

# CS 2750 Machine Learning

## Lecture 22

# Reinforcement learning

Milos Hauskrecht

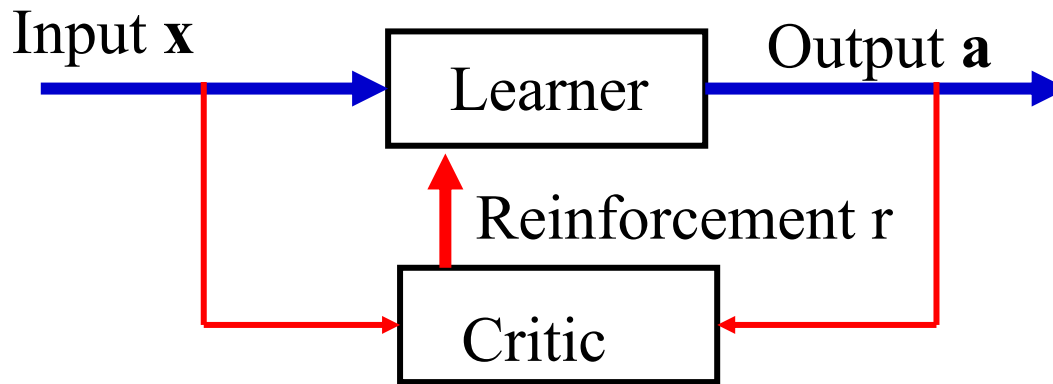
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)

5329 Sennott Square

---

# Reinforcement learning

## Basics:

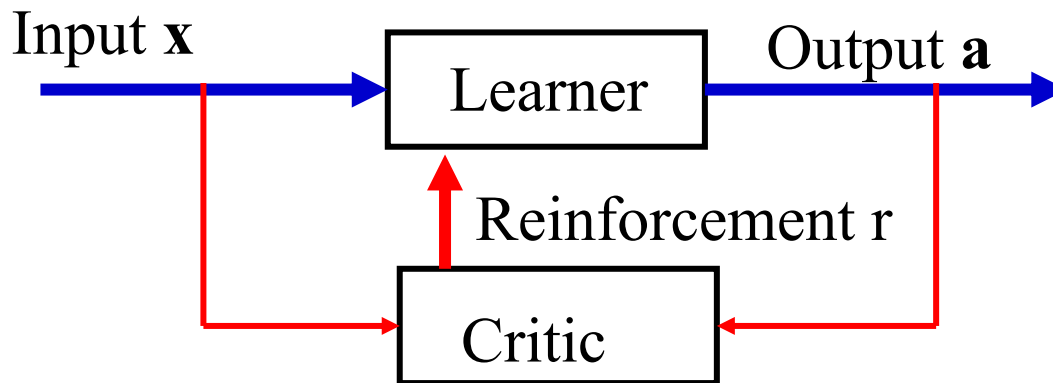


- **Learner interacts with the environment**
    - Receives input with information about the environment (e.g. from sensors)
    - Makes actions that (may) effect the environment
    - Receives a reinforcement signal that provides a feedback on how well it performed
-

# Reinforcement learning

**Objective:** Learn how to act in the environment in order to maximize the reinforcement signal

- The selection of actions should depend on the input
- A policy  $\pi : X \rightarrow A$  maps inputs to actions
- **Goal:** find the optimal policy  $\pi : X \rightarrow A$  that gives the best expected reinforcements



**Example:** learn how to play games (AlphaGo)

---

# Gambling example



- **Game:** 3 biased coins
  - The coin to be tossed is selected randomly from the three coin options. The agent always sees which coin is going to be played next. The agent makes a bet on either a head or a tail with a wage of \$1. If after the coin toss, the outcome agrees with the bet, the agent wins \$1, otherwise it loses \$1

- **RL model:**




- **Input:**  $X$  – a coin chosen for the next toss,
- **Action:**  $A$  – choice of head or tail the agent bets on,
- **Reinforcements:**  $\{1, -1\}$

- **A policy**  $\pi : X \rightarrow A$

**Example:**  $\pi :$ 

Coin1	$\rightarrow$ head
Coin2	$\rightarrow$ tail
Coin3	$\rightarrow$ head

$\pi :$ 

	$\rightarrow$ head
	$\rightarrow$ tail
	$\rightarrow$ head







# Trajectories

- Environment + Agent's actions in time generate State, action, reward trajectories

**Example:** Assume the agent applies the following policy

$$\pi : \left| \begin{array}{l} \text{Coin1} \rightarrow \textit{head} \\ \text{Coin2} \rightarrow \textit{tail} \\ \text{Coin3} \quad \textit{head} \end{array} \right|$$

**One possible SAR trajectory:**

	Step0	Step1	Step2	..	Step k	..
<b>state</b>	Coin2	Coin1	Coin2	..	Coin1	..
<b>action</b>	Tail 	Head 	Tail 	.. 	Head 	.. 
<b>reward</b>	-1	1	1	..	1	..

# Measuring the quality of the policy

- The quality of the policy can be measured in terms of the total rewards received by following the policy

**Example:** Assume the agent applies the following policy

$$\pi : \left| \begin{array}{l} \text{Coin1} \rightarrow \textit{head} \\ \text{Coin2} \rightarrow \textit{tail} \\ \text{Coin3} \quad \textit{head} \end{array} \right|$$

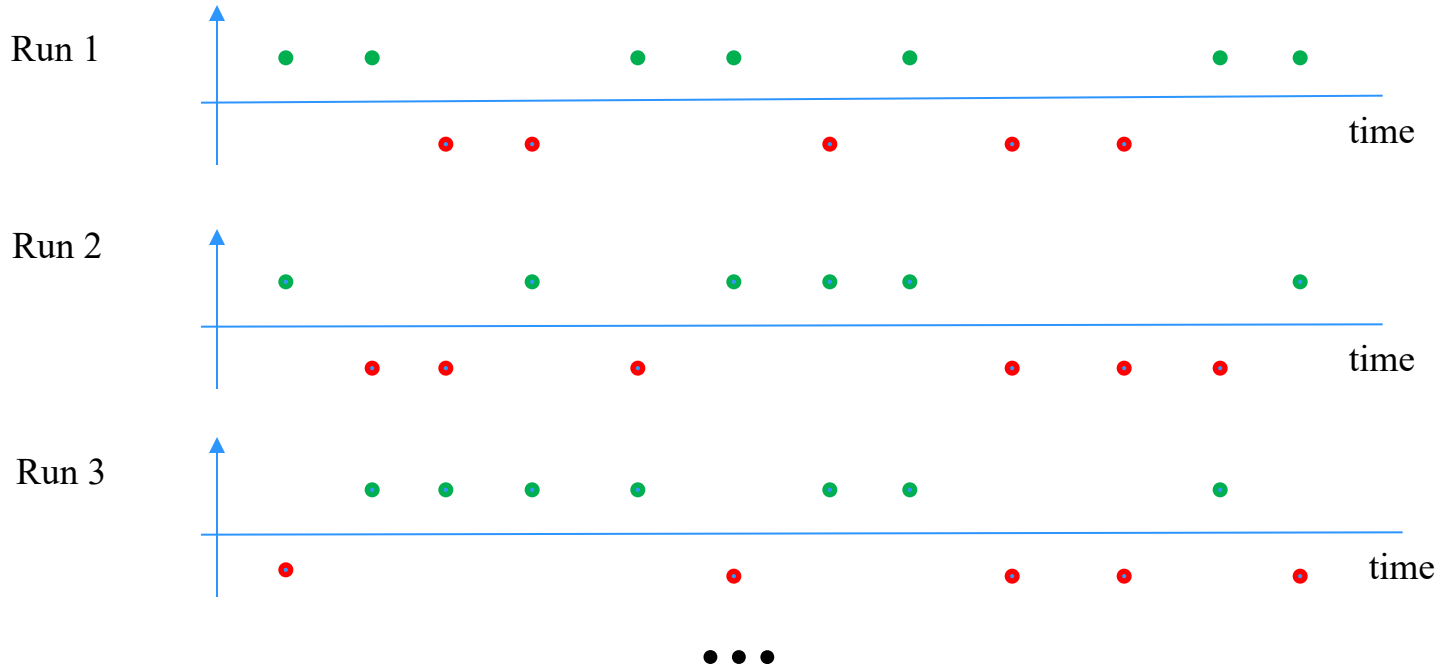
- The total reward for the policy and one SAR trajectory = **sum of rewards for the trajectory**



- But there can be multiple different trajectories the agent may face**
-

# Expected rewards

- Expected rewards for  $\pi : X \rightarrow A$



**A good measure of the quality of policy  $\pi : X \rightarrow A$**

$$E\left(\sum_{t=0}^T r_t\right)$$

Expectation over many possible reward trajectories defined by  $\pi : X \rightarrow A$

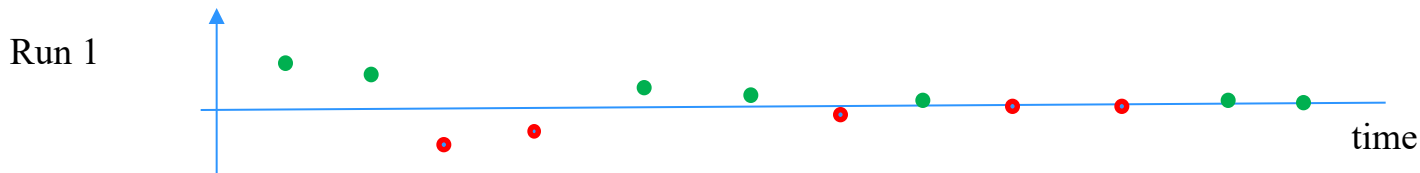
# Expected discounted rewards

- **Expected discounting rewards for  $\pi : X \rightarrow A$**
- **Discounting with  $0 \leq \gamma < 1$**  (future value of money)

No discounting:



Discounting



**Another measure of the quality of policy  $\pi : X \rightarrow A$**

$$E\left(\sum_{t=0}^T \gamma^t r_t\right)$$

Expectation over many possible discounted reward trajectories for  $\pi : X \rightarrow A$

---

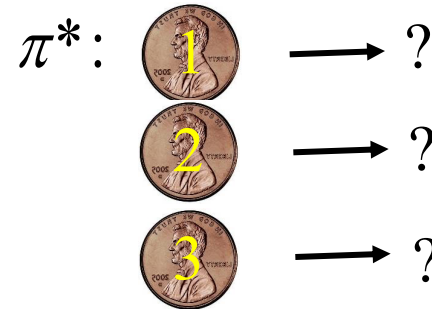


# RL learning objective

**Learning goal:** find the optimal policy

$$\pi^* : X \rightarrow A$$

$$\pi^* : \left| \begin{array}{l} \text{Coin1} \rightarrow ? \\ \text{Coin2} \rightarrow ? \\ \text{Coin3} \rightarrow ? \end{array} \right|$$



**That is, the policy that will maximize the future expected rewards**

$$E\left(\sum_{t=0}^T \gamma^t r_t\right) \quad 0 \leq \gamma < 1$$

a discount factor = present value of money

---

# RL learning: objective functions

- **Objective:**

**Find a policy**  $\pi^* : X \rightarrow A$

That maximizes some combination of future reinforcements (rewards) received over time

- **Valuation models (quantify how good the mapping is):**

- **Finite horizon models**

$$E\left(\sum_{t=0}^T r_t\right)$$

Time horizon:  $T > 0$

$$E\left(\sum_{t=0}^T \gamma^t r_t\right)$$

Discount factor:  $0 \leq \gamma < 1$

- **Infinite horizon discounted model**

$$E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right)$$

Discount factor:  $0 \leq \gamma < 1$

- **Average reward**

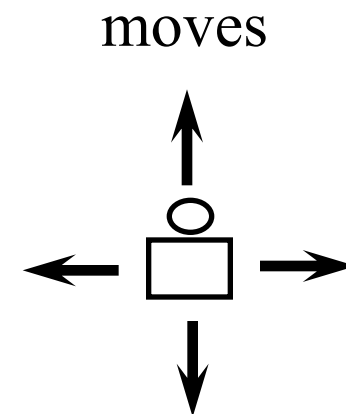
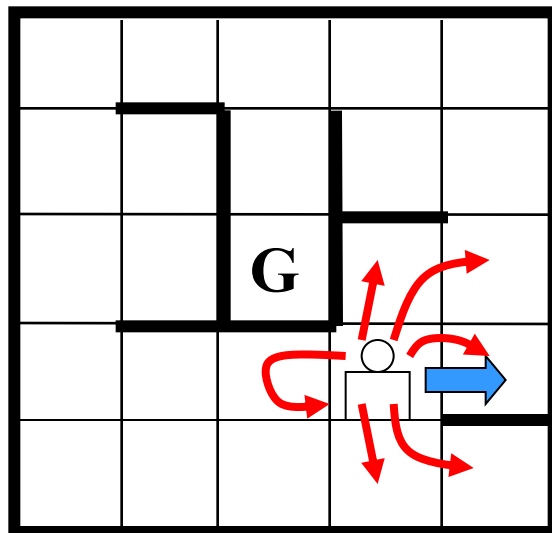
$$\lim_{T \rightarrow \infty} \frac{1}{T} E\left(\sum_{t=0}^T r_t\right)$$

---

# Agent navigation example

- **Agent navigation in the maze:**

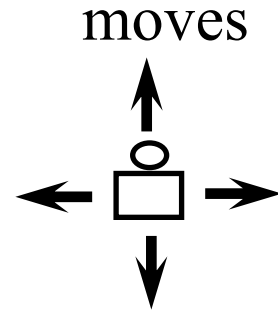
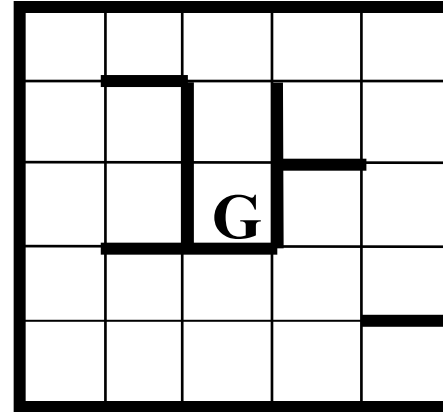
- 4 moves in compass directions
- Effects of moves are stochastic – we may wind up in other than intended location with a non-zero probability
- **Objective:** learn how to reach the goal state in the shortest expected time



# Agent navigation example

- **The RL model:**

- **Input:**  $X$  – a position of an agent
- **Actions:**  $A$  – the next move
- **Reinforcements:**  $R$ 
  - -1 for each move
  - +100 for reaching the goal
- **A policy:**  $\pi : X \rightarrow A$



$$\pi : \left| \begin{array}{l} \text{Position 1} \rightarrow \textit{right} \\ \text{Position 2} \rightarrow \textit{right} \\ \dots \\ \text{Position 25} \rightarrow \textit{left} \end{array} \right|$$

- **Goal:** find the policy maximizing future expected rewards

$$E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) \quad 0 \leq \gamma < 1$$

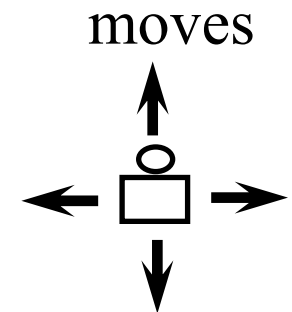
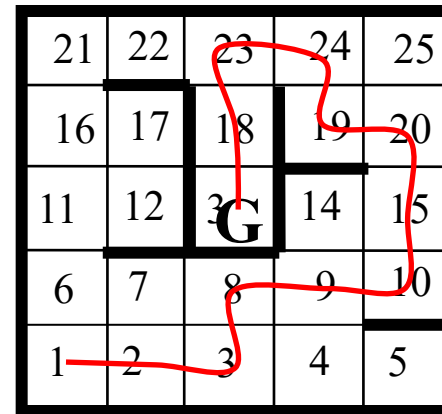
# Agent navigation example

## State, action reward trajectories

- policy

$\pi :$ 

Position 1	→	<i>right</i>
Position 2	→	<i>right</i>
...		
Position 25	→	<i>left</i>



	Step0	Step1	Step2	..	Step k	..
<b>state</b>	Pos1	Pos2	Pos3	..	Pos15	..
<b>action</b>	Right	Right	Up	→	Up	→
<b>reward</b>	-1	-1	-1		-1	

# Effects of actions on the environment

## Effect of actions on the environment

- More specifically on the next input  $\mathbf{x}$  to be seen

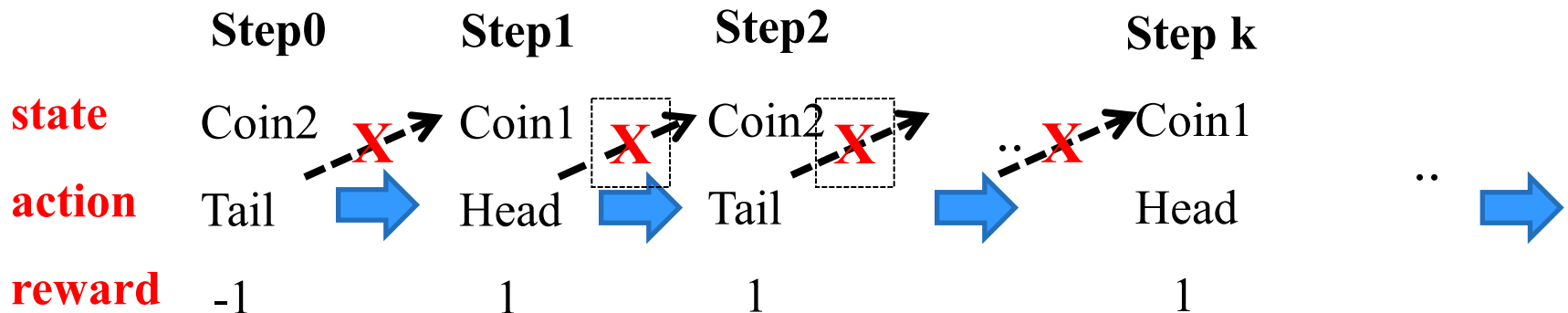
**Case 1. No effect** The distribution over possible  $\mathbf{x}$  is independent of past actions. The rewards received depend only on the current state  $\mathbf{x}$  and the action  $a$  chosen.

- **Reinforcement learning with immediate rewards**

- 3 coin example



**What coin we see next is not affected by our previous action, hence our action does not effect future rewards**



# Effects of actions on the environment

## Effect of actions on the environment

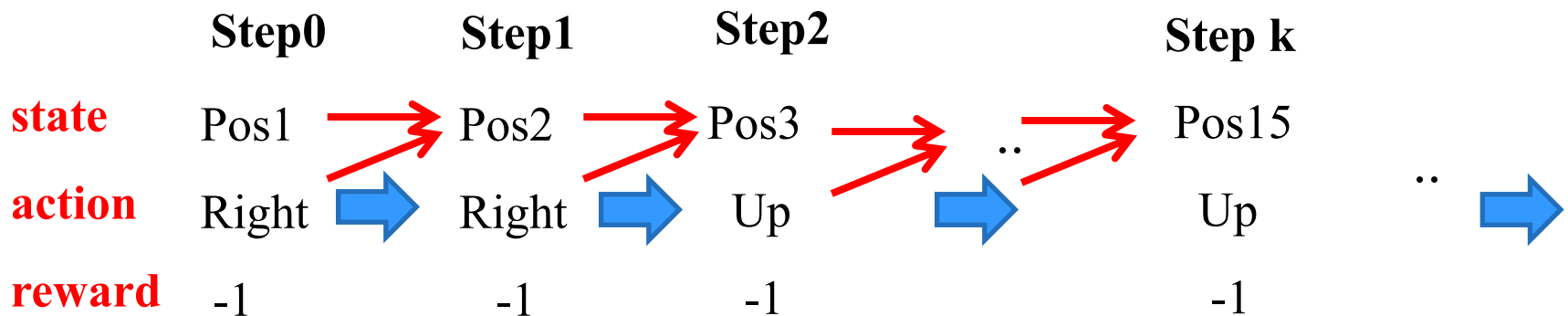
- More specifically on the next input  $\mathbf{x}$  to be seen

**Case 2. Actions may effect the environment** and next inputs  $\mathbf{x}$ . The distribution of  $\mathbf{x}$  can change due to past actions; the rewards related to the action can be seen with some delay.

- Learning with delayed rewards



- **Agent navigation example**; a move action effects next position, and hence more distant future rewards



# RL with immediate rewards

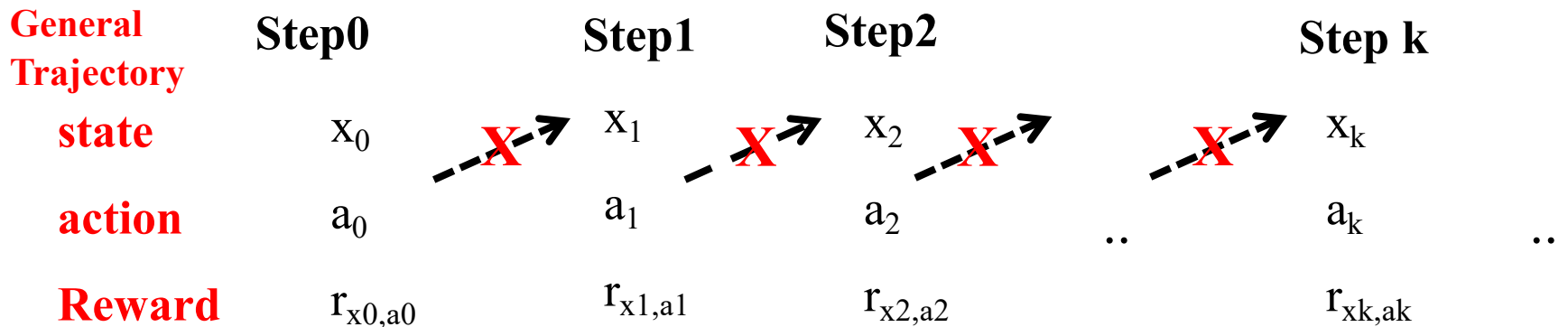


- **Game:** 3 biased coins
    - The coin to be tossed is selected randomly from the three coin options. The agent always sees which coin is going to be played next. The agent makes a bet on either a head or a tail with a wage of \$1. If after the coin toss, the outcome agrees with the bet, the agent wins \$1, otherwise it loses \$1
  - **RL model:**
    - **Input:**  $X$  – a coin chosen for the next toss
    - **Action:**  $A$  – head or tail the agent bets on
    - **Reinforcements:**  $\{1, -1\}$  (\$1 either won or lost)
  - **Learning goal:** find the optimal policy  $\pi^* : X \rightarrow A$   
maximizing the future expected profits over time  
$$E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) \quad 0 \leq \gamma < 1 \quad \text{a discount factor}$$
-



# RL with immediate rewards

- **Expected reward**  $E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) \quad 0 \leq \gamma < 1$
- **Immediate reward case:**
  - Reward depends only on  $\mathbf{x}$  and the action choice
  - The action does not affect the environment and hence future inputs (states) and future rewards:



$$\begin{aligned}
 E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) &= E(r_{x_0,a_0}) + E(\gamma r_{x_1,a_1}) + E(\gamma^2 r_{x_2,a_2}) + \dots E(\gamma^k r_{x_k,a_k}) + \dots \\
 &= E(r_{x_0,a_0}) + \gamma E(r_{x_1,a_1}) + \gamma^2 E(r_{x_2,a_2}) + \dots \gamma^k E(r_{x_k,a_k}) + \dots
 \end{aligned}$$


---

# RL with immediate rewards

## Immediate reward case:

- Reward for input  $\mathbf{x}$  and the action choice  $a$  may vary
- **Expected one-step reward for the input  $\mathbf{x}$  and action  $a$ :**

$$R(\mathbf{x}, a) = E(r_{\mathbf{x}, a})$$

– For the coin bet problem it is:

$$R(\mathbf{x}, a_i) = \sum r(\omega_j | a_i, \mathbf{x}) P(\omega_j | \mathbf{x}, a_i)$$

$\omega_j$  : an outcome of the coin toss  $\mathbf{x}$

$r(\omega_j | a_i, \mathbf{x})$  : reward for an outcome and the bet made on  $\mathbf{x}$

- **Expected one step reward for a policy**  $\pi : X \rightarrow A$

$$R(\mathbf{x}, \pi(x)) = E(r_{\mathbf{x}, \pi(x)})$$

---

# RL with immediate rewards

- **Expected reward**

$$\begin{aligned} E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) &= E(r_{x_0, a_0}) + E(\gamma r_{x_1, a_1}) + E(\gamma^2 r_{x_2, a_2}) + \dots E(\gamma^k r_{x_k, a_k}) + \dots \\ &= E(r_{x_0, a_0}) + \gamma E(r_{x_1, a_1}) + \gamma^2 E(r_{x_2, a_2}) + \dots \gamma^k E(r_{x_k, a_k}) + \dots \end{aligned}$$

- **Optimizing the expected reward**

$$\begin{aligned} \max E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) &= \max E(r_{x_0, a_0}) + \max E(\gamma r_{x_1, a_1}) + \dots \max E(\gamma^k r_{x_k, a_k}) + \dots \\ &= \max E(r_{x_0, a_0}) + \gamma \max E(r_{x_1, a_1}) + \dots \gamma^k \max E(r_{x_k, a_k}) + \dots \\ &= \max_{a_0} R(x_0, a_0) + \gamma \max_{a_1} R(x_1, a_1) + \dots \gamma^k \max_{a_k} R(x_k, a_k) + \dots \end{aligned}$$

**Optimal strategy:**  $\pi^* : X \rightarrow A$

$$\pi^*(\mathbf{x}) = \arg \max_a R(\mathbf{x}, a)$$

# RL with immediate rewards

The optimal choice assumes we know the expected reward

$$R(\mathbf{x}, a)$$

- **Then:**  $\pi^*(\mathbf{x}) = \arg \max_a R(\mathbf{x}, a)$

## Caveats

- **We do not know the expected reward**  $R(\mathbf{x}, a)$ 
    - We need to estimate it using  $\tilde{R}(\mathbf{x}, a)$  from the interactions
  - **We cannot determine the optimal policy if the estimate of the expected reward is not good**
    - We need to try also actions that look suboptimal wrt the current estimates of  $\tilde{R}(\mathbf{x}, a)$
-

# RL with immediate rewards

- **Problem:** In the RL framework we do not know  $R(\mathbf{x}, a)$ 
  - The expected reward for performing action  $a$  at input  $\mathbf{x}$
- **Solution:**
  - For each input  $\mathbf{x}$  try different actions  $a$
  - Estimate  $R(\mathbf{x}, a)$  using the average of observed rewards

$$\tilde{R}(\mathbf{x}, a) = \frac{1}{N_{x,a}} \sum_{i=1}^{N_{x,a}} r_i^{x,a}$$

- Action choice  $\pi(\mathbf{x}) = \arg \max_a \tilde{R}(\mathbf{x}, a)$
  - Accuracy of the estimate: statistics (Hoeffding's bound)
$$P\left(\left|\tilde{R}(\mathbf{x}, a) - R(\mathbf{x}, a)\right| \geq \varepsilon\right) \leq \exp\left[-\frac{2\varepsilon^2 N_{x,a}}{(r_{\max} - r_{\min})^2}\right] \leq \delta$$
  - Number of samples: 
$$N_{x,a} \geq \frac{(r_{\max} - r_{\min})^2}{2\varepsilon^2} \ln \frac{1}{\delta}$$
-

# RL with immediate rewards

- **On-line (stochastic approximation)**

- An alternative way to estimate  $R(\mathbf{x}, a)$

- **Idea:**

- choose action  $a$  for input  $\mathbf{x}$  and observe a reward  $r^{x,a}$
  - Update an estimate in every step  $i$

$$\tilde{R}(\mathbf{x}, a)^{(i)} \leftarrow (1 - \alpha(i))\tilde{R}(\mathbf{x}, a)^{(i-1)} + \alpha(i) r_i^{x,a} \quad \alpha(i) - \text{a learning rate}$$

- **Convergence property:** The approximation converges in the limit for an appropriate learning rate schedule.
- Assume:  $\alpha(n(x, a))$  - is a learning rate for  $n$ th trial of  $(x, a)$  pair
- Then the converge is assured if:

1.  $\sum_{i=1}^{\infty} \alpha(i) = \infty$

2.  $\sum_{i=1}^{\infty} \alpha(i)^2 < \infty$

---

# RL with immediate rewards

- At any step in time  $i$  during the experiment we have estimates of expected rewards for each  $(coin, action)$  pair:

$$\tilde{R}(coin1, head)^{(i)}$$

$$\tilde{R}(coin1, tail)^{(i)}$$

$$\tilde{R}(coin2, head)^{(i)}$$

$$\tilde{R}(coin2, tail)^{(i)}$$

$$\tilde{R}(coin3, head)^{(i)}$$

$$\tilde{R}(coin3, tail)^{(i)}$$

- Assume the next coin to play in step  $(i+1)$  is coin 2 and we pick head as our bet. Then we update  $\tilde{R}(coin2, head)^{(i+1)}$  using the observed reward and one of the update strategy above, and keep the reward estimates for the remaining  $(coin, action)$  pairs unchanged, e.g.  $\tilde{R}(coin2, tail)^{(i+1)} = \tilde{R}(coin2, tail)^{(i)}$
-

# Exploration vs. Exploitation in RL

The (learner) actively interacts with the environment via actions:

- At the beginning the learner does not know anything about the environment
- It gradually gains the experience and learns how to react to the environment

## **Dilemma (exploration-exploitation):**

- After some number of steps, should I select the best current choice (**exploitation**) or try to learn more about the environment (**exploration**)?
  - **Exploitation** may involve the selection of a sub-optimal action and prevent the learning of the optimal choice
  - **Exploration** may spend too much time on trying bad currently suboptimal actions
-



# Exploration vs. Exploitation

- In the RL framework
    - the (learner) actively interacts with the environment and **chooses the action** to play for the current input  $\mathbf{x}$
    - Also at any point in time it has an estimate of  $\tilde{R}(\mathbf{x}, a)$  for any *(input, action)* pair
  - **Dilemma for choosing the action to play for  $\mathbf{x}$ :**
    - Should the learner choose the current best choice of action (exploitation)
$$\hat{\pi}(\mathbf{x}) = \arg \max_{a \in A} \tilde{R}(\mathbf{x}, a)$$
    - Or choose some other action  $a$  which may help to improve its  $\tilde{R}(\mathbf{x}, a)$  estimate (exploration)
- This dilemma is called exploration/exploitation dilemma
- **Different exploration/exploitation strategies exist**
-

# Exploration vs. Exploitation

- **Epsilon greedy exploration:**

- Uses exploration parameter  $0 \leq \varepsilon \leq 1$
- Choose the “current” best choice with probability  $1 - \varepsilon$

$$\hat{\pi}(\mathbf{x}) = \arg \max_{a \in A} \tilde{R}(\mathbf{x}, a)$$

- All other choices are selected with a uniform probability  $\frac{\varepsilon}{|A| - 1}$

## **Advantages:**

- Simple, easy to implement

## **Disadvantages:**

- Exploration more appropriate at the beginning when we do not have good estimates of  $\tilde{R}(\mathbf{x}, a)$
- Exploitation more appropriate later when we have good estimates

# Exploration vs. Exploitation

- **Boltzman exploration**

- The action is chosen randomly but proportionally to its current expected reward estimate
- Can be tuned with a temperature parameter  $T$  to promote exploration or exploitation

- Probability of choosing action  $a$

$$p(a | \mathbf{x}) = \frac{\exp[\tilde{R}(x, a) / T]}{\sum_{a' \in A} \exp[\tilde{R}(x, a') / T]}$$

- **Effect of  $T$ :**

- For high values of  $T$ ,  $p(a | \mathbf{x})$  is uniformly distributed for all actions
  - For low values of  $T$ ,  $p(a | \mathbf{x})$  of the action with the highest value of  $\tilde{R}(\mathbf{x}, a)$  is approaching 1
-