

# Reinforcement Learning with Applications of Intelligent Agents

Dr. Michael Thrun

# Reinforcement Learning (RL)

- Reinforcement learning can solve problems that are sequential and stochastic, e.g.,
  - Healthcare treatment [Håkansson et al., 2020]
  - Drug sensitivity prediction in cancer by ranking prediction algorithms [Daoud & Freisleben et al., 2020]
  - Taxation [Zheng & Socher et al., 2020 & 2021]

## Basic Principles

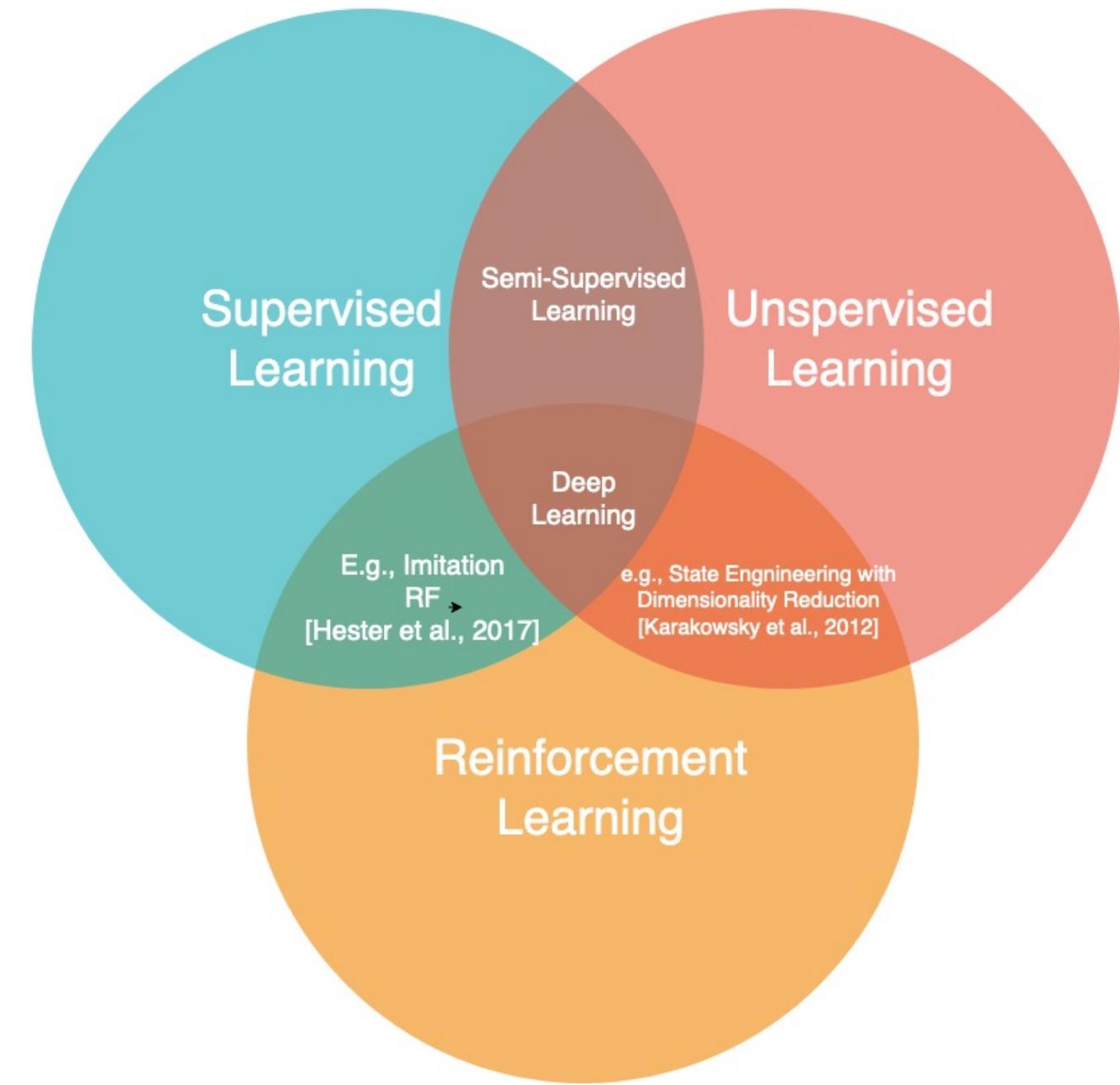
- Agent learns from interactions with environment
  - Feedback is restricted to delayed reward
- => No need for pairs of input and correct output

# Brief History of Reinforcement Learning (RL)

- RL is based on decades of research in AI and animal learning [Sutton & Barto, 1998], e.g.,
  1. Bellmann's optimal control, dynamic programming (**DP**) and markov decision processes (**MDP**), ~1950
  2. Trial and error learning in animals [Woodworth 1938; Pavlov 1927; Alexander Bain (ca. 1850)]
  3. Temporal difference learning (**TD**) [Minsky 1954; Samuel, 1959]
- Chris Watkins [1989 PhD thesis, Learning from Delayed Rewards] integrated these three research area into a new algorithm called Q-learning

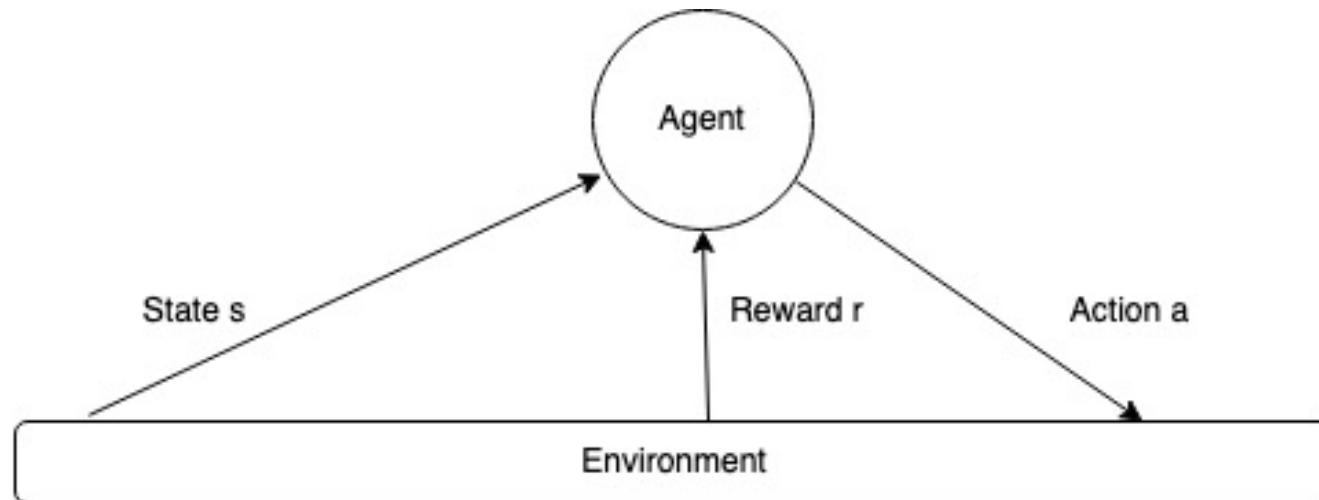
# Breakthroughs in RL

- RL concepts from supervised and unsupervised learning
- RL+MLP reached a level similar to the top three human player in the backgammon [Tesauro, 1995]
- RL+Deep learning (DQN) beats humans on atari games [Mnih et al., 2013 & 2015]
- DQN+supervised learning using human knowledge beats grandmaster in GO [Silver et al., 2016]
- DQN+MC beats its 2016 version in GO [Silver et al., 2017]
- DQN+MC beats grandmaster in chess [Silver et al., 2018]
- DQN beats grandmasters in real time strategy game (starcraft 2) [Wang et al., 2021]



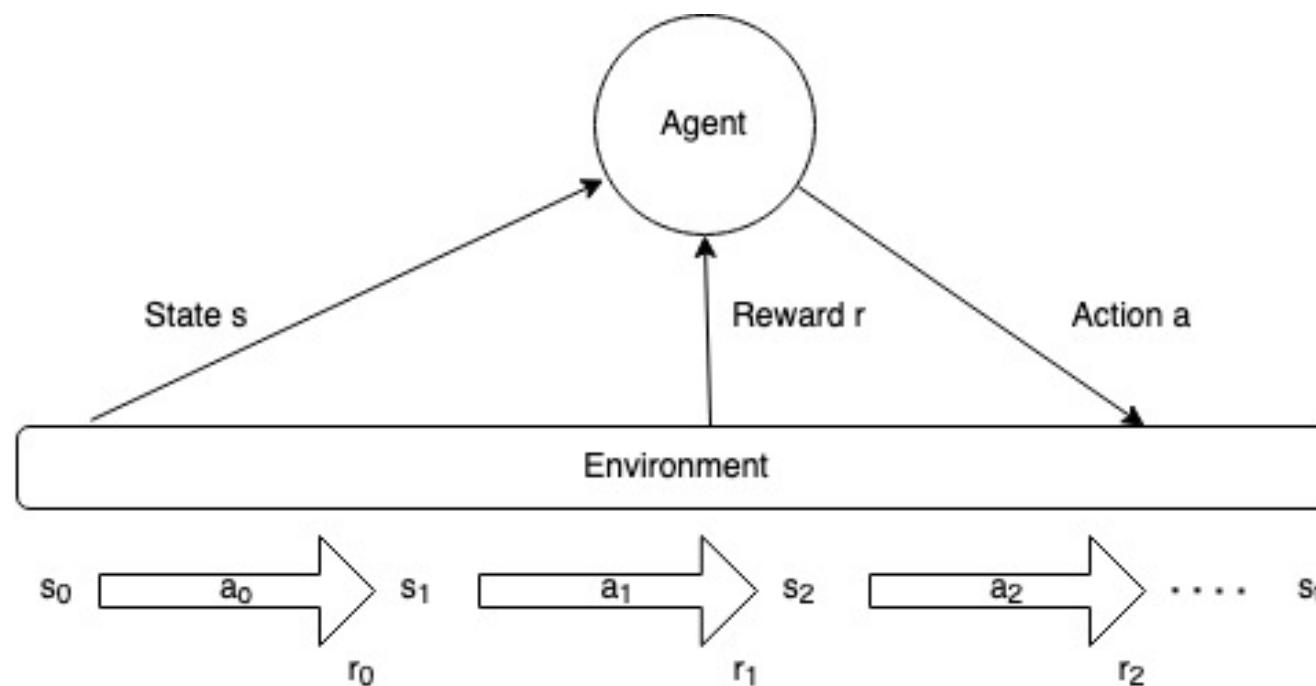
# Environment and Agent

- Environment denotes an abstraction of the world that an agent operates within
  - Either real-life situations are constrained to solvable problems using robotics
  - OR environment is a simulation of a task
    - Defined by the tuple  $(\Gamma, \Omega)$  where  $\Gamma$  is task space,  $\Omega$  is distribution that provides the probability of a task  $\gamma \in \Gamma$  to occur
- Intelligent agent is instantiation of RL algorithm that “intelligently” interacts with the environment and gain experience



# Definition: Rewards R, States S, Actions A, transition probabilities T

- Reward  $r \in R$  is a delayed and often sparse and noisy
- State  $s_i \in S$  is a representation of the environment at step  $i \in \mathbb{N}$  of the agent
- At each state  $s_i$  agent performs action  $a_i \in A$  that changes the state  $s = s_i$  to  $s_{i+1}$
- $T = P(s_{i+1}|s_i)$  are transition probabilities that are unknown in model-free algorithms



# Definition: Markov Decision Process (MDP)

- MDP allows to model *tasks*  $\gamma \in \Gamma$ , i.e., decision making even in situations where outcome are partially random

An MDP is a *task*, i.e., a 4-tuple  $(S, A, R, T)$  that for every state  $s_i$  fulfills the Markov property if

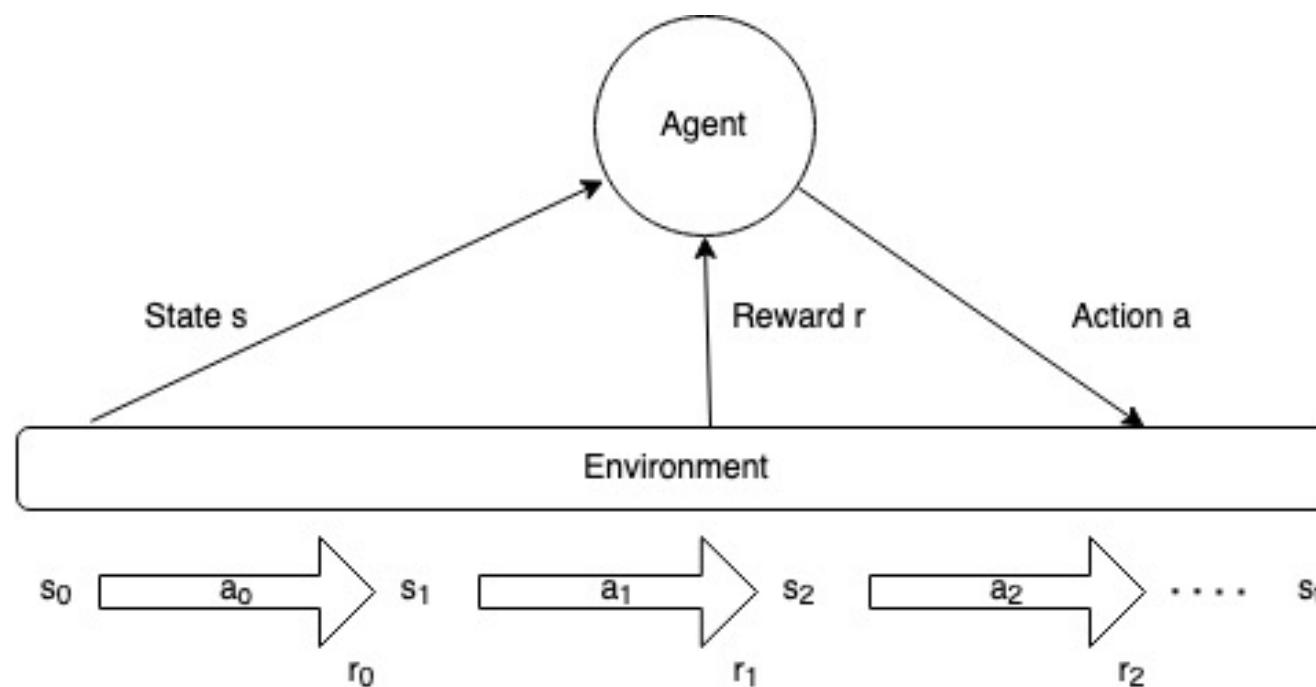
$$P(s_{i+1} | s_i) = P((s_{i+1} | s_1, \dots, s_i))$$

=> Transition from  $s_i$  to  $s_{i+1}$  is entirely independent of the past

=> effects of an action taken in a state depend only on that state but not prior states

# Definitions: Policy $\pi$

- Policy is the probability  $\pi(a|s)$  of an action given a state
  - $\pi: S \rightarrow A$  is *any* function of choosing an action in a given state
  - Fully defines the behavior of an agent (i.e., all possible actions) in environment
- After learning the optimal behavior is denoted with  $\pi^*$



# Two main approaches to solve the control problem of finding $\pi^*$

First

- MDP specifies setup for reinforcement learning

Then

- Policy-based agents will learn policy  $\pi(a|s) = \pi(a|s, \vartheta)$  directly by optimizing parameter  $\vartheta$
- Value-based agents will learn either action-value function  $Q_\pi(s, a)$   
OR state-value function  $V_\pi(s)$

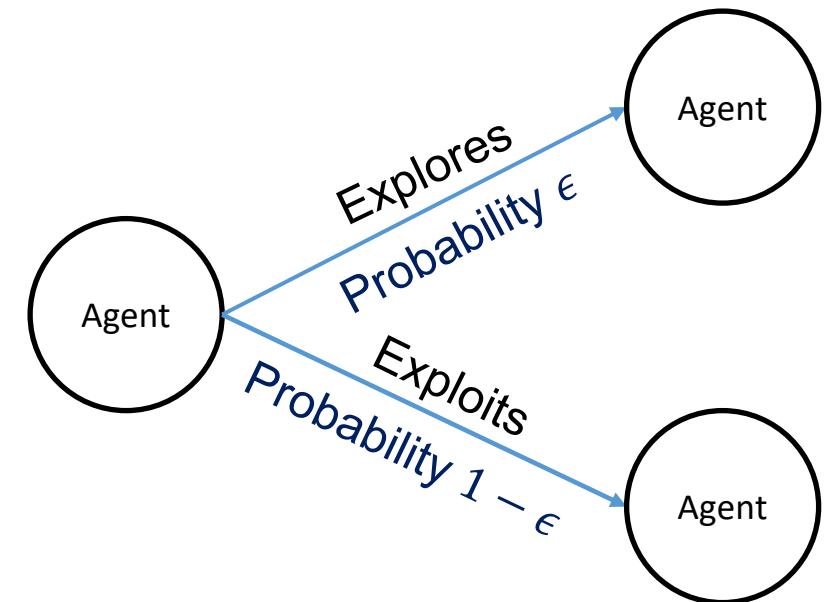
given a behavior policy  $\pi(a|s)$

- What is a simple policy of an agent?

# Try and error policy for learning $\pi^*$ with greedy action selection

## Exploration

- Agent tries new non-optimal action to
    - Learn reward
    - Gain better understanding of environment
- ⇒ Random action with  $\epsilon$ ,  $0 < \epsilon < 1$



## Exploitation

- Use current experience of agent to choose optimal action  $\pi(s|a) = \begin{cases} 1 & \text{if } a = \max Q(s, a) \\ 0 & \text{otherwise} \end{cases}$ 
  - ⇒ Greedy action with  $1 - \epsilon$
  - ⇒ Optimal in the sense that  $Q_\pi$  is maximized for selected  $a$
- Basic improvements are annealing  $\epsilon = \epsilon(\text{epoch})$ , UCB and softmax

# Definition: State-value and Action-value Functions

Let  $G$  be the return and  $E_\pi$  its expected value with regards to the policy of the agent  $\pi$

- Action-value function  $Q_\pi(s, a)$  is defined as

$$Q_\pi(s, a) = E_\pi[G|s, a]$$

- includes all expected rewards until terminal state  $s_T$  given state  $s$  and taken action  $a$
- State value function  $V_\pi(s)$  is defined as

$$V_\pi(s) = E_\pi[G|s]$$

- Gives expected returns starting from the state  $s$  and going to successor states thereafter

# Definition: Return G

- Return is the sum of discounted rewards

$$G = \sum_{i=0}^T \gamma^i r_i$$

- „Rewards received  $i$  steps hence are worth less than rewards received now“ [Watkins & Dayan, 1992]

$$0 < \gamma \leq 1, T < \infty$$

with  $T$  being the terminal step leading to the terminal state  $s_T$

- Goal of the agent is to find a policy that maximizes the expected return and, thus, all rewards

# Optimal Control: Learn $\pi^*$ by Bellman Optimality

- Maximizing Q is equal to maximizing the expected return

$$Q_*(s) = \underset{\pi}{\operatorname{argmax}} Q_{\pi^*}(s, a) = \underset{a_i \in A}{\operatorname{argmax}} E_{\pi}[G|s, a]$$

=> Policy obtains the highest reward if agent chooses action  $a$  with highest expected return

- Optimal policy  $\pi^*$  also maximizes state value function  $V_{\pi}(s)$  with

$$V_*(s) = \underset{\pi}{\operatorname{argmax}} V_{\pi}(s) \quad \forall s$$

=> Value of the state under an optimal policy must be equal to G for next best action of that state [Sutton & Barto, 1998]:

$$V_*(s) = \underset{a_i \in A}{\operatorname{argmax}} Q_{\pi^*}(s, a)$$

# Q-learning [Watkins, 1989]

- Q-learning learns  $Q_*$  directly without modelling transition probabilities

$$Q_*(s) = \underset{\pi}{\operatorname{argmax}} Q_{\pi^*}(s, a) = \underset{a_i \in A}{\operatorname{argmax}} E_{\pi}[G|s]$$

⇒ Q learning learns  $Q_*(s, a)$  using the Bellman optimality

⇒ Agent has to have full experience about future rewards in G

# Algorithm of Q learning

- 1: Input is a policy  $\pi(a|s)$  that uses the action-value function  $Q(s, a)$
- 2: Initialize  $Q(s, a) = 0$ , for all  $s_i \in S, a_i \in A$
- 3: Loop for each episode
- 4:   Initialize environment to provide s
- 5:   do
- 6:     Choose a in state s using  $\pi$  breaking ties randomly
- 7:     Take action a and observe r,  $s_{i+1}$
- 8:     Update  $Q(s, a)$  by Eq. 1
- 9:      $s = s_{i+1}$
- 10:    while s is not terminal

Line 1:  
 $\pi$  is used to decide the probability of which action to take, e.g., greedy action selection

Line 2:  
Initializes look-up table

# Experience of an agent (discrete case)

- Q-learning learns  $Q_*$  directly without modelling transition probabilities or expected rewards of the MDP

$$Q_*(s) = \underset{\pi}{\operatorname{argmax}} Q_{\pi^*}(s, a) = \underset{a_i \in A}{\operatorname{argmax}} E_{\pi}[G|s]$$

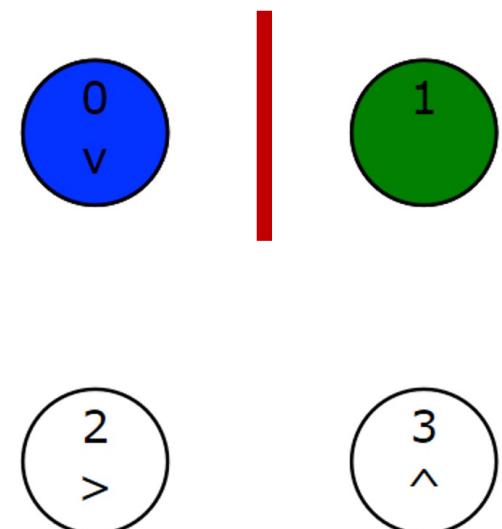
⇒ Q learning learns  $Q_*(s, a)$  using the Bellman optimality

⇒ Agent has to have full experience about future rewards in G

⇒ Experience has to be gained by trying out states/action combinations

⇒ Experience of agent is stored state-action-table

| State/action                     | $a_1=\text{up}$ | $a_2=\text{left}$ | $a_3=\text{down}$ | $a_4=\text{right}$ |
|----------------------------------|-----------------|-------------------|-------------------|--------------------|
| <b>Start <math>s_0</math></b>    | 0               | /                 | 0                 | 0                  |
| <b>Terminal <math>s_1</math></b> | 0               | 0                 | 0                 | /                  |
| $s_2$                            | 0               | 0                 | 0                 | 0                  |
| $s_3$                            | 0               | 0                 | 0                 | 0                  |



# Algorithm of Q learning

- 1: **Input** is a policy  $\pi(a|s)$  that uses the action-value function  $Q(s, a)$
- 2: Initialize  $Q(s, a) = 0$ , for all  $s_i \in S, a_i \in A$
- 3: Loop for each episode
- 4:   Initialize environment to provide  $s$
- 5:   do
- 6:     Choose  $a$  in state  $s$  using  $\pi$  breaking ties randomly
- 7:     Take action  $a$  and observe  $r, s_{i+1}$
- 8:     Update  $Q(s, a)$  by Eq. 1
- 9:      $s = s_{i+1}$
- 10:    while  $s$  is not terminal

Line 3  
Episode is the full experiment in the environment until a terminal state  $s_T$

Line 4  
Sets also agent to the starting state

# Algorithm of Q learning

- 1: **Input** is a policy  $\pi(a|s)$  that uses the action-value function  $Q(s, a)$
- 2: Initialize  $Q(s, a) = 0$ , for all  $s_i \in S, a_i \in A$
- 3: Loop for each episode
- 4:   Initialize environment to provide  $s$
- 5:   do
- 6:     Choose  $a$  in state  $s$  using  $\pi$  breaking ties randomly
- 7:     Take action  $a$  and observe  $r, s_{i+1}$
- 8:     Update  $Q(s, a)$  by Eq. 1
- 9:      $s = s_{i+1}$
- 10:    while  $s$  is not terminal

Line 6:  
Agent decides  
which action to  
take  
Line 7: Agent  
observes effects  
of action

# Algorithm of Q learning

- 1: **Input** is a policy  $\pi(a|s)$  that uses the action-value function  $Q(s, a)$
- 2: Initialize  $Q(s, a) = 0$ , for all  $s_i \in S, a_i \in A$
- 3: Loop for each episode
- 4:   Initialize environment to provide s
- 5:   do
- 6:     Choose a in state s using  $\pi$  breaking ties randomly
- 7:     Take action a and observe r,  $s_{i+1}$
- 8:     Update  $Q(s, a)$  by Eq. 1
- 9:      $s = s_{i+1}$
- 10:    while s is not terminal

Line 8:  
Expected return G is  
calculated as the sum of  
the current reward and  
the expected value of  
the next state and used  
to update Q in lookup  
table  
How?

# Q learning: Value Iteration by Update Rule

- Iterates estimate  $Q(s,a)$  by overwriting previous estimate in look up table
  - For non terminal state

$$Q(s_i, a) = Q(s_i, a) + \alpha \left( r + \gamma * \underset{a_s \in A}{\operatorname{argmax}} Q(s_{i+1}, a_s) - Q(s_i, a) \right) \quad \text{Eq. 1}$$

Learning rate

$$0 < \alpha \leq 1$$

Estimate of future rewards based on current experience

Discount factor

$$0 < \gamma \leq 1$$

$G_{i:i+1}$  one-step estimated expected return (Bellman equation)

- For terminal state  $Q(s_T, a) = r$

TD error  $\delta_i$  for one-step ahead Q-learning at step  $i$  available at time step  $i+1$

# Algorithm of Q learning

- 1: **Input** is a policy  $\pi(a|s)$  that uses the action-value function  $Q(s, a)$
- 2: Initialize  $Q(s, a) = 0$ , for all  $s_i \in S, a_i \in A$
- 3: Loop for each episode
- 4:   Initialize environment to provide s
- 5:   do
- 6:     Choose a in state s using  $\pi$  breaking ties randomly
- 7:     Take action a and observe r,  $s_{i+1}$
- 8:     Update  $Q(s, a)$  by Eq. 1
- 9:      $s = s_{i+1}$
- 10:    while s is not terminal

- Update is performed independently of the policy used by agent
- ⇒ Update is computed using the best action for  $a_{i+1}$

=> **off-policy** algorithm

# Algorithm of SARSA ( $s_i, a_i r_i, s_{i+1}, a_{i+1}$ )

- 1: **Input** is a policy  $\pi(a|s)$  that uses the action-value function  $Q(s, a)$
- 2: Initialize  $Q(s, a) = 0$ , for all  $s_i \in S, a_i \in A$
- 3: Loop for each episode
- 4: Initialize environment to provide  $s$
- 5: do
- 6: Choose  $a$  in state  $s$  using  $\pi$  breaking ties randomly
- 7: Take action  $a$  and observe  $r, s_{i+1}$
- 7b: Choose action  $a_{i+1}$  from  $s_{i+1}$  using  $\pi$
- 8: Update  $Q_{a_{i+1}}(s, a)$
- 9:  $s = s_{i+1}$
- 10: while  $s$  is not terminal

Q-learning is changed & extended by marked lines

- Line 7b:
- Suboptimal action  $a_{i+1}$  is taken
- Line 8:
- Update uses behavior policy to select  $a_{i+1}$
- => **on-policy** algorithm

$$Q(s, a) = Q(s, a) + \alpha(r + \gamma Q(s_{i+1}, a_{i+1}) - Q(s, a))$$

# Q-Learning is TD(0) Learning with Sample Update

Extension of Q-learning can be seen as N-step temporal difference learning (TD)

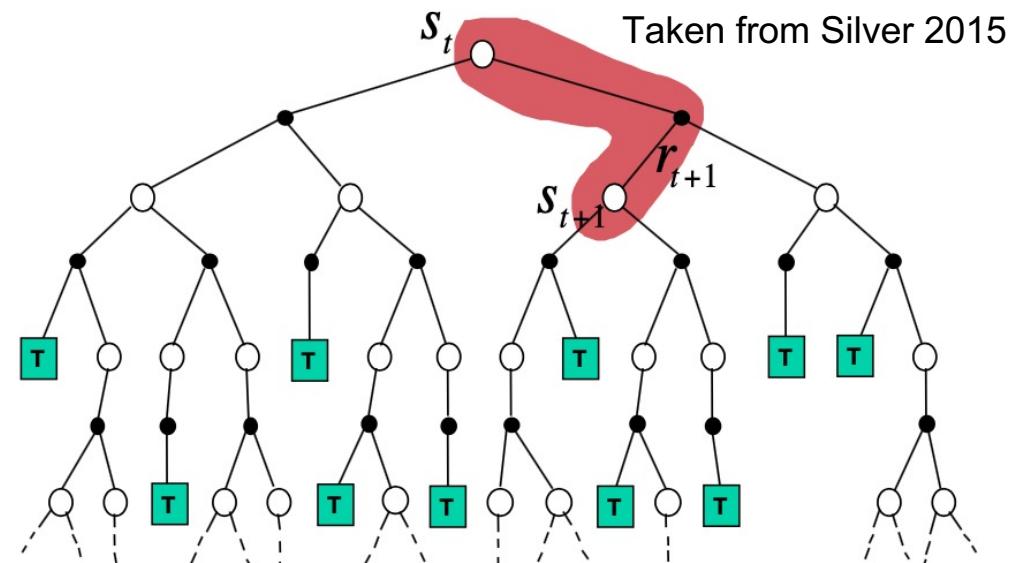
- Agent can buffer the rewards more than one step and take multiple steps forward until updating the lookup table

$$Q_{i:i+1}(s_i, a) = Q(s_i, a) + \alpha \left( r + \gamma * \operatorname{argmax}_{a_s \in A} Q(s_{i+1}, a_s) - Q(s_i, a) \right)$$

$$Q_{i:i+n}(s_i, a) = Q(s_i, a) + \alpha(G_{i:i+n} - Q(s_i, a))$$

$$G_{i:i+n} = r + \gamma r_{i+1} + \dots + \gamma^{n-1} r_{i+n-1} + \gamma^n * \operatorname{argmax}_{a_s \in A} Q(s_{i+n}, a_s)$$

$$0 \leq i \leq T - n$$



# Improvements of Q-Learning

- Replace look-up table with approximate function
  - ⇒ Agent looks at the current state-action pair and predicts the expected value within an regression method
  - ⇒ BUT: state-action pairs tend to be nonlinear and discontinuous
  - ⇒ Use neural networks [Tesauro, 1995]
  - ⇒ Deep Q learning translates raw observations into actions [Mnih et al., 2013]

# Basic Algorithm of Deep Q-learning [Mnih et al., 2013]

- 1: **Input** is a policy  $\pi(a|s)$  that uses the action-value function  $Q(a, s, \vartheta_i)$
- 2: Initialize  $Q(a, s, \vartheta_i)$  with random weights and replay memory  $D_N$
- 3: loop for each episode
- 4:   Initialize environment to provide  $s$  and preprocessed  $\varphi(s)$
- 5:   do
- 6:     Choose  $a$  in state  $s$  using  $\pi$  breaking ties randomly
- a7:     Take action  $a$  and observe  $r_i, s_{i+1}$
- b7:     Store transition  $e_i$  in  $D_N$
- c7:     Sample random minibatch of transition  $e_j \in D_N$
- a8:     Set target  $y$
- b8:**     Update  $Q(a, s, \vartheta_i)$  by  $\nabla_{\vartheta} L(\vartheta_i)$
- 9:      $s = s_{i+1}$
- 10:    while  $s$  is not terminal

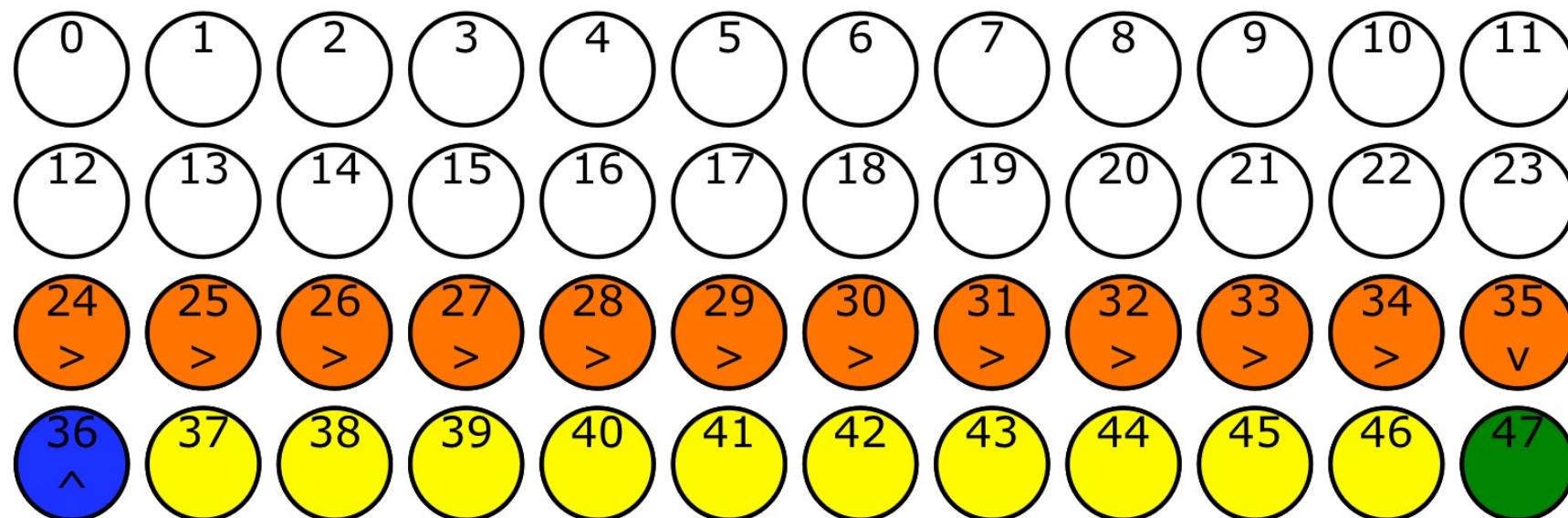
Q-learning is changed & extended by marked lines

Through  $\varphi(s)$  the environment definition is more elaborate to exploit standard deep learning

Q-learning is updated by gradient descent

# Application: Bypass a cliff

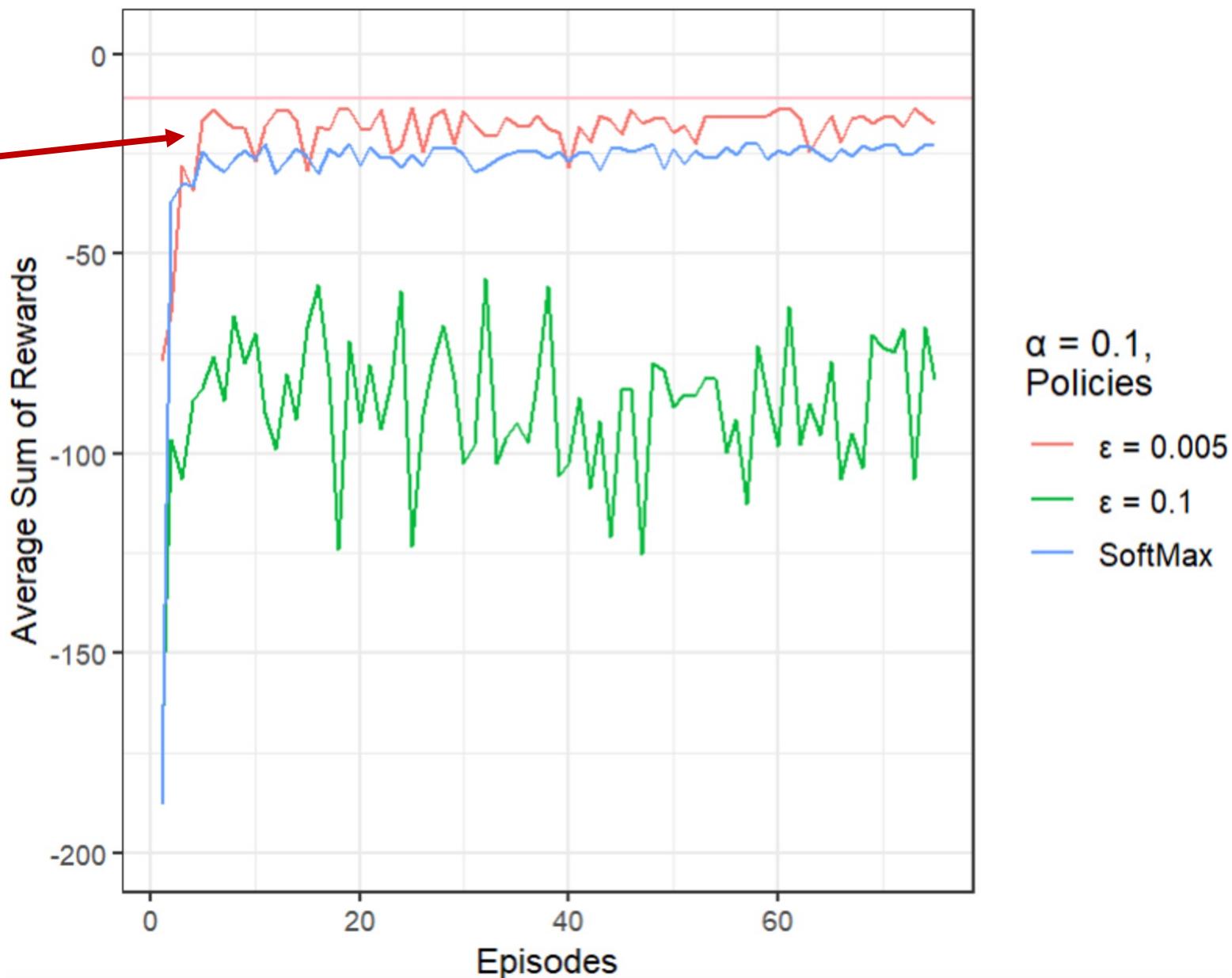
- $s_0, \dots, s_{47}$  states in white with four actions per state: up, down, right left
- Yellow Cliff has reward  $r=-200$  and stepping off the cliff sends agent back to blue start
- No discount with  $\gamma = 1$
- Each step has reward  $r=-1 \Rightarrow$  maximal reward is -13 (orange)
- Green goal has reward  $r=0$
- Approximation of  $\pi^*$  will vary depending on selected policy  $\pi$ , parameters and learning rate



# Typical Learning effects

- Agent converges but never stabilizes („chattering“)
  - => optimal path is near cliff
  - => Chance to be kicked off the cliff by  $\epsilon$

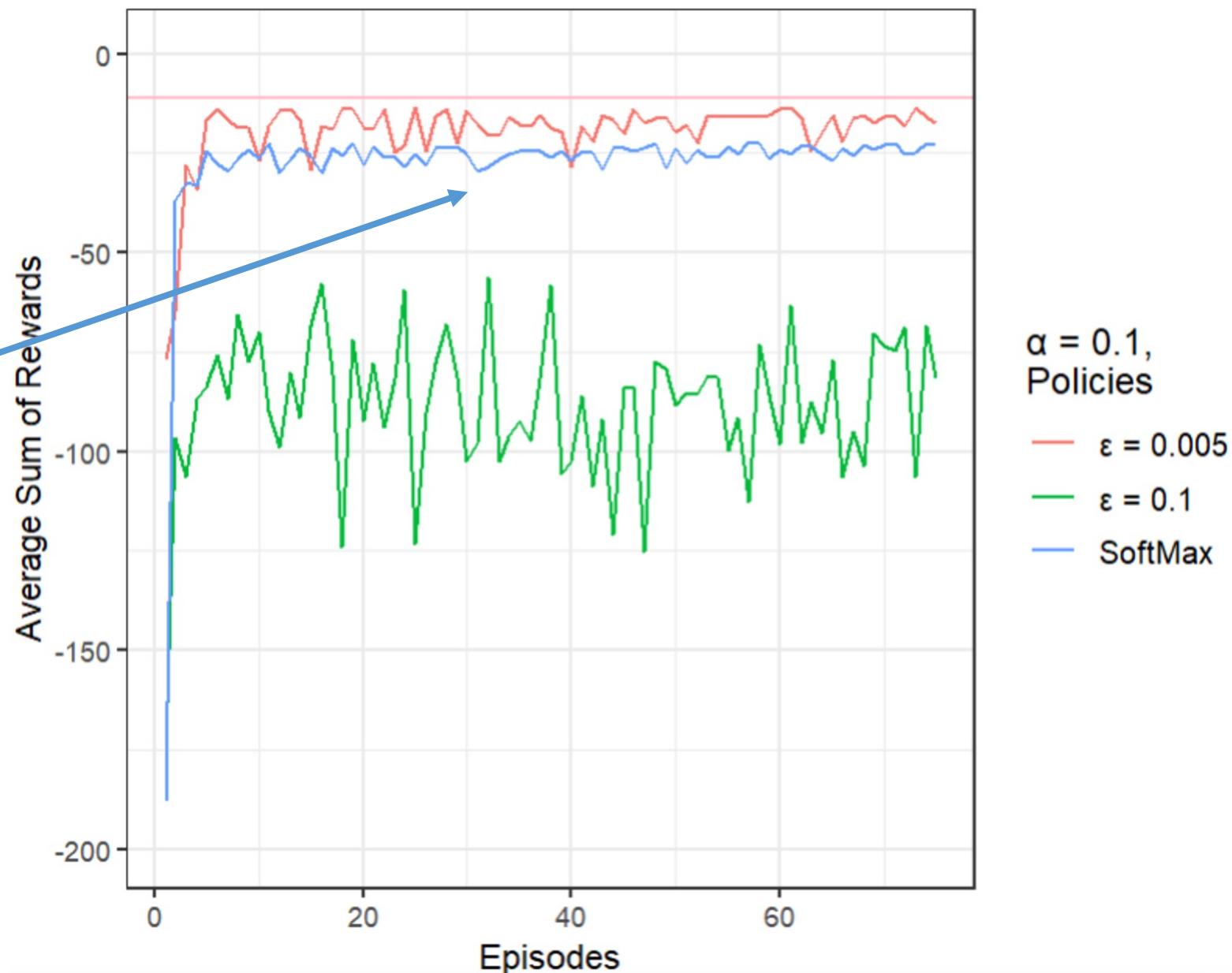
Q Learning Averaged over 100 Trials



# Typical Learning effects

- Agent converges but never stabilizes („chattering“)
  - => optimal path is near cliff
  - => Chance to be kicked off the cliff by  $\epsilon$
- Policy converges (slightly) below optimum
  - ⇒ Possibility to find local optima
  - ⇒ There may be more than one near optimal policy

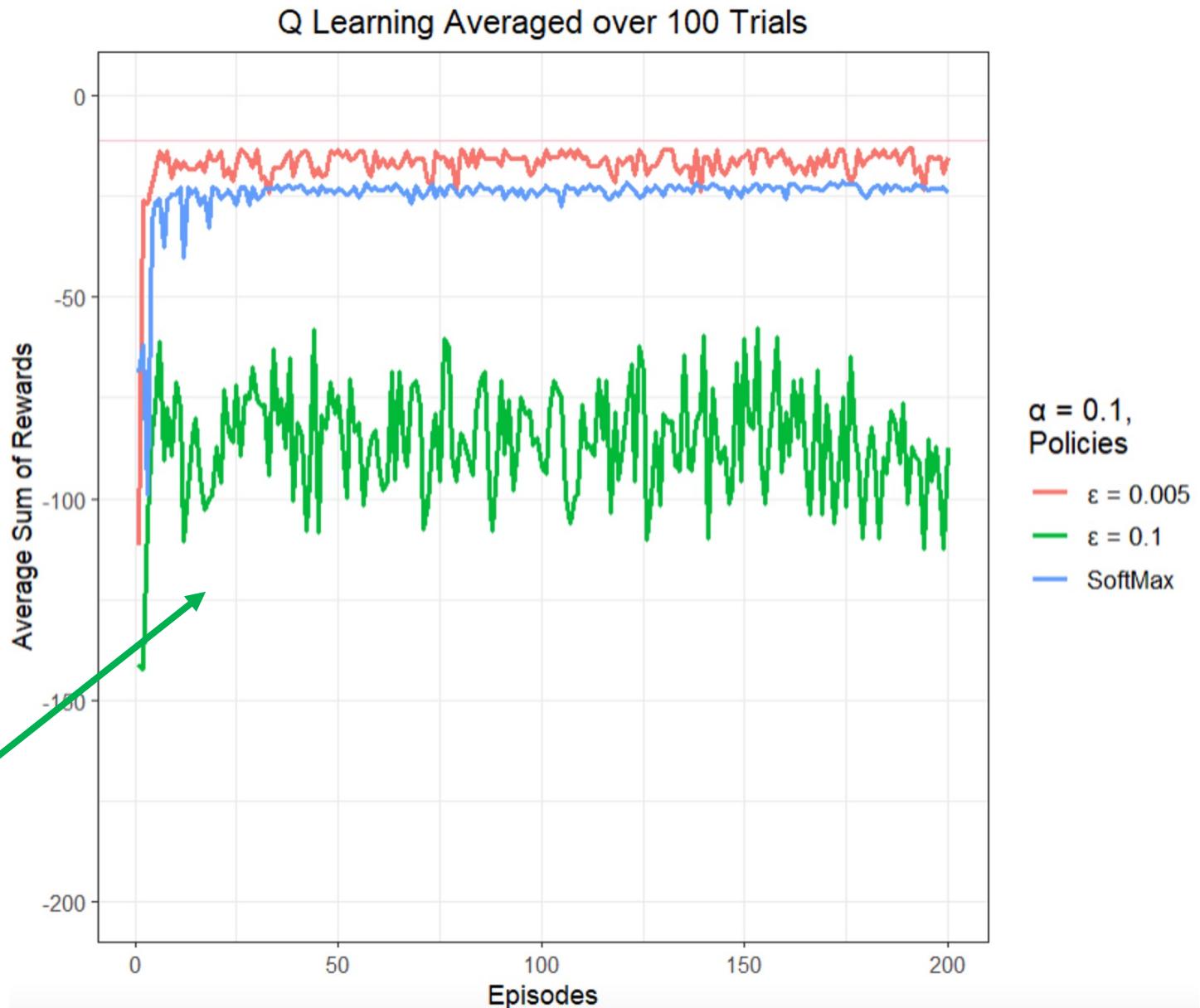
Q Learning Averaged over 100 Trials



# Typical Learning effects

- Agent converges but never stabilizes („chattering“)
  - => optimal path is near cliff
  - => Chance to be kicked off the cliff by  $\epsilon$
- Policy converges (slightly) below optimum
  - ⇒ Possibility to find local optima
  - ⇒ There may be more than one near optimal policy

Policy never converges



# Limitations of Q-Learning I

- Setting „balance“ between exploration and exploitation is challenging
  - Too much exploration, e.g., “agent jumps off the cliff”  
=> “jittered” convergence
  - Too much exploitation, e.g., “long route is taken”  
=> policy converges non-optimally or overfitting
- Value-based methods like Q-learning maximize the exploitation part of RL
  - BUT: Exploration is only ensured by external parameter (e.g.,  $\epsilon$ ) to be set by user
  - One Solution is learning appropriate “balance” via entropy
    - Maximum entropy RL learns by itself how much exploration it needs to learn appropriately, e.g.,
      - Soft Q-learning [Ziebart et al., 2008]

# Demonstrating Overfitting

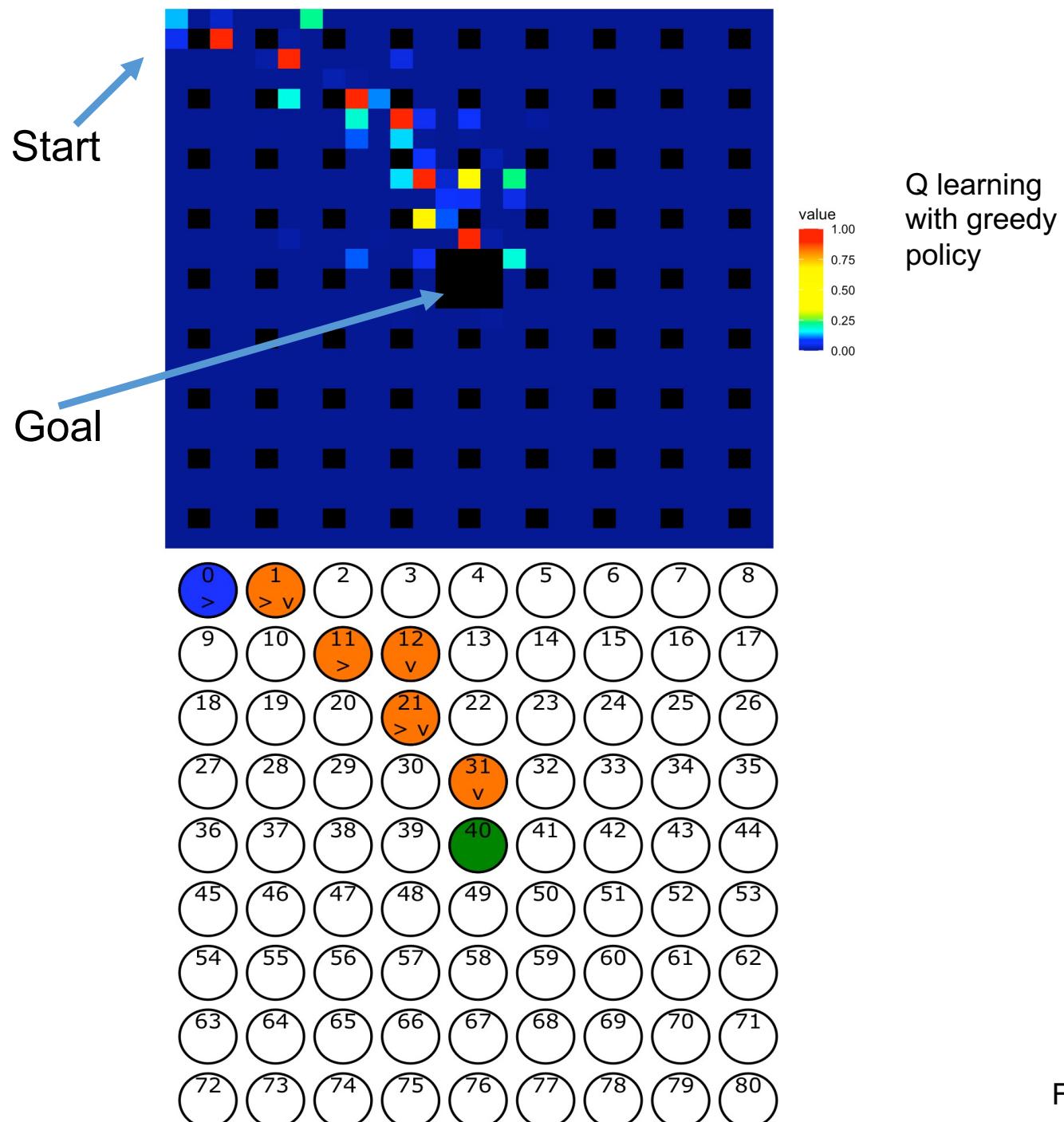
- Grid world with  $\epsilon$  greedy exploration
- 8 discrete actions are allowed
  - Diagonal movements are allowed

Lookup table is visualized as follows

- Black=States
- Color=Value of estimated Q
- Position of Pixel=movement position

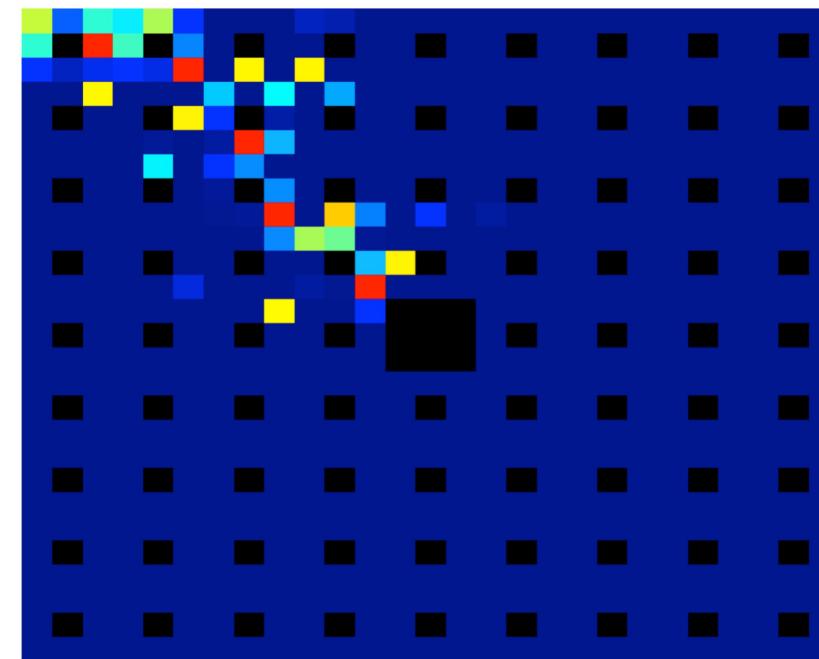
⇒ Blue= actions & states not sampled by agent so far

⇒ Yellow/Red high action states defining the non-optimal trajectory the agents follows

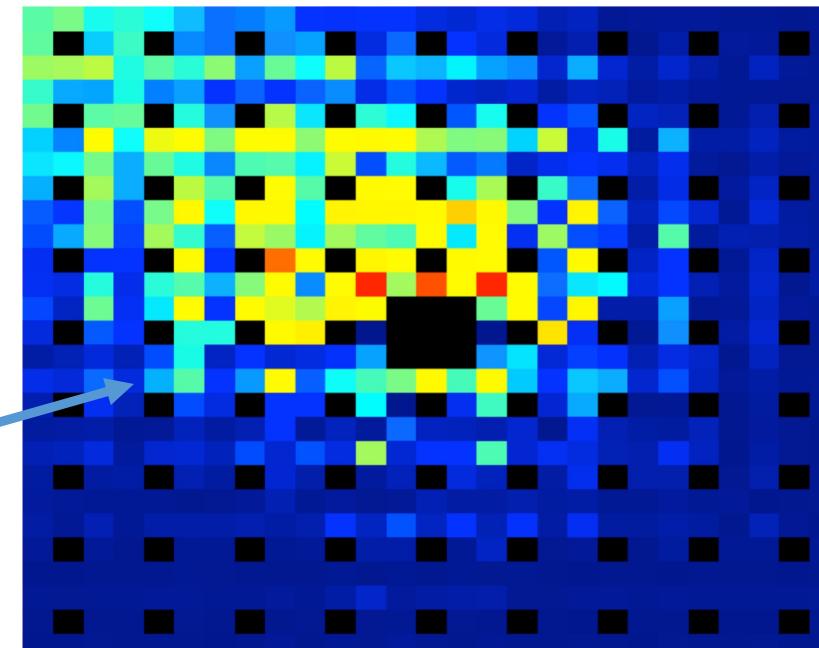


# Demonstrating Overfitting

- Process of accidentally building model that works on very specific circumstances, e.g.
  - Policy  $\pi^*$  has uncharted regions near optimal trajectory
    - Agent did not explore the environment enough
- Modelling actions as a distribution with softmax function results in more explored states
  - In model changes in start position still enable agent to find goal
  - Exploration took 10 times longer yet still left states unexplored
  - Model is more robust



Q learning with greedy policy  
-converged after 10 epochs learned



Q-learning with softmax policy  
-converged after 100 epochs learned

# Limitations of Q-Learning II

- Methods using look-up tables like Q-learning struggle when the state or action space is large
  - Agent has to sample a huge number of states and every action to find a good policy
    - Learning becomes very slow
    - Alternative solutions are, e.g., using information loss [Hoof, Neumann & Peters, 2017] or experience replay [Lazaric et al., 2008]
- Appropriate Parameter settings are unknown and depend on environment, rewards and policy
  - How many steps  $n$  should the  $n$ -step Q-learning take until Q update?
  - Alternatively with eligibility traces the agent can look backward in time, but how to set  $\lambda$ ?
  - How should the learning rate  $\alpha$  be set?

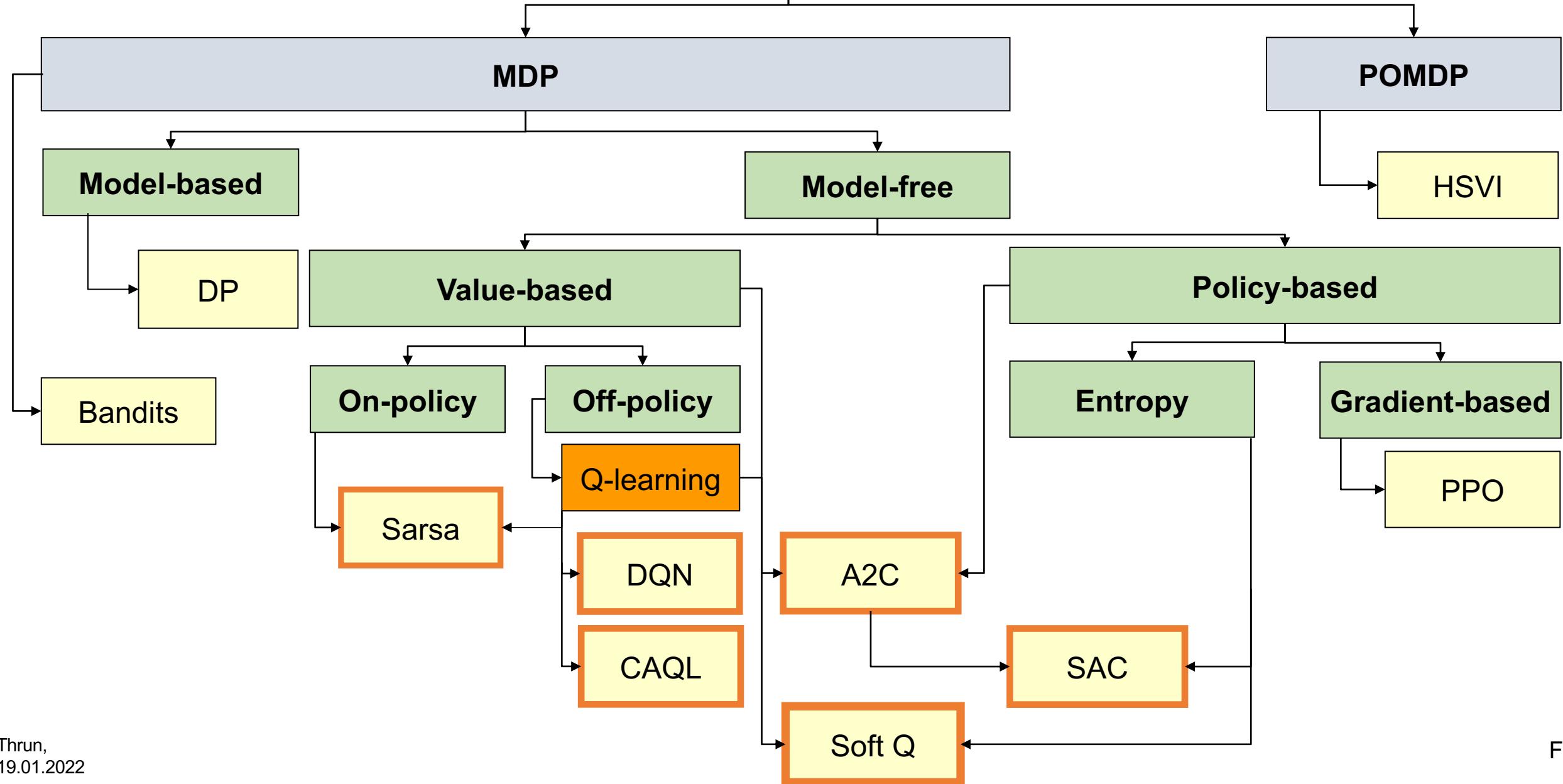
# Q-learning Alternatives

- Major problem can be tackled with an “divide and conquer” approach  
=> Hierarchical RL
  - One agent learns, otherwise use e.g.,
    - Multi-agent RL (MARL)
    - Proximal Policy Optimization (PPO) [Schulman et al., 2015; Schulman et al. 2017]
  - Actions are discrete, otherwise use e.g.,
    - Continuous Action Q-Learning (CAQL) [Ryu et al., 2020]
    - Soft Actor-Critic (SAC) [Haarnajoa et al., 2017]
  - Environment has a terminal state → ends after an episode
    - Alternative: Continuous Tasks for life-long learning
  - If each episode consists only of one time step => multi-armed bandit
  - Expectation value of Q is the moving average
    - For example it is multimodal for cliff-wall scenarios!
- => Distributional RL

# Using Assumptions about MDP to Structure RL types

- Transition probabilities  $T$  are unknown => model-free algorithms
  - Otherwise model-based RL algorithms
- States  $s$  are fully observable, otherwise use
  - Algorithms for partially observable MDP (POMDP)
- Agent learns one policy and algorithm is based on one given policy, otherwise
  - Advantage actor critic (A2C) [Barto,Sutton&Anderson, 1983; Kimura & Kobayashi ,1998]
- States  $S$  are discrete, otherwise use
  - Function approximators with deep learning (DQN) [Mnih et al., 2013]
- Q-learning: Agent learns value-based (i.e.,  $Q$  or  $V$ )
  - Alternative: policy-based algorithms learn  $\pi$  directly
- Policy is deterministic, otherwise
  - Stochastic policy learning

# RF Algorithms



# Summary

- RL problems are sequential contrary to ML
- Policy is responsible for mapping representation of state to action
  - Vast majority of RL literature is dedicated to finding improved ways of defining policies or improving some theoretical guarantees
- Balance exploitation vs. exploration  
⇒ Many algorithms provide elaborate solutions (e.g. SAC)
- Consider varying the environment to help algorithms generalize (not to overfit)
- Some part of the MDP can always be modeled by supervised or unsupervised ML approaches, e.g.,
  - Advantage: Deep learning does not require time-consuming state engineering
  - Disadvantage: There is no proof that RL will converge in this case!

# Practical Challenges

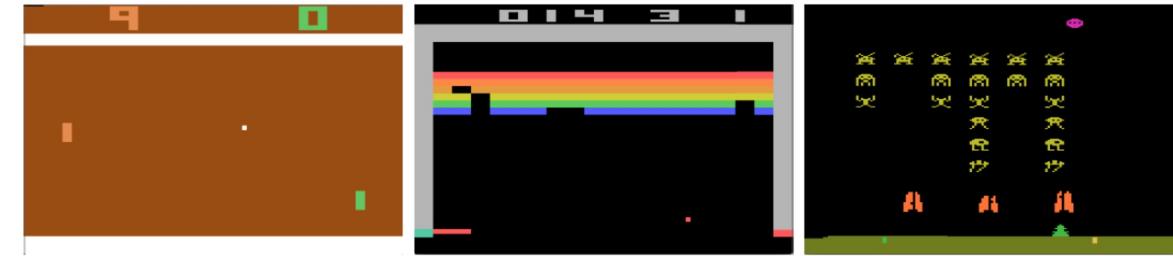
- Learning is the bottleneck of RL because it costs time or/and money  
=> Constrain problem definition (and thus the exploration space) until a valid MDP can be defined
  - Depending on MDP select RL algorithm using Ockham's razor
  - Start simple from straightforward Q-learning and then select “addons”
    - Very vast number of literature is about some combinations of such addons
- BEWARE: RL aims to find a policy that maximizes the expected return irrespective of consequences
  - ⇒ Rewards should be as near as possible to the application (e.g., € in business decisions)
  - ⇒ Hypothesis is that all goals can be described in such a way!

# Sources

Literature, Code & Slide accessible in  
[https://github.com/Mthrun/RLwithAppl\\_IntAgents/](https://github.com/Mthrun/RLwithAppl_IntAgents/)

# Deep Q Learning on Atari Games

- Reward is current game score
- Actions are defined by game with a probability
- Environment is Atari emulator
  - Agent observes image ( $s$ =vector of pixels) from emulator representing the current screen
    - 4 frames are preprocessed  $\varphi(s)$  and combined to one input to Q-function
- $Q(a, s, \vartheta) \approx Q(a, s)$  is the function approximator,  $\vartheta$  are weights of the deep learning network with loss function



$$L_i(\vartheta_i) = E_{\pi} \left[ (y_i - Q(a, s, \vartheta_i))^2 \right]$$

With the target

$$y_i = E_{\pi} \left[ r + \gamma * \operatorname{argmax}_{a_s \in A} Q(s_{i+1}, s, \vartheta_{i-1}) \right]$$

Where weights of previous iteration are fixed

=> BUT Target is in contrast to ML not fixed

# Replay memory

- Problem: Supervised ANN algorithms assume I.I.D
- BUT: RL has highly correlated sequence of states and data distribution changes as agent learns
- Solution: Use replay memory  $D$  [Lohng-ji, 1993] that takes sample of previous transitions  $e$

$$e_i = (\varphi_i, a_i, r_i, \varphi_{i+1})$$

$$D_N = \{e_1, \dots e_N\}, N = 10^6$$

- Sampling technique samples depending on the action, e.g.,
  - Action is move right than samples are dominated by right-hand side
- Mnih et al, 2013 claims that this smoothes out learning and avoids divergence and chattering

# Basic Algorithm of Deep Q-learning

- 1: **Input** is a policy  $\pi(a|s)$  that uses the action-value function  $Q(a, s, \vartheta_i)$
- 2: Initialize  $Q(a, s, \vartheta_i)$  with random weights and replay memory  $D_N$
- 3: loop for each episode
- 4:   Initialize environment to provide  $s$  and preprocessed  $\varphi(s)$
- 5:   do
- 6:     Choose  $a$  in state  $s$  using  $\pi$  breaking ties randomly
- a7:     Take action  $a$  and observe  $r_i, s_{i+1}$
- b7:     Store transition  $e_i$  in  $D_N$
- c7:     Sample random minibatch of transition  $e_j \in D_N$
- a8:     Set target  $y$
- b8:     Update  $Q(a, s, \vartheta_i)$  by  $\nabla_{\vartheta} L(\vartheta_i)$
- 9:      $s = s_{i+1}$
- 10:    while  $s$  is not terminal

Q-learning is changed & extended by marked lines

Through  $\varphi(s)$  the environment definition is more elaborate to exploit standard deep learning

Q-learning is updated by gradient descent

# Basic Algorithm of Deep Q-learning

1: **Input** is a policy  $\pi(a|s)$  that uses the action-value function  $Q(a, s, \vartheta_i)$

**2:** Initialize  $Q(a, s, \vartheta_i)$  with random weights and replay memory  $D_N$

Most important  
knack is to add  
replay memory  
 $D_N$

3: loop for each episode

4: Initialize environment to provide  $s$  and preprocessed  $\varphi(s)$

5: do

6: Choose  $a$  in state  $s$  using  $\pi$  breaking ties randomly

a7: Take action  $a$  and observe  $r_i, s_{i+1}$

Replay memory

- Stores given variables for a specific iteration

**b7:** Store transition  $e_i$  in  $D_N$

**c7:** Sample random minibatch of transition  $e_j \in D_N$

- Takes sample for computing a8-b8

**a8:** Set target  $y$

**b8:** Update  $Q(a, s, \vartheta_i)$  by  $\nabla_{\vartheta} L(\vartheta_i)$

9:  $s = s_{i+1}$

10: while  $s$  is not terminal

# Results Deep Q Learning on Atari Games [Mnih et al., 2013]

- Total average reward by running  $\epsilon$  greedy policy 4 Million frames
- DQN outperforms standard RL with elaborate manual feature engineering for states
- DQN outperforms even human expert in three games

|                        | B. Rider    | Breakout   | Enduro     | Pong      | Q*bert      | Seaquest    | S. Invaders |
|------------------------|-------------|------------|------------|-----------|-------------|-------------|-------------|
| <b>Random</b>          | 354         | 1.2        | 0          | -20.4     | 157         | 110         | 179         |
| <b>Sarsa [3]</b>       | 996         | 5.2        | 129        | -19       | 614         | 665         | 271         |
| <b>Contingency [4]</b> | 1743        | 6          | 159        | -17       | 960         | 723         | 268         |
| <b>DQN</b>             | <b>4092</b> | <b>168</b> | <b>470</b> | <b>20</b> | <b>1952</b> | <b>1705</b> | <b>581</b>  |
| <b>Human</b>           | 7456        | 31         | 368        | -3        | 18900       | 28010       | 3690        |

⇒ Deep learning allows automated feature extraction given a suitable problem for the ANN architecture

- Using 6 prior extensions of Q-learning rainbow DQN extended DQN [Hessel et al., 2017]



# Convergence Proofs

- Q-learning converges [Watkins & Dayan, 1992], so long as
  - All actions are repeatedly sampled in all states
  - Action-values are represented discretely
- Basic n-step Q-learning converges [Tsitsiklis & Van Roy, 1999]
- Given policies of continuous states policies must approximate actions
  - Linear approximations are always convergent [Sutton & Barto, 1998]
  - Non-parametric policies converge [Sutton & Barto, 1998, Hoof, Neumann & Peters, 2017]
  - Non-linear approximations (e.g. neural networks) have no guarantees to converge [Silver 2015, Winder, 2020]
    - In ML, Ben-David et.al., 2019 proofs that learnability is undecidable

# [Silver,2015] claims...

- “TD does not follow the gradient of any objective function”
  - “This is why TD can diverge when off-policy or using non-linear function approximation”
- Gradient TD follows true gradient of projected Bellman error

| On/Off-Policy | Algorithm   | Table Lookup | Linear | Non-Linear | On/Off-Policy | Algorithm       | Table Lookup | Linear | Non-Linear |
|---------------|-------------|--------------|--------|------------|---------------|-----------------|--------------|--------|------------|
| On-Policy     | MC          | ✓            | ✓      | ✓          | On-Policy     | MC              | ✓            | ✓      | ✓          |
|               | TD          | ✓            | ✓      | ✗          |               | TD(0)           | ✓            | ✓      | ✗          |
|               | Gradient TD | ✓            | ✓      | ✓          |               | TD( $\lambda$ ) | ✓            | ✓      | ✗          |
| Off-Policy    | MC          | ✓            | ✓      | ✓          | Off-Policy    | MC              | ✓            | ✓      | ✓          |
|               | TD          | ✓            | ✗      | ✗          |               | TD(0)           | ✓            | ✗      | ✗          |
|               | Gradient TD | ✓            | ✓      | ✓          |               | TD( $\lambda$ ) | ✓            | ✗      | ✗          |

| Algorithm           | Table Lookup | Linear | Non-Linear     |
|---------------------|--------------|--------|----------------|
| Monte-Carlo Control | ✓            | (✓)    | ✗ (chattering) |
| Sarsa               | ✓            | (✓)    | ✗              |
| Q-learning          | ✓            | ✗      | ✗              |
| Gradient Q-learning | ✓            | ✓      | ✗              |

# UCB and Softmax function

## Softmax

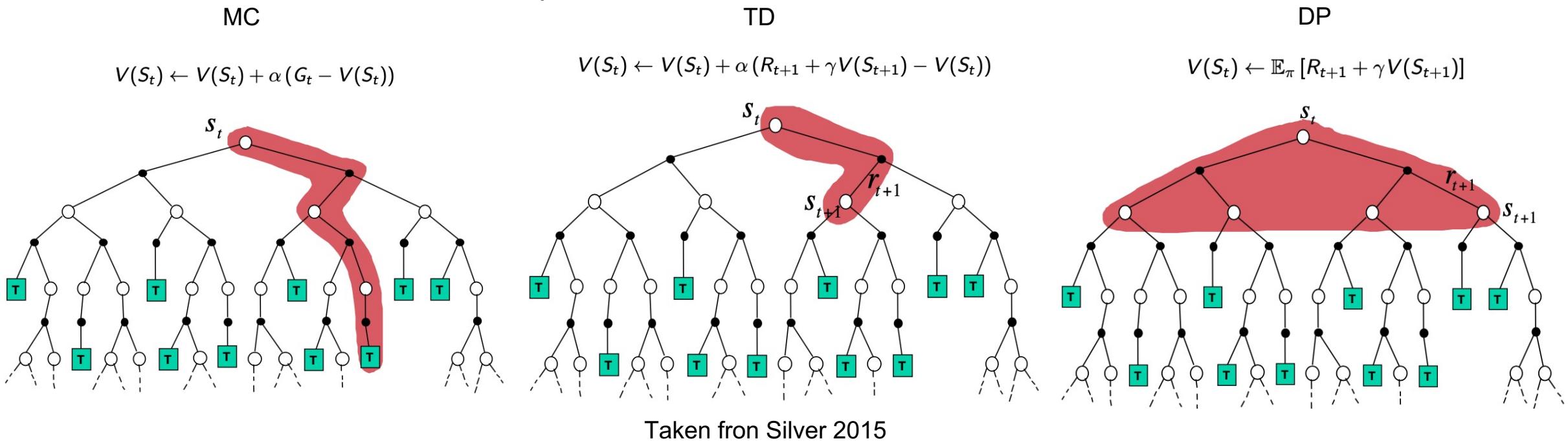
- It is often better to choose the action based on the current estimates of the distribution of rewards
  - Rather than returning a single action, the agent returns probabilities of each action weighted by expected rewards of each state
  - Action with highest preference  $h$  has the highest probability of being selected

$$\pi(s|a, \vartheta) = \frac{e^{h(s,a,\vartheta)}}{\sum_b e^{h(s,b,\vartheta)}}, h \in \mathbb{R} \text{ preference for a state-action pair}$$

- Another option is upper confidence bound (UCB)
  - Agent does not explore randomly but explores specifically states it has not seen before
  - UCB adds bonus to each action for inadequately sampled states

# Monte-Carlo (MC) vs Dynamic Programming DP approaches

- MC is modell-free and samples
  - MC techniques choose a single action but have to simulate full episode
- MC learns directly from experience but needs complete episodes to learn  
=> terminal state required in MDP
- DP is model-based as it requires all transition probabilities and does not sample
  - DP methods look one step ahead and iterate over all actions



# Benchmarking Algorithms – Policy Gradient (PG) vs Q-Learning

- Research showed that soft Q-learning algorithms performed equivalent to policy grading algorithms, e.g., [Schulman et al., 2018]
  - Q learning increases the value of an action
  - Policy gradients increase the probability of selecting an action
    - Entropy-based policy gradients boost the probability of an action proportional to a function  $f(a)$

=> Both solve the same problem

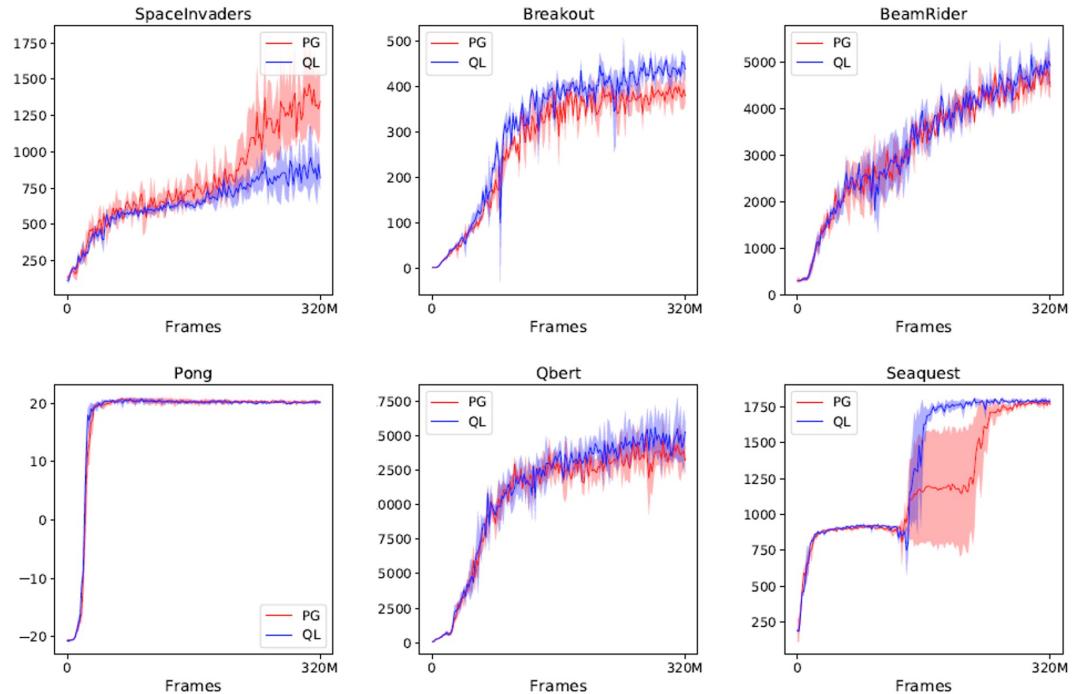


Figure 3: Atari performance with policy gradient vs  $Q$ -learning update rules. Solid lines are average evaluation return over 3 random seeds and shaded area is one standard deviation.

# Benchmarking for continuous control tasks by [Winter, 2020]

- Challenge: Researchers mostly publish plots and Winter had to read the Sum of expected returns
- Result: SAC and PPO perform over a significant number of studies similar

Table 7-1. Performance results from continuous control benchmarks for SAC and PPO

| Environment              | SAC <sup>a</sup> | SAC <sup>b</sup> | SAC <sup>c</sup> | SAC-<br>auto <sup>x,d</sup> | PPO <sup>e</sup> | PPO <sup>f</sup> | SAC-<br>auto <sup>g</sup> | PPO <sup>h</sup> |
|--------------------------|------------------|------------------|------------------|-----------------------------|------------------|------------------|---------------------------|------------------|
| MuJoCo Walker            | 3475             | 3941             | 3868             | 4800                        | 3425             | 3292             |                           |                  |
| MuJoCo Hopper            | 2101             | 3020             | 2922             | 3200                        | 2316             | 2513             |                           |                  |
| MuJoCo Half-cheetah      | 8896             | 6095             | 10994            | 11000                       | 1669             |                  |                           |                  |
| MuJoCo Humanoid          | 4831             |                  | 5723             | 5000                        | 4900             | 806              |                           |                  |
| MuJoCo Ant               | 3250             | 2989             | 3856             | 4100                        |                  |                  |                           |                  |
| MuJoCo Swimmer           |                  |                  | 42               |                             | 111              |                  |                           |                  |
| AntBulletEnv-v0          |                  |                  |                  |                             |                  | 3485             | 2170                      |                  |
| BipedalWalker-v2         |                  |                  |                  |                             |                  | 307              | 266                       |                  |
| BipedalWalkerHardcore-v2 |                  |                  |                  |                             |                  | 101              | 166                       |                  |
| HalfCheetahBulletEnv-v0  |                  |                  |                  |                             |                  | 3331             | 3195                      |                  |
| HopperBulletEnv-v0       |                  |                  |                  |                             |                  | 2438             | 1945                      |                  |
| HumanoidBulletEnv-v0     |                  |                  |                  |                             |                  | 2048             | 1286                      |                  |

| Environment                         | SAC <sup>a</sup> | SAC <sup>b</sup> | SAC <sup>c</sup> | SAC-<br>auto <sup>x,d</sup> | PPO <sup>e</sup> | PPO <sup>f</sup> | SAC-<br>auto <sup>g</sup> | PPO <sup>h</sup> |
|-------------------------------------|------------------|------------------|------------------|-----------------------------|------------------|------------------|---------------------------|------------------|
| InvertedDoublePendulumBulletEnv-v0  |                  |                  |                  |                             |                  |                  |                           | 9357 7703        |
| InvertedPendulumSwingupBulletEnv-v0 |                  |                  |                  |                             |                  |                  |                           | 892 867          |
| LunarLanderContinuous-v2            |                  |                  |                  |                             |                  |                  |                           | 270 128          |
| MountainCarContinuous-v0            |                  |                  |                  |                             |                  |                  |                           | 90 92            |
| Pendulum-v0                         |                  |                  |                  |                             |                  |                  |                           | -160 -168        |
| ReacherBulletEnv-v0                 |                  |                  |                  |                             |                  |                  |                           | 18 18            |
| Walker2DBulletEnv-v0                |                  |                  |                  |                             |                  |                  |                           | 2053 2080        |

[a] Haarnoja, Tuomas, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. "Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor" (<https://oreil.ly/JjjwB>). ArXiv:1801.01290, August.

[b] Wang, Tingwu, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yeming Wen, Eric Langlois, Shunshi Zhang, Guodong Zhang, Pieter Abbeel, and Jimmy Ba. 2019. "Benchmarking Model-Based Reinforcement Learning" (<https://oreil.ly/9ZBoU>). ArXiv: 1907.02057, July.

[c] Wang, Che, and Keith Ross. 2019. "Boosting Soft Actor-Critic: Emphasizing Recent Experience without Forgetting the Past" (<https://oreil.ly/eigQt>). ArXiv:1906.04009, June.

[d] Haarnoja, Tuomas, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, et al. 2019. "Soft Actor-Critic Algorithms and Applications" (<https://oreil.ly/s3apP>). ArXiv:1812.05905, January.

[e] Dhariwal, Prafulla, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. 2017. OpenAI Baselines (<https://oreil.ly/9642Z>). GitHub Repository.

[f] Logan Engstrom et al. 2020. "Implementation Matters in Deep Policy Gradients: A Case Study on PPO and TRPO" (<https://oreil.ly/pOPdK>). ArXiv:2005.12729, May.

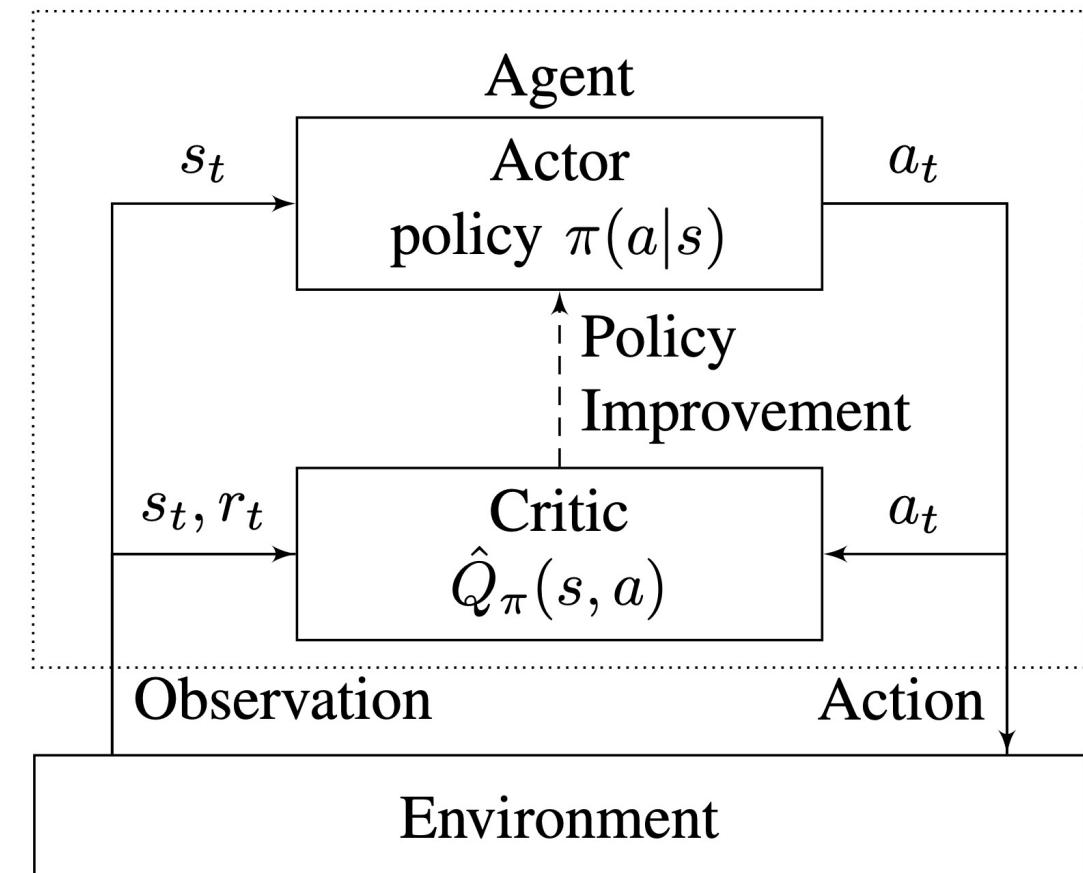
[g] Raffin, Antonin. (2018) 2020. Araffin/RL-Baselines-Zoo. Python. <https://oreil.ly/H3Nc2>

[h] Raffin, Antonin. (2018) 2020. Araffin/RL-Baselines-Zoo. Python. <https://oreil.ly/H3Nc2>

# Policy-based algorithms I: Advantage Actor Critic (A2C)

- Basic idea published in [Barto,Sutton&Anderson, 1983] and refined in [Kimura & Kobayashi ,1998]
  - Actor refers to the learned policy
  - Critic refers to the learned value function with parameters w
  - Policy  $\pi(a|s) = \pi(a|s, \vartheta)$  is a probability function differentiable w.r.t. its parameters  $\vartheta$
- => Continuous actions are easy to implement
- Policy gradient update rule is defined by

$$\vartheta = \vartheta + \alpha \gamma^t B \nabla_{\vartheta} \ln \pi(a|s, \vartheta)$$



Source: [Kämmerer, 2019] cites  
[Kimura & Kobayashi  
,1998]

# Policy-based algorithms I: Advantage Actor Critic (A2C)

- Policy gradient update rule is defined by

$$\vartheta = \vartheta + \alpha_\vartheta \gamma^i B_w \nabla_\vartheta \ln \pi(a|s, \vartheta)$$

- $B = G$

- => reinforce algorithm & no baseline

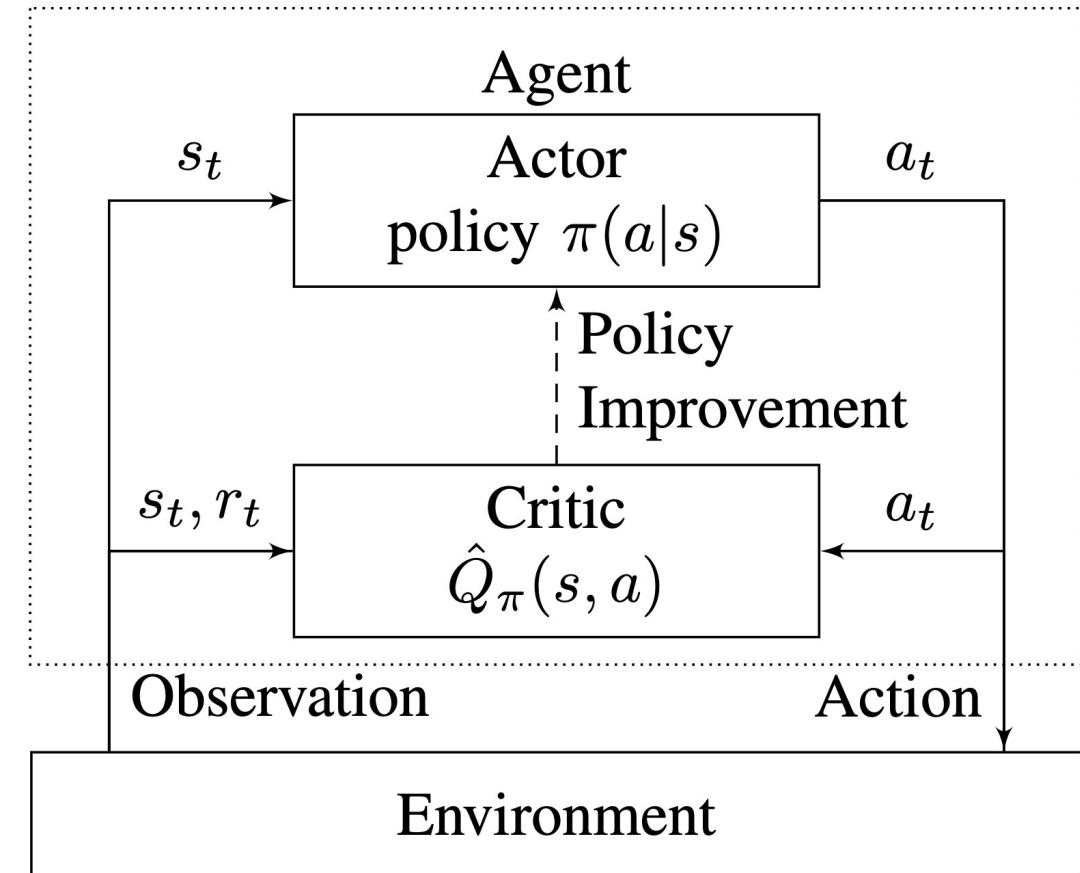
- $B_w = \delta_{t:t+1} = G_{t:t+1} - V(s, w)$

$$G_{t:t+1} = r + \gamma V(s_{t+1}, w)$$

=> reinforce with baseline

- $B_w = \delta_{t:t+1} * I$

=> n-step actor critic (A2C)



Source: [Kämmerer, 2019] cites  
[Kimura & Kobayashi  
, 1998]

# Algorithm of A2C

1: Input is a policy  $\pi(a|s, \vartheta)$  and state value function  $V(s, w)$

2: Initialize  $\vartheta_i$  and  $w$

3: loop for each episode

4:   Initialize environment to provide  $s$  and  $l=1$

5:   do

6:     Choose  $a$  in state  $s$  using  $\pi$  breaking ties randomly

a7:     Take action  $a$  and observe  $r_i, s_{i+1}$

b7:      $\delta_{t:t+n} = G_{t:t+n} - V(s, w)$

c7:      $w = w + \alpha_w \delta_{t:t+n} \nabla_w V(s, w)$

b8:     Update  $\vartheta$  using baseline  $B_w$

c8:      $l = \gamma l$

9:      $s = s_{i+1}$

10:    while  $s$  is not terminal

Q-learning is  
changed &  
extended by  
marked lines

Learns both policy and state-value function on every single step

- Actor updates the policy in the direction suggested by the critic and takes best action based on that policy in b8
- Baseline prediction in b7 criticizes the consequences of the action and is used to update the policy (*critic*)
- Critic alters the probability of actions to be taken by  $\delta$   
=> off-policy algorithm

# Advantage Function $\mathcal{A}$

- Given a predicted state-value function which is the expectation (average over all actions of a state, one can represent the individual actions relative to the average

$$\mathcal{A}_\pi = Q_\pi(s, a) - V_\pi(s, a)$$

- If advantage function  $\mathcal{A}_\pi$  is positive means that action taken by agent is more optimal
- If  $\mathcal{A}_\pi$  is negative, decrease the probability of the action to be taken next time
- For policy-based algorithms using the temporal difference error of Eq. 1

$$\mathcal{A}_\pi = r_{i+1} + \gamma Q_\pi(s_{i+1}, a_{i+1}) - V_\pi(s, a)$$

- When agent is following optimal policy, proof in [Bhatnagar et al., 2007]

$$= r_{i+1} + \gamma V_\pi(s_{i+1}, a_{i+1}) - V_\pi(s, a)$$

# Advantage Actor Critic (A2C)

- Objective function  $J$  is based on the expected rewards and model parameters  $\vartheta$

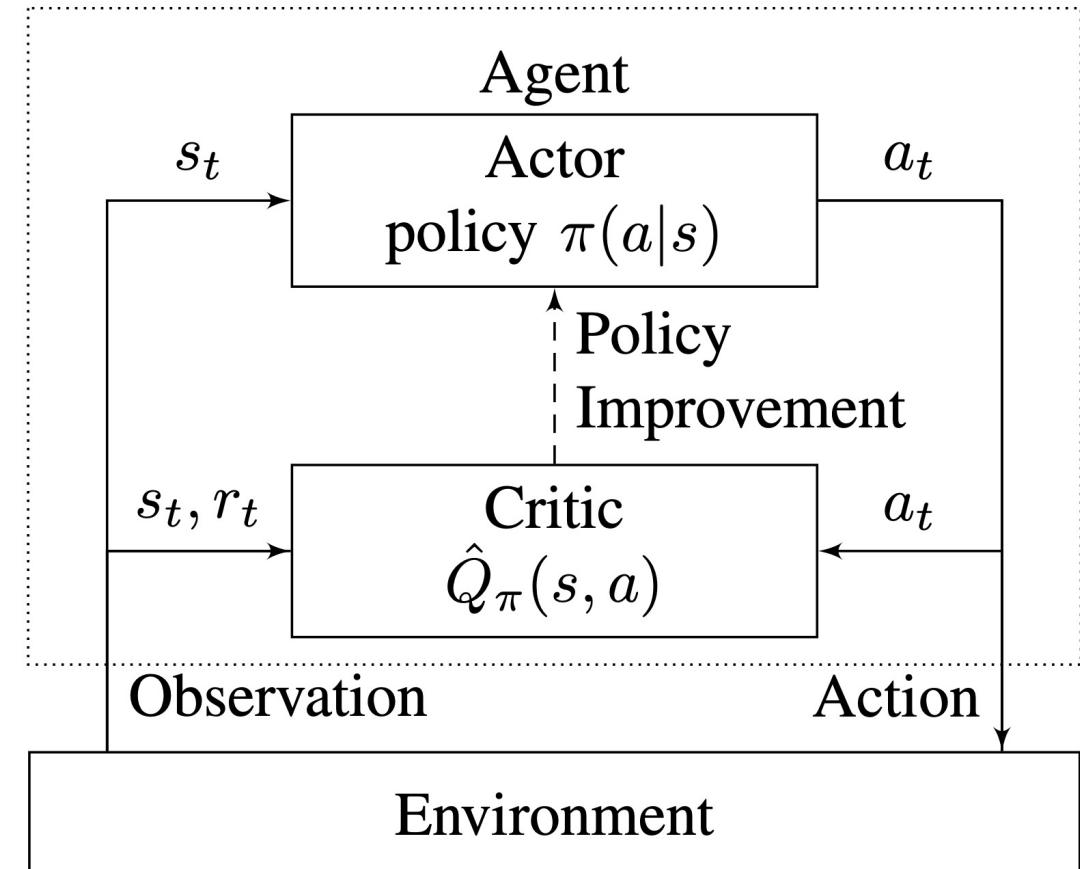
$$J(\vartheta) = E_{\pi}[G|s, a]$$

Parameters are optimized with update rule

$$\vartheta = \vartheta + \alpha(\nabla_{\vartheta} J(\vartheta))$$

With [Sutton&Barto, 1998 p.325]

$$\begin{aligned}\nabla_{\vartheta} J(\vartheta) &= \nabla_{\vartheta} E_{\pi}[G|s, a] = \nabla_{\vartheta} \sum_{a \in A} Q(s, a)\pi(a|s, \vartheta) \\ &= \dots = E_{\pi}[G \nabla_{\vartheta} \ln(\pi(a|s, \vartheta))|s, a]\end{aligned}$$



Source: [Kämmerer, 2019] cites  
[Kimura & Kobayashi  
1998]

# Soft Q-Learning [Ziebart et al., 2008]

- Using Shannon entropy  $H(\pi_\vartheta(s)) = -\sum_{a \in A} \pi_\vartheta(s, a) \log_b \pi_\vartheta(s, a)$  the maximum entropy objective function is

$$\pi * = \operatorname{argmax}_\pi \sum_{s \in S, a \in A} E[r + \beta H(\pi(a|s))]$$

- Large values of hyperparameter  $\beta$  promote nondeterministic actions and hence more exploration
- Beware: in deterministic problems with one optimal trajectory this may not work well
- BUT: Learning a stochastic policy forces the agent to learn **many** optimal solutions to the same problem

$$J(\vartheta) = E_\pi[r + \beta H]$$
$$\pi(s|a) \sim e^{Q_{soft}(s,a)/\beta}$$

$$Q_{soft}(s, a) = E[r + \gamma E[V_{soft}(s_{i+1})]]$$
$$V_{soft}(s) = E[Q(s, a) + \beta \log(\pi(a|s))]$$

- $\beta=0$  is the classical Q-Learning
- Main drawback is that several actions have to be sampled in the next state in order to estimate its current soft V-value, what makes it hard to implement in practice

# Soft Actor-Critic (SAC) algorithm [Haarnaoja et al., 2017]

- Objective Function for critic

$$J_Q(\vartheta) = E_\pi \left[ \left( \frac{1}{2} Q_{soft}(s, a) - (r + \gamma E[V_{soft}(s_{i+1})]) \right)^2 \right]$$

- Objective Function for actor

$$J_V(\varphi) = E_\pi \left[ \beta \log(\pi_\varphi(a|s)) - Q_{soft,\vartheta}(s, a) \right]$$

- Policy parameters can be learned by directly minimizing the expected KL-divergence
- Given,  $Z$  as a partition function to normalize the softmax

$$J(\sigma)_\pi = \mathbb{E}_{s_t \in \mathcal{D}} [D_{KL}(\pi_\theta(s, \cdot) \parallel \frac{\exp Q_\psi(s_t, \cdot)}{Z(s_t)})]$$

---

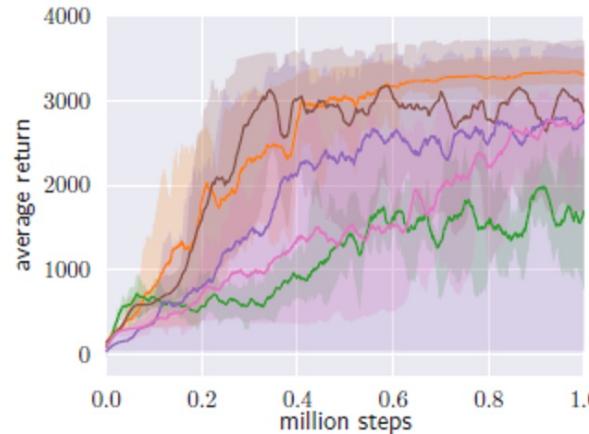
**Algorithm 1** Soft Actor-Critic

---

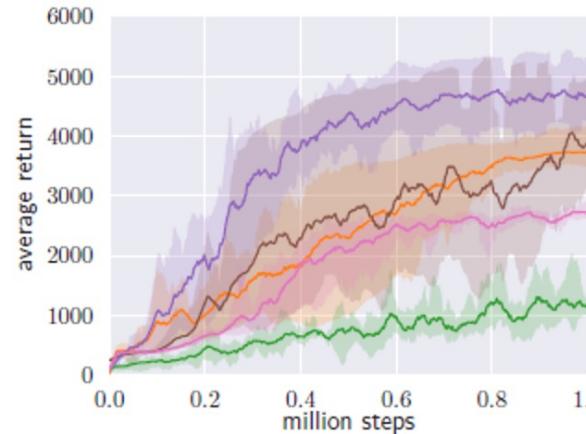
Initialize parameter vectors  $\psi, \bar{\psi}, \theta, \phi$ .  
**for** each iteration **do**  
    **for** each environment step **do**  
         $\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$   
         $\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$   
         $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$   
    **end for**  
    **for** each gradient step **do**  
         $\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$   
         $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$  for  $i \in \{1, 2\}$   
         $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$   
         $\bar{\psi} \leftarrow \tau \psi + (1 - \tau) \bar{\psi}$   
    **end for**  
**end for**

# Soft Actor-Critic (SAC)

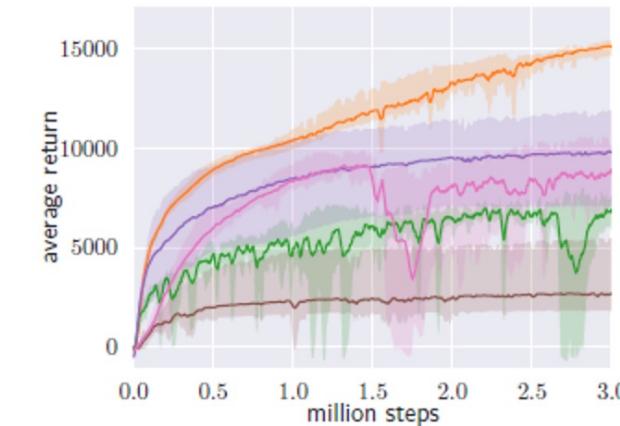
- SAC claims to Outperforms PPO [Haarnaoja et al., 2017]



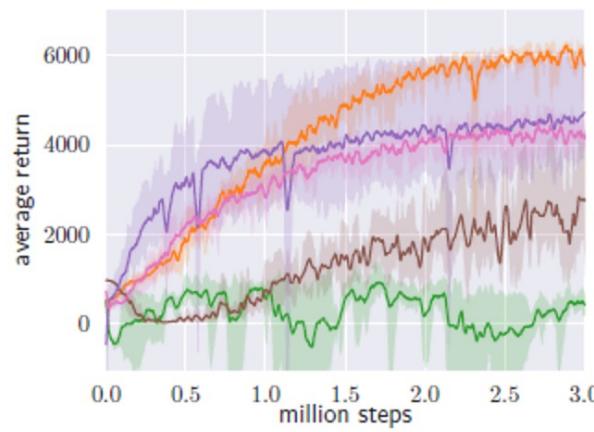
(a) Hopper-v1



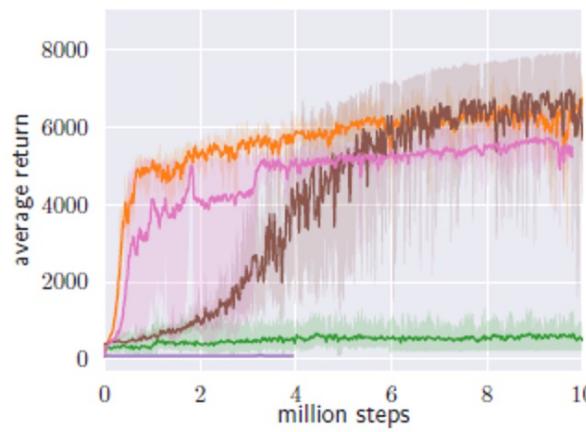
(b) Walker2d-v1



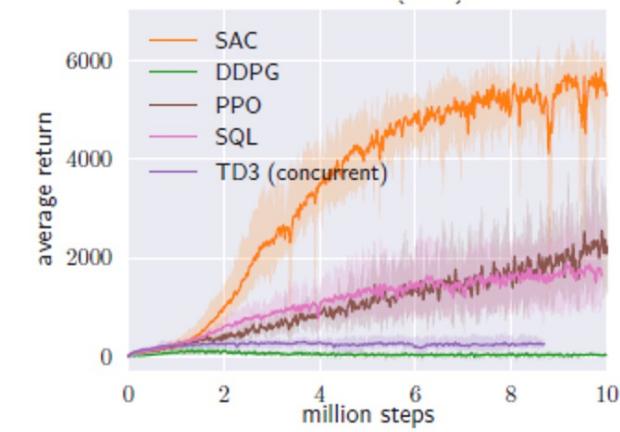
(c) HalfCheetah-v1



(d) Ant-v1



(e) Humanoid-v1



(f) Humanoid (rllab)

Figure 1. Training curves on continuous control benchmarks. Soft actor-critic (yellow) performs consistently across all tasks and outperforming both on-policy and off-policy methods in the most challenging tasks.

## Outlook policy-based algorithms: Proximal Policy Optimization (PPO)

- Published in [Schulman et al., 2015, PMLR] and extended [Schulman et al. 2017 ArXiv]:
- Objective function  $J$  is based on policy and advantage function

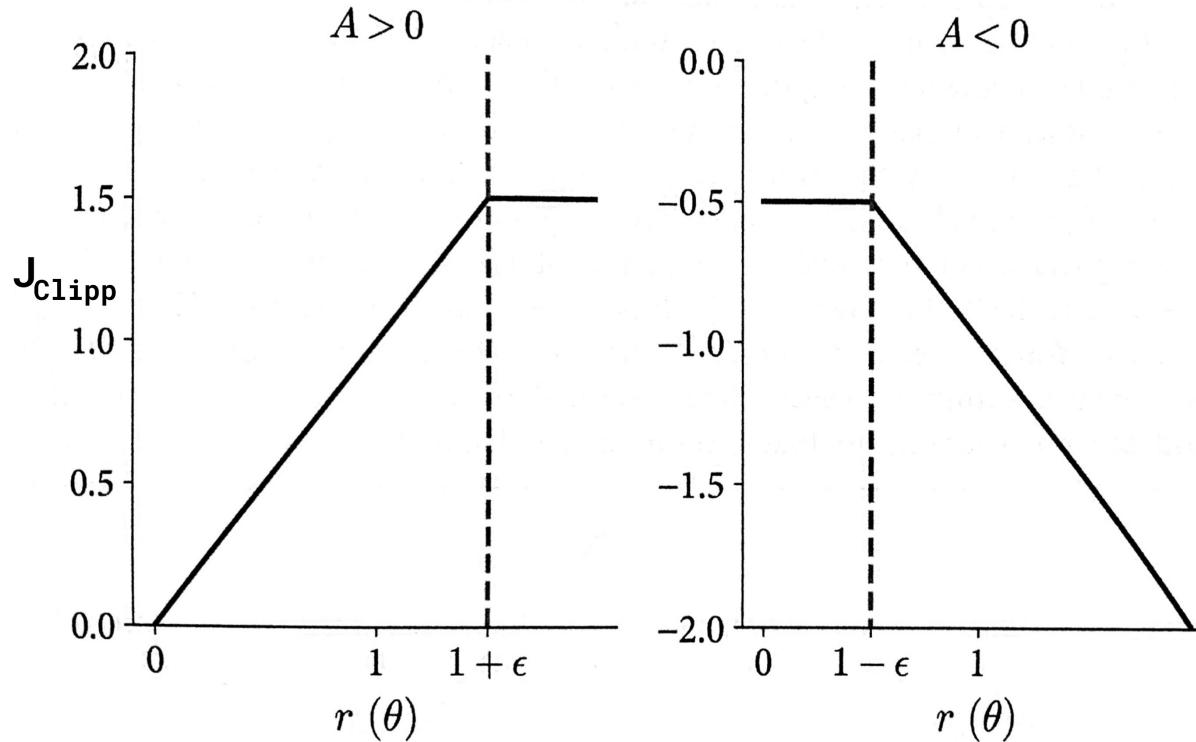
$$J(\vartheta) = E_{\pi} \left[ r((\vartheta) A_{\pi,old}(s, a) \right] \text{ with } r(\vartheta) = \frac{\pi_{\vartheta}(a|s)}{\pi_{\vartheta,old}(a|s)}$$

- Without a limitation on difference between  $\vartheta, old$  and  $\vartheta$ , maximization of  $J(\vartheta)$  would lead to instability with extremely large parameter updates and big policy ratios

$$\Rightarrow J_{clipp}(\theta) = \mathbb{E}[\min(r(\theta)A_{\pi,old}(s, a), clip(r(\theta), 1 - \epsilon, 1 + \epsilon)A_{\pi,old}(s, a))]$$

## Outlook policy-based algorithms II: Proximal Policy Optimization (PPO)

$$J_{clipp}(\theta) = \mathbb{E}[\min(r(\theta)A_{\pi,old}(s,a), clip(r(\theta), 1 - \epsilon, 1 + \epsilon)A_{\pi,old}(s,a))]$$



- Function *clip* clips the ratio to be no more than  $1+\epsilon$  and no less than  $1-\epsilon$ .
  - Objective function of PPO takes the minimum one between original value and clipped version
- => loses motivation for increasing the policy update to extremes for better rewards

# PPO Objective function and algorithm

$$J_{clipp+VF+S}(\theta) = E[J_{clipp} + C_1 L_{VF} + C_2 H(\pi_\theta, s)]$$

With  $C_1 L_{VF} = (V_\theta(s) - V_{target})^2$ , and  $H$  the entropy defined by [Mnih et al., 2016]

- $H$  improves the exploration by discouraging premature convergence to suboptimal deterministic policies
- $C_1, C_2$  hyperparameters
- $V_{target}$  is the prediction target of a neural network

Source: [Schulman et al. 2017 ArXiv]:

---

**Algorithm 1** PPO, Actor-Critic Style

---

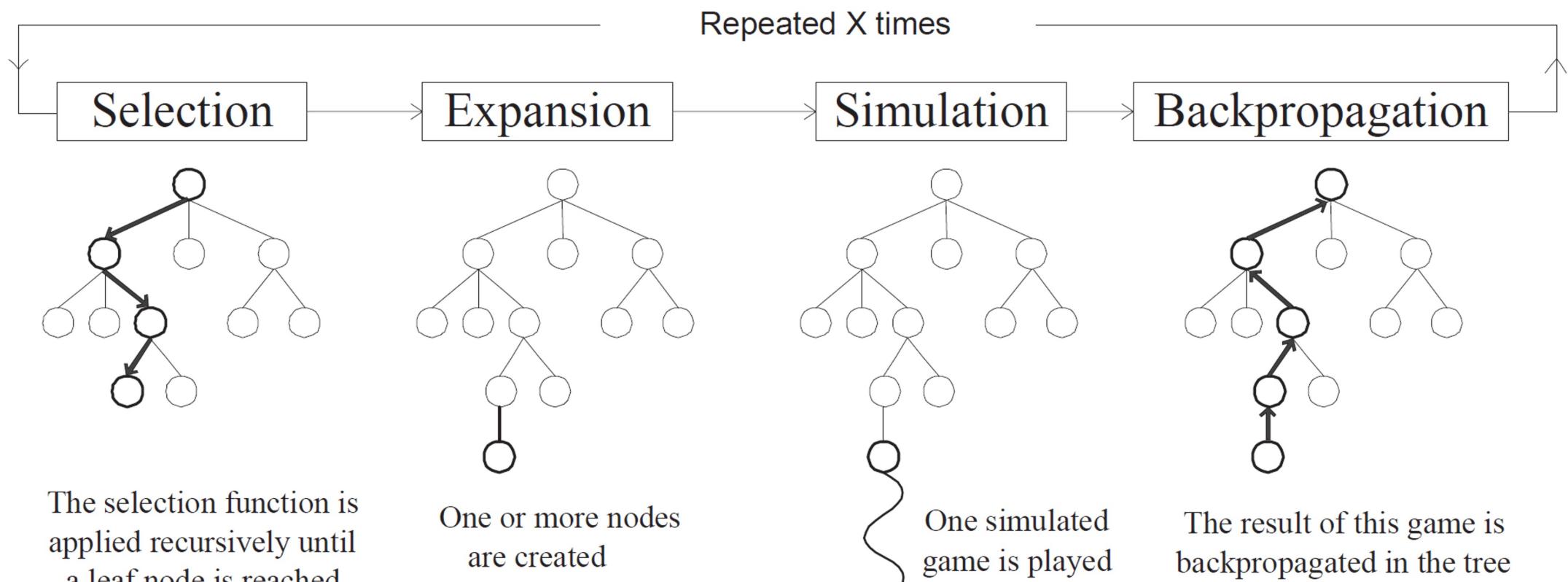
```
for iteration=1, 2, ... do
    for actor=1, 2, ..., N do
        Run policy  $\pi_{\theta_{old}}$  in environment for  $T$  timesteps
        Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
    end for
    Optimize surrogate  $L$  wrt  $\theta$ , with  $K$  epochs and minibatch size  $M \leq NT$ 
     $\theta_{old} \leftarrow \theta$ 
end for
```

# Application of AlphaGO zero has three key components [Silver at al., 2017]

1. Self-play generates training data in each time step  $i$  ( $s_i, \pi_i, z_i$ ),
  - $z_i$  being the game winner from perspective of current step
  - Two player play and are continually evaluated by policy evaluation
  - $(s_i, \pi_i, z_i)$  is taken from best-performing player at the moment
2. Monte-Carlo-Search-Tree (MCTS) output probabilities  $\pi_i = a_{\vartheta_{i-1}}(s_i)$  if playing each move by each player given neural network parameters  $\vartheta$  and root position  $s_i$ 
  - Tree search is restricted to either  $19*19*2=722$  moves or loss of game
3. Training data is used in supervised learning to train parameters  $\vartheta$  from data sample uniformly
  - Minimizing the error between predicted value  $v$  and the self-play winner  $z$
  - Maximize the similarity between the neural network move probabilities  $p$  to the search probabilities  $\pi_i$
$$(p, v) = f(\vartheta) \text{ and } loss = (z-v)^2 - \pi_T \log(p + c\|\vartheta\|^2), c \text{ controlling parameter to prevent overfitting}$$

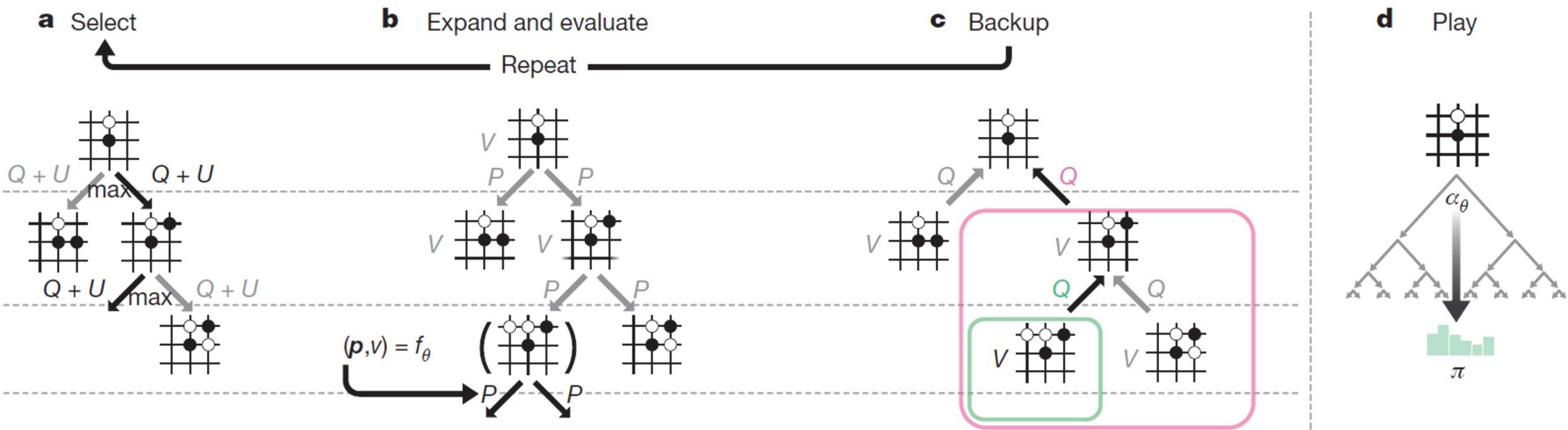
# Monte-Carlo-Search-Tree (MCTS)

- MCTS successively focuses multiple simulations at the current state by extending the initial portions of trajectories that have received high evaluations from earlier simulations
- MCTS executes as many iterations as possible before an action needs to be selected
- MCTS incrementally builds a tree whose root node represents the current state
- Each iteration consists of the four operations: Selection, Expansion, Simulation and Backup/Backpropagation [Chaslot et al., 2008]



# Monte-Carlo-Search-Tree (MCTS) in adapted [Silver at al., 2017]

- Each simulation traverses tree by selecting edge with maximum action value  $Q$ , plus an upper confidence bound  $U$  that depends on stored prior probability  $P$  & visit count  $N$  for that edge which is incremented once traversed
- Leaf node is expanded and the associated position  $s$  is evaluated by neural network  $(P(s, \cdot), V(s)) = f_\theta(s)$ 
  - Vector of  $P$  values are stored in the outgoing edges from  $s$
- Action value  $Q$  is updated to track mean of all evaluations  $V$  in the subtree below that action
- Once the search is complete, search probabilities  $\pi$  are returned, proportional to  $N_1/\tau$ , where  $N$  is the visit count of each move from the root state and  $\tau$  is a parameter controlling temperature



# AlphaGO Zero Results

- Using 64 GPU workers and 19 CPU servers AlphaGOZero learns for 40 days 500.000 games
- Components calculate in parallel using specific computation structures of Google cloud build for performing efficiently on these algorithms
- Thereafter it has enough experience in GO to beat prior AlphaGO(s) that were combined with human experience
  - Prior AlphaGO outperformed grandmaster in GO in 2016 [ilver et al., 2016].

# Analysis of Update Rule

$$\begin{aligned} Q_{t:t+1}(s_i, a) &= Q(s_i, a) + \alpha \left( r + \gamma * \operatorname{argmax}_{a_s \in A} Q(s_{i+1}) - Q(s_i, a) \right) \\ &= (1 - \alpha) Q(s_i, a) + \alpha \left( r + \gamma * \operatorname{argmax}_{a_s \in A} Q(s_{i+1}, a_s) \right) \end{aligned}$$

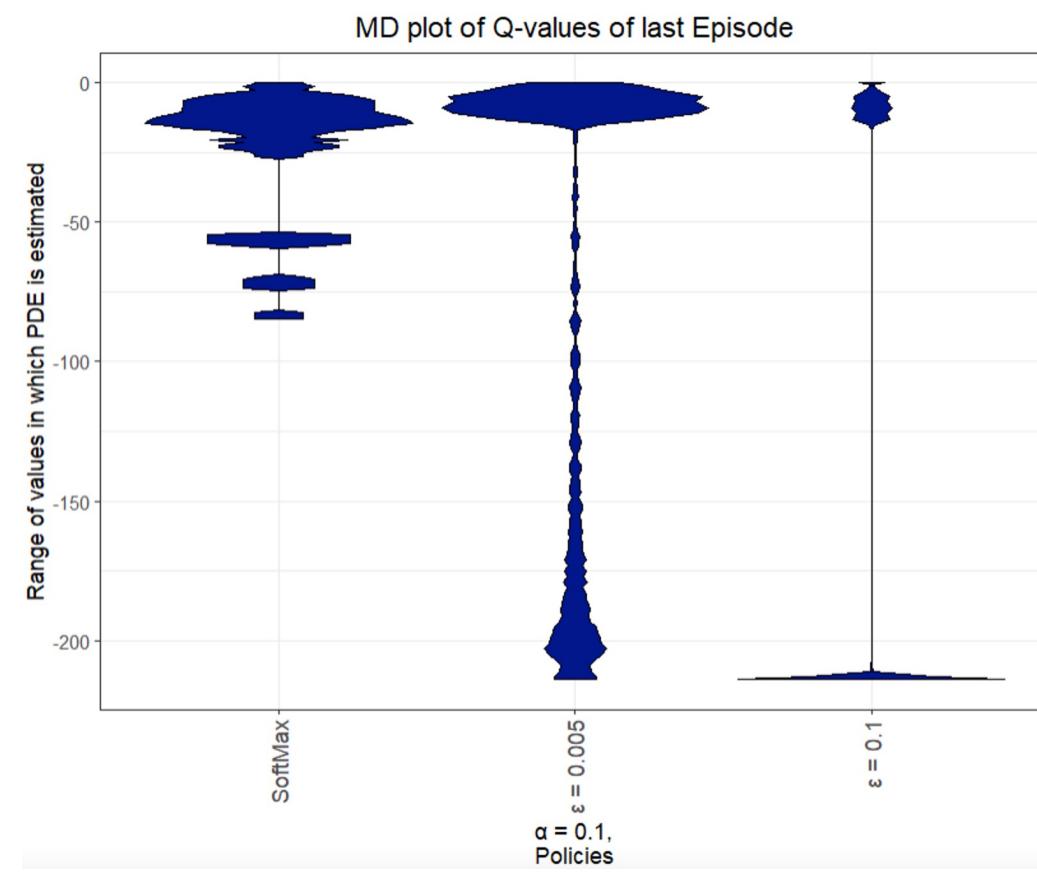
Is the exponential moving average (e.g., see [Awheda/Schwartz,2013])

$$\bar{x}_{n+1} = (1 - \alpha)\bar{x}_n + \alpha x_{n-1}$$

$$Q_*(s) = \operatorname{argmax}_{a_i \in A} E_\pi[G|s]$$

=> Computation of  $Q_*(s)$  assumes approximation by exponential moving average!

Q values of the 1000 Episode show multimodalities  
=> Assumption is not correct



# Improvements of Q-Learning

- In basic algorithm the next best action is chosen  
=> Intial bad update can create loop in which agent keeps making supoptimal decision
  - Solution: Use two look-up tables [Hasselt, 2010]
- Expected returns can be noisy
  - Solution: Buffer the rewards and update action value function after agent has visited the state a certain number of times [Strehl at al.,2006]
- Combine on- and off-policy Q learning by intruducing new parameter  $\sigma$  [Sutton & Barto, 2018]

# Extensions Q Learning II: eligibility trace

- Use tracer to update Q table backwards [Singh & Sutton, 1996]

After line 7:

- Temporarily calculate one step error  $\delta_{t:t+1}$
- Iterate  $z(s, a) = z(s, a) + 1$

Change Line 8 to:

- Add eligibility trace:

For each  $s \in S$  and for each  $a \in S$

$$Q(s_i, a) = Q(s_i, a) + \alpha(G_{t:t+n} - Q(s_i, a)) * z(s, a)$$
$$z(s, a) = \lambda z(s, a), 0 < \lambda < 1$$

# Difference between Sampling Updates and Expected Update

Value estimated

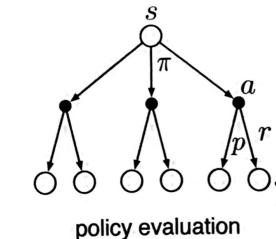
Expected updates (DP)

Sample updates (one-step TD)

In determinisct environment there is no difference

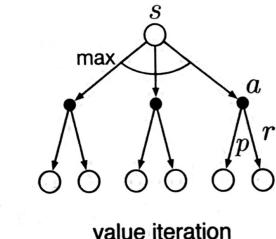
- Sampling Update is used in Q learning with  $Q_{i:i+1}(s_i, a) = Q(s_i, a) + \alpha \left( r + \gamma * \operatorname{argmax}_{a_s \in A} Q(s_{i+1}, a_s) - Q(s_i, a) \right)$
- Expected Updates is used in TD(0) with  $Q_{i:i+1}(s_i, a) = \sum_{s_{i+1}, r} P(s_{i+1}, r | s, a) \left( r + \gamma * \operatorname{argmax}_{a_s \in A} Q(s_{i+1}, a_s) - Q(s_i, a) \right)$
- BUT: In stochastic environment given s,a many next states are possible making sampling update computationally cheaper but influenced by the sampling error
- In general there are 8 backup strategy of updating Q and V possible
  1. Update V or Q
  2. Update on or off-policy
  3. Expected updates or sample updates

$v_\pi(s)$



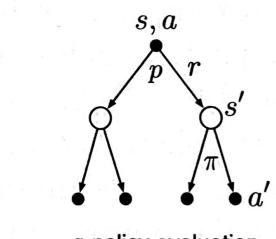
policy evaluation

$v_*(s)$



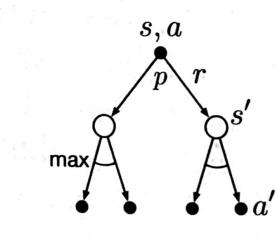
value iteration

$q_\pi(s, a)$



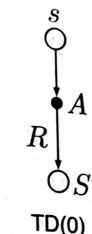
q-policy evaluation

$q_*(s, a)$



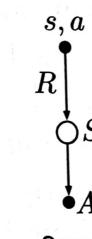
q-value iteration

Source:  
[Sutton & Barto, 1998]

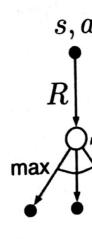


TD(0)

No useful update!



Sarsa



Q-learning

Figure 8.6: Backup diagrams for all the one-step updates considered in this book.

# Bellman Equation

Let  $s = s_t$  be the current state,  $G$  the discounted return then the Bellman equation defines a recursive relationship between values with

$$\begin{aligned} V_\pi(s) &= E_\pi[G_t|s] = E_\pi\left[\sum_{i=0}^T \gamma^i r_i |s\right] \\ &= E_\pi\left[\sum_{i=0}^T \gamma^i r_{t+i} |s\right] = E_\pi\left[r_t + \sum_{i=1}^T \gamma^i r_{t+i} |s\right] \\ &= E_\pi[r_t + \gamma(1 * r_{t+1} + \gamma r_{t+2} + \dots) |s] \\ &= E_\pi\left[r_t + \gamma \sum_{i=0}^T \gamma^i r_{t+1+i} |s\right] \\ &= E_\pi[r_t + \gamma G_{t+1} |s] = \dots = r_t + \gamma * E_\pi[G_{t+1} |s] \\ &= r_t + \gamma * V_\pi(s_{t+1}) \end{aligned}$$

# Bellman Equation

Let  $s = s_t$  be the current state,  $G$  the discounted return then the Bellman equation defines a recursive relationship between values with

$$V_\pi(s) = E_\pi[G_t|s] = r_t + \gamma * V_\pi(s_{t+1})$$

⇒ We have an optimal substructure and overlapping subproblems that can be solved recursively

⇒ Dynamic programming can be used in generalized policy iteration (GPI)

- $\pi_0 \xrightarrow{\text{Evaluation}} V_{\pi_0} \xrightarrow{\text{Iteration}} \pi_1 \xrightarrow{\text{Evaluation}} V_{\pi_1} \xrightarrow{\text{Iteration}} \pi_2 \dots \pi^*$

=> model-based RF algorithms if transition probabilities  $T$  given