

REINFORCEMENT LEARNING

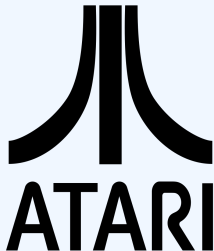
ATARI GAMES

ISABELL LEDERER

[GITHUB.COM/MATHEBELL/RL_ATARI_GAMES](https://github.com/mathebelle/RL_ATARI_GAMES)

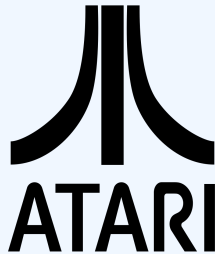
AKNUM MASCHINELLES LERNEN SE

WIEN, 18. DEZEMBER 2018

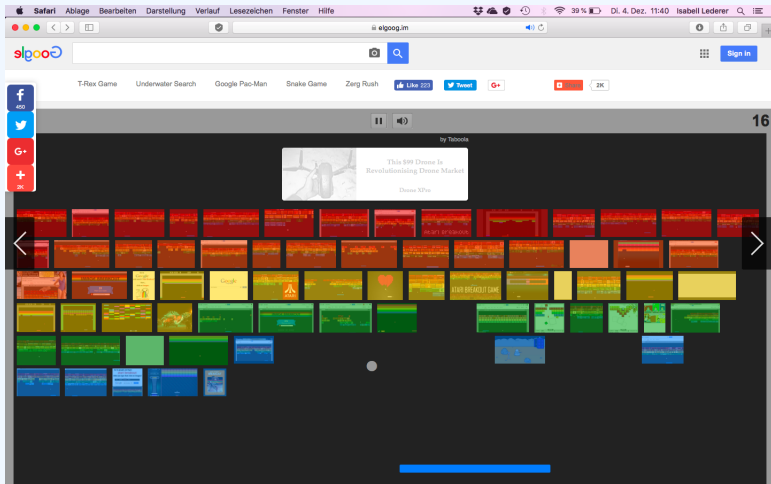


Contents:

1. Introduction
 - 1.1 Atari Games
 - 1.2 New Problem
2. The mathematics behind the algorithm
 - 2.1 Q-learning and DQN
 - 2.2 Policy Gradient Methods and Theorem
3. Implementation
 - 3.1 Preprocessing
 - 3.2 Model architecture
 - 3.3 Pseudocode DQN
4. Results
 - 4.1 Video clip
 - 4.2 Performance



ATARI GAMES

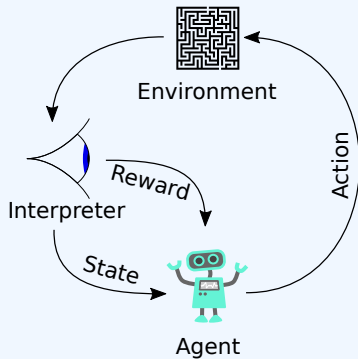


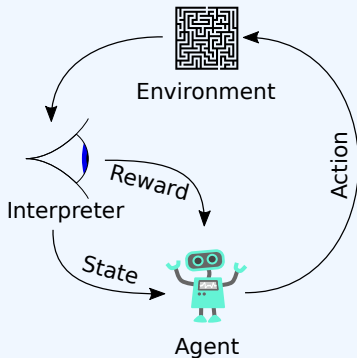
Breakout at <https://elgoog.im/breakout/>

SOME OTHER ATARI GAMES



REINFORCEMENT LEARNING - BASICS





- policy function:

$$\pi(a | s) = \Pr(A_t = a | S_t = s)$$

- new problem

action-value methods vs. policy gradient
methods

methods that approximate an action-value function followed by a policy

- perfect model of environment as Markov Decision Process (MDP)

- perfect model of environment as Markov Decision Process (MDP)
- dynamics function p :

$$p(s', r | s, a) = \Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$

- optimal state-value function v_* and action-value function q_* via Bellmann equations:

$$\begin{aligned}v_*(s) &= \max_a \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \\&= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_*(s')], \\q_*(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\&= \sum_{s', r} p(s', r \mid s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right].\end{aligned}$$

- characteristic for DP: perfect model as a MDP, approximate optimal value function, bootstrapping

- off-policy algorithm that can directly learn from raw experience without a model of the environment's dynamic

- off-policy algorithm that can directly learn from raw experience without a model of the environment's dynamic
- it has been shown that under some assumptions Q converges with probability 1 to q_*

- off-policy algorithm that can directly learn from raw experience without a model of the environment's dynamic
- it has been shown that under some assumptions Q converges with probability 1 to q_*
- Q-learning iteration:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)].$$

Q-learning Pseudocode

Parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose A by ϵ -greedy selection

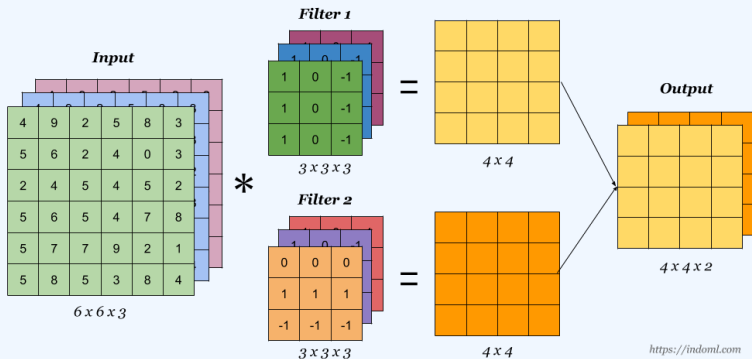
 Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

 until S is terminal

DQN - DEEP Q-NETWORK

- combines the idea of Q-learning with a deep convolutional network



- learn parametrized policy

$$\pi(a | s, \theta) = \Pr(A_t = a | S_t = s, \theta_t = \theta)$$

- learn parametrized policy

$$\pi(a | s, \theta) = \Pr(A_t = a | S_t = s, \theta_t = \theta)$$

- learn the parameter θ based on the gradient of some scalar performance measure $J(\theta)$

- learn parametrized policy

$$\pi(a | s, \theta) = \Pr(A_t = a | S_t = s, \theta_t = \theta)$$

- learn the parameter θ based on the gradient of some scalar performance measure $J(\theta)$
- updating via gradient ascent

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$$

- can be parametrized in any way as long as the policy is differentiable
- to ensure exploration we require that $\pi(a | s, \theta)$ never becomes deterministic
- advantage: can approach a deterministic policy

POLICY GRADIENT METHODS

- can be parametrized in any way as long as the policy is differentiable
- to ensure exploration we require that $\pi(a | s, \theta)$ never becomes deterministic
- advantage: can approach a deterministic policy
- one can proof better performance for policy gradient methods than for action value methods

- define performance as

$$J(\theta) := v_{\pi_\theta}(s_0).$$

- define performance as

$$J(\theta) := v_{\pi_\theta}(s_0).$$

- $\widehat{\nabla J(\theta_t)} = ???$

Theorem

Let $\pi(a | s, \theta)$ be a parameterized policy, $q_\pi(s, a)$ the action-value function under π and $\mu(s)$ the on-policy distribution under π , thus

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi(a | s, \theta).$$

POLICY GRADIENT THEOREM

Proof.

$$\nabla v_{\pi}(s) = \nabla \left[\sum_a \pi(a|s) q_{\pi}(s, a) \right] \quad (1)$$

$$= \sum_a \left[\nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \nabla q_{\pi}(s, a) \right] \quad (2)$$

$$= \sum_a \left[\nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \nabla \left[\sum_{s', r} p(s', r|s, a) (r + v_{\pi}(s')) \right] \right] \quad (3)$$

$$= \sum_a \left[\nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \sum_{s'} p(s'|s, a) \nabla v_{\pi}(s') \right] \quad (4)$$

POLICY GRADIENT THEOREM

$$(1) : v_{\pi}(s) = \sum_a \pi(s|a) q_{\pi}(s, a)$$

(1) \rightarrow (2) : product rule

$$(2) \rightarrow (3) : q_{\pi}(s, a) = \sum_{s', r} p(s', r|s, a) (r + v_{\pi}(s'))$$

$$(3) \rightarrow (4) : p(s'|s, a) = \sum_r p(s', r|s, a)$$

POLICY GRADIENT THEOREM

Proof.

$$\begin{aligned} &= \sum_a \left[\nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \sum_{s'} p(s'|s, a) \right. \\ &\quad \left. \sum_{a'} \left[\nabla \pi(a', s') q_\pi(s', a') + \pi(a'|s') \sum_{s''} p(s''|s', a') \nabla V_\pi(s'') \right] \right] \end{aligned} \quad (5)$$

$$= \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \Pr(s \rightarrow x, k, \pi) \sum_a \nabla \pi(a|x) q_\pi(x, a), \quad (6)$$

where $\Pr(s \rightarrow x, k, \pi)$ is the probability of transitioning from state s to state x in k steps under policy π .

$\eta(s)$ denote the number of time steps spent, on average, in s .

On-policy distribution is then $\mu(s) = \frac{\eta(s)}{\sum_{s'} \eta(s')}$.

POLICY GRADIENT THEOREM

Proof.

$$\nabla J(\theta) = \nabla v_{\pi}(s_0) \quad (7)$$

$$= \sum_s \left(\sum_{k=0}^{\infty} \Pr(s_0 \rightarrow s, k, \pi) \right) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \quad (8)$$

$$= \sum_s \eta(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \quad (9)$$

$$= \sum_{s'} \eta(s') \sum_s \frac{\eta(s)}{\sum_{s'} \eta(s')} \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \quad (10)$$

$$= \sum_{s'} \eta(s') \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \quad (11)$$

$$\propto \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \quad (12)$$

□





OpenAI

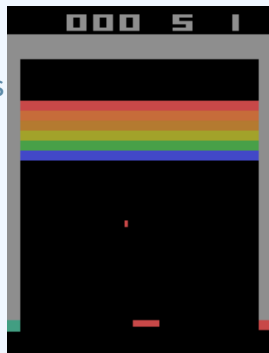
- OpenAI's Gym is a toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Pinball.



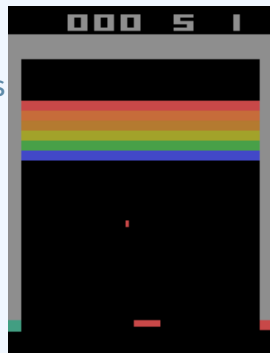
OpenAI

- OpenAI's Gym is a toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Pinball.
- Arcade Learning Environment (ALE)

- Input: 3 Atari frames in 210×160 pixels in RGB color



- Input: 3 Atari frames in 210×160 pixels in RGB color
- a lot of data for memory: memorize recent 1,000,000 of $(S_t, A_t, R_t, S_{t+1}, \text{isTerminal})$



- Rescaling, e.g. from 210×160 pixels to 84×84

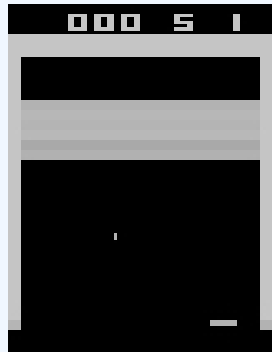
- Rescaling, e.g. from 210×160 pixels to 84×84
- Greyscaling



- Rescaling, e.g. from 210×160 pixels to 84×84
- Greyscaling



- Rescaling, e.g. from 210×160 pixels to 84×84
- Greyscaling



- adequate collection for memory I can sample on
- deque perfect for needs, but bad runtime
- implement own RingBuffer?

MODEL ARCHITECTURE OF DQN

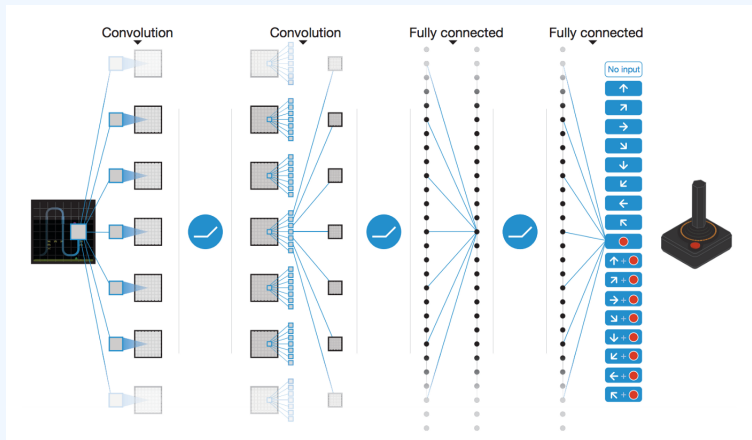
```
#Create Model
ATARI_SHAPE = (105, 80,3)

# Inputs
frames_input = keras.layers.Input(shape = ATARI_SHAPE, name = 'frames')
actions_input = keras.layers.Input(shape = (n_actions,), name = 'actions')

# "The first hidden layer convolves 16 8x8 filters with stride 4 with the input image
# and applies a rectifier nonlinearity."
conv_1 = keras.layers.convolutional.Conv2D(16, (8, 8), strides = 4, activation = 'relu')(frames_input)
# "The second hidden layer convolves 32 4x4 filters with stride 2, again followed by a
# rectifier nonlinearity."
conv_2 = keras.layers.convolutional.Conv2D(32,(4, 4),strides = 2, activation = 'relu')(conv_1)
# Flattening the second convolutional layer.
conv_flattened = keras.layers.core.Flatten()(conv_2)
# "The final hidden layer is fully-connected and consists of 256 rectifier units."
hidden = keras.layers.Dense(256, activation = 'relu')(conv_flattened)
# "The output layer is a fully-connected linear layer with a single output for each
# valid action."
output = keras.layers.Dense(n_actions)(hidden)
filtered_output = keras.layers.multiply([output, actions_input]) # get corresponding Q val

model = keras.models.Model(input = [frames_input, actions_input], output = filtered_output)
# parameters from paper
optimizer = keras.optimizers.RMSprop(lr = 0.00025, rho = 0.95, epsilon = 0.01)
model.compile(optimizer, loss='mse')
```

MODEL ARCHITECTURE OF DQN



Implementation with keras/ tensorflow

Pseudocode

```
env = gym.make('BreakoutDeterministic-v4')
```

```
Create model, initialize memory  $D$ 
```

```
Loop for each episode:
```

```
     $s$  = env.reset() and preprocess  $s$ 
```

```
    env.render()
```

```
    Loop for each step of episode:
```

```
        Choose  $\epsilon$ , choose with probability  $\epsilon$  a random action,  
        otherwise choose best action predicted by model
```

```
         $s', r, is\_done, info$  = env.step( $a$ )
```

```
        preprocess  $s'$  and append ( $s, a, s', is\_done$ ) to  $D$ 
```

```
        sample batch from  $D$ , predict  $Q_{t+1}$ , update  $Q_t$ 
```

```
        fit batch of ( $s, a$ ) with  $Q_t$  to model
```

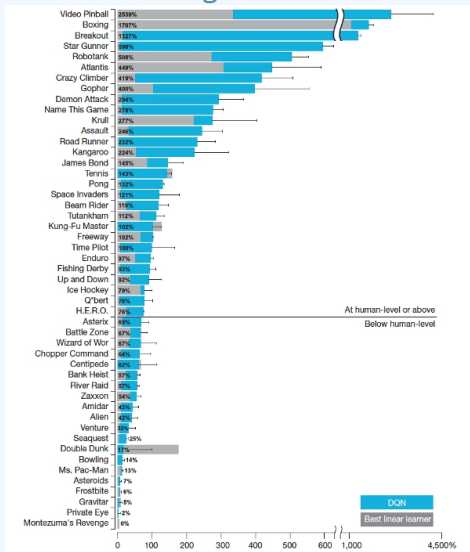
```
    env.render()
```

```
until  $S$  is terminal
```

Video of Playing Atari with Deep Reinforcement Learning:
<https://www.youtube.com/watch?v=V1eYniJ0Rnk>

RESULTS

Performance of DQN agent on various Atari games



QUESTIONS?



R. S. SUTTON, A. G. BARTO.

REINFORCEMENT LEARNING - AN INTRODUCTION.

The MIT Press, 2018.



GOOGLE DEEPMIND.

HUMAN-LEVEL CONTROL THROUGH DEEP REINFORCEMENT LEARNING.

Nature, doi:10.1038/nature14236, 2015.