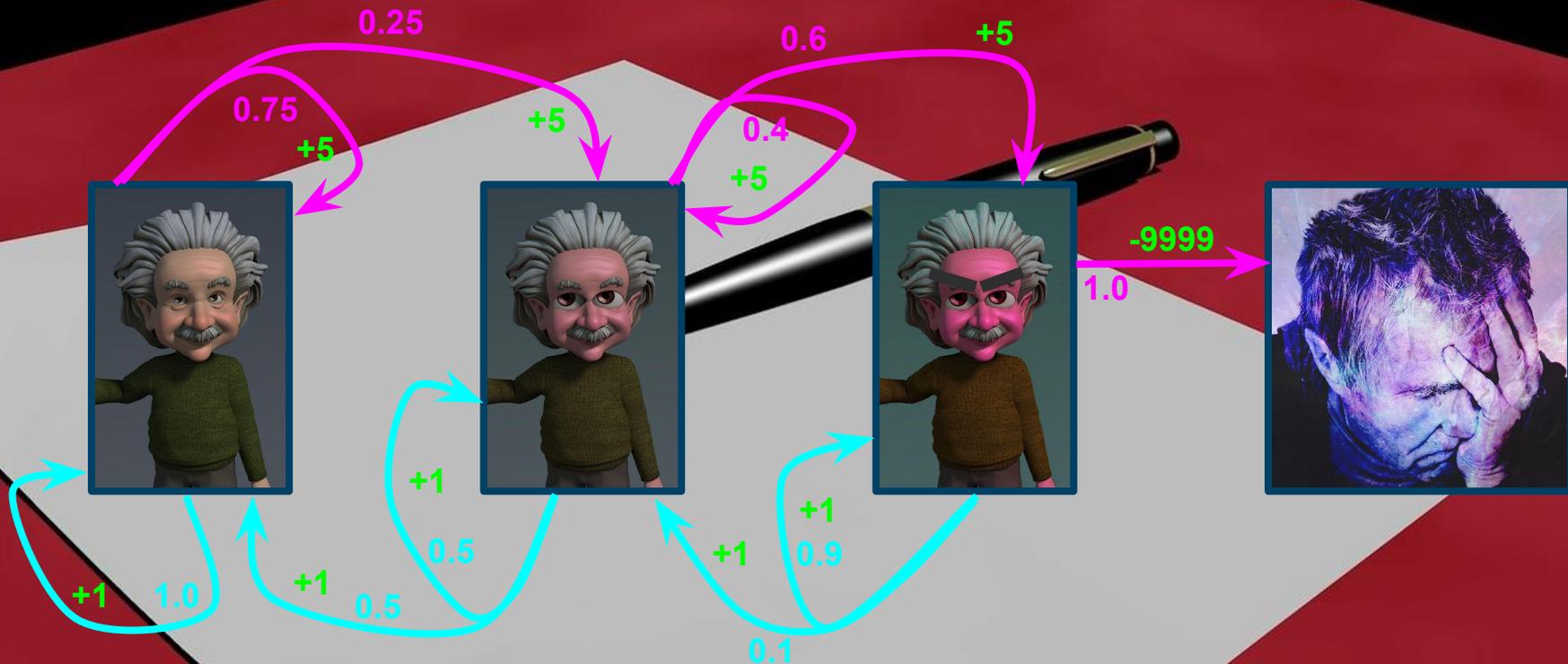


Einführung ins Reinforcement Learning

~ Patrick Dammann ~
~ Ist künstliche Intelligenz gefährlich? ~



how to cheat in an exam



markov decision process

the given problem:

- set of **states** $S = \{s_1, \dots, s_n\}$
- set of **actions** $A = \{a_1, \dots, a_m\}$
- **transition** between states via actions (and randomness)
- **rewards** for transitions
- markov property is given

$S =$



$A =$



(transition probability) (reward)
0.6 +5



markov decision process ~ cont.

$$\pi(\text{img}) = \rightarrow$$


what we want:

- maximize rewards (just a value)
- a **policy** π^* that defines the best action for every state
- $\pi: S \rightarrow A$

another example

- 4x4 **states** (visualized as 2d grid)
- 4 **actions** (north, south, east, west)
- chance to take random orthogonal direction (e.g. 10%)
- invalid movement results in transition to same state
- negative reward for moving
- terminal fields (colored) end game on any action, giving noted reward

	0	1	2	3
0				
1		X	-100	+100
2				
3	S	X	-100	-100

some notation

	0	1	2	3
0				
1		X	 -100	+100
2				
3	S	X	-100	-100

- $T(s, a, s')$ # transition
 - **probability** of getting into state s' when using action a in state s
 - $T(s_{0,2}, a_E, s_{0,3}) = 0.8$ # go E, as planned
 - $T(s_{0,2}, a_E, s_{0,2}) = 0.1$ # go N, bump against wall
 - $T(s_{0,2}, a_E, s_{1,2}) = 0.1$ # go S
 - $T(s_{0,2}, a_E, s_{0,1}) = 0.0$ # go W
- $R(s, a, s')$ # reward
 - **reward** for getting into state s' when using action a in state s
 - $R(s_{1,2}, _, _) = -100$ # lose
 - $R(s_{2,0}, a_N, s_{1,0}) = -2$ # moved
 - $R(s_{1,3}, _, _) = +100$ # win
- $\gamma \in [0,1]$:= discount factor
 - in timestep t , rewards are worth $\gamma^t \cdot R(s,a,s')$
 - makes sooner rewards worth more
 - $\gamma = 1$: don't care when rewards are achieved
 - $\gamma = 0$: only care about immediate rewards

the V-values

$V^*(s)$ is the **estimated reward** when **starting in s**, taking the **optimal action** and continue to **act optimally**.

- $V^*(s) = \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma \cdot V^*(s')]$
- $V^*(s) = \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma \cdot V^*(s')]$
- $V^*(s) = \max_a \sum_{s'} T(s,a,s') [\text{est. reward when taking action } a, \text{ landing in state } s' \text{ and continue acting optimally}]$
- $V^*(s) = \max_a \sum_{s'} [\text{est. reward when taking action } a, \text{ landing in state } s' \text{ and continue acting optimally, weighted by probability}]$
- $V^*(s) = \max_a [\text{est. reward when taking action } a \text{ and continue acting optimally}]$
- $V^*(s) = [\text{est. reward when taking } \text{best action} \text{ and continue acting optimally}]$

value iteration

$V^k(s)$ is the estimated reward when starting in s , taking the optimal action and continue to act optimally with only **k timesteps left**.

- find V^* via bottom up, iterative approach
- V^1 is known (by problem definition)
- calculate V^{k+1} via information from V^k

$$V^{k+1}(s) = \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma \cdot V^k(s')]$$

- for $k \rightarrow \infty$: $V^k \rightarrow V^*$

	0	1	2	3
0				
1		X	-100	+100
2				
3	S	X	-100	-100

value iteration ~ example

$$V^k(s) = \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma V^{k-1}(s')]$$

$$\gamma = 0.9$$

$$r = -2 \text{ (moving reward)}$$

$$p(\text{random}) = 0.2$$

$$\Rightarrow T(s, a_N, s_N) = 0.8$$

$$\Rightarrow T(s, a_N, s_W) = 0.1$$

$$\Rightarrow T(s, a_N, s_E) = 0.1$$

	0	1	2	3
0	-2	-2	-2	-2
1	-2	0	-100	+100
2	-2	-2	-2	-2
3	-2	0	-100	-100

value iteration ~ example steps

	0	1	2	3
0	-2	-2	-2	2
1	-2	0	-100	+100
2	-2	-2	-2	-2
3	-2	0	-100	-100



	0	1	2	3
0	-3.8	-3.8	-3.8	69.6
1	-3.8	0	-100	+100
2	-3.8	-3.8	-21.4	69.6
3	-3.8	0	-100	-100



	0	1	2	3
0	-5.4	-5.4	47.5	69.3
1	-5.4	0	-100	+100
2	-5.4	-5.4	30.1	66.1
3	-5.4	0	-100	-100

V^1 immediate reward for best option

V^2 probability weighted average over immediate rewards and discounted reward from there for best action

$V^3 = \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma \cdot V^2(s')]$

live demo



value iteration

policy extraction

what we want:

- maximize rewards
- a **policy** π^* that defines the best action for every state
- $\pi: S \rightarrow A$

- assume we have V^* , how to get π^* ?
- simulate **one timestep for every action**, take best action

$$V^*(s) = \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma \cdot V^*(s')]$$

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma \cdot V^*(s')]$$

policy extraction ~ example

- value iteration might give approximations for V^* converged to machine ϵ
- policy extraction then generates the optimal policy π^*
- we now have the **perfect action** for **every state** in a game with **unlimited time**

	0	1	2	3
0	49.5	59.1	70.1	82.6
1	41.0	0	-100	+100
2	33.6	28.2	34.9	76.3
3	27.0	0	-100	-100



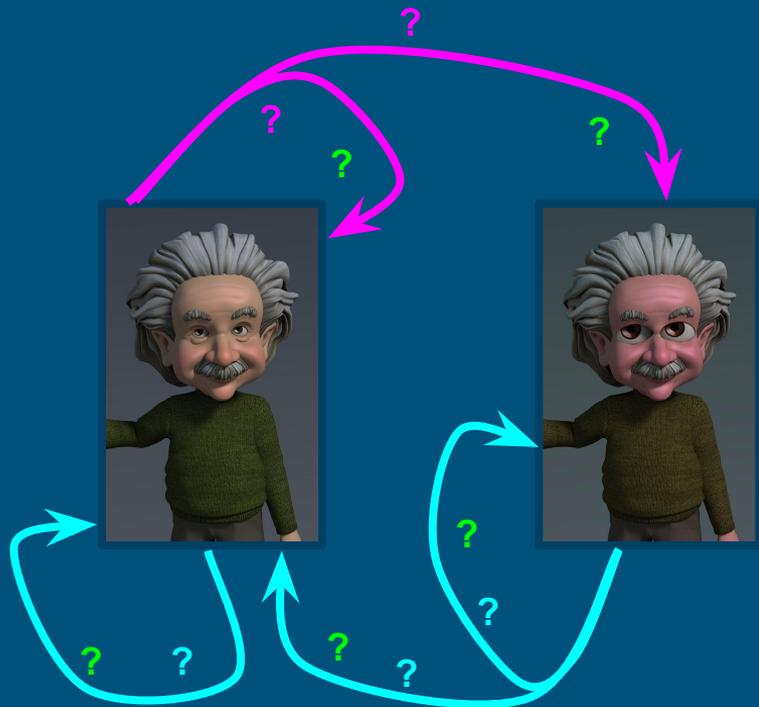
	0	1	2	3
0	⇒	⇒	⇒	↓
1	↑	X	*	*
2	↑	⇒	⇒	↑
3	↑	X	*	*

live demo



policy extraction

a slightly different problem

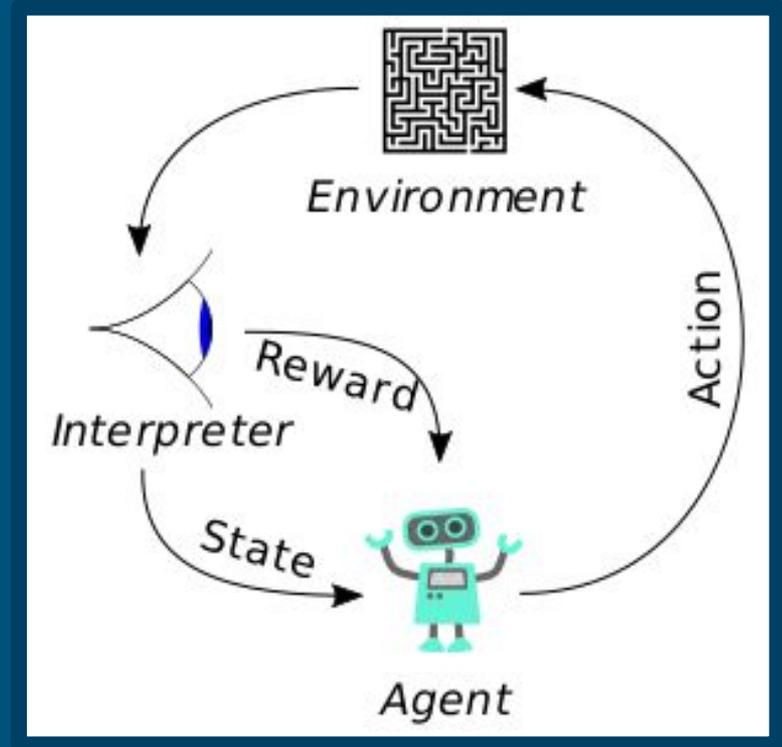


- assume having an MDP
- set of action $A = \{a_1, \dots, a_n\}$ is known
- current state s is known

$T(s,a,s')$ and $R(s,a,s')$ are **not** known and must be **determined by trial and error**

⇒ the agent must actively explore the environment

reinforcement learning



model based approach



$T(s,a,s')$ & $R(s,a,s')$



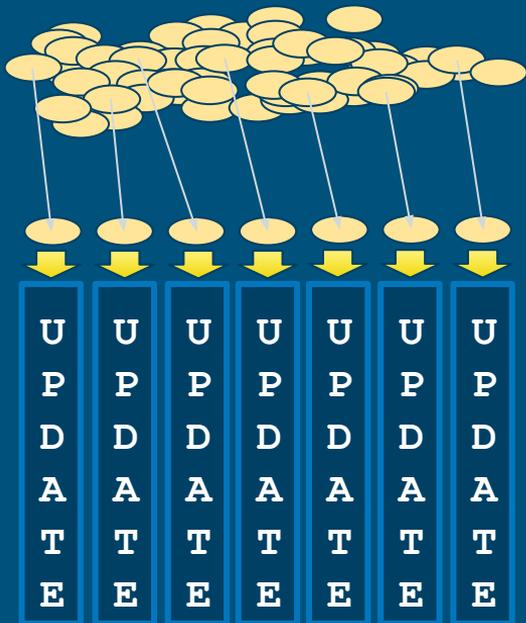
normal MDP

- **approximate** $T(s,a,s')$ and $R(s,a,s')$
 - by collecting as many samples as possible
- **solve** MDP
 - e.g. via value iteration and policy extraction

- problematic, since the agent often can't move "freely"
- requires **huge amount** of samples

Why not learn V directly, without a model?

temporal difference learning



- initialize V randomly

1. take action based on your policy
2. **update V based on your experience**
(only for state you came from)
3. **update policy**
4. if terminated, go to start state
5. go to 1.

(2. and 3. can be made batchwise every n actions)

td learning ~ update V based on your experience

Took action a in state s .

Landed in state s' gaining reward r .

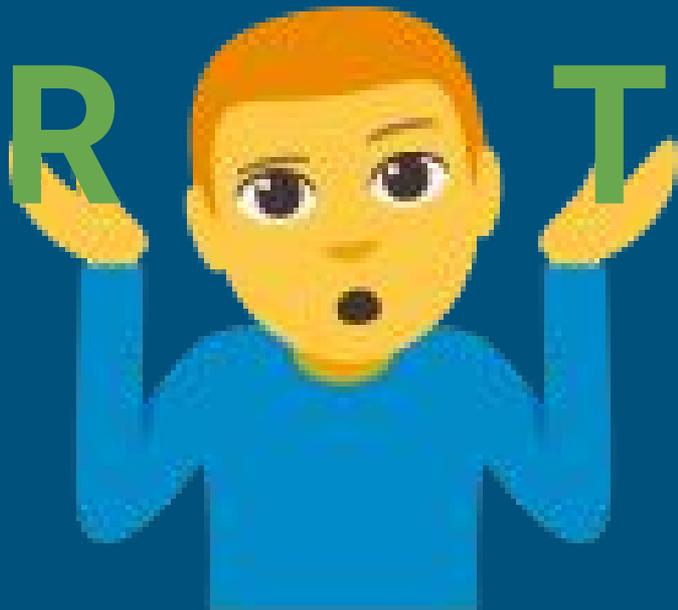
$$V(s) \leftarrow V(s) + \alpha(r + \gamma \cdot V(s') - V(s))$$

(α : learning rate)

- adjust $V(s)$ a little into the direction of the **sample**
- let α decay over time
- converges to V^* under certain circumstances



td learning ~ update policy



$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma \cdot V^*(s')]$$

- can't use policy extraction, since R and T are not present
- other methods not that trivial

Can we directly learn a policy?

the Q-values

$V^*(s)$ is the estimated reward when starting in s , taking the optimal action and continue to act optimally.

$Q^*(s,a)$ is the estimated reward when starting in s , taking the action a and continue to act optimally.

⇒ $|A|$ values per state instead of one

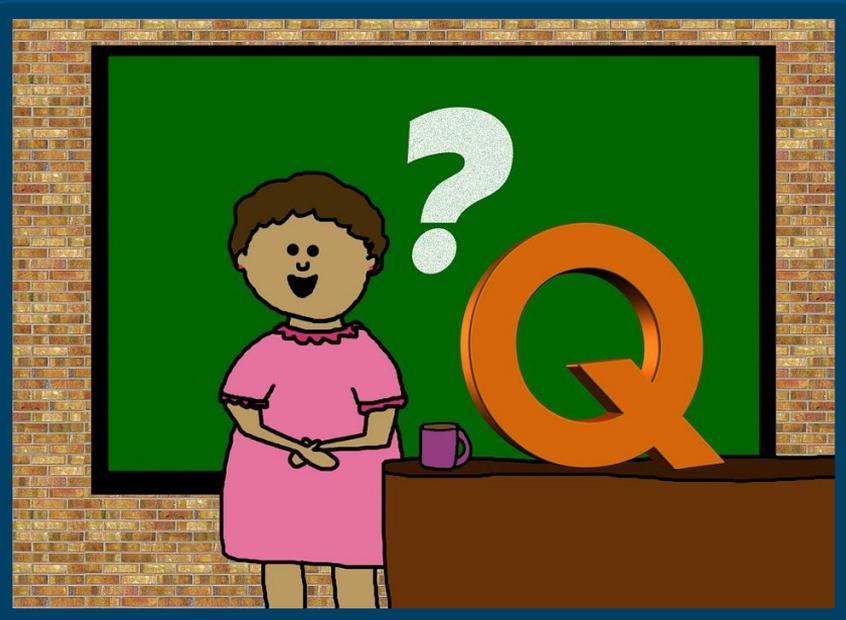
$$V(s) = \max_a Q(s,a)$$

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

	0	1	2	3
0	$\begin{matrix} 2.1 & 3.7 \\ 1.7 & 8.6 \end{matrix}$	$\begin{matrix} 3.4 & 9.3 \\ 3.4 & 9.3 \end{matrix}$	$\begin{matrix} 3.8 & 5.5 \\ 4.8 & 5.5 \end{matrix}$	$\begin{matrix} 3.8 & 8.2 \\ 5.7 & 8.2 \end{matrix}$
1	$\begin{matrix} 2.4 & 3.5 \\ 3.4 & 3.5 \end{matrix}$	X	-100	+100
2	$\begin{matrix} 3.2 & 3.3 \\ 3.2 & 3.3 \end{matrix}$	$\begin{matrix} 2.9 & -5 \\ 2.9 & -5 \end{matrix}$	$\begin{matrix} 0.4 & -15 \\ 0.4 & -15 \end{matrix}$	$\begin{matrix} 0 & -3 \\ -11 & -3 \end{matrix}$
3	$\begin{matrix} 3.0 & 2.2 \\ 2.2 & 2.2 \end{matrix}$	X	-100	-100

How to learn those Q-values?

Q-learning



- temporal difference learning on Q-values
- ⊖ needs **more samples** to converge
- ⊕ **policy** is learned **implicitly**
- ⊕ no V values needed

$$V(s) \leftarrow V(s) + \alpha(r + \gamma \cdot V(s') - V(s))$$

$$Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \cdot \mathbf{V}(s') - Q(s,a))$$

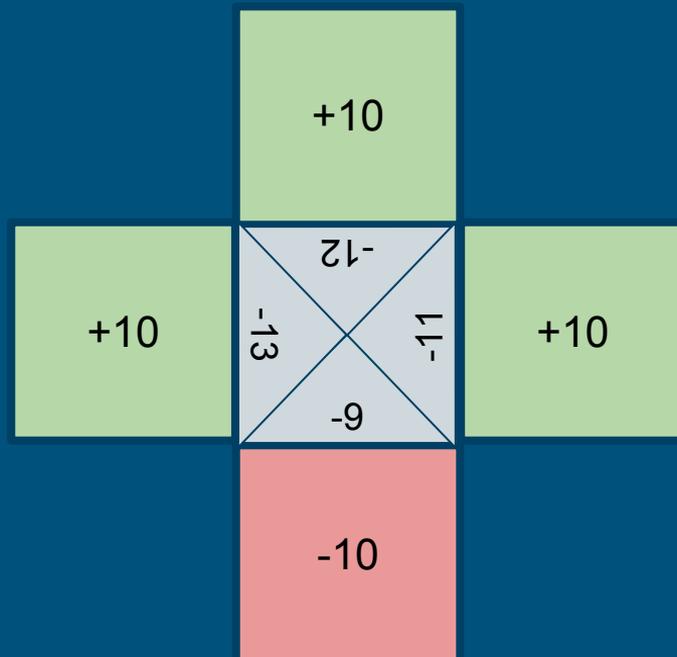
$$Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \cdot \mathbf{\max}_{a'} Q(s',a') - Q(s,a))$$

where to go?

imagine this situation:

- all Q-values are initialized with some random value
- **worst action a** has initially **highest Q-value**
- other actions have initially lower Q-values than the optimal Q-value of **a**

- ⇒ the agent will never try the other actions
- ⇒ we need to motivate him doing so



exploration

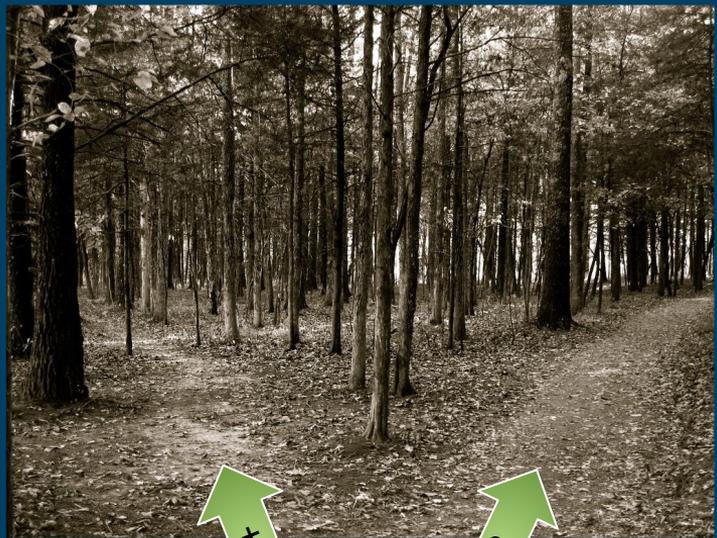
optimistic initial conditions

- estimate maximum Q-values
- **initialize all Q-values higher**
- updates will decrease Q-values
(since Q-learning converges)
- agent will prefer other action in later iterations

ϵ -greedy exploration

- new hyperparameter $0 \leq \epsilon \leq 1$
(this can change during training)
- agent will perform a **random action** with a **chance of ϵ**
- at test time, ϵ is usually set to 0

exploration ~ cont.



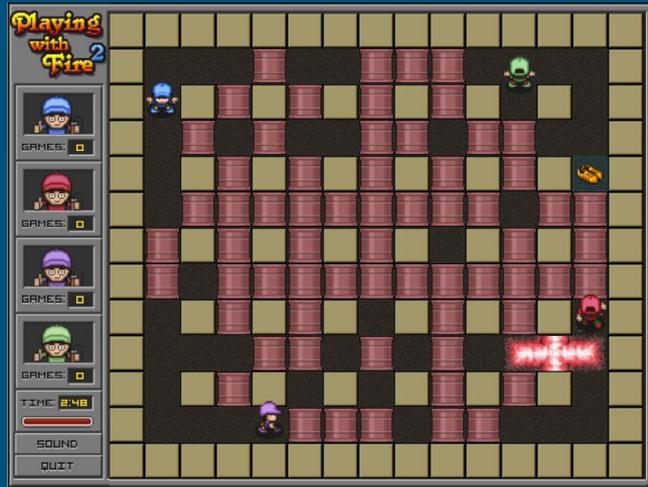
exploration function

- artificially boost Q-values of state-actions that were not used frequently
- choose an exploration function $E(s,a,n)$ with $E \rightarrow 0$ for $n \rightarrow \infty$ (n : number of state-action uses)
- add this function during Q-value updates

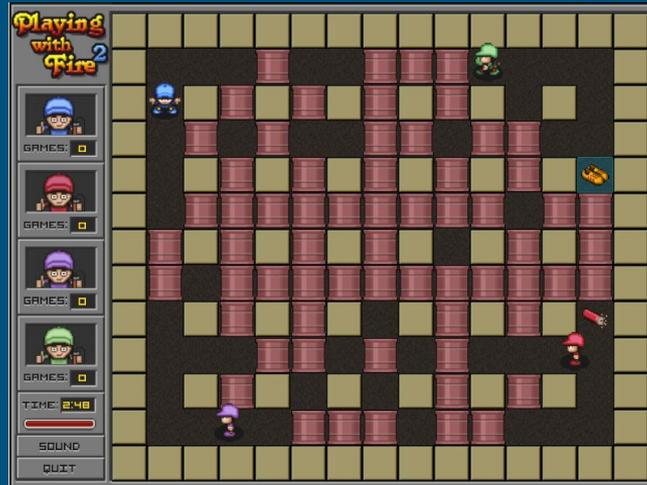
e.g.: $E(s,a,n) = k/n$, $k \in \mathbb{Q}$

$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma \cdot \max_{a'} (Q(s',a') + E(s',a',n_{s,a'})) - Q(s,a)]$$

the state space and its redundancies



- these states are **completely different** for our agent
- the **optimal action** here is (most likely) **the same**



state features

- generate **powerful features** from states
- use features to generate **continuous Q-function** instead of lookup table

$$Q(s,a) = w_1 \cdot f_1(s,a) + w_2 \cdot f_2(s,a) + \dots + w_n \cdot f_n(s,a)$$

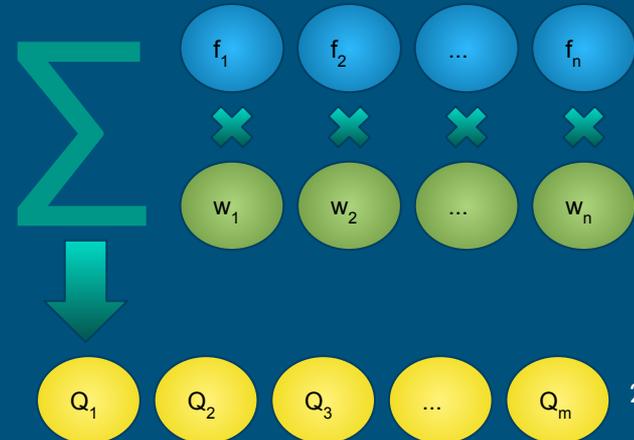
- during training: **tweak weights instead of entries**

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot \langle error \rangle$$

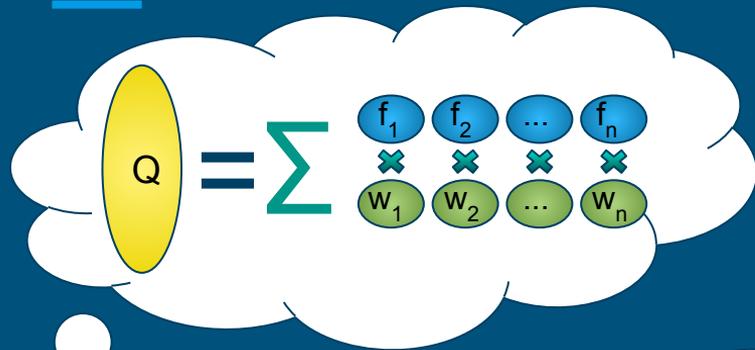
$$w_i \leftarrow w_i + \alpha \cdot \langle error \rangle \cdot f_i(s,a)$$



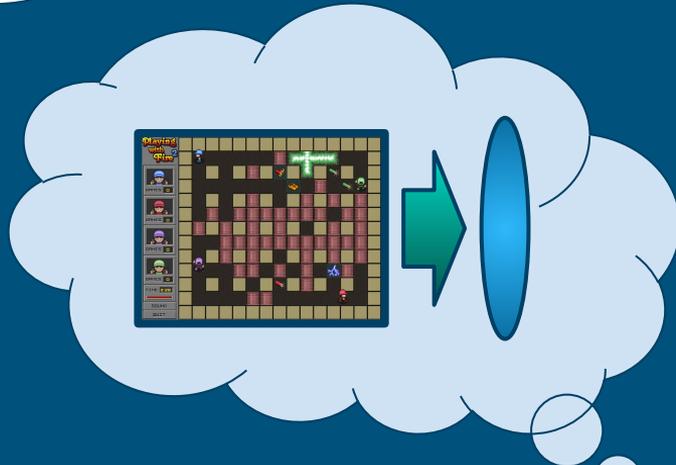
feature
generation



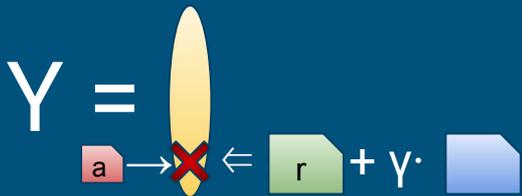
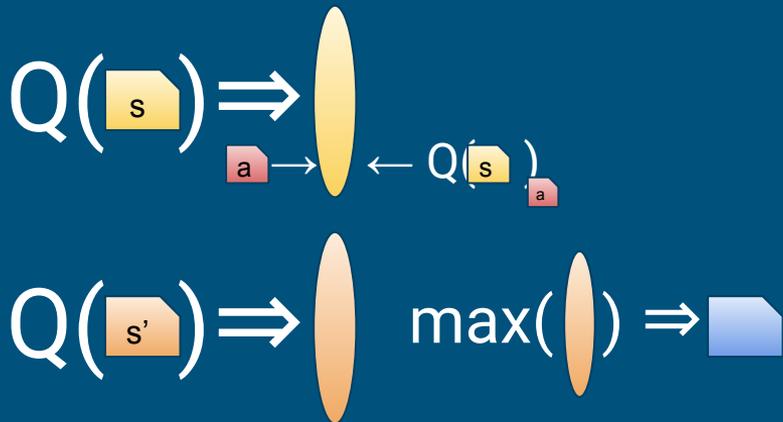
does this sound familiar..?



$$w_i \leftarrow w_i + \alpha \cdot \langle \text{error} \rangle \cdot f_i(s,a)$$



deep Q-learning



- Q-function is approximated by **neural network**
 - input: state
 - output: vector w/ Q-values for all actions
- CNNs allow use of pixel data (game screens, camera) as input
- train with the same samples as in normal Q-learning (s,a,r,s')
- output “label” for training contains:
 - $r + \gamma \cdot \max_a Q(s')_a$ for a action taken
 - $Q(s)_b$ for all other actions $b \in A$
(\Rightarrow no error)

sources

- AI Course CS188 (University Berkley)
 - <http://ai.berkeley.edu/home.html>
 - https://www.youtube.com/channel/UCB4_W1V-KfwpTLxH9jG1_iA/videos
- Harmon & Harmon: Reinforcement Learning: A Tutorial
 - <http://web.cs.iastate.edu/~honavar/rltut.pdf>
- qlearning4k (deep Q-learning framework)
 - <https://github.com/farizrahman4u/qlearning4k>

The End.

Thank you for your attention!

Any Questions?
