

BITS F464: Machine Learning

29

Reinforcement Learning (RL)



Dr. Kamlesh Tiwari

Assistant Professor, Department of CSIS,
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

April 05, 2021

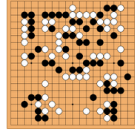
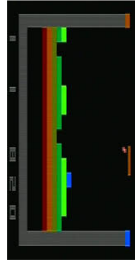
ONLINE (Campus @ BITS-Pilani Jan-May 2021)

<http://katiwari.in/ml>

Learning A Game

Consider playing game of Atari, Chess or Go.

- They all requires intuition, creative and strategic thinking



DeepBlue won 2 chess game over Gary Kasparov (3D+1L) in 1997
AlphaGo in March 2016, beat Lee Sedol a professional GO player.

What one want to learn?

- Rule of game, or
- Tactic/skill (policy) to play

Reinforcement Learning



- Reinforcement learning works on generate and test methodology
- An agent automatically enumerate solution and compute reward/penalty
- Task of agent is to learn from indirect, delayed reward, to choose sequences of actions that produce the greatest cumulative reward
- Agent may have no prior knowledge of the effects of its actions on the environment
- Neither supervised nor unsupervised, it is trial and error
- Reinforcement learning algorithms are related to dynamic programming
- Example includes learning to control a mobile robot, learning to optimize operations in factories, and learning to play board games

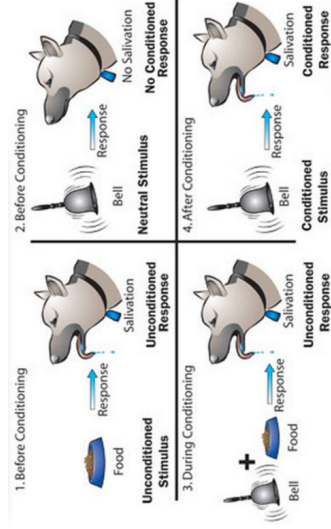
Trial and Error

- In 1956 **Arthur Samuel** wanted to develop a Checkers playing program that could beat him
- Idea was to let the computer play lot of games against itself and learn
- In 1962 the computer won over human player
- Father of ML



He defined ML as a field of study that gives computers the ability to learn without being explicitly programmed.

Reinforcement Learning

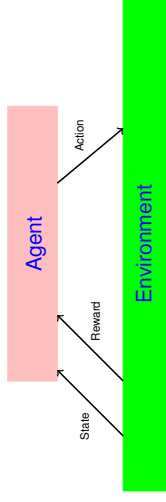


Rewards can lead to learning¹

¹ Image from: <http://www.krigolson.com/reinforcement-learning.html>

State, Action, and Reward

Problem of learning a control policy to maximize cumulative reward is very general



$$S_0 \xrightarrow{a_0} S_1 \xrightarrow{a_1} S_2 \xrightarrow{a_2} S_3 \xrightarrow{a_3} \dots$$

Goal is to choose a policy that maximizes

$$r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \dots$$

where $0 \leq \gamma < 1$

Reinforcement Learning

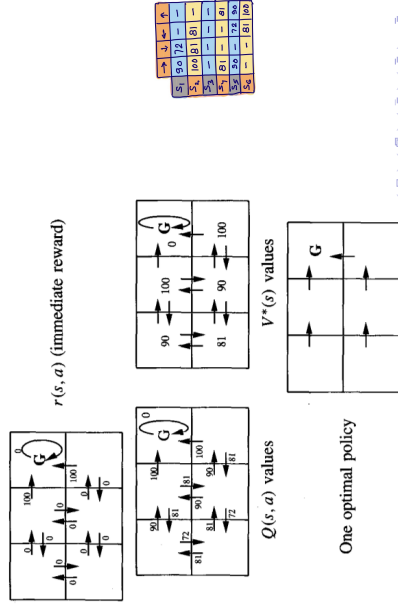
- Target function to be learned is a **policy** $\pi : S \rightarrow A$ where S and A are set of states and actions. Important points are
 - Delayed reward: $r(s, \pi(s))$ is unavailable
 - Exploration: agent influences the distribution of training examples
 - Partially observable states: forward camera cannot see behind
 - Life-long learning: using previous experience to learn new tasks
- Assumption is *Markov decision process* (MDP). Therefore, reward $r_t = r(s_t, a_t)$ and next state $s_{t+1} = \delta(s_t, a_t)$
- Discounted cumulative reward** of a policy π from initial state s_1 is

$$V^\pi(s_1) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$$

here $0 \leq \gamma < 1$ is a constant determining relative value of delayed versus immediate rewards

An Example

Let $\gamma = 0.9$



Q Learning

- One can choose globally optimal action sequences by reacting repeatedly to the local values of Q for the current state
- Value of Q can be found through iterative approximation
- As $V^*(s) = \max_a Q(s, a)$ it is easy to see

$$Q(s, a) = r(s, a) + \gamma \max_{a'} Q(\delta(s, a), a')$$

- In terms of iterative approximation

$$\hat{Q}(s, a) = r + \gamma \max_{a'} \hat{Q}(s', a')$$

where immediate reward $r = r(s, a)$ and $s' = \delta(s, a)$

Other Reward Function

- Finite horizon reward

$$\sum_{i=0}^h r_{t+i}$$
- Average reward

$$\lim_{h \rightarrow \infty} \frac{1}{h} \sum_{i=0}^h r_{t+i}$$

Agent's learning task is

to determine optimal policy

$$\pi^* = \operatorname{argmax}_{\pi} V^{\pi}(s), \forall s$$

$V^{\pi^*}(s)$ is also represented as $V^*(s)$ for simplicity

Q Learning

- Specifies a method for an agent to learn an optimal policy
- Agent should prefer state s_1 over state s_2 when $V^*(s_1) > V^*(s_2)$
- Optimal action in state s is

$$\pi^*(s) = \operatorname{argmax}_a (r(s, a) + \gamma V^*(\delta(s, a)))$$

problem is that r and δ are not known

- Define a function $Q(s, a)$ representing maximum discounted cumulative reward that can be achieved starting from state s and applying action a as the first action

$$Q(s, a) = r(s, a) + \gamma V^*(\delta(s, a))$$

- Then

$$\pi^*(s) = \operatorname{argmax}_a Q(s, a)$$

Q Learning Algorithm

Algorithm for Q Learning can be specified as below

Algorithm 1: Q Learning

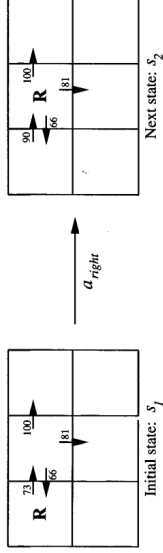
- For each s and a pair: **initialize** table entry $\hat{Q}(s, a) = 0$
- Observe the current state s
- while True do**
- Select an action a and **execute**
- Receive immediate response/reward r
- Observe **new state** s'
- Update** the $\hat{Q}(s, a) = r + \gamma \max_{a'} \hat{Q}(s', a')$
- $s = s'$

S	a	1	2
1	1	0	0
1	2	0	0
2	1	0	0
2	2	0	0

- Produces the, Q Table
- Guess expected running time?

Q Learning Example

Consider a single action for the example below (let $\gamma = 0.9$)



$$\begin{aligned} \hat{Q}(s_1, a_{right}) &= r + \gamma \max_a \hat{Q}(s_2, a) \\ &= 0 + 0.9 \max_a \{66, 81, 100\} \\ &= 90 \end{aligned}$$

Does it converge (contd)

- Maximum error in \hat{Q}

$$\Delta_n = \max_{s,a} |\hat{Q}(s, a) - Q(s, a)|$$

is maximum cell wise difference ³

	→	↓	←	↑
s_1	90	72	-	-
s_2	100	81	81	-
s_3	-	-	-	-
s_4	81	-	-	81
s_5	90	-	72	90
s_6	-	-	81	100

	→	↓	←	↑
s_1	81	55	-	-
s_2	98	71	31	-
s_3	-	-	-	-
s_4	72	-	-	41
s_5	58	-	66	80
s_6	-	-	56	72

³ $a_n = 50$ in the above example

Experimentation Strategies

- Can we select the action that maximizes $\hat{Q}(s, a)$
NO, as one **have to visit each state infinitely often** for convergence
- For some $k > 0$

$$Pr(a_i | s) = \frac{k \hat{Q}(s, a_i)}{\sum_j k \hat{Q}(s, a_j)}$$

Does it converge?

- Yes, under certain assumptions
 - Deterministic MDP
 - Rewards are bounded $|r(s, a)| < c$
 - Agent visits every possible state-action pair infinitely often
- Note that $\hat{Q}(s, a)$ does never decreases in any iteration

$$\hat{Q}_n(s, a) \leq \hat{Q}_{n+1}(s, a)$$

- And is upper bounded by $Q(s, a)$ i.e.

$$0 \leq \hat{Q}_n(s, a) \leq Q(s, a)$$
- Consider consecutive intervals during which each state-action transition occurs at least once²

² Random Walk takes $n \log n$ time to traverse fully connected graph with n nodes

Does it converge (contd)

- Let Δ_n be the maximum error in \hat{Q} i.e.

$$\Delta_n = \max_{s,a} |\hat{Q}(s, a) - Q(s, a)|$$

- For any $\hat{Q}_n(s, a)$ that is updated on iteration $n + 1$, the magnitude of the error in the revised estimate $\hat{Q}_{n+1}(s, a)$ is

$$\begin{aligned} |\hat{Q}_{n+1}(s, a) - Q(s, a)| &= |(r + \gamma \max_{a'} \hat{Q}_n(s', a')) - (r + \gamma \max_{a'} Q(s', a))| \\ &= \gamma |\max_{a'} \hat{Q}_n(s', a') - \max_{a'} Q(s', a)| \\ &\leq \gamma \max_{s', a'} |\hat{Q}_n(s', a') - Q(s', a)| \\ &\leq \gamma \max_{s', a'} |\hat{Q}_n(s'', a) - Q(s'', a)| \\ &= \gamma \Delta_n \end{aligned}$$

- Starting from Δ_0 error reduces to $\gamma \Delta_0, \gamma^2 \Delta_0, \dots, \gamma^n \Delta_0, \dots, 0$

Therefore, \hat{Q} converges to Q after infinitely many steps ...

Nondeterministic Rewards and Actions

- Reward and transition function may have probabilistic nature
- Policy π is then based on expected value. Discounted cumulative reward is taken as $V^\pi(s_t) = E[\sum_{j=0}^{\infty} \gamma^j r_{t+j}]$
- Definition of Q becomes

$$\begin{aligned} Q(s, a) &= E[r(s, a) + \gamma V^*(\delta(s, a))] = E[r(s, a)] + \gamma E[V^*(\delta(s, a))] \\ &= E[r(s, a)] + \gamma \sum_{s'} P(s' | s, a) V^*(s') \\ &= E[r(s, a)] + \gamma \sum_{s'} P(s' | s, a) \max_{a'} Q(s', a) \end{aligned}$$

- Earlier training rule (for the deterministic) fails to converge here. Therefore, decaying weight average is taken as

$$\hat{Q}_n(s, a) = (1 - \alpha_n) \hat{Q}_{n-1}(s, a) + \alpha_n [r + \gamma \max_{a'} \hat{Q}_{n-1}(s', a)]$$

where $\alpha_n = \frac{1}{1 + \text{visits}_n(s, a)}$ where $\text{visits}_n(s, a)$ is number of times (s, a) is visited till n th iteration.

There is a problem

- The method is not scalable.
- As one have to compute $Q(s, a)$ for every pair and infinitely many times.
- It is huge problem as the state space for typical games is large.

What is the solution?

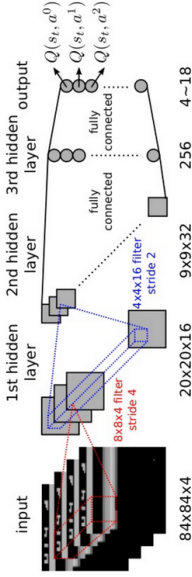
We can use a function approximator, that would estimate the action value function $Q(s, a)$.

- One such an approximator could be a Neural Network.
- When it is a Deep Neural Network, we call it a deep Q-learning.⁴

⁴ <https://arxiv.org/pdf/1511.07289.pdf> (Mnih, Heess, Silver, Kavukcuoglu, Wainwright, Schaul, Peters, and Schmidhuber, 2015)

Deep Reinforcement Learning for Atari Game

The network that we are going to use is this⁵



- Preprocessing: Grayscale conversion and scale down
- Network takes a stack of four frames as an input to exploit spatial relationships in images
- Convolution layer uses ELU (exponential linear unit) activation⁶
- Experience Replay is used to make more efficient use of observed experience

⁵ <https://leonardobarbosa.com/2015/07/01/identical-intelligence-content-deep-ql-learning.html>

⁶ <https://arxiv.org/pdf/1511.07289.pdf>

Thank You!

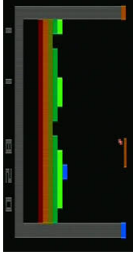
Thank you very much for your attention!

Queries ?

(Reference⁷)

⁷ Book - *Machine Learning*, ch-9/13, Tom Mitchell.

Deep Reinforcement Learning for Atari Game



- Objective is to complete the game with highest score
- States are the raw pixel values in the game screen
- Actions are moving controls left/right/up/down
- Reward is the increase/decrease of score in each timestamps

Deep Q-Learning with Experience Replay

Algorithm 1 Deep Q-learning with Experience Replay

```
Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1,  $M$  do
  Initialise sequence  $s_1 = \{s_t\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$ 
  for  $t = 1, T$  do
    With probability  $\epsilon$  select a random action  $a_t$ 
    otherwise select  $a_t = \arg \max_a Q^*(\phi(s_t), a; \theta)$ 
    Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
    Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
    Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ 
    Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ 
    Set  $\theta_j = \begin{cases} r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a_j; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ 
    Perform a gradient descent step on  $(\theta_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3
  end for
end for
```