} J(\tau_{i,k})}]}$$

$$\delta \theta_{t_i} = \sum_{k=1}^K [P(\tau_{i,k}) \mathbf{M}_{t_i,k} \epsilon_{t_i,k}]$$

- **Input:** DMP with initial parameters θ , cost function J
- While (cost not converged)

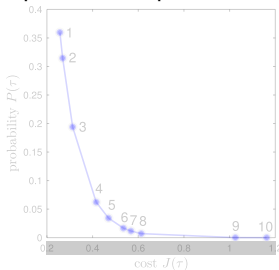
Explore

sample exploration vectors
execute DMP
determine cost



Update

compute prob. from cost
prob.-weighted averaging
parameter update



$$\frac{1}{\tau} \ddot{x}_t = \alpha(\beta(g - x_t) - \dot{x}_t) + \mathbf{g}_t^T(\theta + \epsilon_{t,k})$$

$$J(\tau_i) = \phi_{t_N} + \int_{t_i}^{t_N} (q_t + \frac{1}{2} \theta_t^T \mathbf{R} \theta_t) dt$$

$$\epsilon_{t,k} \sim \mathcal{N}(0, \Sigma^\epsilon)$$

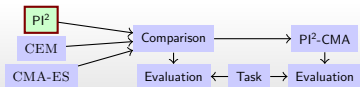
$$\theta_k = \theta + \epsilon_k$$

$$\theta \leftarrow \theta + \delta \theta$$

$$P(\tau_{i,k}) = \frac{e^{-\frac{1}{\lambda} J(\tau_{i,k})}}{\sum_{k=1}^K [e^{-\frac{1}{\lambda} J(\tau_{i,k})}]}$$

$$\delta \theta_{t_i} = \sum_{k=1}^K [P(\tau_{i,k}) \mathbf{M}_{t_i,k} \epsilon_{t_i,k}]$$

PI²- Algorithm



- Advantages

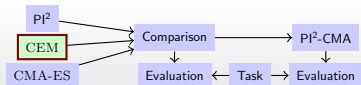
- No gradient \Rightarrow Deals with discontinuous noisy cost functions
- Update δ within convex hull of $\epsilon_{k=1\dots K} \Rightarrow$ Safe update rule
- Arbitrary cost functions
- Model-free
- Fast convergence
- Only one open parameter: magnitude of exploration ($\epsilon_{t,k} \sim \mathcal{N}(0, \Sigma)$)

- Disadvantage

- No global convergence guarantees...
- Robotics: This is where imitation comes in! $\pi(\theta^{imit}) \approx \pi(\theta^*)$

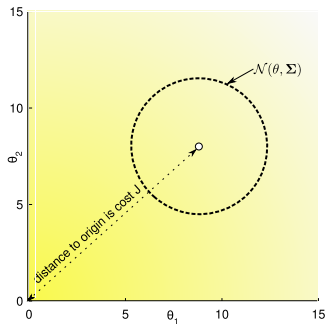
- Applied to very high-dimensional, complex tasks...

Cross-Entropy Method (CEM)



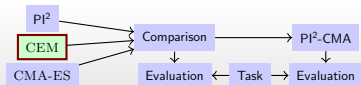
$$\mathcal{N}(\theta, \Sigma)$$

Probability



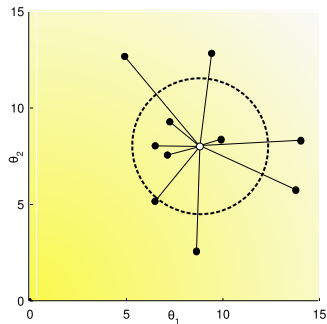
Cost

Cross-Entropy Method (CEM)



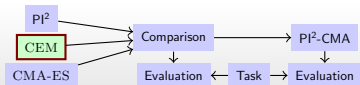
$$\theta_{k=1\dots K} \sim \mathcal{N}(\theta, \Sigma)$$

Probability



Cost

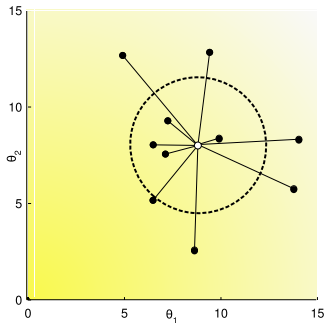
Cross-Entropy Method (CEM)



$$\theta_{k=1\dots K} \sim \mathcal{N}(\theta, \Sigma)$$

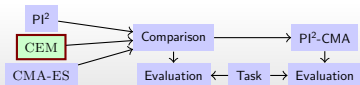
$$\forall k J_k = J(\theta_k)$$

Probability



Cost

Cross-Entropy Method (CEM)

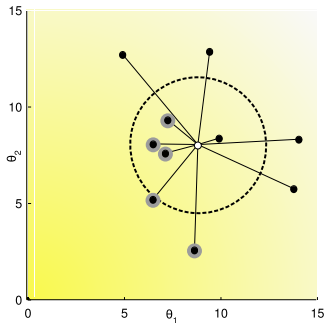


$$\theta_{k=1\dots K} \sim \mathcal{N}(\theta, \Sigma)$$

$$\forall k J_k = J(\theta_k)$$

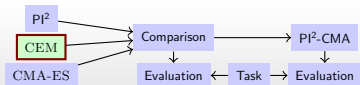
$$\theta_{k=1\dots K} \leftarrow \text{sort } \theta_{k=1\dots K} \text{ w.r.t } J_{k=1\dots K}$$

Probability



Cost

Cross-Entropy Method (CEM)



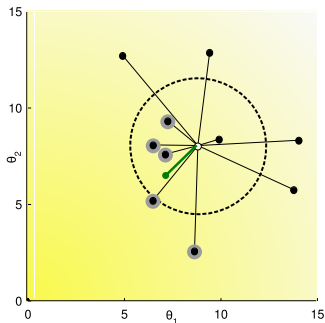
$$\theta_{k=1\dots K} \sim \mathcal{N}(\theta, \Sigma)$$

$$\forall k J_k = J(\theta_k)$$

$$\theta_{k=1\dots K} \leftarrow \text{sort } \theta_{k=1\dots K} \text{ w.r.t } J_{k=1\dots K}$$

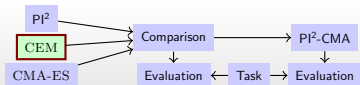
$$\theta^{new} = \sum_{k=1}^{K_e} \frac{1}{K_e} \theta_k$$

Probability



Cost

Cross-Entropy Method (CEM)



$$\theta_{k=1\dots K} \sim \mathcal{N}(\theta, \Sigma)$$

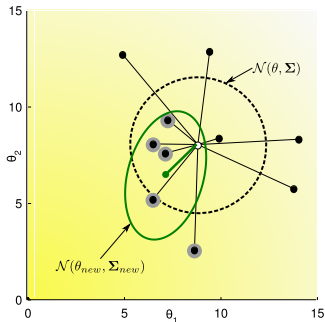
$$\forall k J_k = J(\theta_k)$$

$$\theta_{k=1\dots K} \leftarrow \text{sort } \theta_{k=1\dots K} \text{ w.r.t } J_{k=1\dots K}$$

$$\theta^{new} = \sum_{k=1}^{K_e} \frac{1}{K_e} \theta_k$$

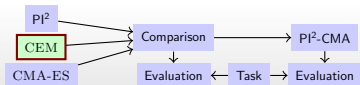
$$\Sigma^{new} = \sum_{k=1}^{K_e} \frac{1}{K_e} (\theta_k - \theta)(\theta_k - \theta)^T$$

Probability



Cost

Cross-Entropy Method (CEM)



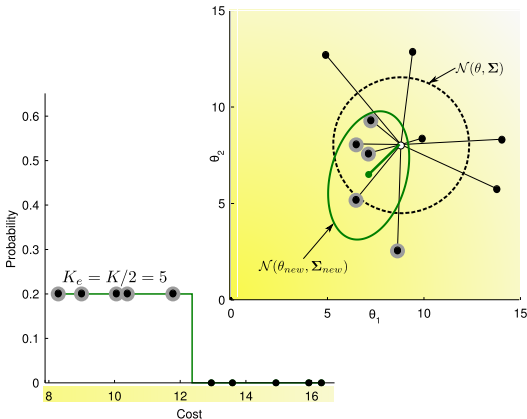
$$\theta_{k=1\dots K} \sim \mathcal{N}(\theta, \Sigma)$$

$$\forall k J_k = J(\theta_k)$$

$$\theta_{k=1\dots K} \leftarrow \text{sort } \theta_{k=1\dots K} \text{ w.r.t } J_{k=1\dots K}$$

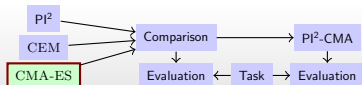
$$\theta^{new} = \sum_{k=1}^{K_e} \frac{1}{K_e} \theta_k$$

$$\Sigma^{new} = \sum_{k=1}^{K_e} \frac{1}{K_e} (\theta_k - \theta)(\theta_k - \theta)^T$$



- This algorithm can be interpreted as performing reward-weighted averaging

CMA-ES



- Covariance Matrix Adaptation - Evolutionary Strategy
- Like CEM, but
 - Different mapping from cost to probability
 - More sophisticated method for updating covariance matrix:

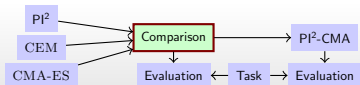
$$p_{\sigma} \leftarrow (1 - c_{\sigma}) p_{\sigma} + \sqrt{c_{\sigma}(2 - c_{\sigma})} \mu_P \Sigma^{-1} \frac{\theta^{new} - \theta}{\sigma} \quad (1)$$

$$\sigma_{new} = \sigma \times \exp \left(\frac{c_{\sigma}}{d_{\sigma}} \left(\frac{\|p_{\sigma}\|}{E\|\mathcal{N}(0, I)\|} - 1 \right) \right) \quad (2)$$

$$p_{\Sigma} \leftarrow (1 - c_{\Sigma}) p_{\Sigma} + h_{\sigma} \sqrt{c_{\Sigma}(2 - c_{\Sigma})} \mu_P \frac{\theta^{new} - \theta}{\sigma} \quad (3)$$

$$\Sigma^{new} = (1 - c_1 - c_{\mu}) \Sigma + c_1 (p_{\Sigma} p_{\Sigma}^T + \delta (h_{\sigma}) \Sigma) + c_{\mu} \sum_{k=1}^{K_e} P_k (\theta_k - \theta) (\theta_k - \theta)^T \quad (4)$$

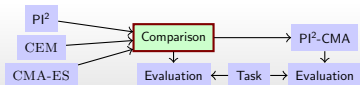
PI²/CEM/CMA-ES- Similarities



- PI²/CEM/CMA-ES are all based on
 - Exploration: sample from a Gaussian
 - Parameter update: Reward-Weighted Averaging

$$\theta_k \sim \mathcal{N}(\theta, \Sigma)$$
$$\theta^{new} = \sum_{k=1}^K P_k \theta_k$$

PI²/CEM/CMA-ES- Similarities

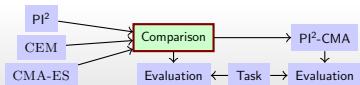


- PI²/CEM/CMA-ES are all based on
 - Exploration: sample from a Gaussian
 - Parameter update: Reward-Weighted Averaging

$$\theta_k \sim \mathcal{N}(\theta, \Sigma)$$
$$\theta^{new} = \sum_{k=1}^K P_k \theta_k$$

Similarity striking, as algorithms derived from very different principles!

PI²/CEM/CMA-ES- Similarities



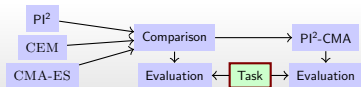
- PI²/CEM/CMA-ES are all based on
 - Exploration: sample from a Gaussian
 - Parameter update: Reward-Weighted Averaging

$$\theta_k \sim \mathcal{N}(\theta, \Sigma)$$
$$\theta^{new} = \sum_{k=1}^K P_k \theta_k$$

Similarity striking, as algorithms derived from very different principles!

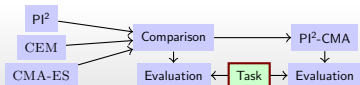
- CEM is a special case of CMA-ES: proof in paper.
(maybe this was already known?)

PI²/CEM/CMA-ES- Differences

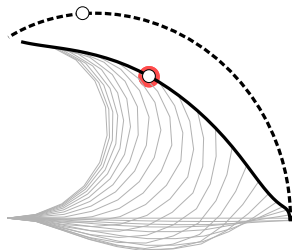
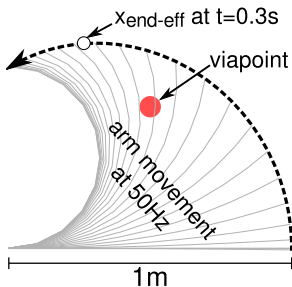


- Also some differences
- Evaluated on following task:

PI²/CEM/CMA-ES- Differences

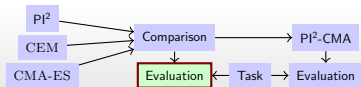


- Also some differences
- Evaluated on following task:



$$J(\tau_t) = \delta(t - 0.3) \cdot ((x_t - 0.5)^2 + (y_t - 0.5)^2) + \frac{\sum_{d=1}^D (D + 1 - d)(\ddot{a}_t)^2}{\sum_{d=1}^D (D + 1 - d)} \quad (5)$$

PI²/CEM/CMA-ES- Differences



Exploration Noise

PI²

$$\theta_{k,t} \sim \mathcal{N}(\theta, \Sigma)$$

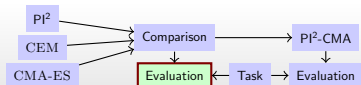
CEM

$$\theta_k \sim \mathcal{N}(\theta, \Sigma)$$

CMA-ES

$$\theta_k \sim \mathcal{N}(\theta, \Sigma)$$

PI²/CEM/CMA-ES- Differences



Exploration Noise

PI²

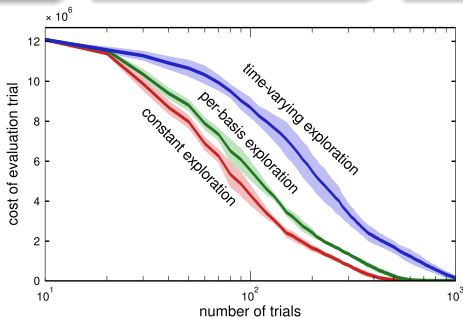
$$\theta_{k,t} \sim \mathcal{N}(\theta, \Sigma)$$

CEM

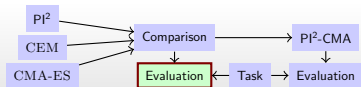
$$\theta_k \sim \mathcal{N}(\theta, \Sigma)$$

CMA-ES

$$\theta_k \sim \mathcal{N}(\theta, \Sigma)$$

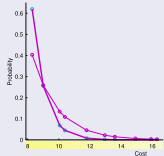


PI²/CEM/CMA-ES- Differences

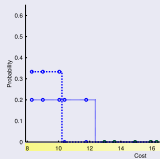


Eliteness

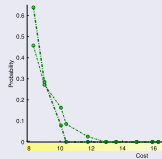
PI²



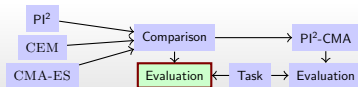
CEM



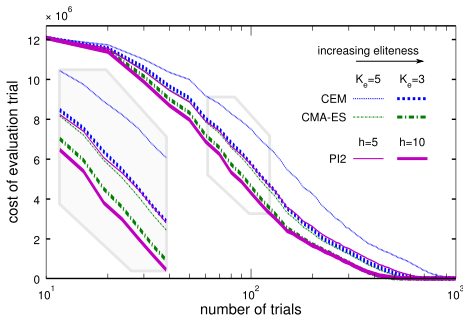
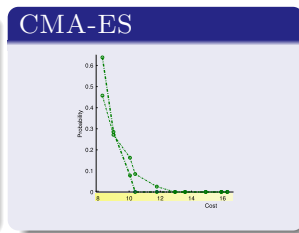
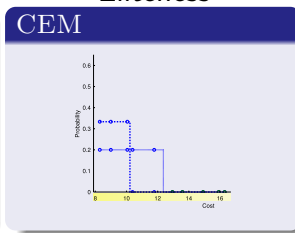
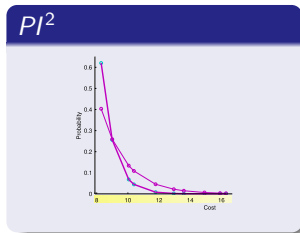
CMA-ES



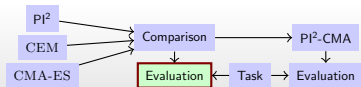
PI²/CEM/CMA-ES- Differences



Eliteness



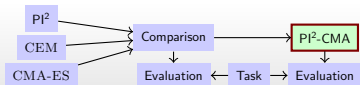
PI²/CEM/CMA-ES- Differences



Covariance Matrix Updating

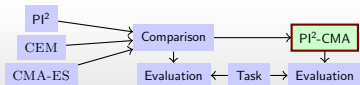
| PI ² | CEM | CMA-ES |
|-----------------|--------------|---------------------|
| No | Yes (simple) | Yes (sophisticated) |

PI²-CMA

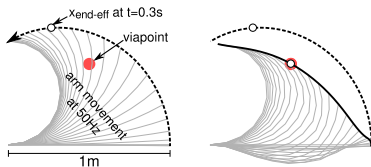


- Suggests a new algorithm: PI²-CMA
 - Constant exploration noise
 - Eliteness measure from PI²
 - Covariance matrix updating from CEM/CMA-ES

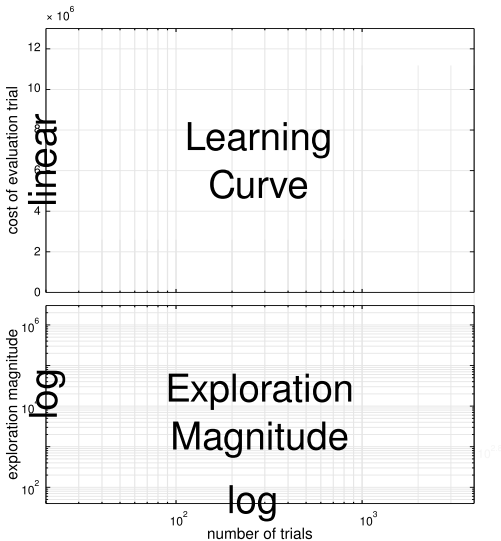
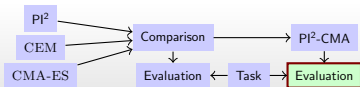
PI²-CMA



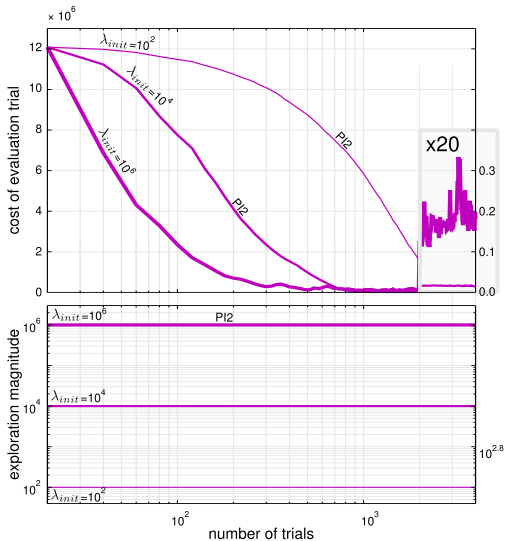
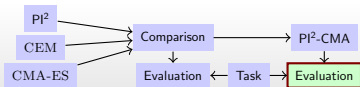
- Suggests a new algorithm: PI²-CMA
 - Constant exploration noise
 - Eliteness measure from PI²
 - Covariance matrix updating from CEM/CMA-ES
- Basically PI², but with adaptive exploration
 - In PI², exploration magnitude must be tuned by hand
 - Next evaluation demonstrates advantages of PI²-CMA



PI²-CMA- Adaptive Exploration

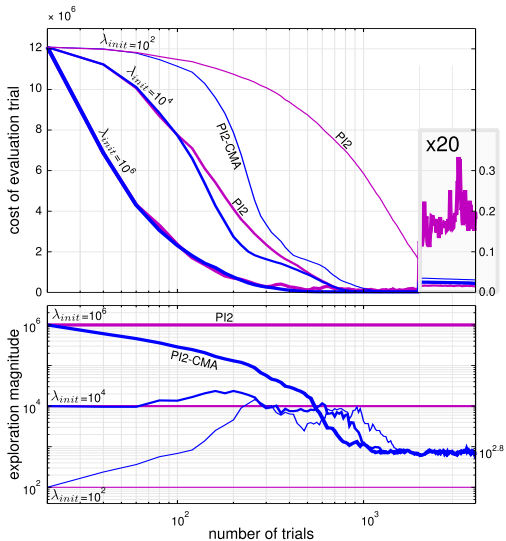
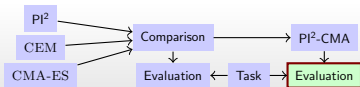


PI²-CMA- Adaptive Exploration

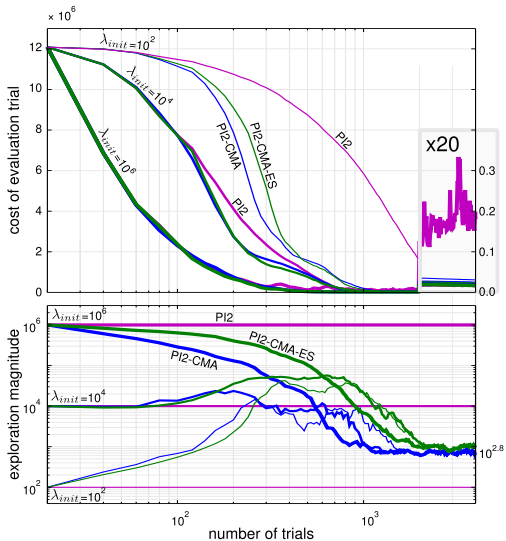
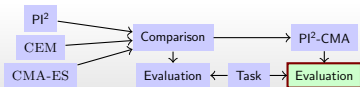


⇒ Exploration magnitude influences convergence speed and exploitation

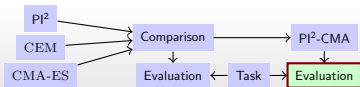
PI²-CMA- Adaptive Exploration



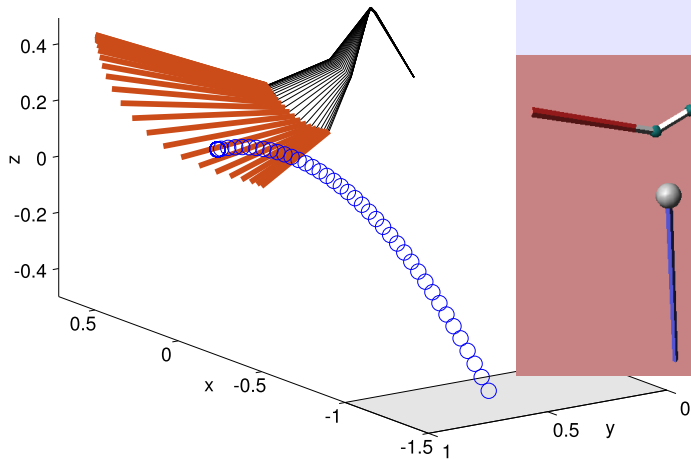
PI²-CMA- Adaptive Exploration



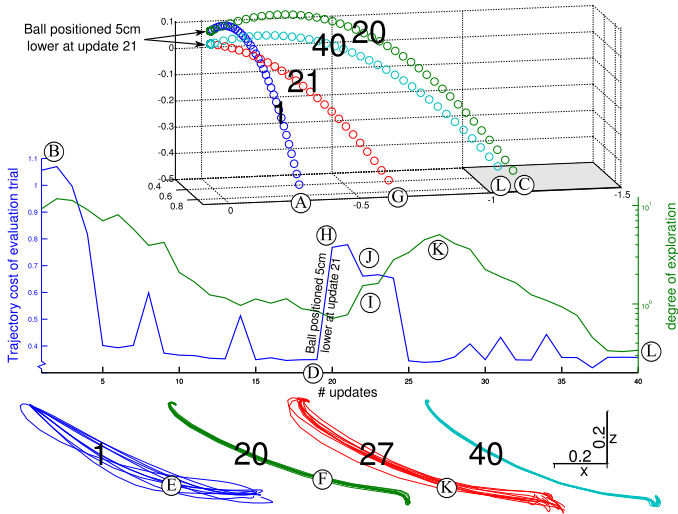
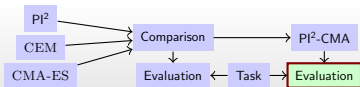
PI²-CMA- Adaptive Exploration



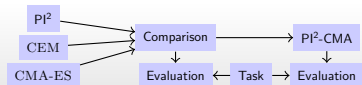
More recent results (not in paper)



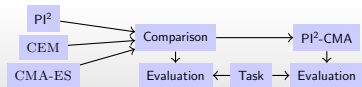
PI²-CMA- Adaptive Exploration



Summary

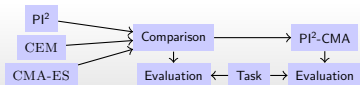


- PI²/CEM/CMA-ES have identical update rules: reward-weighted averaging
- Apply covariance matrix adaptation (as in CEM/CMA-ES) to PI²
- Novel algorithm PI²-CMA
 - With adaptive exploration (other algorithmic parameters trivial to tune)
- Future work
 - Further analysis and theoretical validation
 - Evaluation on real robots



Thank you for your attention!
Questions?

Bibliography



Stulp, F. and Schaal, S. (2011).

Hierarchical reinforcement learning with motion primitives.

In 11th IEEE-RAS International Conference on Humanoid Robots.



Stulp, F., Theodorou, E., Buchli, J., and Schaal, S. (2011a).

Learning to grasp under uncertainty.

In Proceedings of the IEEE International Conference on Robotics and Automation (ICRA).



Stulp, F., Theodorou, E., Kalakrishnan, M., Pastor, P., Righetti, L., and Schaal, S. (2011b).

Learning motion primitive goals for robust manipulation.

In International Conference on Intelligent Robots and Systems (IROS).