

# Planning and Optimal Control

## 6. Reinforcement Learning

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)  
Institute for Computer Science  
University of Hildesheim, Germany

# Syllabus

Tue. 24.10.	(1)	1. Markov Models
Tue. 31.10.	—	— <i>Luther Day</i> —
Tue. 7.11	(2)	2. Hidden Markov Models
Tue. 14.11.	(3)	2b. (ctd.)
Tue. 21.11.	(4)	3. State Space Models
Tue. 28.11.	(5)	3b. (ctd.)
Tue. 5.12.	(6)	4. Markov Random Fields
Tue. 12.12.	(7)	4b. (ctd.)
Tue. 19.12.	(8)	4c. (ctd.)
Tue. 26.12.	—	— <i>Christmas Break</i> —
Tue. 9.1.	(9)	5. Markov Decision Processes
Tue. 16.1.	(10)	5b. (ctd.)
Tue. 23.1.	(11)	6. Reinforcement Learning
Tue. 30.1.	(12)	6b. (ctd.)
Tue. 6.2.	(13)	7. Reinforcement Learning for Games
next year		8. Partially Observable Markov Decision Processes

# Outline

1. Introduction
2. Monte Carlo Methods
3. Learning the Value Function: TD(0)
4. Learning the Action Value Function: SARSA
5. Learning the Optimal Action Value Function: Q-Learning
6. Eligibility Traces
7. R-Learning
8. Model-Based Reinforcement Learning

# Outline

1. Introduction
2. Monte Carlo Methods
3. Learning the Value Function: TD(0)
4. Learning the Action Value Function: SARSA
5. Learning the Optimal Action Value Function: Q-Learning
6. Eligibility Traces
7. R-Learning
8. Model-Based Reinforcement Learning

# Policy Inference vs. Policy Learning

## ▶ **Markov Decision Problem:**

- ▶ transition model  $p$  and reward model  $r$  are known.
- ▶ compute an optimal policy  
(**optimal policy inference**).

## ▶ **Reinforcement Learning:**

- ▶ transition model  $p$  and reward model  $r$  are **not** known.
- ▶ learn an optimal policy from state/action/reward sequence data  
(**optimal policy learning**).

# Problems & Sampling MDP Data

## a. **state/action/reward Markov process learning:**

- ▶ learn the
    - ▶ **transition model**  $P^\pi$  of the state/action/reward Markov process or the
    - ▶ **value function**  $V^\pi$
- of an MDP under a fixed pre-existing policy  $\pi$ .
- ▶ called **generating policy**, **explorative policy** or **sampling policy**.
  - ▶  $\pi$  does not depend on past samples and current estimates.
  - ▶ e.g., choose actions in each state uniformly at random

## b. **MDP learning:**

- ▶ learn the **transition model**  $p$  and **reward model**  $r$  of an MDP.
- ▶ **passive sampling**: sample with a generating policy that
  - ▶ does not depend on past samples and current estimates, but should
  - ▶ guarantee to explore the whole state/action space.
- ▶ **active sampling**: use a generating policy that
  - ▶ selects actions s.t. informative samples are created.
    - e.g., uncertain transitions, uncertain rewards.

# Problems & Sampling MDP Data (2/2)

## c. **reinforcement learning:**

- ▶ learn an **optimal policy**  $\pi^*$ ,  
without knowing the MDP  $(p, r)$ .
  - ▶ i.e., a policy with maximal value function.
- ▶ the generating policy has to balance
  - ▶ **exploration** of actions with uncertain effects (on transitions and rewards) and
  - ▶ **exploitation** of actions with likely best valueto ensure focus on valuable actions.

## d. **joint process learning and control:**

- ▶ create state/action/reward sequences with maximal value.
  - ▶ i.e., learn and execute a policy that overall will lead to maximal value at the same time.
- ▶ such a policy also needs to balance exploration and exploitation to ensure it realizes a high value.

# Reinforcement Learning Approaches

## 1. **model-based reinforcement learning:**

- ▶ given state/action/reward sequences, learn the MDP model  $(p, r)$ ,
- ▶ afterwards use a policy optimization algorithm on the learnt model to find an optimal policy.

## 2. **direct reinforcement learning:**

- ▶ given state/action/reward sequences, learn the
  - ▶ optimal policy  $\pi^*$  or even only the
  - ▶ value function  $V^*$  or the
  - ▶ action value function  $Q^*$  of the optimal policy.



# Generating Policies

- ▶ **uniform at random policy:**

$$\pi_{\text{uniform}}(s, a) := \frac{1}{|A|}$$

- ▶ maximally explores all actions.

- ▶ **greedy policy:**

$$\pi_{\text{greedy}}(s; \hat{Q}) := \arg \max_{a \in A} \hat{Q}(s, a)$$

- ▶ maximally exploits current estimates; does not guarantee exploration.

- ▶  **$\epsilon$ -greedy policy:**

$$\pi(s, a; \hat{Q}, \epsilon) := (1 - \epsilon) \pi_{\text{greedy}}(s, a; \hat{Q}) + \epsilon \pi_{\text{uniform}}(s, a), \quad \epsilon \in [0, 1]$$

- ▶ **Boltzmann at random policy:**

$$\pi(s, a; \hat{Q}, \tau) := \frac{e^{\hat{Q}(s, a)/\tau}}{\sum_{a' \in A} e^{\hat{Q}(s, a')/\tau}}, \quad \tau \rightarrow 0$$

Note:  $\pi_{\text{greedy}}(s, a; \hat{Q}) = \mathbb{I}(a = \arg \max_{a' \in A} \hat{Q}(s, a'))$

# Outline

1. Introduction
- 2. Monte Carlo Methods**
3. Learning the Value Function: TD(0)
4. Learning the Action Value Function: SARSA
5. Learning the Optimal Action Value Function: Q-Learning
6. Eligibility Traces
7. R-Learning
8. Model-Based Reinforcement Learning

# Observed and Estimated Values

**data:**  $N$  state/action/reward sequences sampled from an MDP with unknown transition/reward model under a policy  $\pi$ :

$$\mathcal{D} := \{((s_{n,t}, a_{n,t}, r_{n,t})_{t=0:T_n-1}, s_{T_n}) \mid n = 1 : N\} \subseteq (\mathcal{S} \times \mathcal{A} \times \mathbb{R})^* \times \mathcal{S}$$

**observed values:**

$$V_{n,t} := \text{value}(r_{n,t:T_n-1}), \quad n = 1 : N, t = 0 : T_n - 1$$

e.g., for the discounted criterion:

$$:= \sum_{t'=t}^{T_n-1} \gamma^{t'-t} r_{n,t'}$$

**occurrences of state  $s \in \mathcal{S}$**  in data  $\mathcal{D}$ :

$$I_s := \{(n, t) \mid n = 1 : N, t = 0 : T_n - 1, s_{n,t} = s\}$$

**estimated value of state  $s \in \mathcal{S}$ :**

$$\hat{V}_s := \frac{1}{|I_s|} \sum_{(n,t) \in I_s} V_{n,t}$$

# Observed Values?

- ▶ in general, values cannot be observed.
  - ▶ as they may depend on infinitely many future rewards.
- ▶ to observe values, one needs to ensure/assume that all sequences terminate in a state of true value zero:

$$V_n, T_n = 0 \quad n = 1 : N$$

- ▶ e.g., if the state Markov chain has an absorbing state  $s$  with reward 0,

$$\forall s' \exists t : p(s_t = s \mid s_0 = s') = 1, \quad p(s \mid s) = 1, \quad \forall a : r(s, a) = 0$$

all sequences will be essentially finite (**terminal state**).

# Recursive Estimation of Values

observed state values:

$$V_{s,k} := V_{n_k, t_k} \quad \text{for } I_s = \{(n_1, t_1), (n_2, t_2), \dots, (n_K, t_K)\}$$

estimated state values from first  $k$  occurrences:

$$\begin{aligned}\hat{V}_{s,k} &:= \frac{1}{k} \sum_{j=1}^k V_{s,j} \\ &= \hat{V}_{s,k-1} + \frac{1}{k} (V_{s,k} - \hat{V}_{s,k-1})\end{aligned}$$

$$\hat{V}_{s,0} := 0$$

- ▶ does not have to store all values  $V_{s,k}$  for all  $k$ ,
- ▶ but just two values:  $\hat{V}_{s,k}$  (for current  $k$ ) and  $k$  itself.  
**(online update of the mean).**

# Recursive Estimation of Values

for

$$\hat{V}_{s,k} = \hat{V}_{s,k-1} + \frac{1}{k}(V_{s,k} - \hat{V}_{s,k-1})$$

and even with

$$\hat{V}_{s,k} = \hat{V}_{s,k-1} + \alpha_k(V_{s,k} - \hat{V}_{s,k-1}), \quad \alpha_k > 0, \alpha_k \rightarrow 0$$

we will get

$$\lim_{k \rightarrow \infty} \hat{V}_{s,k} = V^\pi(s)$$

# Unbiased Estimation from First Visit only

- ▶ if a state is visited twice in a sequence, two estimates of its value enter the overall mean (**all visits**).
  - ▶ which are not independent and introduce a bias.
- ▶ fix: retain only the first occurrence (**first visit**)
  - ▶ unbiased
  - ▶ experimentally said to yield somewhat smaller squared error.

## Monte Carlo method update rule:

$$\hat{V}_{s_t} = \hat{V}_{s_t} + \alpha(r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{T-t} r_T - \hat{V}_{s_t}), \quad t = 1 : T$$

- ▶ update values of all visited states after completion of each sequence  $s, a, r, T, \alpha := s_n, a_n, r_n, T_n, \alpha_n$ .
  - ▶ as it has to compute the observed value from all future rewards

# Outline

1. Introduction
2. Monte Carlo Methods
- 3. Learning the Value Function: TD(0)**
4. Learning the Action Value Function: SARSA
5. Learning the Optimal Action Value Function: Q-Learning
6. Eligibility Traces
7. R-Learning
8. Model-Based Reinforcement Learning



# Temporal Differences

Monte Carlo: update after completion of a sequence:

$$\hat{V}_{s_t} = \hat{V}_{s_t} + \alpha(r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{T-t} r_T - \hat{V}_{s_t}), \quad t = 1 : T$$

Temporal differences:

$$\begin{aligned} \hat{V}_{s_t} &= \hat{V}_{s_t} + \alpha(r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{T-t} r_T - \hat{V}_{s_t}) \\ &= \hat{V}_{s_t} + \alpha \left( \begin{aligned} &(r_t + \gamma \hat{V}_{s_{t+1}} - \hat{V}_{s_t}) \\ &+ \gamma (r_{t+1} + \gamma \hat{V}_{s_{t+2}} - \hat{V}_{s_{t+1}}) \\ &+ \gamma^2 (r_{t+2} + \gamma \hat{V}_{s_{t+3}} - \hat{V}_{s_{t+2}}) \\ &\vdots \\ &+ \gamma^{T-t} (r_{T-1} + \gamma \hat{V}_{s_T} - \hat{V}_{s_{T-1}}) \end{aligned} \right) \\ &= \hat{V}_{s_t} + \alpha(\delta_t + \gamma \delta_{t+1} + \gamma^2 \delta_{t+2} + \dots + \gamma^{T-t} \delta_{T-1}) \\ \delta_t &:= r_t + \gamma \hat{V}_{s_{t+1}} - \hat{V}_{s_t} \end{aligned}$$

Note: Still assuming the sequence terminates in a value 0 state:  $\hat{V}_{s_T} = V_{s_T} = 0$ .

## Temporal Difference update rule TD(0):

$$\hat{V}_{s_t} = \hat{V}_{s_t} + \alpha(r_t + \gamma \hat{V}_{s_{t+1}} - \hat{V}_{s_t})$$

- ▶ update values of a visited state immediately after action is taken, reward is received and next state has been determined.
- ▶ TD(0) estimates only the value function.
  - ▶ the value function of the generating policy.
  - ▶ if the transition model is known, one easily can estimate a policy from its value function.
  - ▶ but without access to the transition model, the value function alone does not allow to do so.

# Outline

1. Introduction
2. Monte Carlo Methods
3. Learning the Value Function: TD(0)
- 4. Learning the Action Value Function: SARSA**
5. Learning the Optimal Action Value Function: Q-Learning
6. Eligibility Traces
7. R-Learning
8. Model-Based Reinforcement Learning

# Action Value Function

Given

- ▶ an MDP  $(p, r)$ ,
- ▶ a value aggregation function  $\text{value}(r_1, r_2, \dots)$ 
  - ▶ e.g., for discounted criterion:  $\text{value}(r_1, r_2, \dots) := r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$
- ▶ a policy  $\pi$ :

**value function** (review):

$$V^\pi(s) := \mathbb{E}(\text{value}(r_0, r_1, \dots, r_t, \dots) \mid s_0 = s, \pi), \quad s \in \mathcal{S}$$

**action value function:**

$$Q^\pi(s, a) := \mathbb{E}(\text{value}(r_0, r_1, \dots, r_t, \dots) \mid s_0 = s, a_0 = a, \pi)$$

# Optimal Action Value Function

action value function:

$$Q^\pi(s, a) := \mathbb{E}(\text{value}(r_0, r_1, \dots, r_t, \dots) \mid s_0 = s, a_0 = a, \pi)$$

Bellman equation for optimal action value function:

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s'} p(s' \mid s, a) \max_{a'} Q^*(s', a'), \quad s \in S, a \in A$$

Optimal value function from optimal action value function:

$$V^*(s) = \max_a Q^*(s, a), \quad s \in S$$

Optimal policy from optimal action value function:

$$\pi^*(s) = \arg \max_a Q^*(s, a), \quad s \in S$$

## Temporal difference update rule for action value function (SARSA):

$$\hat{Q}_{s_t, a_t} = \hat{Q}_{s_t, a_t} + \alpha(r_t + \gamma \hat{Q}_{s_{t+1}, a_{t+1}} - \hat{Q}_{s_t, a_t})$$

- ▶ update values of a visited state after an action is taken, a reward is received and the next state and **next action has been determined**.
  - ▶ requires  $(s_t, a_t, r_t, s_{t+1}, a_{t+1})$ , therefore name SARSA.
- ▶ SARSA estimates the action value function, thus allows to infer the generating policy (“**on policy**”).

# SARSA Algorithm

1 **learn-action-value-discounted-sarsa**( $S, A, \gamma, s_{\text{term}}, N, \pi, q_0, \alpha$ ):

2  $\hat{Q} := (q_0)_{s \in S, a \in A}$

3 for  $n := 1, \dots, N$ :

4  $s' := \text{new\_process}()$

5  $a' := \pi(s')$

6 while  $s' \neq s_{\text{term}}$ :

7  $s := s'$

8  $a := a'$

9  $(r, s') := \text{execute\_action}(s, a)$

10  $a' := \pi(s')$

11  $\hat{Q}_{s,a} := \hat{Q}_{s,a} + \alpha(r + \gamma \hat{Q}_{s',a'} - \hat{Q}_{s,a})$

12 for  $s \in S$ :

13  $\pi_s^* := \arg \max_{a \in A} Q(s, a)$

14 return  $\hat{Q}, \pi^*$

where

- ▶  $s_{\text{term}}$  **terminal state** with zero reward.
- ▶  $\pi$  **generating policy**
- ▶  $q_0 \in \mathbb{R}$  initial value of all state/action pairs.
- ▶  $\alpha_k$  **learning rate** for update step  $k$ , e.g.,  $\alpha_k := 1/k$ .
- ▶ `new_process()` sets up a new process.
- ▶ `execute_action(s, a)` executes action  $a$  in process in state  $s$ .

# Outline

1. Introduction
2. Monte Carlo Methods
3. Learning the Value Function: TD(0)
4. Learning the Action Value Function: SARSA
5. Learning the Optimal Action Value Function: Q-Learning
6. Eligibility Traces
7. R-Learning
8. Model-Based Reinforcement Learning



# Update Rule

SARSA update rule:

$$\hat{Q}_{s_t, a_t} = \hat{Q}_{s_t, a_t} + \alpha(r_t + \gamma \hat{V}_{s_{t+1}, a_{t+1}} - \hat{V}_{s_t, a_t})$$

**Q-Learning update rule:**

$$\hat{Q}_{s_t, a_t} = \hat{Q}_{s_t, a_t} + \alpha(r_t + \gamma \max_a \hat{Q}_{s_{t+1}, a} - \hat{Q}_{s_t, a_t})$$

- ▶ does not require next action  $a_{t+1}$
- ▶ does not learn the action value function of the generating policy !  
 (“**off policy**”)
- ▶ learns the action value function of the policy that
  - ▶ decides in time  $t$  as the generating policy and
  - ▶ in time  $t + 1$  optimally,  
— based on current estimates of the action value function.

# Convergence

If

- ▶ there are finite many states and actions,
- ▶ the generating policy visits each state/action pair an infinite number of times,
- ▶ the learning rates are slowly diminishing  
( $\sum_k \alpha_k(s, a) = \infty$ ,  $\sum_k \alpha_k(s, a)^2 < \infty$ ) and
- ▶  $\gamma < 1$  or  
(if  $\gamma = 1$ ) there exists an absorbing state with zero reward for any policy.

then the Q-learning estimates converge almost surely to the optimal action value function  $Q^*$ .

Note: E.g., learning rates  $\alpha_k(s, a) = 1/k$ .

# Q-Learning

1 **learn-opt-policy-discounted-q-learning**( $S, A, \gamma, s_{\text{term}}, N, \pi, q_0, \alpha$ ):

2  $\hat{Q} := (q_0)_{s \in S, a \in A}$

3 for  $n := 1, \dots, N$ :

4  $s := \text{new\_process}()$

5 while  $s \neq s_{\text{term}}$ :

6  $a := \pi(s)$

7  $(r, s') := \text{execute\_action}(s, a)$

8  $\hat{Q}_{s,a} := \hat{Q}_{s,a} + \alpha_n(r + \gamma \max_{a' \in A} \hat{Q}_{s',a'} - \hat{Q}_{s,a})$

9  $s := s'$

10 for  $s \in S$ :

11  $\pi_s^* := \arg \max_{a \in A} Q(s, a)$

12 return  $\hat{Q}, \pi^*$

where

- ▶  $s_{\text{term}}$  **terminal state** with zero reward.
- ▶  $\pi$  **generating policy**
- ▶  $q_0 \in \mathbb{R}$  initial value of all state/action pairs.
- ▶  $\alpha_k$  **learning rate** for update step  $k$ , e.g.,  $\alpha_k := 1/k$ .
- ▶ `new_process()` sets up a new process.
- ▶ `execute_action(s, a)` executes action  $a$  in process in state  $s$ .

# Action Value Models $Q$

- ▶ from an ML perspective, the action value function  $Q$  can be understood as a **regression model**:

$$Q : S \times A \rightarrow \mathbb{R}$$

- ▶ with predictors  $s$  and  $a$
- ▶ for nominal states and actions:
  - ▶ constant model: just a number for every pair  $(s, a)$ .
  - ▶ a factorization model likely would make better use of the data?
- ▶ for **structured states and actions**:
  - ▶ a proper regression model for values regressed on state and action properties.
- ▶ as for using  $Q$  to choose actions,  $Q(s, a)$  for all competing actions  $a$  have to be computed, the model could be represented as:

$$Q : S \rightarrow \mathbb{R}^A$$

- ▶ with structured input and
- ▶ with structured output

# Action Value Models $Q$ / Learning

- ▶ Updating the action value model  $Q$ :

- ▶ constant model:

$$8 \quad \hat{Q}_{s,a} := \hat{Q}_{s,a} + \alpha_n (r + \gamma \max_{a' \in A} \hat{Q}_{s',a'} - \hat{Q}_{s,a})$$

- ▶ general model:

$$8 \quad \mathcal{D} := \mathcal{D} \cup \{(s, a, r + \gamma \max_{a' \in A} \hat{Q}_{s',a'})\}$$

$$9 \quad \hat{Q} := \text{update-model}(\hat{Q}, \mathcal{D})$$

- ▶ if the generating policy  $\pi$  depends on  $\hat{Q}$ , learning can be accelerated by discounting older data.
  - ▶ e.g., with a discounting case weight.
  - ▶ or simply be forgetting / removing.

# Outline

1. Introduction
2. Monte Carlo Methods
3. Learning the Value Function: TD(0)
4. Learning the Action Value Function: SARSA
5. Learning the Optimal Action Value Function: Q-Learning
- 6. Eligibility Traces**
7. R-Learning
8. Model-Based Reinforcement Learning

# SARSA Algorithm with Eligibility Traces: SARSA( $\lambda$ )

1 **learn-action-value-discounted-sarsa-eleg**( $S, A, \gamma, s_{\text{term}}, K, \pi, q_0, \alpha, \lambda$ ):

2  $\hat{Q} := (q_0)_{s \in S, a \in A}$

3 for  $k := 1, \dots, K$ :

4  $N(s, a) := 0 \quad \forall a \in A, s \in S$

5  $s' := \text{new\_process}()$

6  $a' := \pi(s')$

7 while  $s' \neq s_{\text{term}}$ :

8  $s := s'$

9  $a := a'$

10  $N(s, a) := N(s, a) + 1$

11  $N(\tilde{s}, \tilde{a}) := \lambda \gamma N(\tilde{s}, \tilde{a}) \quad \forall \tilde{a} \in A, \tilde{s} \in S$

12  $(r, s') := \text{execute\_action}(s, a)$

13  $a' := \pi(s')$

14  $\delta := r + \gamma \hat{Q}_{s', a'} - \hat{Q}_{s, a}$

15  $\hat{Q}_{\tilde{s}, \tilde{a}} := \hat{Q}_{\tilde{s}, \tilde{a}} + \alpha N(\tilde{s}, \tilde{a}) \delta, \quad \forall \tilde{a} \in A, \tilde{s} \in S$

16 for  $s \in S$ :

17  $\pi_s^* := \arg \max_{a \in A} Q(s, a)$

where

- ▶  $s_{\text{term}}$  **terminal state** with zero reward.
- ▶  $\pi$  **generating policy**
- ▶  $q_0 \in \mathbb{R}$  initial value of all state/action pairs.
- ▶  $\alpha_k$  **learning rate** for update step  $k$ , e.g.,  $\alpha_k := 1/k$ .
- ▶  $\lambda \in [0, 1]$  **eligibility discount factor**

# Outline

1. Introduction
2. Monte Carlo Methods
3. Learning the Value Function: TD(0)
4. Learning the Action Value Function: SARSA
5. Learning the Optimal Action Value Function: Q-Learning
6. Eligibility Traces
- 7. R-Learning**
8. Model-Based Reinforcement Learning



...

# Outline

1. Introduction
2. Monte Carlo Methods
3. Learning the Value Function: TD(0)
4. Learning the Action Value Function: SARSA
5. Learning the Optimal Action Value Function: Q-Learning
6. Eligibility Traces
7. R-Learning
- 8. Model-Based Reinforcement Learning**

...

# Further Readings

- ▶ Reinforcement Learning:
  - ▶ Olivier Sigaud, Frederick Garcia (2010): *Reinforcement Learning*, ch. 1 in Sigaud and Buffet [2010].

# References

Olivier Sigaud and Olivier Buffet, editors. *Markov Decision Processes in Artificial Intelligence*. Wiley, 2010.