

OpenAI: Hide and Seek Game

Multiple learning agents split in two teams: hiders and seekers



[go to web](#)

- Observations: other agents and blocks
- Actions: move, shift blocks
- Rewards to hiders: -1 if any hider seen by a seeker, +1 if all hiders hidden. Seekers get opposite rewards.

Observations = States

How does the example fall into our framework? The *observation* set X is *finite*. Observations in RL are also called **states**.

The initial state x_1 is distributed by $P(x_1)$ independently of y_1 (so y_1 is irrelevant). For $k > 1$, states are distributed as

$$x_k \stackrel{c}{\sim} P(x_k | x_{k-1}, y_k) \quad (1)$$

which is a special case of **(??)**. As the distribution is the same for all $k > 1$, we will also use the index-less form

$$P(x' | x, y) \quad (2)$$

where x' is the state immediately following state x .

Terminal States

A state $x \in X$ is **terminal** if

$$P(x' | x, y) = P(x')$$

Intuitively, when a terminal state x is reached, the environment is 'restarted', i.e. the next state x' is sampled independently of x, y from the same distribution as the initial state x_1 .

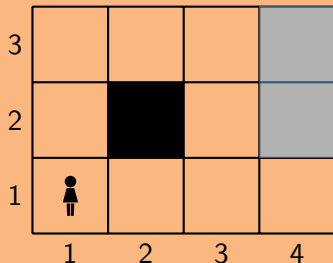
An **episode** is the span of interaction

- starting with x_1 or a state immediately following a terminal state
- ending with the next terminal state.

Episodes model e.g. repeated games where terminal states may be 'game lost' or 'game won'.

Example: Grid World - States

Consider a simpler example with a single agent:





X: 11 possible agent's positions on the grid.

$x_1 = (1, 1)$ with probability 1

$(4, 2), (4, 3)$ are *terminal*, always followed by state $(1, 1)$

Example: Grid World - Actions

Actions $Y = \{ \text{left, right, up, down} \}$

 0.1	0.8		
0.1		0.8	
	0.1		0.1

Uncertain effects: prescribed direction with 0.8 probability,
each of perpendicular directions with 0.1 probability

Bouncing: if new state out of grid, agent stays put

State Transitions under Policy

Due to (1), the probability of a sequence of states x_1, x_2, \dots given actions y_1, y_2, \dots is

$$P(x_1, x_2, \dots | y_1, y_2, \dots) = P(x_1)P(x_2 | x_1, y_2)P(x_3 | x_2, y_3) \dots \quad (3)$$

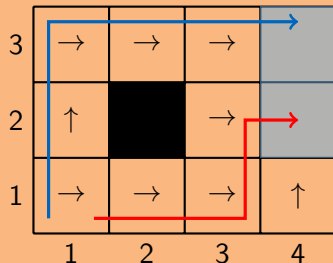
Define for $k > 1$

$$P^\pi(x_k | x_{k-1}) = P(x_k | x_{k-1}, \pi(x_{k-1}))$$

where π is an agent's policy. So (3) can be written as

$$P^\pi(x_1, x_2, \dots) = P(x_1)P^\pi(x_2|x_1)P^\pi(x_3|x_2) \dots \quad (4)$$

Example: State Transitions under Policy



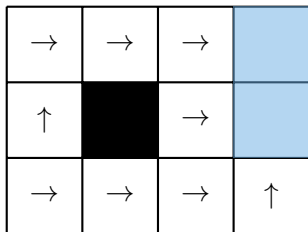
- Arrows indicate policy π for all states except terminal (irrelevant)
- Probability of blue path and red path:

$$P((1,1)) \cdot P^\pi((1,2) | (1,1)) \cdot \dots = 1 \cdot 0.1 \cdot 0.8 \cdot 0.8 \cdot 0.8 \cdot 0.8$$

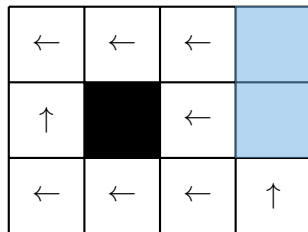
$$P((1,1)) \cdot P^\pi((2,1) | (1,1)) \cdot \dots = 1 \cdot 0.8 \cdot 0.8 \cdot 0.1 \cdot 0.8$$

Proper Policy

A **proper policy** is a policy such that from each state $x \in X$ there is a sequence of state transitions that has non-zero probability and that ends a terminal state.



proper



not proper

Rewards are instant, i.e., as in (??) but since x_k (1) depends on the same conditions x_{k-1}, y_k , we can assume instead a functional dependence:

$$r_k = r(x_k) \quad (5)$$

(This is because r_k can be formally made a part of x_k and 'extracted' by the function above.)

Rewards are thus a *function* of states.

Example: Grid World - Rewards

r_s	r_s	r_s	1
r_s		r_s	-1
r_s	r_s	r_s	r_s

Terminal states : 'win' (1) and 'lose' (-1) rewards

Non-terminal states : a small punishment r_s (e.g., $r_s = -0.05$) for every step taken

State Utility under a Proper Policy

For a $\gamma \in [0; 1]$, define the **utility** $U^\pi(x)$ **of state** $x \in X$ **under a proper policy** π as

$$U^\pi(x) = \begin{cases} r(x) & \text{if } x \text{ is terminal} \\ r(x) + \gamma \sum_{x' \in X} P^\pi(x' | x) U^\pi(x') & \text{otherwise} \end{cases} \quad (6)$$

$U^\pi(x)$ is thus recursively defined as the expected sum of γ -discounted rewards in an episode starting in state x . Since an episode has a finite length, $U^\pi(x)$ converges even with $\gamma = 1$ (i.e., if γ dropped).

This is similar to U^π (??) but U^π is the expected sum of γ -discounted rewards in the entire interaction consisting of an infinite succession of episodes.

State Utility under a Proper Policy (cont'd)

Given π , $U^\pi(x)$ for all $x \in X$ can be calculated by solving $|X|$ linear equations. Example, abbreviating $U^\pi((i,j))$ as u_{ij} :

3	→	→	→	
2	↑		→	
1	→	→	→	↑
	1	2	3	4

$$u_{43} = 1 \text{ (terminal)}$$

$$u_{42} = -1 \text{ (terminal)}$$

$$u_{41} = r_s + \gamma(0.8 \cdot u_{42} + 0.1 \cdot u_{31} + 0.1 \cdot u_{41})$$

etc.

Theorem 1

For any $x, x' \in X$, $U^\pi(x)$ and $U^\pi(x')$ are maximized by the same policies:

$$\bar{\pi} = \arg \max_{\pi} U^\pi(x) = \arg \max_{\pi} U^\pi(x')$$

and if $\gamma < 1$, these policies also maximize U^π :

$$\bar{\pi} = \arg \max_{\pi} U^\pi$$

Intuitive proof: the best action to take in any state x evidently does not depend on how x was reached - whether going from x or from x' . Thus $\arg \max_{\pi} U^\pi(x) = \arg \max_{\pi} U^\pi(x')$. Since $\bar{\pi}$ is optimal for each starting state, i.e., for each episode, the discounted sum U^π over all episodes is also maximized by $\bar{\pi}$. Hence the last equality.

State Utility and Policy Iteration

$\bar{\pi}$ is called the **optimal policy** and the **utility of state** $x \in X$ is defined as the utility of x under the optimal policy, i.e.

$$U(x) = U^{\bar{\pi}}(x) \quad (7)$$

An optimal policy can be computed by **policy iteration**, which is similar in spirit to the EM algorithm. Start with a random initial policy π and iterate

- 1 $U(x) := U^{\pi}(x)$
- 2 compute a new policy π such that

$$\pi(x) := \arg \max_y \sum_{x' \in X} P(x' | x, y) U(x') \quad (8)$$

until the policy does not change.

Value Iteration

As an alternative to policy iteration, the following method can be used to calculate $U(x)$ and $\bar{\pi}$:

- 1 Calculate $U(x)$, $x \in X$ as a solution to the set of Bellman equations

$$U(x) = r(x) + \gamma \max_{y \in Y} \sum_{x' \in X} P(x' | x, y) U(x') \quad (9)$$

one for each non-terminal $x \in X$ (for a terminal x , $U(x) = r(x)$).
This is called **value iteration**.

- 2 Calculate $\bar{\pi}$ by (8).

Unlike in policy iteration, the two steps are not *iterated*. However, (9) is not linear so an iterative algorithm is needed to find its solution.

Example: Grid World - Optimal Policy and Utilities

-0.04	-0.04	-0.04	1
-0.04		-0.04	-1
-0.04	-0.04	-0.04	-0.04

rewards

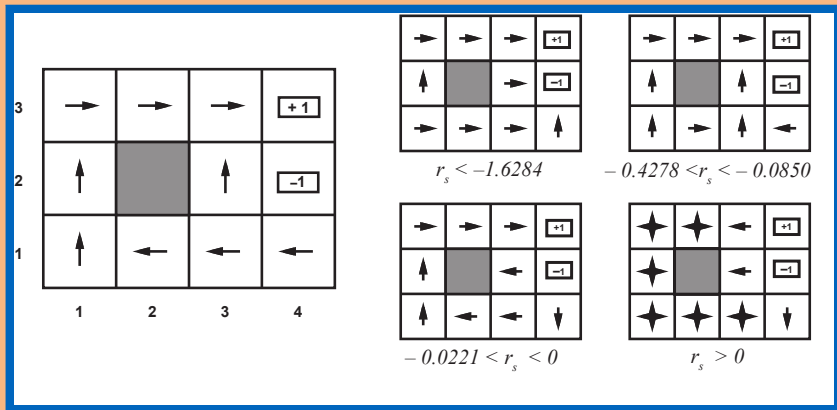
→	→	→	
↑		↑	
↑	←	←	←

optimal policy

0.812	0.868	0.918	1
0.762		0.660	-1
0.705	0.655	0.611	0.388

state utilities

Example: Optimal Policies Under Different Rewards



Optimal policy for $\gamma = 1$ and $r_s = -0.04$ for non-terminal states.

Optimal policies for $\gamma = 1$ and different ranges of r_s for non-terminal states.

Reinforcement Learning

The goal of RL is to find a good (high-utility) policy in an *unknown* environment, i.e., the agent does not know P (2) and r (5).

After each action y following a state x , the agent observes the new state x' and receives $r(x')$.

From these (x', x, y) samples, it can learn an estimate $p(x' | x, y)$ of $P(x' | x, y)$.

- for small $|X|, |Y|$, p can be just a 3-dimensional array of relative frequencies
- larger $|X|, |Y|$ may require a parametric approach (such as a Bayes Network if independence structure available)

Learning $r(x)$ is easier - just store the reward for each state x visited

Exploration vs. Exploitation

What actions should the agent take while collecting the (x', x, y) and r samples? Consider two options:

- 1 *Random* actions. Reason: a diverse sample will be collected, agent will not get stuck looping over a subset of states.
- 2 *Greedy* actions, i.e. optimal w.r.t. the current estimates

$$\pi(x) = \arg \max_y \sum_{x' \in X} p(x' | x, y) \hat{U}(x') \quad (10)$$

Reason: the agent's goal is to maximize rewards.

This is the exploration-exploitation dilemma we have already investigated (although with a different estimated quantity).

$\hat{U}(x')$ is the estimate of $U(x')$ computed with p instead of P .

Greedy Approach

→	→	→	1
↑		↑	-1
↑	←	←	←

Optimal policy
($r_s = -0.04$)

→	→	→	1
↓		↑	-1
→	→	↑	↓

Policy learned with
greedy actions (10)

After each greedy action (10), the agent re-computes π by policy iteration using $p(x' | x, y)$ estimated from collected samples of x', x, y, r .

This process converges into a *non-optimal* policy.

GLIE Approach

To balance btw. exploration and exploitation, we can adopt this strategy, i.e. instead of (10), take a random action with probability approaching 0.

This approach is called **GLIE** = *greedy in the limit of infinite exploration*.

A popular instance of GLIE is to take action y in state x with probability

$$\frac{e^{Q(x,y)/\tau}}{\sum_{y' \in Y} e^{Q(x,y')/\tau}} \quad (11)$$

where

$$Q(x,y) = \sum_{x' \in X} p(x' | x,y) \hat{U}(x') \quad (12)$$

and the *temperature* $\tau \rightarrow 0$ with $k \rightarrow \infty$.