# Policy Gradient Methods

**Milan Straka**

📅 **November 09, 2020**

Charles University in Prague
Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics

Instead of predicting expected returns, we could train the method to directly predict the policy

$$\pi(a|s; \boldsymbol{\theta}).$$
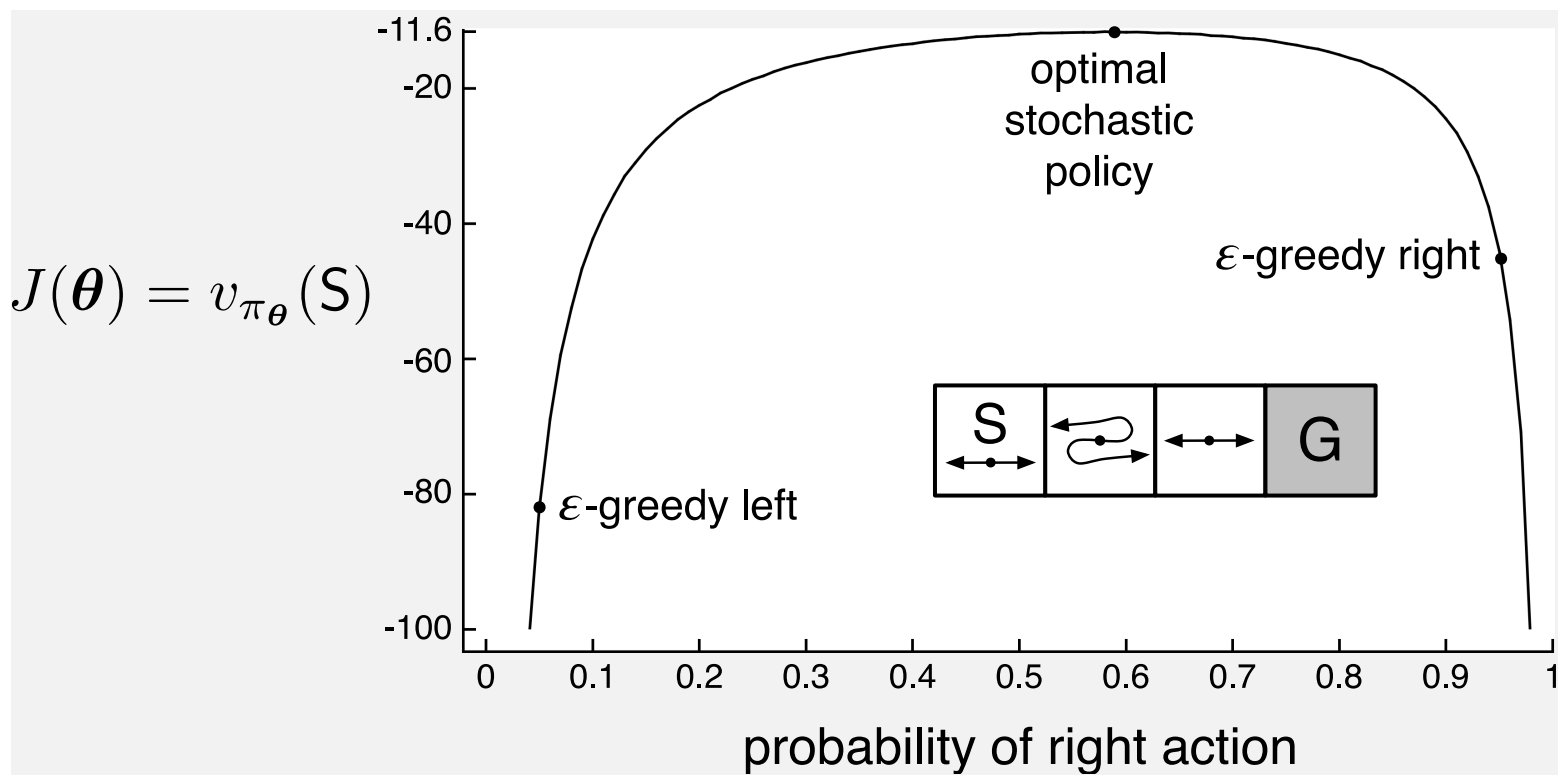
Obtaining the full distribution over all actions would also allow us to sample the actions according to the distribution $\pi$ instead of just $\varepsilon$-greedy sampling.

However, to train the network, we maximize the expected return $v_\pi(s)$ and to that account we need to compute its *gradient* $\nabla_{\boldsymbol{\theta}} v_\pi(s)$.

In addition to discarding $\varepsilon$-greedy action selection, policy gradient methods allow producing policies which are by nature stochastic, as in card games with imperfect information, while the action-value methods have no natural way of finding stochastic policies (distributional RL might be of some use though).



$$J(\boldsymbol{\theta}) = v_{\pi_{\boldsymbol{\theta}}}(\mathrm{S})$$

*Example 13.1 of "Reinforcement Learning: An Introduction, Second Edition".*

Let $\pi(a|s; \boldsymbol{\theta})$ be a parametrized policy. We denote the initial state distribution as $h(s)$ and the on-policy distribution under $\pi$ as $\mu(s)$. Let also $J(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \mathbb{E}_{s \sim h} v_\pi(s)$.

Then

$$\nabla_{\boldsymbol{\theta}} v_\pi(s) \propto \sum_{s' \in \mathcal{S}} P(s \to \ldots \to s' | \pi) \sum_{a \in \mathcal{A}} q_\pi(s', a) \nabla_{\boldsymbol{\theta}} \pi(a|s'; \boldsymbol{\theta})$$

and

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \propto \sum_{s \in \mathcal{S}} \mu(s) \sum_{a \in \mathcal{A}} q_\pi(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s; \boldsymbol{\theta}),$$

where $P(s \to \ldots \to s' | \pi)$ is the probability of getting to state $s'$ when starting from state $s$, after any number of 0, 1, … steps. The $\gamma$ parameter should be treated as a form of termination, i.e., $P(s \to \ldots \to s' | \pi) \propto \sum_{k=0}^{\infty} \gamma^k P(s \to s' \text{ in } k \text{ steps} | \pi)$.

$$\nabla v_\pi(s) = \nabla\Big[\sum_a \pi(a|s;\boldsymbol{\theta})q_\pi(s,a)\Big]$$

$$= \sum_a \Big[\nabla\pi(a|s;\boldsymbol{\theta})q_\pi(s,a) + \pi(a|s;\boldsymbol{\theta})\nabla q_\pi(s,a)\Big]$$

$$= \sum_a \Big[\nabla\pi(a|s;\boldsymbol{\theta})q_\pi(s,a) + \pi(a|s;\boldsymbol{\theta})\nabla\big(\sum_{s'} p(s'|s,a)(r + \gamma v_\pi(s'))\big)\Big]$$

$$= \sum_a \Big[\nabla\pi(a|s;\boldsymbol{\theta})q_\pi(s,a) + \gamma\pi(a|s;\boldsymbol{\theta})\big(\sum_{s'} p(s'|s,a)\nabla v_\pi(s')\big)\Big]$$

We now expand $v_\pi(s')$.

$$= \sum_a \Big[\nabla\pi(a|s;\boldsymbol{\theta})q_\pi(s,a) + \gamma\pi(a|s;\boldsymbol{\theta})\Big(\sum_{s'} p(s'|s,a)\Big($$
$$\sum_{a'}\Big[\nabla\pi(a'|s';\boldsymbol{\theta})q_\pi(s',a') + \gamma\pi(a'|s';\boldsymbol{\theta})\big(\sum_{s''} p(s''|s',a')\nabla v_\pi(s'')\big)\Big]\Big)\Big)\Big]$$

Continuing to expand all $v_\pi(s'')$, we obtain the following:

$$\nabla v_\pi(s) = \sum_{s'\in\mathcal{S}}\sum_{k=0}^{\infty}\gamma^k P(s\to s' \text{ in } k \text{ steps } |\pi)\sum_{a\in\mathcal{A}}q_\pi(s',a)\nabla_{\boldsymbol{\theta}}\pi(a|s';\boldsymbol{\theta}).$$

To finish the proof of the first part, recall that

$$\sum_{k=0}^{\infty} \gamma^k P(s \to s' \text{ in } k \text{ steps} | \pi) \propto P(s \to \ldots \to s' | \pi).$$

For the second part, we know that

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{s \sim h} \nabla_{\boldsymbol{\theta}} v_{\pi}(s) \propto \mathbb{E}_{s \sim h} \sum_{s' \in \mathcal{S}} P(s \to \ldots \to s' | \pi) \sum_{a \in \mathcal{A}} q_{\pi}(s', a) \nabla_{\boldsymbol{\theta}} \pi(a|s'; \boldsymbol{\theta}),$$

therefore using the fact that $\mu(s') = \mathbb{E}_{s \sim h} P(s \to \ldots \to s' | \pi)$ we get

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \propto \sum_{s \in \mathcal{S}} \mu(s) \sum_{a \in \mathcal{A}} q_{\pi}(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s; \boldsymbol{\theta}).$$

The REINFORCE algorithm (Williams, 1992) uses directly the policy gradient theorem, minimizing $-J(\boldsymbol{\theta}) \overset{\text{def}}{=} -\mathbb{E}_{s \sim h} v_\pi(s)$. The loss gradient is then

$$\nabla_{\boldsymbol{\theta}} - J(\boldsymbol{\theta}) \propto -\sum_{s \in \mathcal{S}} \mu(s) \sum_{a \in \mathcal{A}} q_\pi(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s; \boldsymbol{\theta}) = -\mathbb{E}_{s \sim \mu} \sum_{a \in \mathcal{A}} q_\pi(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s; \boldsymbol{\theta}).$$

However, the sum over all actions is problematic. Instead, we rewrite it to an expectation which we can estimate by sampling:

$$\nabla_{\boldsymbol{\theta}} - J(\boldsymbol{\theta}) \propto \mathbb{E}_{s \sim \mu} \mathbb{E}_{a \sim \pi} q_\pi(s, a) \nabla_{\boldsymbol{\theta}} - \ln \pi(a|s; \boldsymbol{\theta}),$$

where we used the fact that

$$\nabla_{\boldsymbol{\theta}} \ln \pi(a|s; \boldsymbol{\theta}) = \frac{1}{\pi(a|s; \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} \pi(a|s; \boldsymbol{\theta}).$$

REINFORCE therefore minimizes the loss

$$\mathbb{E}_{s \sim \mu} \mathbb{E}_{a \sim \pi} q_\pi(s, a) \nabla_{\boldsymbol{\theta}} - \ln \pi(a|s; \boldsymbol{\theta}),$$

estimating the $q_\pi(s, a)$ by a single sample.

Note that the loss is just a weighted variant of negative log likelihood (NLL), where the sampled actions play a role of gold labels and are weighted according to their return.

---

**REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for $\pi_*$**

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$
Algorithm parameter: step size $\alpha > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):
    Generate an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \boldsymbol{\theta})$
    Loop for each step of the episode $t = 0, 1, \ldots, T - 1$:
        $G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k$              $(G_t)$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \quad \alpha \, G \nabla \ln \pi(A_t|S_t, \boldsymbol{\theta})$

*Modified from Algorithm 13.3 of "Reinforcement Learning: An Introduction, Second Edition" by removing γ̂t from the update of θ.*

---

In the proof, we assumed $\gamma$ is used as a form of termination in the definition of the on-policy distribution.

However, even when discounting is used during training (to guarantee convergence even for very long episodes), evaluation is often performed without discounting.

Consequently, the distribution $\mu$ used in the REINFORCE algorithm is almost always the unterminated (undiscounted) on-policy distribution (I am not aware of any implementation or paper that would use it), so that we learn even in states that are far from the beginning of an episode.

The returns can be arbitrary – better-than-average and worse-than-average returns cannot be recognized from the absolute value of the return.

Hopefully, we can generalize the policy gradient theorem using a baseline $b(s)$ to

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \propto \sum_{s \in \mathcal{S}} \mu(s) \sum_{a \in \mathcal{A}} \big(q_\pi(s, a) - b(s)\big) \nabla_{\boldsymbol{\theta}} \pi(a|s; \boldsymbol{\theta}).$$

The baseline $b(s)$ can be a function or even a random variable, as long as it does not depend on $a$, because

$$\sum_a b(s) \nabla_{\boldsymbol{\theta}} \pi(a|s; \boldsymbol{\theta}) = b(s) \sum_a \nabla_{\boldsymbol{\theta}} \pi(a|s; \boldsymbol{\theta}) = b(s) \nabla_{\boldsymbol{\theta}} \sum_a \pi(a|s; \boldsymbol{\theta}) = b(s) \nabla 1 = 0.$$

A good choice for $b(s)$ is $v_\pi(s)$, which can be shown to minimize variance of the estimator. Such baseline reminds centering of returns, given that

$$v_\pi(s) = \mathbb{E}_{a\sim\pi}q_\pi(s, a).$$

Then, better-than-average returns are positive and worse-than-average returns are negative. The resulting $q_\pi(s, a) - v_\pi(s)$ function is also called an **advantage** function

$$a_\pi(s, a) \overset{\text{def}}{=} q_\pi(s, a) - v_\pi(s).$$

Of course, the $v_\pi(s)$ baseline can be only approximated. If neural networks are used to estimate $\pi(a|s; \boldsymbol{\theta})$, then some part of the network is usually shared between the policy and value function estimation, which is trained using mean square error of the predicted and observed return.

**REINFORCE with Baseline (episodic), for estimating $\pi_{\boldsymbol{\theta}} \approx \pi_*$**

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$
Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$
Algorithm parameters: step sizes $\alpha^{\boldsymbol{\theta}} > 0$, $\alpha^{\mathbf{w}} > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):
　　Generate an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \boldsymbol{\theta})$
　　Loop for each step of the episode $t = 0, 1, \ldots, T-1$:
　　　　$G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k$ 　　　　　　　　　　　　　　　　　 $(G_t)$
　　　　$\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$
　　　　$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S_t, \mathbf{w})$
　　　　$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \quad \alpha^{\boldsymbol{\theta}} \delta \nabla \ln \pi(A_t|S_t, \boldsymbol{\theta})$

*Modified from Algorithm 13.4 of "Reinforcement Learning: An Introduction, Second Edition" by removing γ^t from the update of θ.*
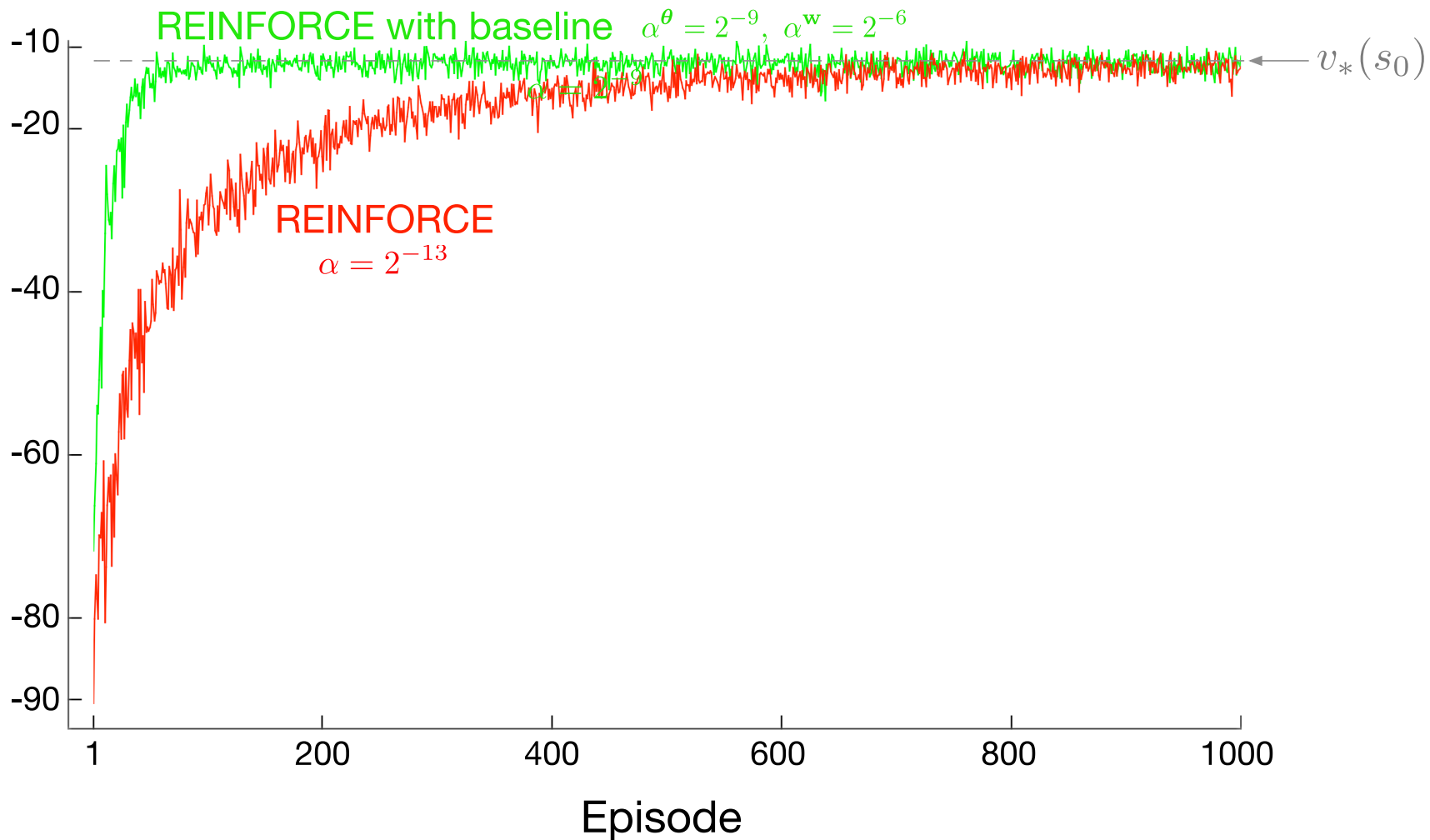
Figure 13.2 of "Reinforcement Learning: An Introduction, Second Edition".

# Actor-Critic

It is possible to combine the policy gradient methods and temporal difference methods, creating a family of algorithms usually called *actor-critic* methods.

The idea is straightforward – instead of estimating the episode return using the whole episode rewards, we can use $n$-step temporal difference estimation.

**One-step Actor–Critic (episodic), for estimating $\pi_{\boldsymbol{\theta}} \approx \pi_*$**

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$
Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$
Parameters: step sizes $\alpha^{\boldsymbol{\theta}} > 0$, $\alpha^{\mathbf{w}} > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g., to $\mathbf{0}$)
Loop forever (for each episode):
    Initialize $S$ (first state of episode)

    Loop while $S$ is not terminal (for each time step):
        $A \sim \pi(\cdot|S, \boldsymbol{\theta})$
        Take action $A$, observe $S', R$
        $\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$         (if $S'$ is terminal, then $\hat{v}(S', \mathbf{w}) \doteq 0$)
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S, \mathbf{w})$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \quad \alpha^{\boldsymbol{\theta}} \delta \nabla \ln \pi(A|S, \boldsymbol{\theta})$
        $S \leftarrow S'$

*Modified from Algorithm 13.5 of "Reinforcement Learning: An Introduction, Second Edition" by removing I.*

A 2015 paper from Volodymyr Mnih et al., the same group as DQN.

The authors propose an asynchronous framework, where multiple workers share one neural network, each training using either an off-line or on-line RL algorithm.

They compare 1-step Q-learning, 1-step Sarsa, $n$-step Q-learning and A3C (an *asynchronous advantage actor-critic* method). For A3C, they compare a version with and without LSTM.

The authors also introduce *entropy regularization term* $-\beta H(\pi(s; \boldsymbol{\theta}))$ to the loss to support exploration and discourage premature convergence (they use $\beta = 0.01$).

---

**Algorithm 1** Asynchronous one-step Q-learning - pseudocode for each actor-learner thread.

---

// *Assume global shared $\theta$, $\theta^-$, and counter $T = 0$.*
Initialize thread step counter $t \leftarrow 0$
Initialize target network weights $\theta^- \leftarrow \theta$
Initialize network gradients $d\theta \leftarrow 0$
Get initial state $s$
**repeat**
    Take action $a$ with $\epsilon$-greedy policy based on $Q(s, a; \theta)$
    Receive new state $s'$ and reward $r$
    $y = \begin{cases} r & \text{for terminal } s' \\ r + \gamma \max_{a'} Q(s', a'; \theta^-) & \text{for non-terminal } s' \end{cases}$
    Accumulate gradients wrt $\theta$: $d\theta \leftarrow d\theta + \frac{\partial(y - Q(s,a;\theta))^2}{\partial \theta}$
    $s = s'$
    $T \leftarrow T + 1$ and $t \leftarrow t + 1$
    **if** $T \mod I_{target} == 0$ **then**
        Update the target network $\theta^- \leftarrow \theta$
    **end if**
    **if** $t \mod I_{AsyncUpdate} == 0$ or $s$ is terminal **then**
        Perform asynchronous update of $\theta$ using $d\theta$.
        Clear gradients $d\theta \leftarrow 0$.
    **end if**
**until** $T > T_{max}$

---

*Algorithm 1 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.*

---

**Algorithm S2** Asynchronous n-step Q-learning - pseudocode for each actor-learner thread.

---

*// Assume global shared parameter vector $\theta$.*
*// Assume global shared target parameter vector $\theta^-$.*
*// Assume global shared counter $T = 0$.*
Initialize thread step counter $t \leftarrow 1$
Initialize target network parameters $\theta^- \leftarrow \theta$
Initialize thread-specific parameters $\theta' = \theta$
Initialize network gradients $d\theta \leftarrow 0$
**repeat**
    Clear gradients $d\theta \leftarrow 0$
    Synchronize thread-specific parameters $\theta' = \theta$
    $t_{start} = t$
    Get state $s_t$
    **repeat**
        Take action $a_t$ according to the $\epsilon$-greedy policy based on $Q(s_t, a; \theta')$
        Receive reward $r_t$ and new state $s_{t+1}$
        $t \leftarrow t + 1$
        $T \leftarrow T + 1$
    **until** terminal $s_t$ **or** $t - t_{start} == t_{max}$
    $R = \begin{cases} 0 & \text{for terminal } s_t \\ \max_a Q(s_t, a; \theta^-) & \text{for non-terminal } s_t \end{cases}$
    **for** $i \in \{t - 1, \ldots, t_{start}\}$ **do**
        $R \leftarrow r_i + \gamma R$
        Accumulate gradients wrt $\theta'$: $d\theta \leftarrow d\theta + \frac{\partial \left( R - Q(s_i, a_i; \theta') \right)^2}{\partial \theta'}$
    **end for**
    Perform asynchronous update of $\theta$ using $d\theta$.
    **if** $T \mod I_{target} == 0$ **then**
        $\theta^- \leftarrow \theta$
    **end if**
**until** $T > T_{max}$

---

*Algorithm S2 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.*

---

**Algorithm S3** Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

---

// *Assume global shared parameter vectors $\theta$ and $\theta_v$ and global shared counter $T = 0$*
// *Assume thread-specific parameter vectors $\theta'$ and $\theta'_v$*
Initialize thread step counter $t \leftarrow 1$
**repeat**
    Reset gradients: $d\theta \leftarrow 0$ and $d\theta_v \leftarrow 0$.
    Synchronize thread-specific parameters $\theta' = \theta$ and $\theta'_v = \theta_v$
    $t_{start} = t$
    Get state $s_t$
    **repeat**
        Perform $a_t$ according to policy $\pi(a_t|s_t; \theta')$
        Receive reward $r_t$ and new state $s_{t+1}$
        $t \leftarrow t + 1$
        $T \leftarrow T + 1$
    **until** terminal $s_t$ **or** $t - t_{start} == t_{max}$
    $R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t \end{cases}$// Bootstrap from last state
    **for** $i \in \{t - 1, \ldots, t_{start}\}$ **do**
        $R \leftarrow r_i + \gamma R$
        Accumulate gradients wrt $\theta'$: $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$
        Accumulate gradients wrt $\theta'_v$: $d\theta_v \leftarrow d\theta_v + \partial \left(R - V(s_i; \theta'_v)\right)^2 / \partial \theta'_v$
    **end for**
    Perform asynchronous update of $\theta$ using $d\theta$ and of $\theta_v$ using $d\theta_v$.
**until** $T > T_{max}$

---

*Algorithm S3 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.*

All methods performed updates every 5 actions ($t_{\mathrm{max}} = I_{\mathrm{AsyncUpdate}} = 5$), updating the target network each $40\,000$ frames.

The Atari inputs were processed as in DQN, using also action repeat 4.

The network architecture is: 16 filters $8 \times 8$ stride 4, 32 filters $4 \times 4$ stride 2, followed by a fully connected layer with 256 units. All hidden layers apply a ReLU non-linearity. Values and/or action values were then generated from the (same) last hidden layer.

The LSTM methods utilized a 256-unit LSTM cell after the dense hidden layer.

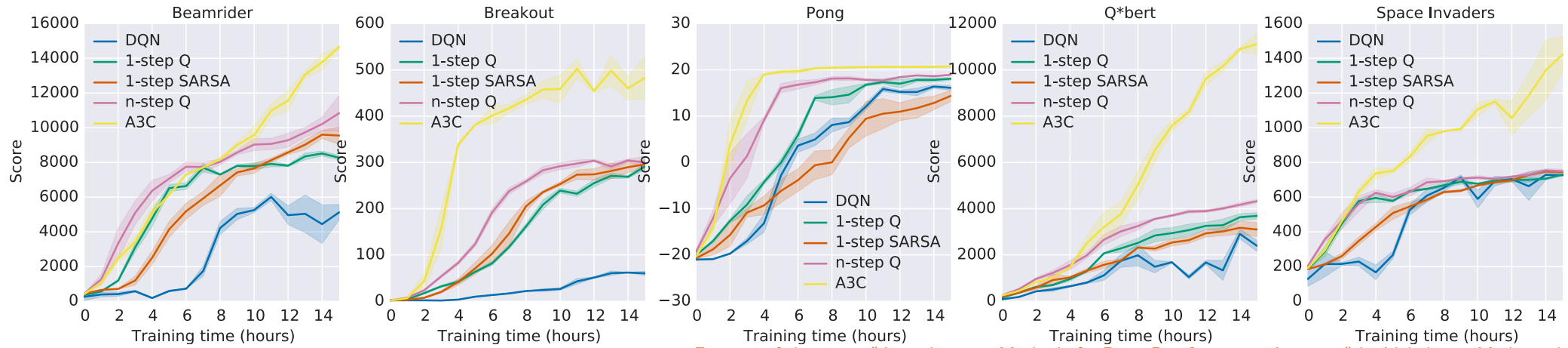All experiments used a discount factor of $\gamma = 0.99$ and used RMSProp with momentum decay factor of $0.99$.

Figure 1 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.

| Method | Training Time | Mean | Median |
|---|---|---|---|
| DQN | 8 days on GPU | 121.9% | 47.5% |
| Gorila | 4 days, 100 machines | 215.2% | 71.3% |
| D-DQN | 8 days on GPU | 332.9% | 110.9% |
| Dueling D-DQN | 8 days on GPU | 343.8% | 117.1% |
| Prioritized DQN | 8 days on GPU | 463.6% | 127.6% |
| A3C, FF | 1 day on CPU | 344.1% | 68.2% |
| A3C, FF | 4 days on CPU | 496.8% | 116.6% |
| A3C, LSTM | 4 days on CPU | 623.0% | 112.6% |

Table 1 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.

| Method | Number of threads | | | | |
|---|---|---|---|---|---|
| | 1 | 2 | 4 | 8 | 16 |
| 1-step Q | 1.0 | **3.0** | **6.3** | **13.3** | **24.1** |
| 1-step SARSA | 1.0 | **2.8** | **5.9** | **13.1** | **22.1** |
| n-step Q | 1.0 | **2.7** | **5.9** | **10.7** | **17.2** |
| A3C | 1.0 | 2.1 | 3.7 | 6.9 | 12.5 |

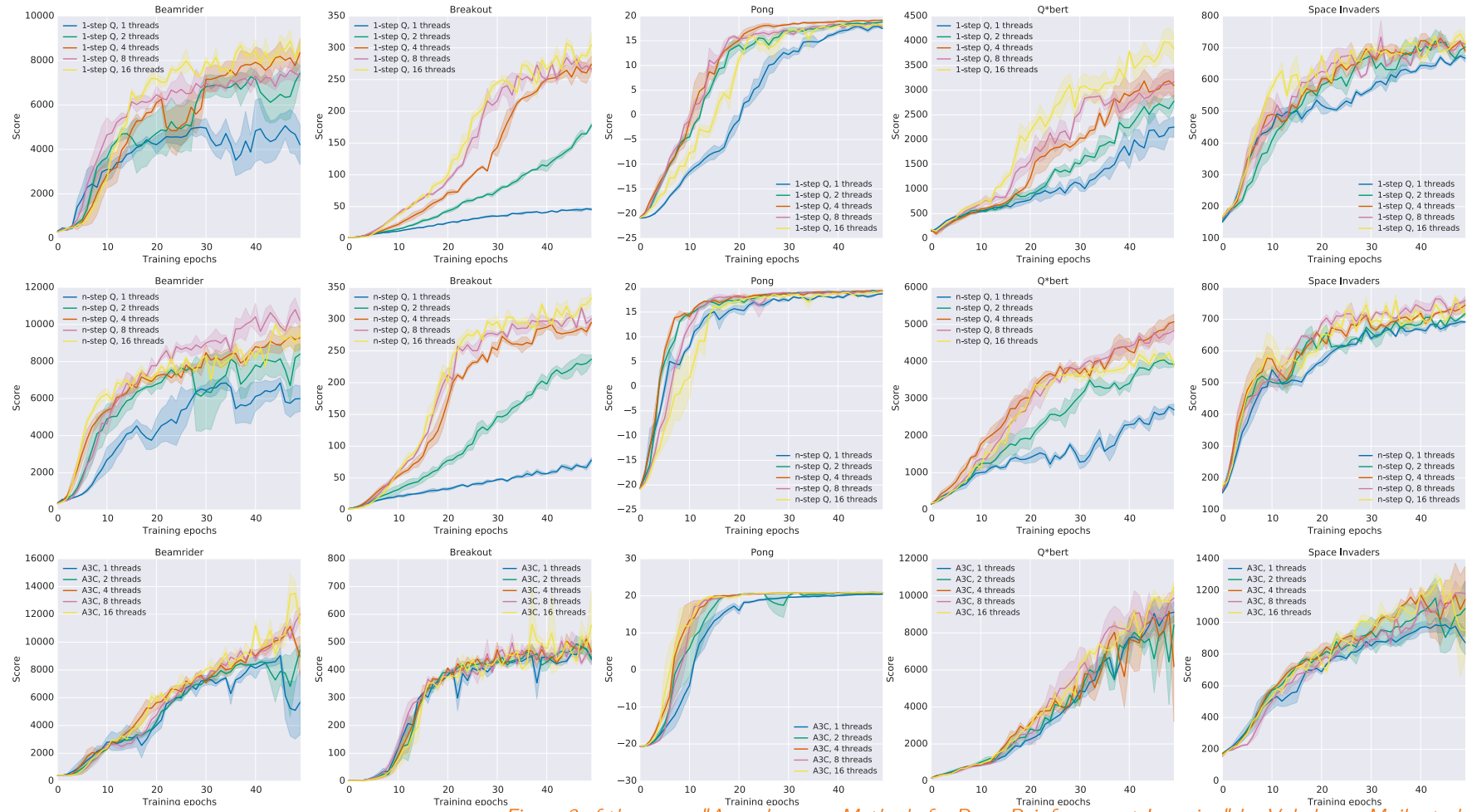Table 2 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.

*Figure 3 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.*
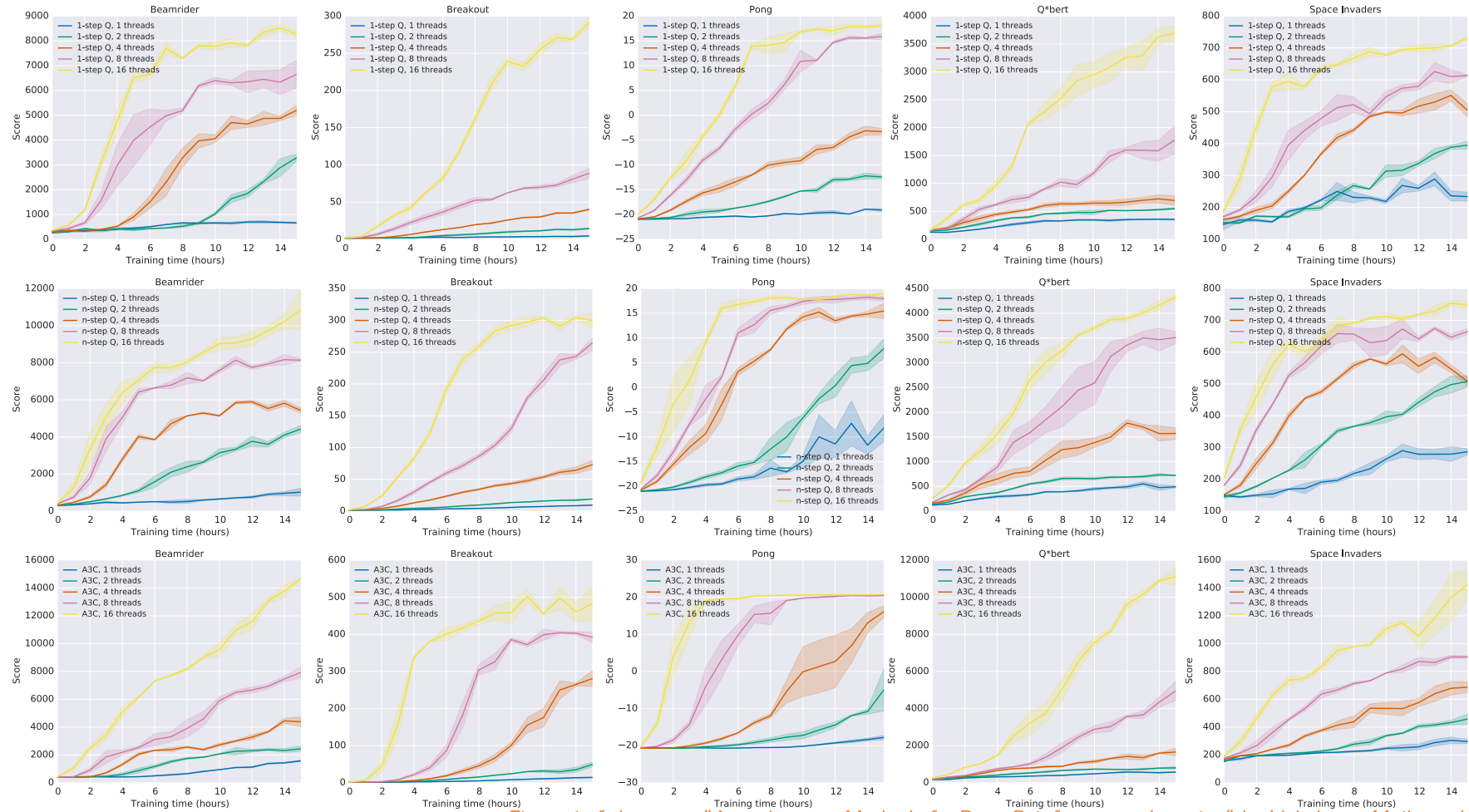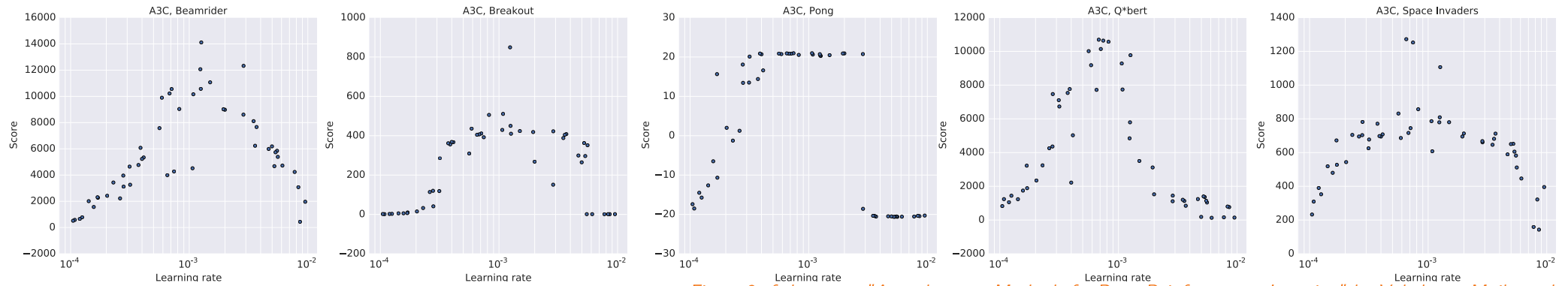
*Figure 4 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.*

Figure 2 of the paper "Asynchronous Methods for Deep Reinforcement Learning" by Volodymyr Mnih et al.