

Faster Deep Reinforcement Learning with Slower Online Network

Kavosh Asadi, Rasool Fakoor, Omer Gottesman, Taesup Kim, Michael L. Littma, Alexander J.
Smola

Proximal Bellman Operator

$$(T_{c,f}^\pi)^n = \arg \min_{v'} \|v' - (T^\pi)^n\|_2^2 + \frac{1}{c} D_f(v', v), \text{ где}$$

$$D_f(v', v) = f(v') - f(v) - \langle \nabla f(v), v' - v \rangle$$

Если $D_f(v', v) = \|v' - v\|_2^2$, то

$$(T_{c,f}^\pi)^n = \arg \min_{v'} \|v' - (T^\pi)^n\|_2^2 + \frac{1}{c} \|v' - v\|_2^2$$

Теорема. $T_{c,f}^*$ - сжимающий оператор с неподвижной точкой v^*

DQN PRO

$$\omega_{t+1} \leftarrow \arg \min_{\omega} h(\omega_t, \omega) + \frac{1}{2\tilde{c}} \|\omega - \omega_t\|_2^2, \text{ где}$$

$$h(\theta, \omega) = \hat{E}_{\langle s, a, r, s' \rangle} \left[\left(r + \gamma \max_{a'} \hat{Q}(s', a'; \theta) - \hat{Q}(s, a; \omega) \right)^2 \right]$$

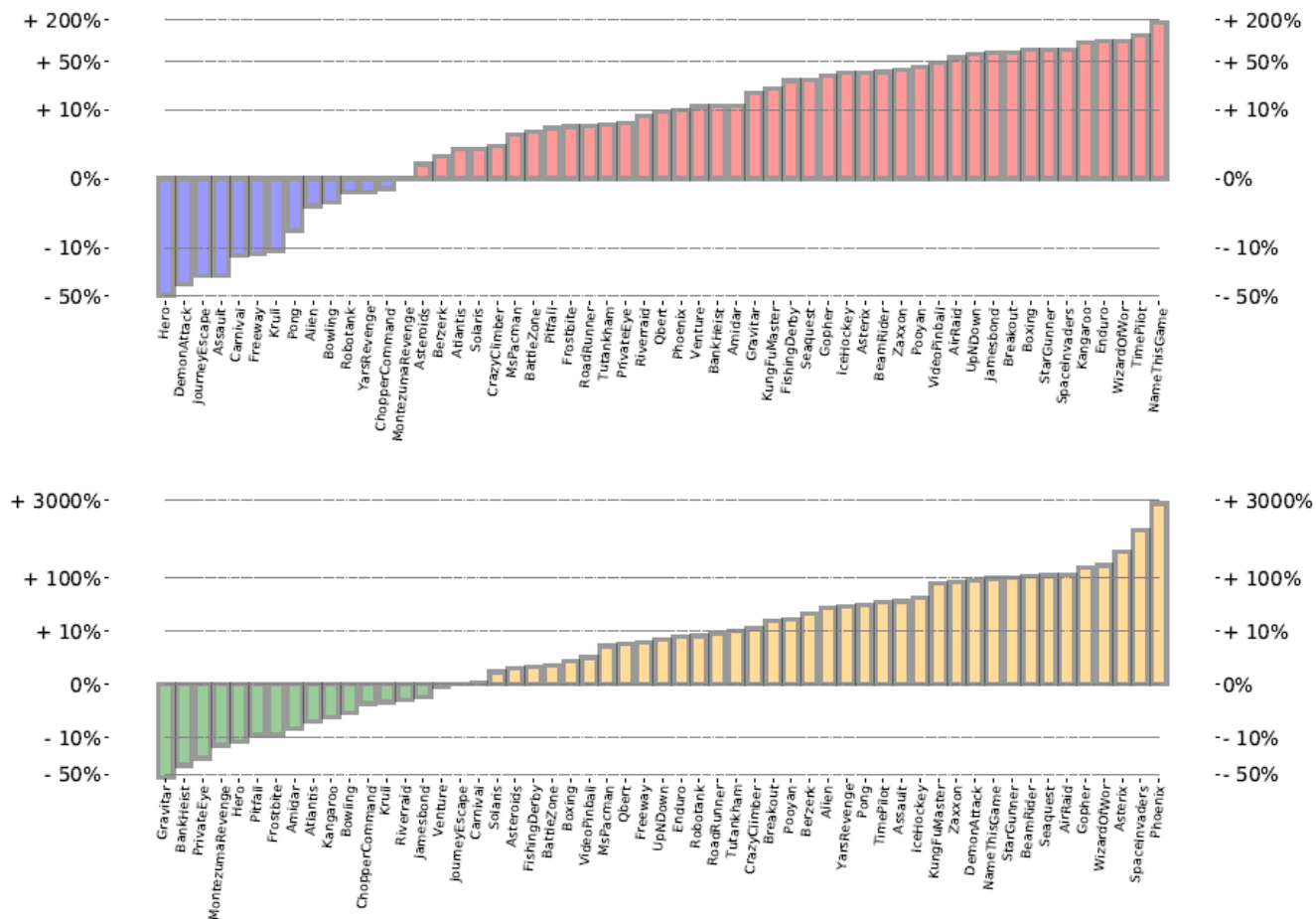
ИЛИ

$$\omega_{t+1} \leftarrow \left(1 - \frac{\alpha}{\tilde{c}}\right) \omega + \frac{\alpha}{\tilde{c}} \omega_t - \alpha \nabla_2 h(\omega_t, \omega)$$

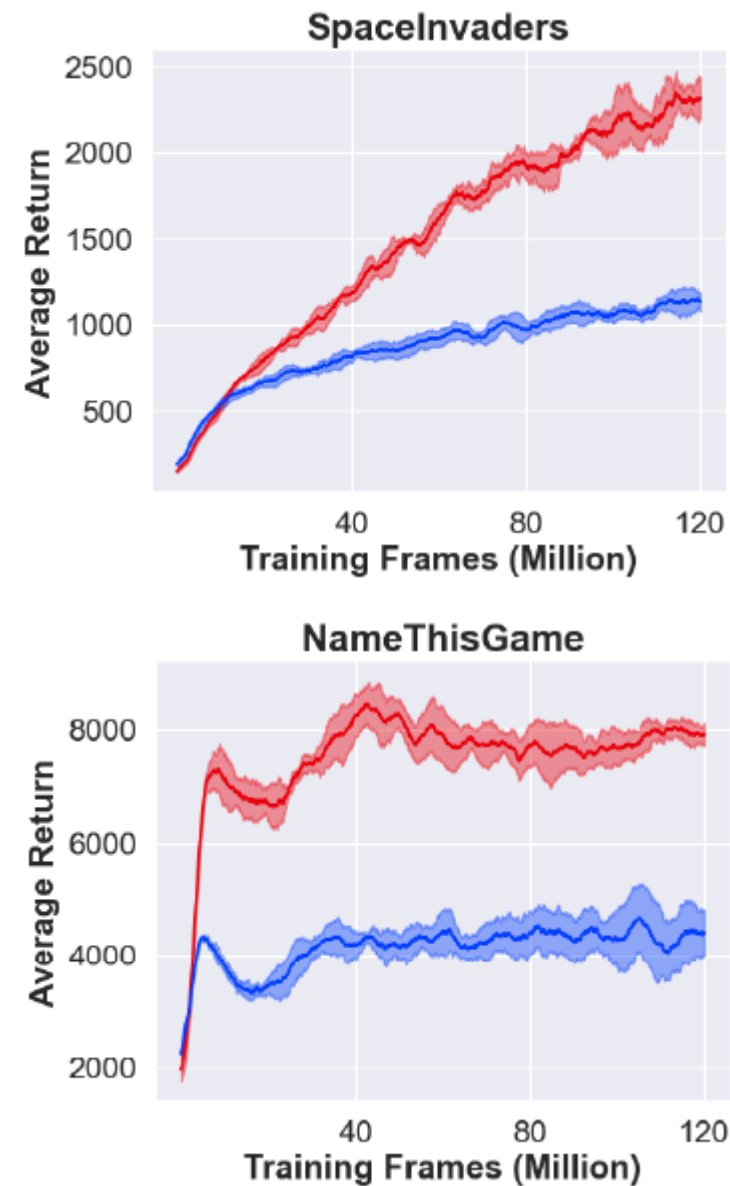
Algorithm 1 DQN with Proximal Iteration (DQN Pro)

```
1: Initialize  $\theta$ ,  $N$ ,  $period$ , replay buffer  $\mathcal{D}$ ,  $\alpha$ , and  $\tilde{c}$ 
2:  $s \leftarrow \text{env.reset}()$ ,  $w \leftarrow \theta$ ,  $\text{numUpdates} \leftarrow 0$ 
3: repeat
4:    $a \sim \epsilon\text{-greedy}(Q(s, \cdot; w))$ 
5:    $s', r \leftarrow \text{env.step}(s, a)$ 
6:   add  $\langle s, a, r, s' \rangle$  to  $\mathcal{D}$ 
7:   if  $s'$  is terminal then
8:      $s \leftarrow \text{env.reset}()$ 
9:   end if
10:  for  $n$  in  $\{1, \dots, N\}$  do
11:    sample  $\mathcal{B} = \{\langle s, a, r, s' \rangle\}$ , compute  $\nabla_w h(w)$ 
12:     $w \leftarrow (1 - (\alpha/\tilde{c}))w + (\alpha/\tilde{c})\theta - \alpha \nabla_w h(w)$ 
13:     $\text{numUpdates} \leftarrow \text{numUpdates} + 1$ 
14:    if  $\text{numUpdates} \% period = 0$  then
15:       $\theta \leftarrow w$ 
16:    end if
17:  end for
18: until convergence
```

Результаты исследователей

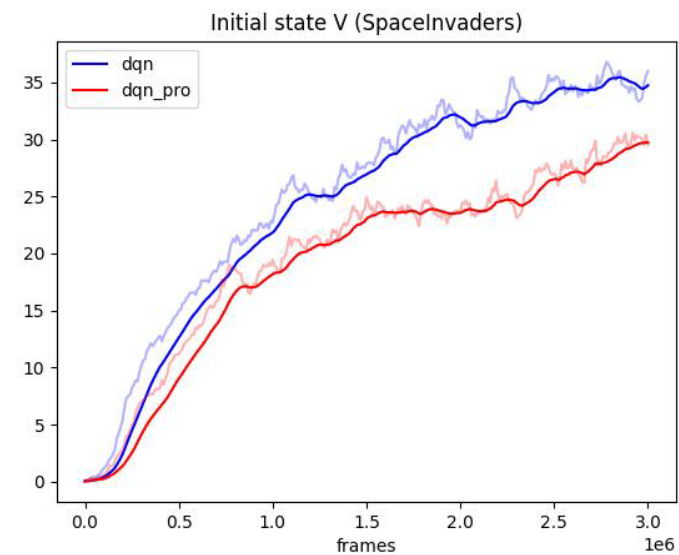
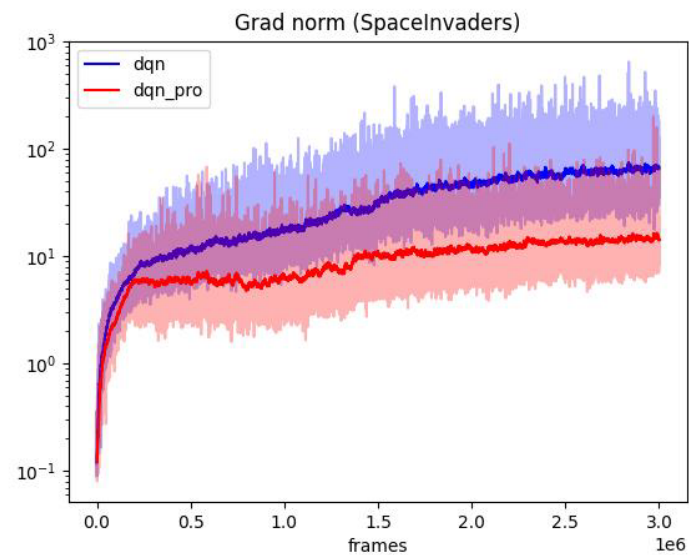
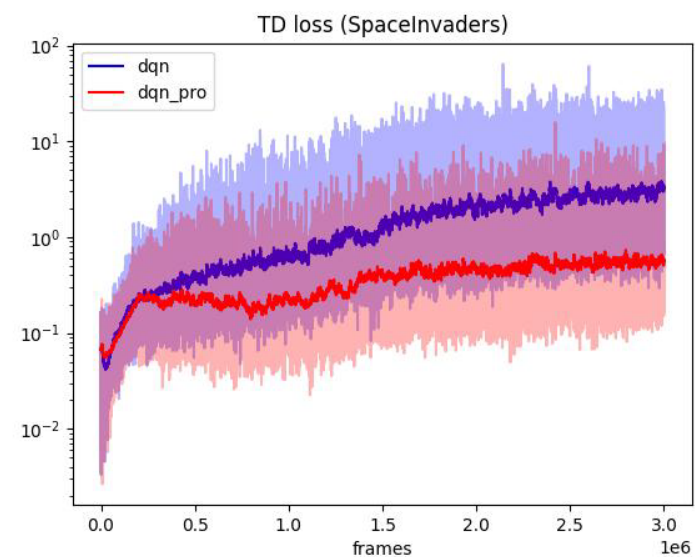
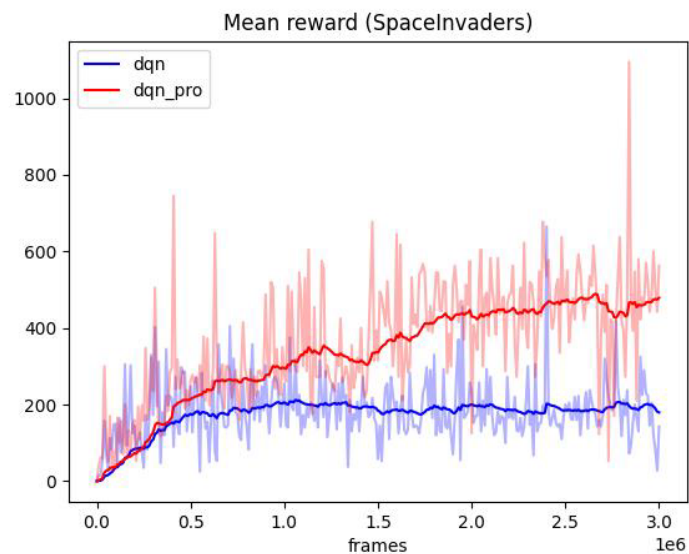
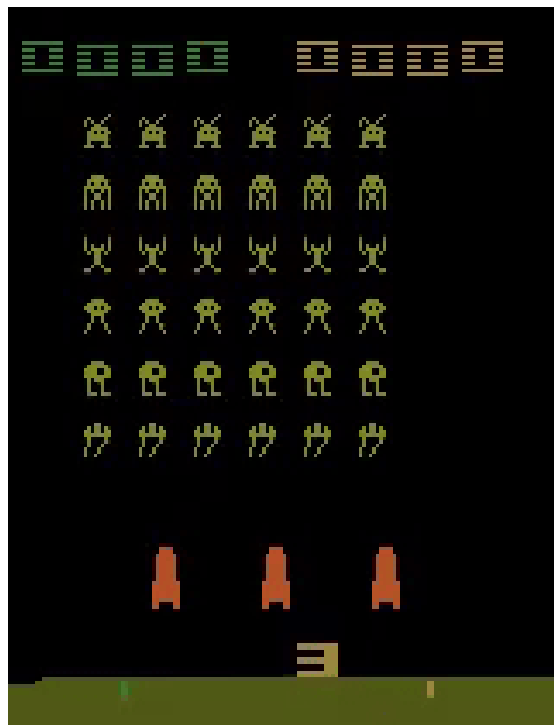


Финальное преимущество DQN Pro над DQN (вверху) и Rainbow Pro над Rainbow (внизу)

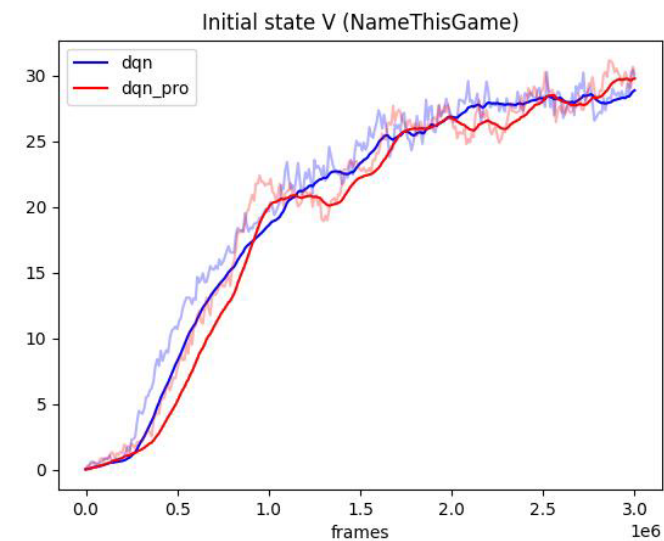
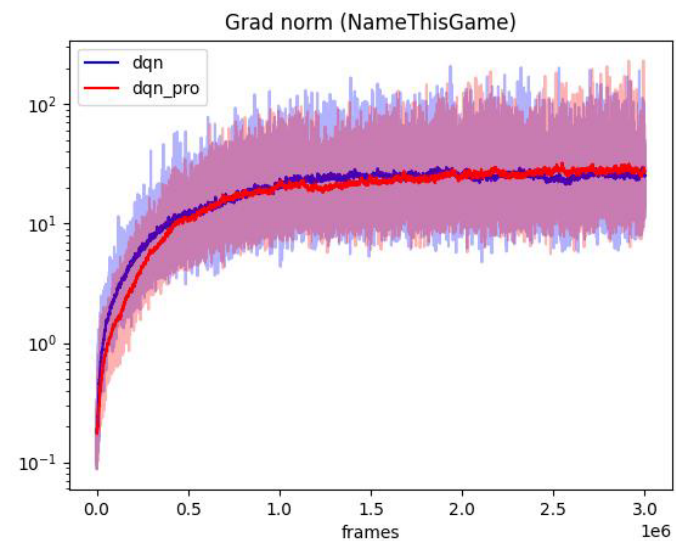
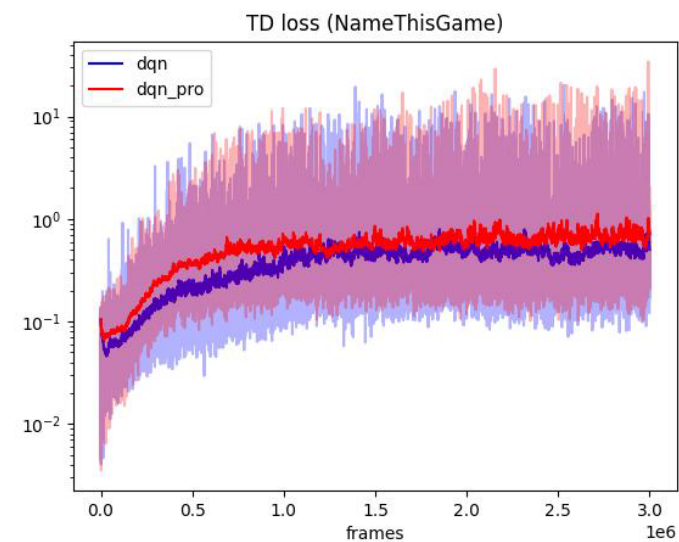
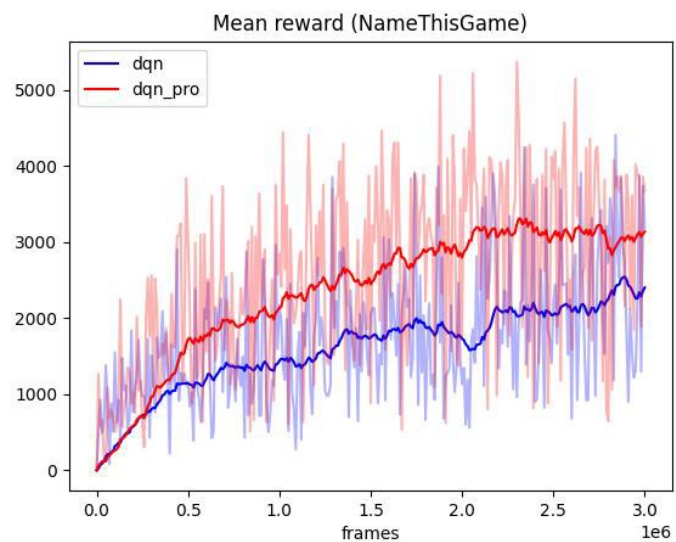
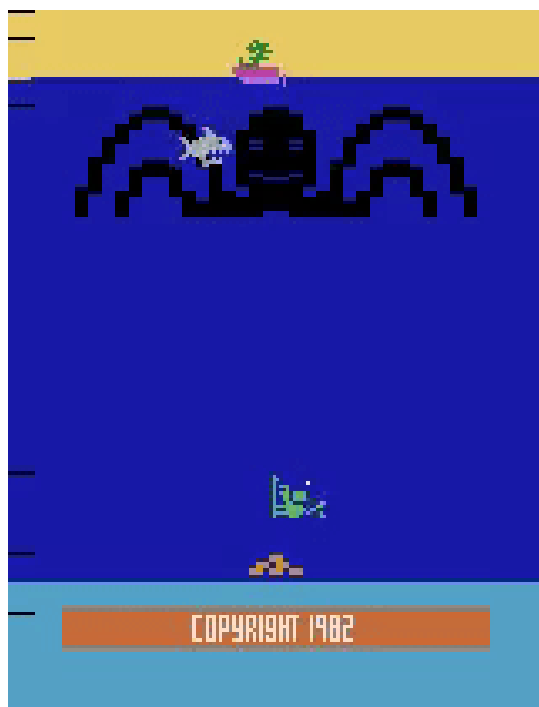


DQN Pro – красный, DQN - синий

SpaceInvaders



NameThisGame



Выводы

- DQN Pro позволяет добиться заметно большей награды по сравнению с DQN (но для этого по-прежнему нужно долго его обучать)
- Менее склонен к переоценке значения V -функции

Спасибо за внимание!