

Sapienza University of Rome

Master in Artificial Intelligence and Robotics  
Master in Engineering in Computer Science

## Machine Learning

A.Y. 2019/2020

Prof. L. Iocchi, F. Patrizi, V. Ntouskos

## 17. Markov Decision Processes and Reinforcement Learning

L. Iocchi, F. Patrizi, V. Ntouskos

# Overview

- Dynamic systems
- Markov Decision Processes
- Q-Learning
- Experimentation strategies
- Evaluation

## References

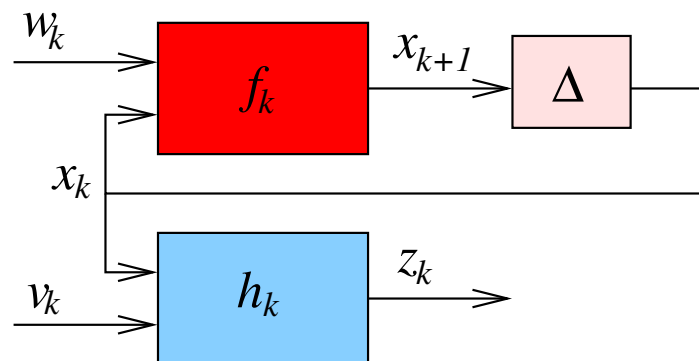
Richard S. Sutton and Andrew G. Barto.

Reinforcement Learning: An Introduction (2nd edition).

On-line: <http://incompleteideas.net/book/the-book.html>

## Dynamic System

The classical view of a dynamic system



$x$ : state

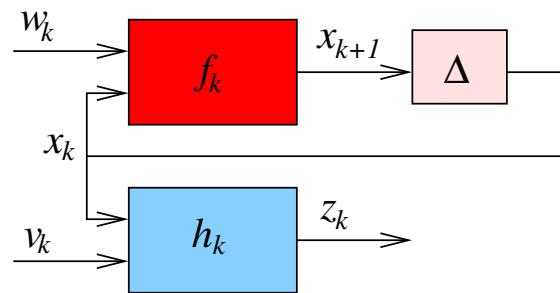
$z$ : observations

$w, v$ : noise

$f$ : state transition model

$h$ : observation model

# Reasoning vs. Learning in Dynamic Systems



**Reasoning:** given the model  $(f, h)$  and the current state  $x_k$ , predict the future  $(x_{k+T}, z_{k+T})$ .

**Learning:** given past experience  $(z_{0:k})$ , determine the model  $(f, h)$ .

## State of a Dynamic System

The state  $x$  encodes:

- all the past knowledge needed to predict the future
- the knowledge gathered through operation
- the knowledge needed to pursue the goal

Examples:

- configuration of a board game
- configuration of robot devices
- screenshot of a video-game

## Observability of the state

When the state is fully observable, the decision making problem for an agent is to decide which *action* must be executed in a given *state*.

The agent has to compute the function

$$\pi : \mathbf{X} \mapsto A$$

When the model of the system is not known, the agent has to *learn* the function  $\pi$ .

## Supervised vs. Reinforcement Learning

### Supervised Learning

Learning a function  $f : X \rightarrow Y$ , given

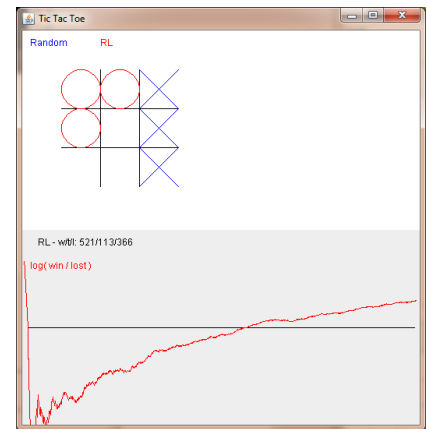
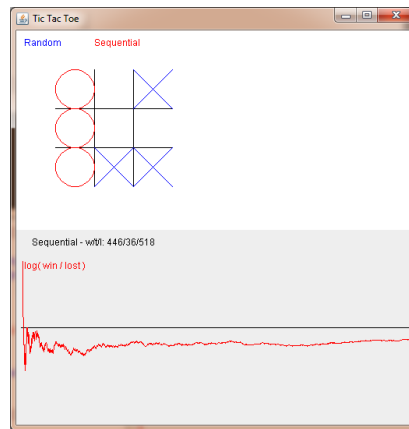
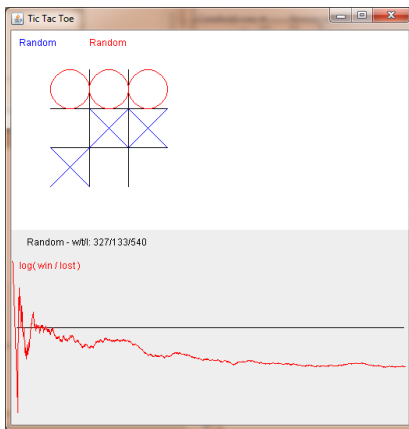
$$D = \{\langle x_i, y_i \rangle\}$$

### Reinforcement Learning

Learning a behavior function  $\pi : \mathbf{X} \rightarrow A$ , given

$$D = \{\langle \mathbf{x}_1, a_1, r_1, \dots, \mathbf{x}_n, a_n, r_n \rangle^{(i)}\}$$

# RL Example: Tic Tac Toe



# RL Example: Humanoid Walk

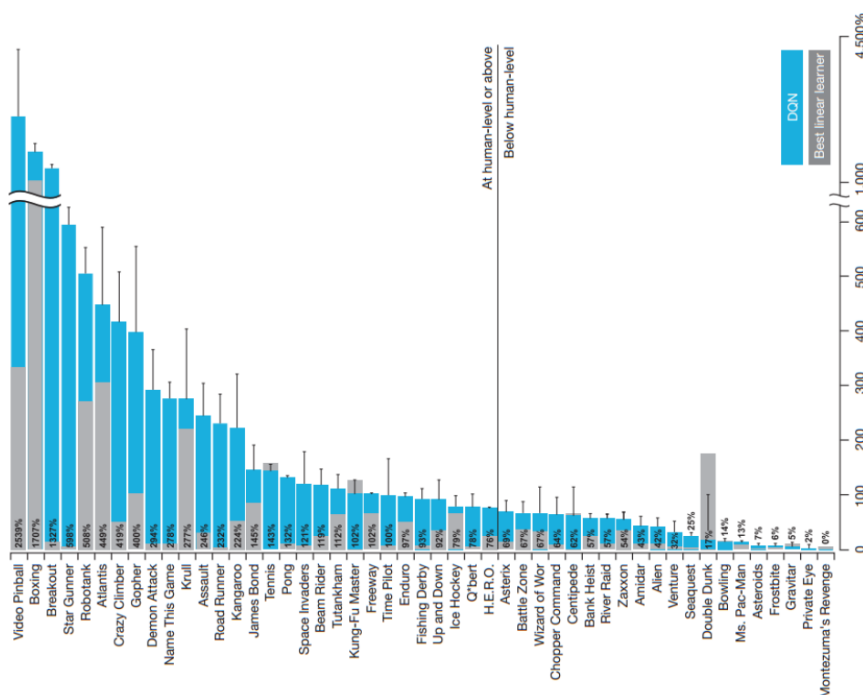


# RL Example: Controlling an Helicopter



## Deep Reinforcement Learning: ATARI games

Deep Q-Network (DQN) by Deep Mind (Google)



Human-level control through Deep Reinforcement Learning, Nature (2015)

# Dynamic System Representation

**X**: set of states

- explicit discrete and finite representation  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
- continuous representation  $\mathbf{X} = F(\dots)$  (state function)
- probabilistic representation  $P(\mathbf{X})$  (probabilistic state function)

**A**: set of actions

- explicit discrete and finite representation  $\mathbf{A} = \{a_1, \dots, a_m\}$
- continuous representation  $\mathbf{A} = U(\dots)$  (control function)

$\delta$ : transition function

- deterministic / non-deterministic / probabilistic

**Z**: set of observations

- explicit discrete and finite representation  $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_k\}$
- continuous representation  $\mathbf{Z} = Z(\dots)$  (observation function)
- probabilistic representation  $P(\mathbf{Z})$  (probabilistic observation function)

## Markov property

*Markov property*

- Once the current state is known, the evolution of the dynamic system does not depend on the history of states, actions and observations.
- The current state contains all the information needed to predict the future.
- Future states are conditionally independent of past states and past observations given the current state.
- The knowledge about the current state makes past, present and future observations statistically independent.

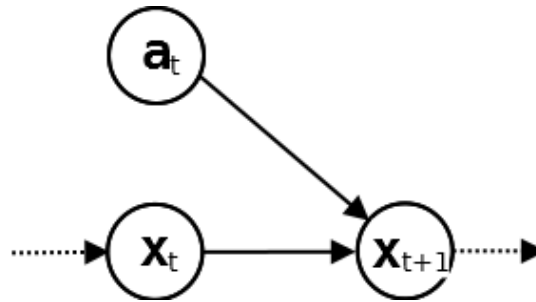
*Markov process* is a process that has the Markov property.

# Markov Decision Processes (MDP)

Markov processes for decision making.

States are fully observable, no need of observations.

Graphical model



# Markov Decision Processes (MDP)

Deterministic transitions

$$MDP = \langle \mathbf{X}, \mathbf{A}, \delta, r \rangle$$

- $\mathbf{X}$  is a finite set of states
- $\mathbf{A}$  is a finite set of actions
- $\delta : \mathbf{X} \times \mathbf{A} \rightarrow \mathbf{X}$  is a transition function
- $r : \mathbf{X} \times \mathbf{A} \rightarrow \mathbb{R}$  is a reward function

Markov property:  $\mathbf{x}_{t+1} = \delta(\mathbf{x}_t, a_t)$  and  $r_t = r(\mathbf{x}_t, a_t)$

Sometimes, the reward function is defined as  $r : \mathbf{X} \rightarrow \mathbb{R}$



# Markov Decision Processes (MDP)

## Non-deterministic transitions

$$MDP = \langle \mathbf{X}, \mathbf{A}, \delta, r \rangle$$

- $\mathbf{X}$  is a finite set of states
- $\mathbf{A}$  is a finite set of actions
- $\delta : \mathbf{X} \times \mathbf{A} \rightarrow 2^{\mathbf{X}}$  is a transition function
- $r : \mathbf{X} \times \mathbf{A} \times \mathbf{X} \rightarrow \mathbb{R}$  is a reward function

# Markov Decision Processes (MDP)

## Stochastic transitions

$$MDP = \langle \mathbf{X}, \mathbf{A}, \delta, r \rangle$$

- $\mathbf{X}$  is a finite set of states
- $\mathbf{A}$  is a finite set of actions
- $P(\mathbf{x}'|\mathbf{x}, \mathbf{a})$  is a probability distribution over transitions
- $r : \mathbf{X} \times \mathbf{A} \times \mathbf{X} \rightarrow \mathbb{R}$  is a reward function

# Full Observability in MDP

States are fully observable.

In presence of non-deterministic or stochastic actions, the state resulting from the execution of an action is not known *before* the execution of the action, but it can be fully observed *after* its execution.

## MDP Solution Concept

Given an MDP, we want to find an optimal policy.

Policy is a function

$$\pi : \mathbf{X} \rightarrow \mathbf{A}$$

For each state  $\mathbf{x} \in \mathbf{X}$ ,  $\pi(\mathbf{x}) \in \mathbf{A}$  is the *optimal* action to be executed in such state.

*optimality* = maximize the cumulative reward.

Optimality is defined with respect to maximizing the (expected value of the) cumulative discounted reward.

$$V^\pi(\mathbf{x}_1) \equiv E[\bar{r}_1 + \gamma \bar{r}_2 + \gamma^2 \bar{r}_3 + \dots]$$

where  $\bar{r}_t = r(\mathbf{x}_t, a_t, \mathbf{x}_{t+1})$ ,  $a_t = \pi(\mathbf{x}_t)$ , and  $\gamma \in [0, 1]$  is the discount factor for future rewards.

Optimal policy:  $\pi^* \equiv \operatorname{argmax}_\pi V^\pi(\mathbf{x}), \forall \mathbf{x} \in \mathbf{X}$

## Value function

Deterministic case

$$V^\pi(\mathbf{x}) \equiv r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$$

Non-deterministic/stochastic case:

$$V^\pi(\mathbf{x}) \equiv E[r_1 + \gamma r_2 + \gamma^2 r_3 + \dots]$$

## Optimal policy

$\pi^*$  is an **optimal policy** iff for any other policy  $\pi$

$$V^{\pi^*}(\mathbf{x}) \geq V^{\pi}(\mathbf{x}), \forall \mathbf{x}$$

For infinite horizon problems, a stationary MDP always has an optimal stationary policy.

## Reasoning and Learning in MDP

Problem: MDP  $\langle \mathbf{X}, \mathbf{A}, \delta, r \rangle$

Solution: Policy  $\pi : \mathbf{X} \rightarrow \mathbf{A}$

If the MDP  $\langle \mathbf{X}, \mathbf{A}, \delta, r \rangle$  is completely known  $\rightarrow$  reasoning or planning

If the MDP  $\langle \mathbf{X}, \mathbf{A}, \delta, r \rangle$  is not completely known  $\rightarrow$  learning

Note: Simple examples of reasoning in MDP can be modeled as a search problem and solved using standard search algorithm (e.g.,  $A^*$ ).

# One-state Markov Decision Processes (MDP)

$$MDP = \langle \{\mathbf{x}_0\}, \mathbf{A}, \delta, r \rangle$$

- $\mathbf{x}_0$  unique state
- $\mathbf{A}$  finite set of actions
- $\delta(\mathbf{x}_0, a_i) = \mathbf{x}_0, \forall a_i \in \mathbf{A}$  transition function
- $r(\mathbf{x}_0, a_i, \mathbf{x}_0) = r(a_i)$  reward function

Optimal policy:  $\pi^*(\mathbf{x}_0) = a_i$

## Deterministic One-state MDP

If  $r(a_i)$  is **deterministic** and **known**, then

Optimal policy:  $\pi^*(\mathbf{x}_0) = \operatorname{argmax}_{a_i \in \mathbf{A}} r(a_i)$

What if reward function  $r$  is unknown?

## Deterministic One-state MDP

If  $r(a_i)$  is **deterministic** and **unknown**, then

*Algorithm:*

- ① for each  $a_i \in \mathbf{A}$ 
  - **execute**  $a_i$  and **collect** reward  $r_{(i)}$
- ② Optimal policy:  $\pi^*(\mathbf{x}_0) = a_i$ , with  $i = \operatorname{argmax}_{i=1 \dots |\mathbf{A}|} r_{(i)}$

*Note:* exactly  $|\mathbf{A}|$  iterations are needed.

## Non-Deterministic One-state MDP

If  $r(a_i)$  is **non-deterministic** and **known**, then

Optimal policy:  $\pi^*(\mathbf{x}_0) = \operatorname{argmax}_{a_i \in \mathbf{A}} E[r(a_i)]$

*Example:*

If  $r(a_i) = \mathcal{N}(\mu_i, \sigma_i)$ , then

$\pi^*(\mathbf{x}_0) = a_i$ , with  $i = \operatorname{argmax}_{i=1 \dots |\mathbf{A}|} \mu_i$

# Non-Deterministic One-state MDP

If  $r(a_i)$  is **non-deterministic** and **unknown**, then

*Algorithm:*

- ① Initialize a data structure  $\Theta$
- ② For each time  $t = 1, \dots, T$  (until termination condition)
  - **choose** an action  $a_{(t)} \in \mathbf{A}$
  - **execute**  $a_{(t)}$  and **collect** reward  $r_{(t)}$
  - Update the data structure  $\Theta$
- ③ Optimal policy:  $\pi^*(\mathbf{x}_0) = \dots$ , according to the data structure  $\Theta$

*Note:* **many** iterations ( $T \gg |\mathbf{A}|$ ) are needed.

# Non-Deterministic One-state MDP

*Example:*

If  $r(a_i)$  is **non-deterministic** and **unknown** and  $r(a_i) = \mathcal{N}(\mu_i, \sigma_i)$ , then

*Algorithm:*

- ① Initialize  $\Theta_{(0)}[i] \leftarrow 0$  and  $c[i] \leftarrow 0$ ,  $i = 1 \dots |\mathbf{A}|$
- ② For each time  $t = 1, \dots, T$  (until termination condition)
  - **choose** an index  $\hat{i}$  for action  $a_{(t)} = a_{\hat{i}} \in \mathbf{A}$
  - **execute**  $a_{(t)}$  and **collect** reward  $r_{(t)}$
  - increment  $c[\hat{i}]$
  - update  $\Theta_{(t)}[\hat{i}] \leftarrow \frac{1}{c[\hat{i}]} (r_{(t)} + (c[\hat{i}] - 1)\Theta_{(t-1)}[\hat{i}])$
- ③ Optimal policy:  $\pi^*(\mathbf{x}_0) = a_i$ , with  $i = \operatorname{argmax}_{i=1 \dots |\mathbf{A}|} \Theta_{(T)}[i]$

# Learning with Markov Decision Processes

Given an agent accomplishing a task according to an MDP  $\langle \mathbf{X}, \mathbf{A}, \delta, r \rangle$ , for which functions  $\delta$  and  $r$  are **unknown** to the agent,

determine the optimal policy  $\pi^*$

Note: This is not a supervised learning approach!

- Target function is  $\pi : \mathbf{X} \rightarrow \mathbf{A}$
- but we do not have training examples  $\{(\mathbf{x}_{(i)}, \pi(\mathbf{x}_{(i)}))\}$
- training examples are in the form  $\langle (\mathbf{x}_{(1)}, a_{(1)}, r_{(1)}), \dots, (\mathbf{x}_{(t)}, a_{(t)}, r_{(t)}) \rangle$

## Agent's Learning Task

Since  $\delta$  and  $r$  are not known, the agent cannot predict the effect of its actions. But it can execute them and then observe the outcome.

The learning task is thus performed by repeating these steps:

- **choose** an action
- **execute** the chosen action
- **observe** the resulting new state
- **collect** the reward



# Approaches to Learning with MDP

- Value iteration  
(estimate the Value function and then compute  $\pi$ )
- Policy iteration  
(estimate directly  $\pi$ )

## Learning through value iteration

The agent could learn the value function  $V^{\pi^*}(\mathbf{x})$  (written as  $V^*(\mathbf{x})$ )

From which it could determine the optimal policy:

$$\pi^*(\mathbf{x}) = \operatorname{argmax}_{a \in \mathbf{A}} [r(\mathbf{x}, a) + \gamma V^*(\delta(\mathbf{x}, a))]$$

However, this policy cannot be computed in this way because  $\delta$  and  $r$  are not known.

## Q Function

$Q^\pi(\mathbf{x}, a)$ : expected value when executing  $a$  in the state  $\mathbf{x}$  and then act according to  $\pi$ .

$$Q^\pi(\mathbf{x}, a) \equiv r(\mathbf{x}, a) + \gamma V^\pi(\delta(\mathbf{x}, a))$$

$$Q^\pi(\mathbf{x}, a) \equiv \sum_{\mathbf{x}'} P(\mathbf{x}'|\mathbf{x}, a)(r(\mathbf{x}, a, \mathbf{x}') + \gamma V^\pi(\mathbf{x}'))$$

If the agent learns  $Q$ , then it can determine the optimal policy without knowing  $\delta$  and  $r$ .

$$\pi^*(\mathbf{x}) = \operatorname{argmax}_{a \in \mathbf{A}} Q(\mathbf{x}, a)$$

## Training Rule to Learn $Q$

Deterministic case:

$$Q(\mathbf{x}_t, a_t) = r(\mathbf{x}_t, a_t) + \gamma \max_{a' \in \mathbf{A}} Q(\mathbf{x}_{t+1}, a')$$

Let  $\hat{Q}$  denote learner's current approximation to  $Q$ .

Training rule:

$$\hat{Q}(\mathbf{x}, a) \leftarrow \bar{r} + \gamma \max_{a'} \hat{Q}(\mathbf{x}', a')$$

where  $\bar{r}$  is the immediate reward and  $\mathbf{x}'$  is the state resulting from applying action  $a$  in state  $\mathbf{x}$ .

# Q Learning Algorithm for Deterministic MDPs

- ① for each  $\mathbf{x}, a$  initialize table entry  $\hat{Q}_{(0)}(\mathbf{x}, a) \leftarrow 0$
- ② observe current state  $\mathbf{x}$
- ③ for each time  $t = 1, \dots, T$  (until termination condition)
  - **choose** an action  $a$
  - **execute** the action  $a$
  - **observe** the new state  $\mathbf{x}'$
  - **collect** the immediate reward  $\bar{r}$
  - update the table entry for  $\hat{Q}(\mathbf{x}, a)$  as follows:

$$\hat{Q}_{(t)}(\mathbf{x}, a) \leftarrow \bar{r} + \gamma \max_{a' \in \mathbf{A}} \hat{Q}_{(t-1)}(\mathbf{x}', a')$$

- $\mathbf{x} \leftarrow \mathbf{x}'$

- ④ Optimal policy:  $\pi^*(\mathbf{x}) = \operatorname{argmax}_{a \in \mathbf{A}} \hat{Q}_{(T)}(\mathbf{x}, a)$

Note: not using  $\delta$  and  $r$ , but just observing new state  $\mathbf{x}'$  and immediate reward  $\bar{r}$ , after the execution of the chosen action.

## Experimentation Strategies

How actions are chosen by the agents?

*Exploitation*: select action  $a$  that maximizes  $\hat{Q}(\mathbf{x}, a)$

*Exploration*: select random action  $a$  (with low value of  $\hat{Q}(\mathbf{x}, a)$ )

### $\epsilon$ -greedy strategy

Given,  $0 \leq \epsilon \leq 1$ ,

select a random action with probability  $\epsilon$

select the best action with probability  $1 - \epsilon$

$\epsilon$  can decrease over time (first exploration, then exploitation).

# Experimentation Strategies

## soft-max strategy

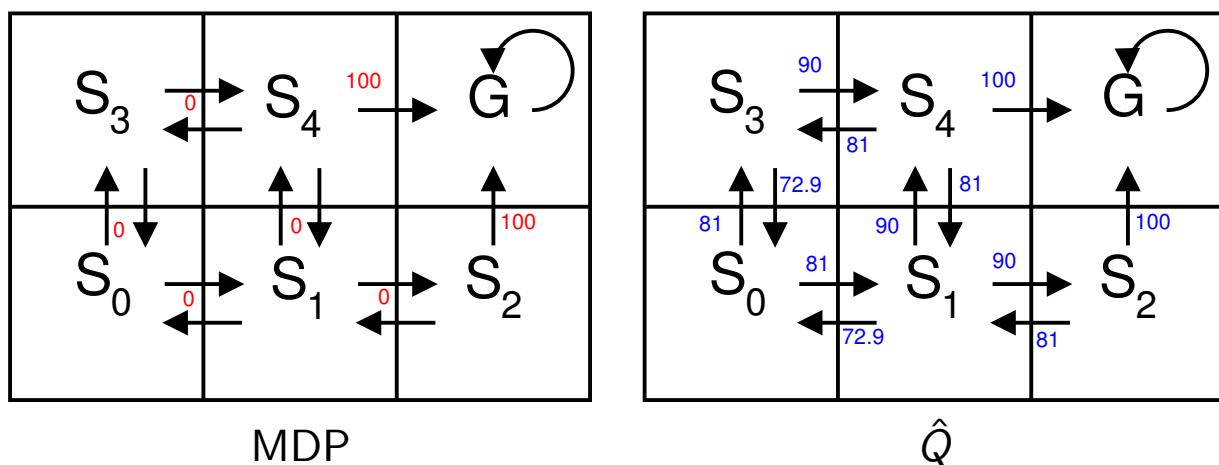
actions with higher  $\hat{Q}$  values are assigned higher probabilities, but every action is assigned a non-zero probability.

$$P(a_i|\mathbf{x}) = \frac{k^{\hat{Q}(\mathbf{x}, a_i)}}{\sum_j k^{\hat{Q}(\mathbf{x}, a_j)}}$$

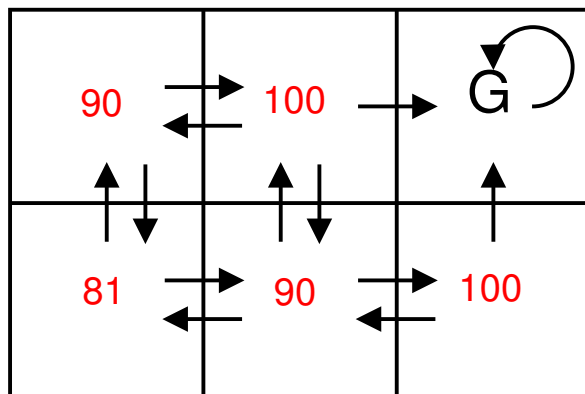
$k > 0$  determines how strongly the selection favors actions with high  $\hat{Q}$  values.

$k$  may increase over time (first exploration, then exploitation).

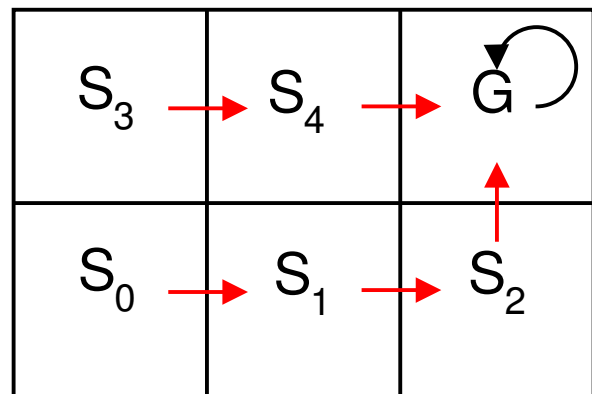
## Example: Grid World



## Example: Grid World

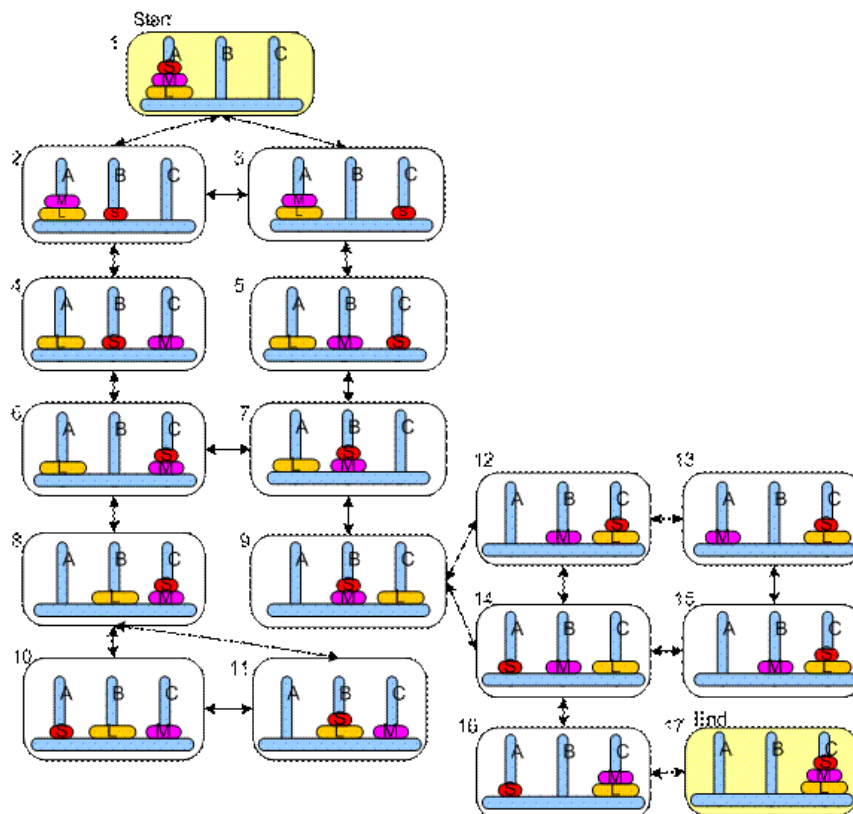


$V^*(\mathbf{x})$  values



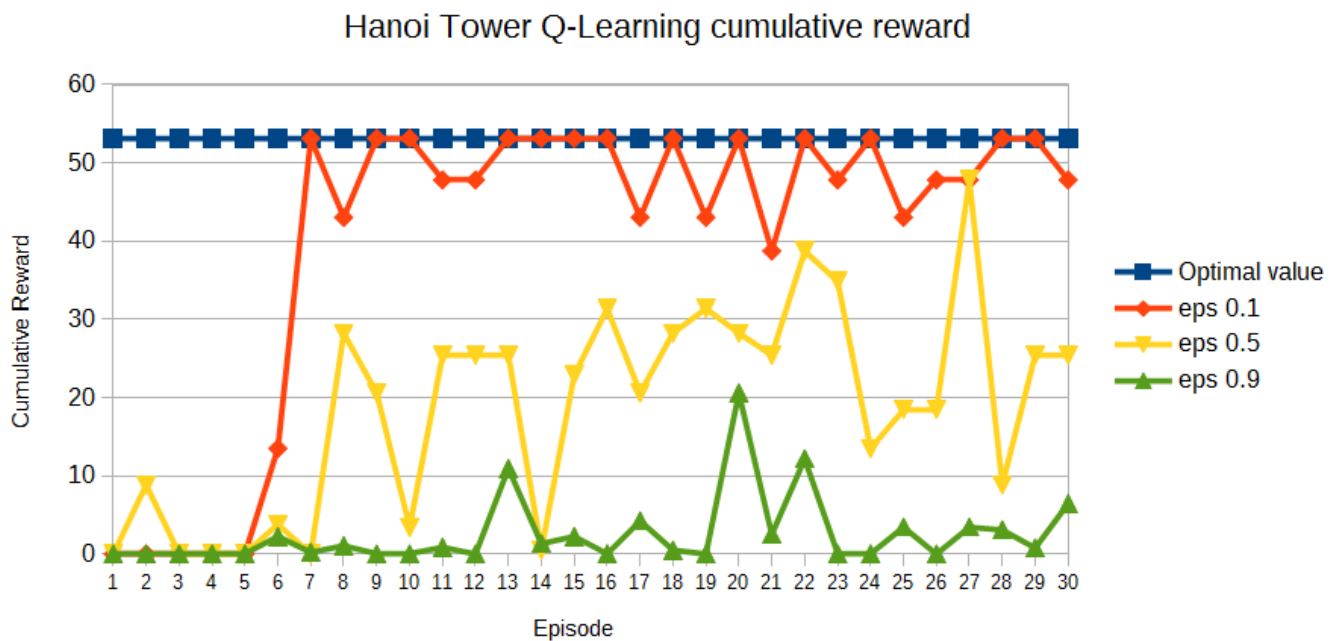
One optimal policy

## Example: Hanoi Tower



# Evaluating Reinforcement Learning Agents

Evaluation of RL agents is usually performed through the cumulative reward gained over time.



# Evaluating Reinforcement Learning Agents

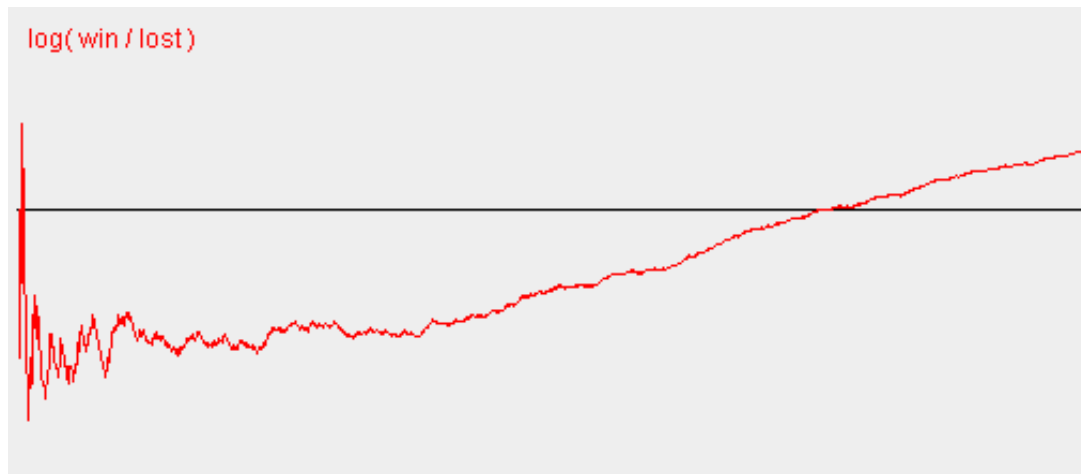
Cumulative reward plot may be very noisy (due to exploration phases).  
A better approach could be:

Repeat until termination condition:

- ① Execute  $k$  steps of learning
- ② Evaluate the current policy  $\pi_k$  (average and stddev of cumulative reward obtained in  $d$  runs with no exploration)

# Evaluating Reinforcement Learning Agents

Domain-specific performance metrics can also be used.



Average of all the results obtained during the learning process.