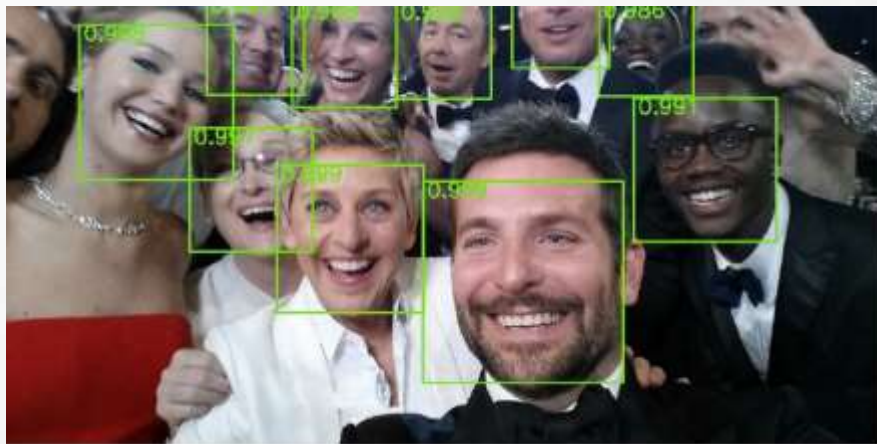
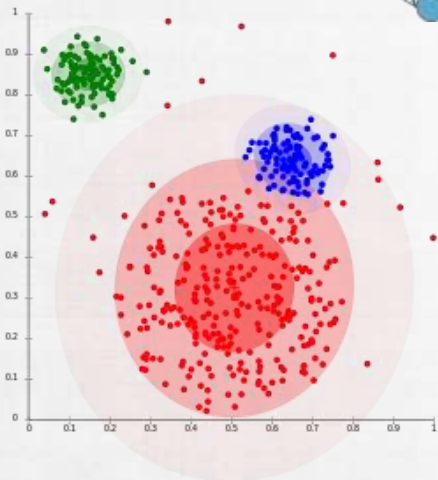
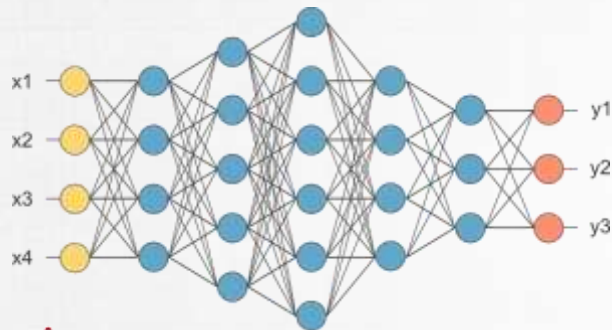

Reinforcement Learning

A short introduction (powered by *Sutton and Barto, 2012*)

10.01.2019

Severin Angerpointner
LMU Munich

AI in different areas



How do organisms learn?

Classification: microscopic model

panda



not a panda



Neural Networks, Clustering, ...

Behaviour (policy): "black box"

?

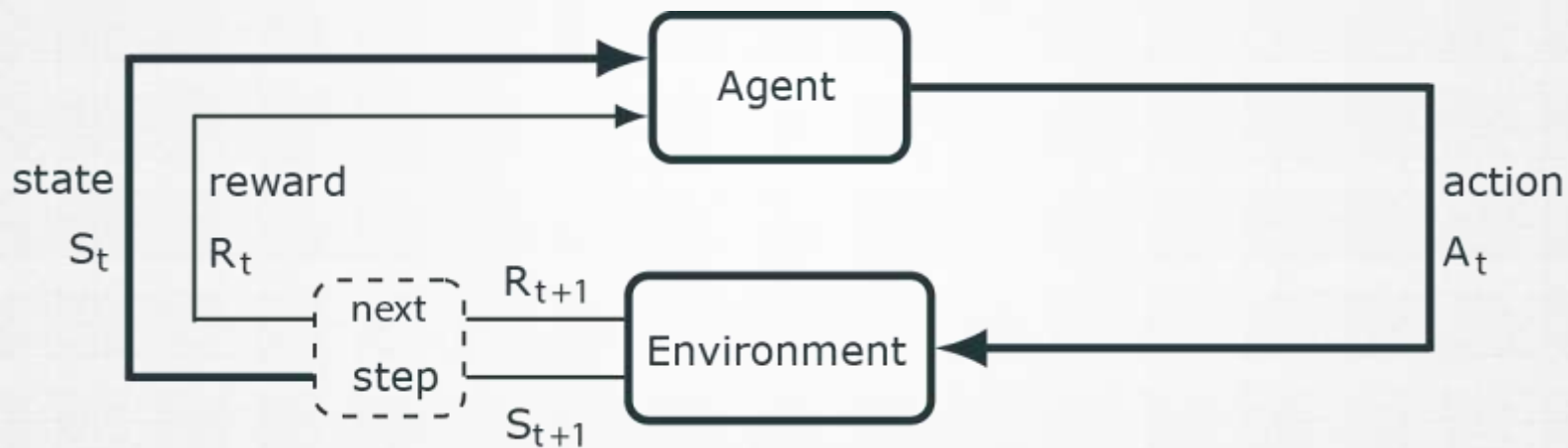
approach
panda

flee

Reinforcement Learning, ...

Dynamic-Programming
Agent Policy
On-Policy Return Off-Policy
Environment Action
Bellman-Equation
Greedy Exploitation Exploration
Reinforcement
Value Learning State
Non-greedy SARSA
Q-learning
Reward
Temporal-Difference

The RL Problem



$t \in \{0, 1, \dots, T\}$

(episodic tasks: $T < \infty$)

Basic Not(at)ions

- States: $S_t \in \mathcal{S}$
- Actions: $A_t \in \mathcal{A}(S_t)$
- Rewards: $R_t \in \mathbb{R}$ (i.g. randomly distributed)

"input"

- Policies: $\pi : \mathcal{S} \rightarrow \text{PDF}[\mathcal{A}]$

$$s \mapsto \pi(a | s)$$

(Probability to select action a being in state s)

"output"

Examples (blackboard)

- Pole Balancing

- Gridworld

Reward vs. Return vs. Value

Learner should achieve an overall *goal* (not just immediate reward)

- Return: $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T \quad \gamma \in [0, 1]$

- State Value: $v_\pi(s) = \mathbb{E}_\pi [G_t | S_t = s]$

State-Action Value

Alternative quantity: $q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t \mid S_t = s, A_t = a]$

(Expected return after choosing action a in state s and then following policy π)

+ No search for best action necessary

- Higher computational cost

Bellman Equation

$$q_{\pi}(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a, s') + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right]$$

$$p(s' | s, a) = p(S_{t+1} = s' | S_t = s, A_t = a)$$

$$r(s, a, s') = \mathbb{E}_{\pi} [R_{t+1} | S_t = s, A_t = a, S_{t+1} = s']$$

Strictly applies only to *Markov Decision Processes (MDP)*!

Optimal Policies

There is at least one "optimal" policy π_* , i.e.:

$$\forall s \in \mathcal{S} \quad \forall a \in \mathcal{A}(s) : \quad q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

(M.L. Puterman, "Markov Decision Processes", 2016)

→ π_* = "Always choose action with highest value" (greedy)

→ But we don't know corresponding q

General RL Algorithms

RL problem solved by finding q_* , with:

$$\pi_*(a_* | s) \neq 0 \quad \text{only for} \quad a_* = \arg \max_a q_*(s, a)$$

Apply Bellman equation:

$$q_*(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a, s') + \gamma \max_{a'} q_*(s', a') \right]$$

Unique solution exists for finite MDP!

Common RL Methods

Dynamic
Programming

Monte
Carlo

Temporal Difference-/
Q-Learning

Q(λ)-Learning



Dynamic Programming

Example: Policy Iteration

- 1) Use Bellman eq. iteratively to update v for given policy
- 2) Find better policy by selecting argmax of q as action

+ Guaranteed convergence (finite MDP)

- Computationally expensive
- Complete knowledge of MDP required

Policy Iteration

1. Initialization

$v(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Repeat

$\Delta \leftarrow 0$

For each $s \in \mathcal{S}$:

$temp \leftarrow v(s)$

$v(s) \leftarrow \sum_{s'} p(s'|s, \pi(s)) [r(s, \pi(s), s') + \gamma v(s')]$

$\Delta \leftarrow \max(\Delta, |temp - v(s)|)$

until $\Delta < \theta$ (a small positive number)

3. Policy Improvement

$policy-stable \leftarrow true$

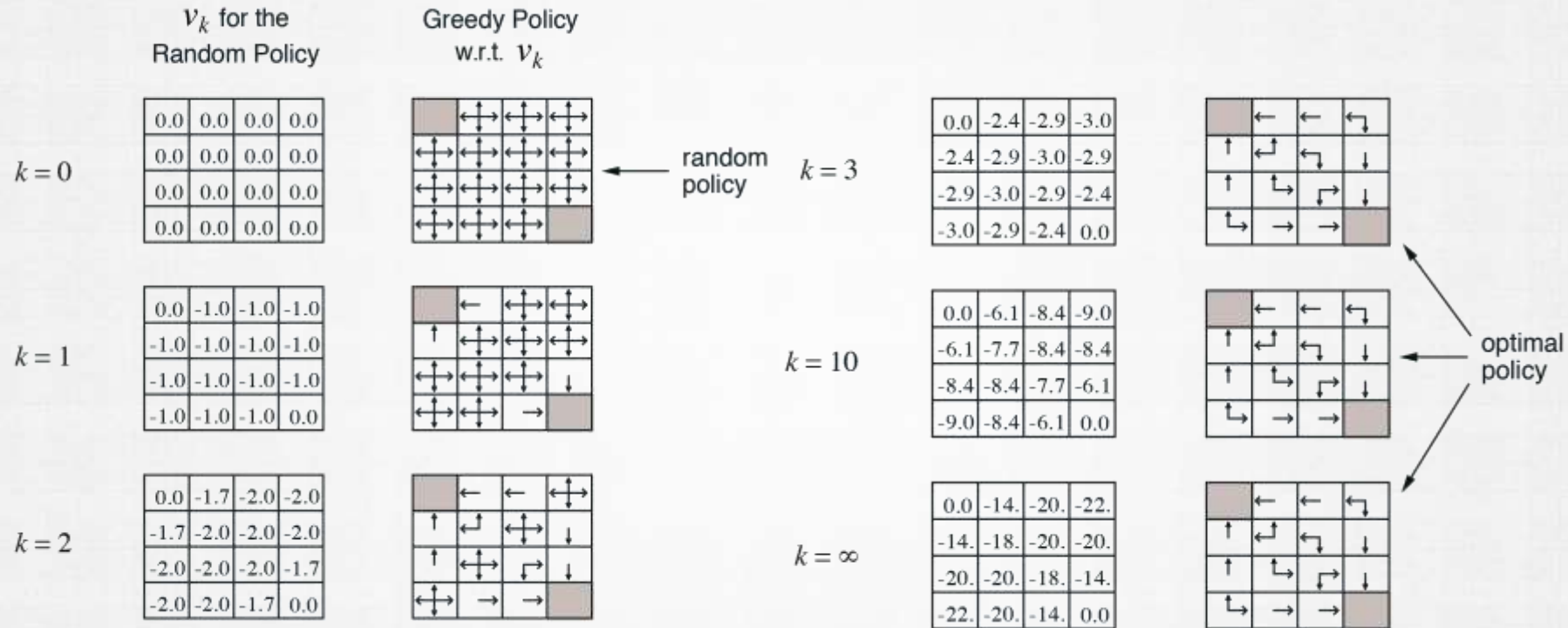
For each $s \in \mathcal{S}$:

$temp \leftarrow \pi(s)$

$\pi(s) \leftarrow \arg \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$

If $temp \neq \pi(s)$, then $policy-stable \leftarrow false$

If $policy-stable$, then stop and return v and π ; else go to 2



Monte Carlo

Suppose only sample of MDP known, not full process

- 1) Approximate value functions empirically
 - 2) Improve policy similar to DP
- $$q_{\pi}(s, a) \leftarrow \frac{1}{N} \sum_{m=1}^N G_t^{(m)}$$

+ Requires only sample returns/episodes

- Maintaining exploration
- Can only update after each episode

On-/Off-Policy

- On-Policy: follow *and* evaluate same policy π (as before)
 - Off-Policy: follow *behaviour policy* π / evaluate *estimation policy* π'
- Can choose exploring (*soft*) policy to sample whole state space

Examples: ϵ -greedy, Softmax, ...

Temporal Difference Learning

General idea (combine DP and MC):

Gradually update q towards optimum

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \delta_t$$

learning rate

error/temporal difference
= "target - current value"

e.g. update towards G after full episode (MC): $\delta_t = G_t^{(m)} - Q(S_t, A_t)$

One Step TD Learning

On-Policy: SARSA $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$

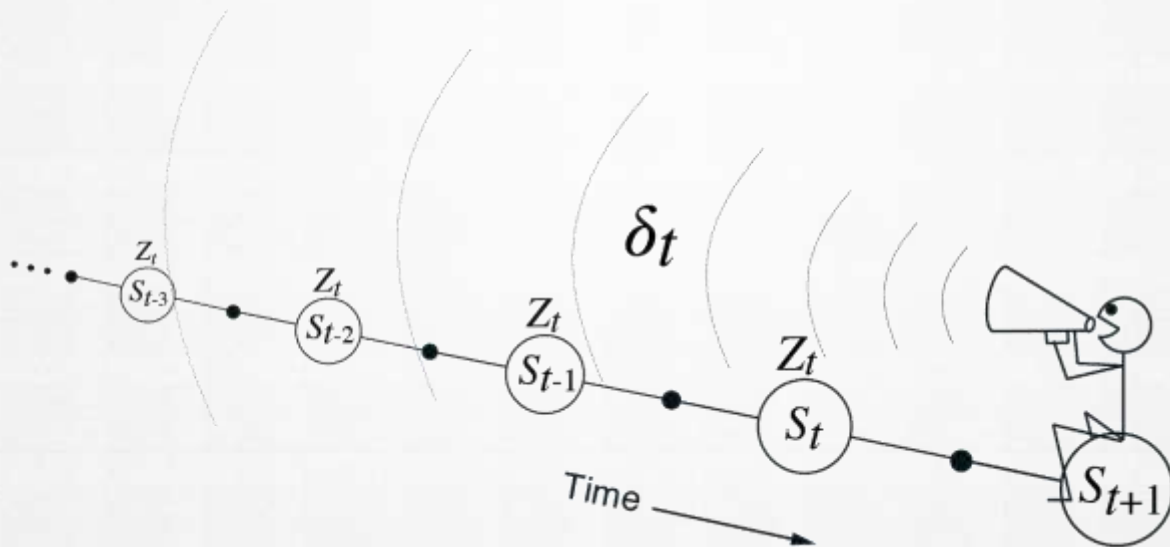
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

Off-Policy: Q-Learning

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

Eligibility Traces

Update more than one previously visited states
→ Compromise between full MC and one step TD



TD(λ) Algorithms

New update rule: $\forall s \in \mathcal{S} \quad \forall a \in \mathcal{A}(s) :$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta_t Z_t(s, a)$$

With eligibility trace:

$$Z_t(s, a) = \gamma \lambda Z_{t-1}(s, a) + \delta_{S_t, s} \delta_{A_t, a}$$

Sarsa(λ)

Initialize $Q(s, a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat (for each episode):

$Z(s, a) = 0$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Initialize S, A

Repeat (for each step of episode):

Take action A , observe R, S'

Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)

$\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$

$Z(S, A) \leftarrow Z(S, A) + 1$

For all $s \in \mathcal{S}, a \in \mathcal{A}(s)$:

$Q(s, a) \leftarrow Q(s, a) + \alpha \delta Z(s, a)$

$Z(s, a) \leftarrow \gamma \lambda Z(s, a)$

$S \leftarrow S'; A \leftarrow A'$

until S is terminal

Summary

- 1) Model specific tasks with **state** and **action** spaces
- 2) Define **goal** via reward function
- 3) RL problem: find optimal **policy/value** function
- 4) Methods: systematic policy improvement (**DP**), learning from experience (**MC, TD, Q**)

*References

- <https://qph.fs.quoracdn.net/main-qimg-330e8b2941bc0164211bbdc7d5c693f3>
- <https://de.wikipedia.org/wiki/Datei:AlphaGo.svg>
- <https://de.wikipedia.org/wiki/Clusteranalyse#/media/File:EM-Gaussian-data.svg>
- <https://www.dailydot.com/debug/face-detection-algorithm-image-search/>
- <http://m.koreatimes.co.kr/pages/article.asp?newsIdx=260722>

- Any other graphics are taken from *Sutton and Barto, 2012*