# ELEC-E8125 Reinforcement Learning Actor-critic methods

Ville Kyrki

19.10.2021

# Motivation

- Policy gradient (PG) methods may be often ineffective in terms of requiring lots (and lots and lots) of data because of high variance of gradient estimates.
  - Similar to MC approaches for value function estimation.

- Temporal difference (TD) approaches have smaller variance compared to MC but they cannot handle stochastic policies or continuous action spaces like PG.

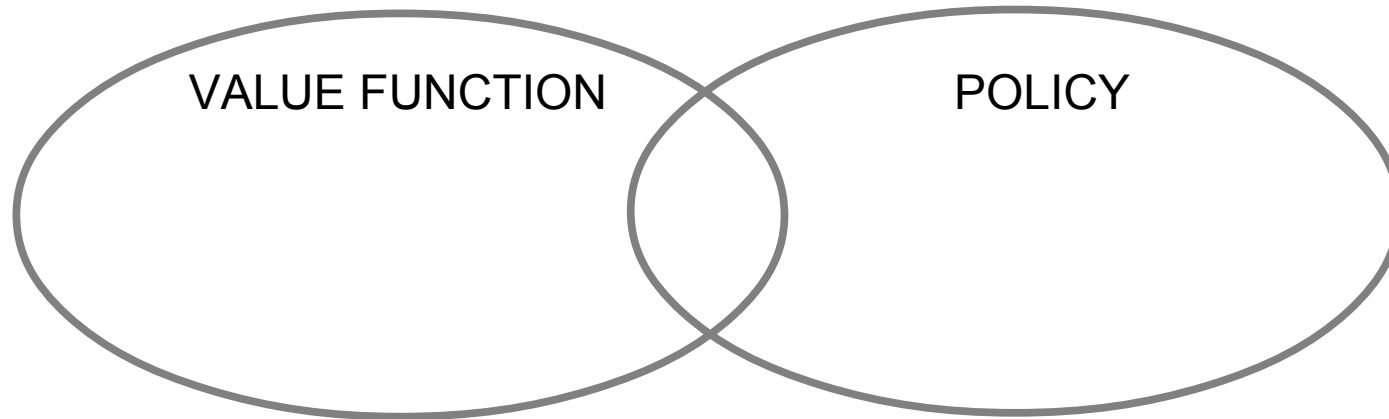- Can we combine PG with something like TD?

# **Today**

- Combining policy gradient with value functions
  → actor-critic methods

# Learning goals

- Build basis for understanding recent approaches combining policy learning and value estimation.

# Value-based vs policy-based RL



Value-based
· Learnt value function.
· Implicit policy.
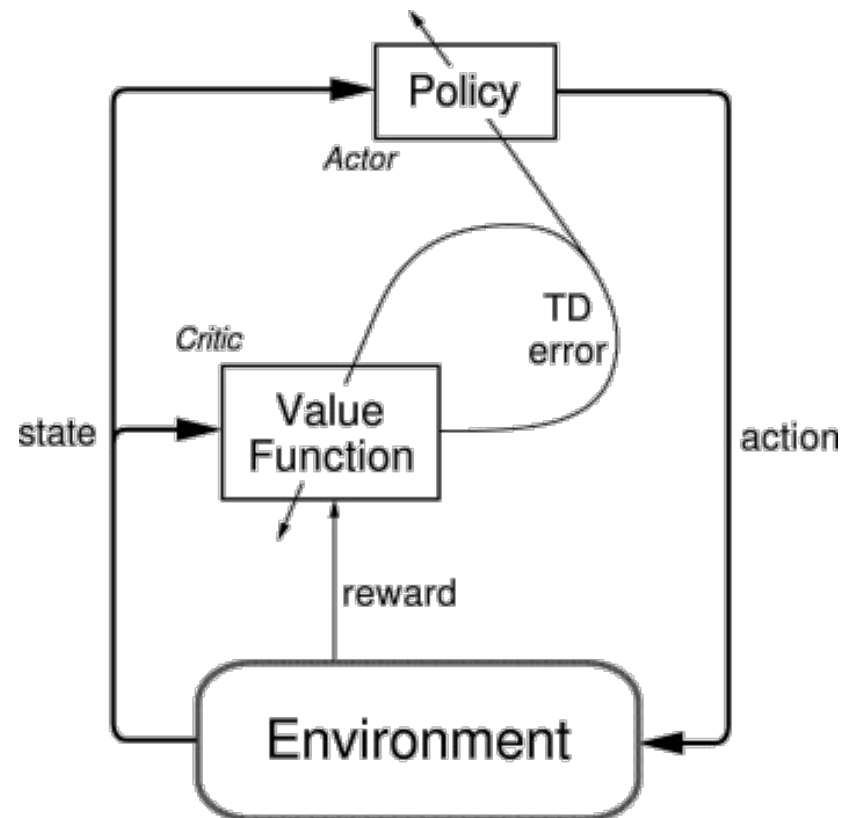
**Actor-critic**
· **Learnt value function.**
· **Learnt policy.**

Policy-based
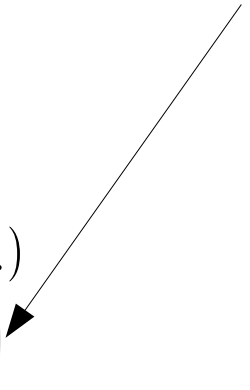· No value function.
· Learnt policy.

# Actor-critic approach – overview

- *Critic* estimates value function.

- *Actor* updates policy in direction of critic.

- For policy gradient, critic estimates value function.
  - See previous lectures.

# Policy gradient – recap

REINFORCE

1. Run policy, collect $\{\boldsymbol{\tau}_i\}$ $\qquad \boldsymbol{\tau_i}=(\boldsymbol{s_0^i}, \boldsymbol{a_0^i}, r_0^i, \boldsymbol{s_1^i}, \boldsymbol{a_1^i}, r_1^i, \ldots)$

2. $\nabla_\theta R(\theta) \approx \dfrac{1}{J}\sum_{i=1}^{J}\left(\sum_{t=0}^{T}\nabla_{\boldsymbol{\theta}}\log\pi_{\boldsymbol{\theta}}(\boldsymbol{a_t^i}|\boldsymbol{s_t^i})\left(\sum_{t'=t}^{T}r(\boldsymbol{s_{t'}^i}, \boldsymbol{a_{t'}^i})\right)\right)$

3. $\theta \leftarrow \theta + \alpha\nabla_{\boldsymbol{\theta}}R(\theta)$

**Aalto University**
**School of Electrical**
**Engineering**

# Policy gradient – recap

REINFORCE

1. Run policy, collect $\{\tau_i\}$ $\quad \tau_i = (s_0^i, a_0^i, r_0^i, s_1^i, a_1^i, r_1^i, \ldots)$

2. $\nabla_\theta R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta (a_t^i | s_t^i) \left( \underbrace{\sum_{t'=t}^{T} r(s_{t'}^i, a_{t'}^i)}_{} \right) \right)$

3. $\theta \leftarrow \theta + \alpha \nabla_\theta R(\theta)$

What's this?
Does it look familiar?

–

# Policy gradient – recap

REINFORCE

1. Run policy, collect $\{\tau_i\}$    $\tau_i = (s_0^i, a_0^i, r_0^i, s_1^i, a_1^i, r_1^i, \dots)$

2. $\nabla_\theta R(\theta) \approx \dfrac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \underbrace{\left( \sum_{t'=t}^{T} r(s_{t'}^i, a_{t'}^i) \right)} \right)$

3. $\theta \leftarrow \theta + \alpha \nabla_\theta R(\theta)$

What's this?
Does it look familiar?

$$Q_\pi(s_t, a_t) = \sum_{t=t'}^{T} E\left[ r(s_{t'}^i, a_{t'}^i) | s_t, a_t \right]$$

Q is true expected reward, unlike the estimate in step 2.
This would reduce variance of the gradient estimate.

$$\nabla_\theta R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) Q(s_t^i, a_t^i) \right)$$

$$\nabla_{\boldsymbol{\theta}} R(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}}\left[\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau})(R(\boldsymbol{\tau}) - \boxed{b})\right]$$

# Remember the baselines?

$$\nabla_{\boldsymbol{\theta}} R(\theta) \approx \frac{1}{J}\sum_{i=1}^{J}\left(\sum_{t=0}^{T}\nabla_{\boldsymbol{\theta}}\log \pi_{\boldsymbol{\theta}}(\boldsymbol{a}_t^i|\boldsymbol{s}_t^i)(Q(\boldsymbol{s}_t^i, \boldsymbol{a}_t^i) - b)\right)$$

Average is a good baseline: $b_t = \dfrac{1}{J}\sum_{i=1}^{J} Q(s_t^i, a_t^i)$

But what does the average mean here?

**A"** Aalto University
School of Electrical
Engineering

$$\nabla_\theta R(\boldsymbol{\theta}) = E_\theta \left[ \nabla_\theta \log p_\theta(\boldsymbol{\tau})(R(\boldsymbol{\tau}) - \boxed{b}) \right]$$

# Remember the baselines?

$$\nabla_\theta R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(\boldsymbol{a}_t^i | \boldsymbol{s}_t^i)(Q(\boldsymbol{s_t^i}, \boldsymbol{a_t^i}) - b) \right)$$

Average is a good baseline: $b_t = \frac{1}{J} \sum_{i=1}^{J} Q(s_t^i, a_t^i)$

But what does the average mean here?

*b* approximates the state value function *V(x)!*

$$\nabla_\theta R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(\boldsymbol{a}_t^i | \boldsymbol{s}_t^i)(Q(\boldsymbol{s_t^i}, \boldsymbol{a_t^i}) - V(\boldsymbol{s_t^i})) \right)$$

$$\nabla_{\theta} R(\theta) = E_{\theta} \left[ \nabla_{\theta} \log p_{\theta}(\tau)(R(\tau) - \boxed{b}) \right]$$

# Remember the baselines?

$$\nabla_{\theta} R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i)(Q(s_t^i, a_t^i) - b) \right)$$

Average is a good baseline: $b_t = \frac{1}{J} \sum_{i=1}^{J} Q(s_t^i, a_t^i)$

But what does the average mean here?

*b* approximates the state value function *V(s)!*

How much better is taking action *a* compared to average?

$$\nabla_{\theta} R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i)(\underbrace{Q(s_t^i, a_t^i) - V(s_t^i)}) \right)$$

*advantage* function $A(s_t^i, a_t^i)$

$$\nabla_{\theta} R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) A(s_t^i, a_t^i) \right)$$

# Determining the advantage

How to find a good estimate for this?

$$\nabla_{\theta} R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) A(s_t^i, a_t^i) \right)$$

Estimate $Q$, $V$, or $A$?

$V$ has the fewest parameters, so let's estimate it (from data).
But how to then get $A$?

$$A(s_t, a_t) = Q(s_t, a_t) - V(s_t)$$

$$Q(s_t, a_t) = r(s_t, a_t) + \gamma E_{s_{t+1} \sim \pi(s_{t+1}|s_t, a_t)}[V(s_{t+1})]$$

$$A(s_t, a_t) \approx r_{(s_t, a_t)} + \gamma V(s_{t+1}) - V(s_t)$$

Thus, knowing $V$ allows approximating $A$.

How to fit $V$?

Does this look familiar?

# Fitting value functions (mostly recap)

- Episodic batch fitting: (1) gather data, (2) fit (least squares) over gathered data.
- Data = state-value pairs $\left\{ \left( s_t^i, \underbrace{\sum_{t'=t}^{T} r_{t'}^i}_{y_t^i} \right) \right\}$

- Requires episodic environments to get the value.
- Fitting criterion $L(\phi) = \sum_i \left\| V_\phi(s_i) - y_i \right\|^2$

Any parametric function approximator

But what about non-episodic?
What do we do then?

# Fitting value functions (mostly recap)

- Non-episodic batch fitting: (1) gather data, (2) fit (least squares) over gathered data.

- Data = state-value pairs $\left\{ \left( s_t^i, \underbrace{r_t^i + V\left( s_{t+1}^i \right)}_{y_t^i} \right) \right\}$

- Identical fitting criterion

$$L(\phi) = \sum_i \left\| V_\phi(s_i) - y_i \right\|^2$$

Any parametric function approximator

# Wrap-up: A batch TD actor critic

Batch actor-critic

1. Run policy, collect $\{\tau_i\}$    $\tau_i = (s_0^i, a_0^i, r_0^i, s_1^i, a_1^i, r_1^i, \dots)$

2. Fit $V_\phi(s_t)$

3. Evaluate $A(s_t, a_t) \approx r_{(s_t, a_t)} + \gamma V(s_{t+1}) - V(s_t)$

4. Evaluate $\nabla_\theta R(\theta) \approx \frac{1}{J} \sum_{i=1}^{J} \left( \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) A(s_t^i, a_t^i) \right)$

5. Update $\theta \leftarrow \theta + \alpha \nabla_\theta R(\theta)$

6. Repeat

What about discount?

# An on-line TD actor critic (with discount)

learning rate

On-line actor-critic

1. Take action $a = \pi(a|s)$

From lecture 4!

2. Update $V_\phi(s_t)$ using $\phi \leftarrow \phi + \beta \left( r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t) \right) \nabla_\phi V_\phi(s_t)$

3. Evaluate $A(s_t, a_t) \approx r_{(s_t, a_t)} + \gamma V(s_{t+1}) - V(s_t)$

4. Evaluate $\nabla_\theta R(\theta) \approx \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) A(s_t^i, a_t^i)$

5. Update $\theta \leftarrow \theta + \alpha \nabla_\theta R(\theta)$

In practice, even this works best in batches (decreases variance in gradient estimates).

Note: TD estimate can be biased.

# Challenge: Gradient step sizes

$$\theta \leftarrow \theta + \alpha \nabla_{\boldsymbol{\theta}} R(\theta)$$

Gradient step size affects convergence (speed) greatly but is difficult to set.

Incorrect step size may lead to divergence or slow convergence.

How to guarantee policy improvement?

# Reformulating policy gradient through surrogate advantage

- How to predict performance of updated policy (since we do not have data about it yet)?

- Surrogate advantage $R^{IS}_{\theta_{old}}(\theta)$ approximates performance difference between previous and updated policies

$$R^{IS}_{\theta_{old}}(\theta) = E_{\tau \sim \pi_{\theta_{old}}} \left[ \frac{\pi_{\theta}(\boldsymbol{a_t}|\boldsymbol{s_t})}{\pi_{\theta_{old}}(\boldsymbol{a_t}|\boldsymbol{s_t})} A^{\pi_{\theta_{old}}}(\boldsymbol{s_t}, \boldsymbol{a_t}) \right]$$

See the importance sampling in effect!

Can we find a lower bound for this?
Yes, using KL-divergence.

# Bounding surrogate advantage

Result: Policy is guaranteed to improve by optimizing

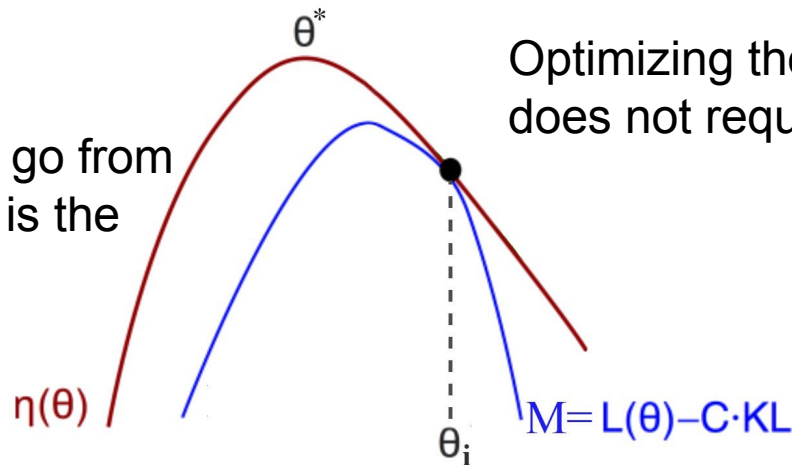$$max_{\theta}\left(R_{\theta_{old}}^{IS}(\theta) - c\,D_{KL}^{max}(\theta_{old}, \theta)\right)$$

where

known constant

$$D_{KL}^{max}(\theta_{old}, \theta)$$

is the maximum Kullback-Leibler divergence between the policies.

$\theta^*$

Optimizing the lower bound function does not require step size!

Intuition: The further you go from current policy, the larger is the penalty.

$\eta(\theta)$

$\theta_i$

$M = L(\theta) - C \cdot KL$

In practice leads to slow convergence, not easy to optimize.

# Trust region policy optimization (Schulman et al. 2015)

Instead of lower bound, optimize surrogate advantage and constrain KL-divergence:

$$max_\theta R^{IS}_{\theta_{old}}(\theta)$$

such that

$$\bar{D}_{KL}(\theta_{old}, \theta) \equiv E_{\tau \sim \pi_{\theta_{old}}}\left[D_{KL}\left(\pi_\theta(.|s), \pi_{\theta_{old}}(.|s)\right)\right] \leq \delta$$

Intuition: Limit the policy parameter change such that the actions do not change too much in the relevant part of state space.

In practice, this is still costly to compute and instead of using that constraint, the constraint is approximated (details in the paper).

**A"** Aalto University
School of Electrical
Engineering

Next: a simple and practical way to implement the same idea (and it even works well usually).
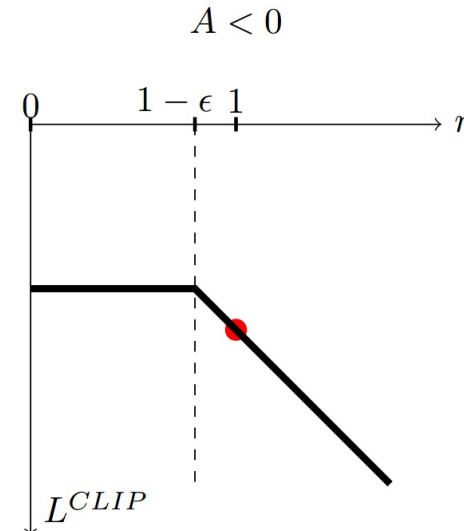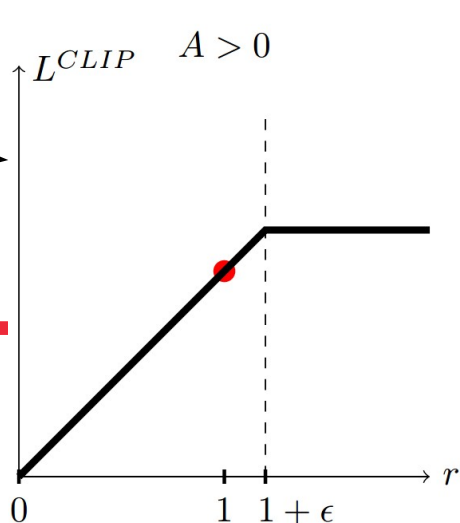
# Proximal policy optimization (Schulman et al. 2017)

Remember the surrogate advantage?

$$R^{IS}_{\theta_{old}}(\theta) = E_{\tau \sim \pi_{\theta_{old}}}\left[ \underbrace{\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}}_{g_t(\theta)} A^{\pi_{\theta_{old}}}(s_t, a_t) \right]$$

Optimize instead

$$L^{CLIP}(\theta) = E_{\tau \sim \pi_{\theta_{old}}}\left[ min\left( g_t(\theta) A, clip\left( g_t(\theta), 1-\epsilon, 1+\epsilon \right) A \right) \right]$$

Looks horrible, look at the figure instead.
In practice: Limit how much benefit there is for changes.

# Proximal policy optimization (Schulman et al. 2017)
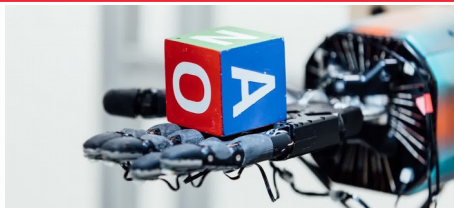
Other variants possible

**Algorithm**: PPO

**for** i = 1, 2, … do

Run policy, collect trajectories $\{\tau_i\}$  $\tau_i = (s_0^i, a_0^i, r_0^i, s_1^i, a_1^i, r_1^i, \ldots)$

Compute advantage estimates $A(s_t, a_t) \approx r_{(s_t, a_t)} + \gamma V(s_{t+1}) - V(s_t)$
using current value function $V_\phi(s_t)$

Update policy by maximizing $L^{CLIP}(\theta)$ for K epochs of stochastic gradient ascent

Fit $V_\phi(s_t)$ by minimizing $L(\phi) \equiv \sum_i \left\| V_\phi(s_i) - y_i \right\|^2$ using gradient descent

PPO is a standard baseline at the moment.

# Recent successful algorithms

- Deep Deterministic Policy Gradient (DDPG)            2016
    - Q-function learning + deterministic policy for continuous action spaces
    - Off-policy


- Soft Actor Critic (SAC)            2018
    - Q-function learning + stochastic policy for continuous action spaces
    - Off-policy

# Summary

- Actor-critic approaches allow addressing continuing problems and continuous action spaces.

- They may also learn faster than policy gradient because variance of policy gradient estimate is reduced.

- TRPO/PPO aim to control extent of policy update steps to avoid oscillation/divergence due to large updates.

Aalto University
School of Electrical
Engineering

# Next: Optimal control – Toward model-based RL

- Even with a critic, policy-based approaches often require huge amounts of data.

- Can we somehow benefit even more from earlier experiences?

- Next week: Lecturer changes to Joni Pajarinen.

**Aalto University**
School of Electrical Engineering