

Natural Actor Critic

Do you have a subtitle?

If so, write it here

First Author · Second Author

Received: date / Accepted: date

F. Author
first address
Tel.: +123-45-678910
Fax: +123-45-678910
E-mail: fauthor@example.com

S. Author
second address

1 Paper

Natural Actor Critic:

- Main Version [6].
- 2nd Version: Natural Actor-Critic in Neurocomputing [4].
- 3rd version: RL of motor skills with policy gradients in NN [5].

Must read paper to understand basics by Jan:

- Policy Evaluation with TD [3].

Recommended by Jan:

- Incremental NAC algorithms [1].
- Jan said that a paper from C. Dann is very important. Did he mean Policy Evaluation with TD by Dann or did he mean a second paper?

Research:

- Comparison of four natural gradient algorithms (co-author Sutton) [2].

2 Meetings & Notes

Meetings:

- 12.12.18: Notes from Jan can be found in “.\\Notes Jan 12.12.18”

3 Introduction

- Steepest ascent direction of performance object with respect to any metric $M(\theta)$: $M(\theta)^{-1} \nabla_{\theta} J(\mu_{\theta})$
- The natural gradient is the steepest ascent direction with respect to the Fisher information metric $M_{\pi}(\theta) = E_{s \sim \rho^{\pi}, a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s)^T \nabla_{\theta} \log \pi_{\theta}(a|s)]$
- For deterministic policies: $M_{\mu}(\theta) = E_{s \sim \rho^{\mu}} [\nabla_{\theta} \mu_{\theta}(s) \nabla_{\theta} \mu_{\theta}(s)^T]$
 - Limiting case of the Fisher information metric: policy variance reduced to zero
- Combining DPG theorem with compatible function approximation gives $\nabla_{\theta} J(\mu_{\theta}) = E_{s \sim \rho^{\mu}} [\nabla_{\theta} \mu_{\theta}(s) \nabla_{\theta} \mu_{\theta}(s)^T] w$ so steepest ascent direction reduces to $M_{\mu}(\theta)^{-1} \nabla_{\theta} J_{\beta}(\mu_{\theta}) = w$

4 Ideas

4.1 Fitted NAC

$$\pi(a|s) = p(a|s, \Theta) = \mathcal{N}(a|\mu = NN_{\Theta}(s), \sigma) \quad (1)$$

$$f(s, a) = \log p(a|s, \Theta)^T w \quad (2)$$

$$f_V(s) = NN_V(s) \quad (3)$$

$$\text{"fitted"} \quad (4)$$

$$\min_{V_t, W_t} (r(s, a) + \gamma f_{V_{t+1}}(s') - f_{V_t}(s) + f_{W_t}(s, a))^2 \quad (5)$$

4.2 Ideas of TRPO

4.3 Other Extentions

- stochastic
- minibatches
- importance sampling:

$$V_{\Theta} J = \sum \mu(s) \pi'(a|s) \nabla \log \pi'(a|s) Q^{\pi'}(a|s) \quad (6)$$

$$\approx \frac{1}{N} \sum \nabla \log \pi'(a|s) Q(s, a) = g(\Theta) \quad (7)$$

5 Episodic NAC

Important to understand beforehand: In episodic NAC, our system of equations has one equation per trajectory and not one equation per action as in the normal NAC algorithm.

First we start by adding together the advantage function across an episode e where we made N steps.

$$A(s, a) = r(s, a) + \gamma V(s') - V(s) \quad (8)$$

$$\gamma A(s', a') = \gamma r(s', a') + \gamma^2 V(s'') - \gamma V(s') \quad (9)$$

$$A(s, a) + \gamma A(s', a') = r(s, a) + \gamma r(s', a') + \gamma^2 V(s'') - V(s) \quad (10)$$

$$\sum_{i=0}^N \gamma^i A(s_i, a_i) = \sum_{i=0}^N \gamma^i r(s_i, a_i) + \gamma^N V(S_{N+1}) - V(S_0) \quad (11)$$

If we assume $\gamma \neq 1$, we can remove the term $\gamma^N V(S_{N+1})$, because in the limit the term becomes zero ($\gamma^N \rightarrow 0$). Additionally, if we assume that we always

start in the same start S_0 , we can write $V(S_0)$ as our cost function J because it will exactly sum up the expected Reward/cost of our problem.

$$\Rightarrow \sum_{i=0}^N \gamma^i A(s_i, a_i) = \sum_{i=0}^N \gamma^i r(s_i, a_i) - J \quad (12)$$

Now we can plug in the parametrized gradient descent for the advantage function. That this works and is indeed the same has been proven by reference. Additionally we bring the cost J to the other side of the equation.

$$\Rightarrow \sum_{i=0}^N \gamma^i \nabla_{\Theta} [\log \pi(a_i | s_i)^T] \cdot w + 1 \cdot J = \sum_{i=0}^N \gamma^i r(s_i, a_i) \quad (13)$$

Let's do some rewriting. We define the following two terms:

$$\Phi_e = \left[\sum_{i=0}^N \gamma^i \nabla_{\Theta} [\log \pi(a_i | s_i)^T], 1 \right] \quad (14)$$

$$R_e = \sum_{i=0}^N \gamma^i r(s_i, a_i) \quad (15)$$

This let's us rewrite equation 13 as:

$$\Phi_e \cdot \begin{bmatrix} w \\ J \end{bmatrix} = R_e \quad (16)$$

An easy way to solve this system of equations is by taking the pseudo inverse of Φ_e .

$$\begin{bmatrix} w \\ J \end{bmatrix} = (\Phi_e^T \Phi_e)^{-1} \Phi_e^T R_e \quad (17)$$

Algorithm 1 Episodic Natural Actor-Critic Algorithm (eNAC)

Require: $n \geq 0 \vee x \neq 0$

Ensure: $y = x^n$

for $u = 1, 2, 3, \dots$ **do**

for $e = 1, 2, 3, \dots$ **do**

Execute Rollout: Draw initial state $s_0 \sim p(s_0)$

for $t = 1, 2, 3, \dots, N$ **do**

 Draw action $u_t \sim \pi(a_t | s_t)$, observe next state $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$, and reward $r_t = r(s_t, a_t)$.

end for

end for

end for

References

1. Shalabh Bhatnagar, Mohammad Ghavamzadeh, Mark Lee, and Richard S Sutton. Incremental natural actor-critic algorithms. In *Advances in neural information processing systems*, pages 105–112, 2008.
2. Shalabh Bhatnagar, Richard S Sutton, Mohammad Ghavamzadeh, and Mark Lee. Natural actor-critic algorithms. *Automatica*, 45(11):2471–2482, 2009.
3. Christoph Dann, Gerhard Neumann, and Jan Peters. Policy evaluation with temporal differences: A survey and comparison. *The Journal of Machine Learning Research*, 15(1):809–883, 2014.
4. Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, 71(7-9):1180–1190, 2008.
5. Jan Peters and Stefan Schaal. Reinforcement learning of motor skills with policy gradients. *Neural networks*, 21(4):682–697, 2008.
6. Jan Peters, Sethu Vijayakumar, and Stefan Schaal. Natural actor-critic. In *European Conference on Machine Learning*, pages 280–291. Springer, 2005.