From Policy Gradient to Actor-Critic methods The Policy Search problem

Olivier Sigaud

Sorbonne Université http://people.isir.upmc.fr/sigaud



The five routes to deep RL



Five different ways to come to Deep RL

2 / 14

DES SYSTÈMES INTELLIDENTS ET DE ROBOTIQ

Э

The Policy Search route



- The favorite route of roboticists
- Central question: difference between PG with baseline and Actor-Critic

Marc P. Deisenroth, Gerhard Neumann, Jan Peters, et al. A survey on policy search for robotics. Foundations and Trends® in Robotics, 2(1-2):1-142, 2013



Policy search: general intuitions



Example: a (cheap) tennis ball collector



- A robot without a ball sensor
- Travels on a tennis court based on a parametrized controller
- Performance: number of balls collected in a given time
- Just depends on robot trajectories and ball positions



Influence of policy parameters



- Controller parameters: proba of turn per time step, travelling speed
- How do the parameters influence the performance?
- Policy search: find the optimal policy parameters



Two sources of stochasticity



- From the environment: position of the balls
- From the policy, if it is stochastic
- \blacktriangleright The performance can vary a lot \rightarrow need to repeat
- Tuning parameters can be hard



Policy search: formalization



The policy search problem: formalization



τ_i is a robot trajectory
R(τ_i) is the corresponding return
π_θ is the parametrized policy of the robot

- ▶ We want to optimize $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}}[R(\tau)]$, the global utility function
- We tune policy parameters θ , thus the goal is to find

$$\boldsymbol{\theta}^* = \operatorname*{argmax}_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{\tau} P(\tau | \boldsymbol{\theta}) R(\tau)$$
(1)

イロト イヨト イヨト イヨト



Deisenroth, M. P., Neumann, G., Peters, J., et al. (2013) A survey on policy search for robotics. Foundations and Trends® Robotics, 2(1-2):1-142



Direct Policy Search is black box optimization



- $J(\boldsymbol{\theta})$ is the performance over policy parameters
- Choose a θ
- Generate trajectories τ_{θ}
- Get the return $J(\boldsymbol{\theta})$ of these trajectories
- Look for a better θ , repeat

DPS uses $(\theta, J(\theta))$ pairs and directly looks for θ with the highest $J(\theta)$



(Truly) Random Search



- Select θ_i randomly
- ► Evaluate $J(\theta_i)$
- If $J(\boldsymbol{\theta}_i)$ is the best so far, keep $\boldsymbol{\theta}_i$

・ロト ・回ト ・ヨト ・ヨト

- ▶ Loop until $J(\boldsymbol{\theta}_i) > target$
- Of course, this is not efficient if the space of θ is large
- General "blind" algorithm, no assumption on $J(\boldsymbol{\theta})$
- We can do better if $J(\theta)$ shows some local regularity

Sigaud, O. & Stulp, F. (2019) Policy search in continuous action domains: an overview. Neural Networks, 113:28-40

Direct policy search

Locality assumption: The function is locally smooth, good solutions are close to each other



Variation - selection

- Variation selection: Perform well chosen variations, evaluate them
- Variations generally controlled using a multivariate Gaussian



Gradient ascent



- Gradient ascent: Following the gradient from analytical knowledge
- ▶ Issue: in general, the function $J(\theta)$ is unknown
- How can we apply gradient ascent without knowing the function?
- The answer is the Policy Gradient Theorem



From Policy Gradient to Actor-Critic methods Policy search: formalization

Any question?



Send mail to: Olivier.Sigaud@upmc.fr





Deisenroth, M. P., Neumann, G., Peters, J., et al. (2013).

A survey on policy search for robotics. Foundations and $Trends(\mathbb{R})$ in Robotics, 2(1-2):1-142.



Sigaud, O. and Stulp, F. (2019).

Policy search in continuous action domains: an overview. Neural Networks, 113:28-40.

