

RL Notes:

## Module 2 Learning Action Values

- Sample-average method

→ The value of action is the expected reward

$$q_t(a) \doteq E[R_t | A_t = a]$$

- a fallen prior to +

$$\sum_{i=1}^{t+1} R_i$$

$$Q_t(a) = \frac{\sum_{i=1}^{t+1} R_i}{t+1}$$

- Incremental update rule

$$B_{n+1} = \frac{1}{n} \sum_{i=1}^n R_i$$

$$= Q_n + \frac{1}{n} (R_n - Q_n)$$

$$\alpha_n \in [0, 1]$$

$$Q_n = (1-\alpha) Q_{n-1} + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

- Greedy action selection

↳ exploit current knowledge

↳ explore to gain knowledge

## Exploration-Exploitation Dilemma

- Estimated incrementally

- General update rule to solve non-stationary bandit prob

## Trade off

- Exploration - exploitation trade off

when to?

→ Exploration - improve knowledge for long-term benefit

$$q(a), N(a), q_*(a)$$

→ Exploitation - exploit knowledge for short-term benefit

- Epsilon greedy

$$A_t \leftarrow \begin{cases} \underset{a}{\operatorname{argmax}} Q_t(a) & \text{with prob } 1-\epsilon \\ a \sim \text{Uniform}\{a_1, \dots, a_{|C|}\} & \text{with prob } \epsilon \end{cases}$$

→ the 10-armed Testbed

Expectation across possible stochastic outcomes

⇒ To Balance exploration and exploitation

## Optimistic Initial Values

- How optimistic initial values encourage?

$$Q_{0,1} \leftarrow Q_w + \alpha (R_w - Q_w)$$

- Describe limitations of OIV

- Drive early exploration
- Not well-suited for non-stationary problems
- May not know max reward value

Upper-Confidence Bound (UCB), Action Selection

→ uses uncertainty in estimates to drive exploration

better than

$\epsilon$  greedy

$$A_t = \arg \max \left[ Q_t(a) + C \sqrt{\frac{1}{N_t(a)}} \right]$$

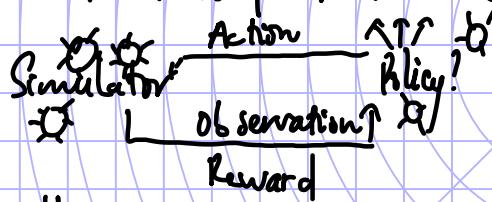
↑  
exploit      ↑  
explore

$t$ : timesteps

$N_t(a)$ : times action  $a$  taken

→ the  $Q_t$  values for these actions are normally distributed with mean zero and standard deviation one

- Contextual Bandits for Real World Reinforcement Learning



How to align?

→ Mind the gap:

Large simulator / reality divergence

How do you RWRL?

↳ realworld-based RL

⇒ Shift your priorities

Temporal Credit Assignment ↓ Generalization ↑

Control environment ↓ Environment controls ↑

Computational efficiency ↓ Statistical efficiency ↑

State ↓

Features ↑

Learning ↓

Evaluation ↑

Last policy ↓

Every Policy ↑

1995 EXP4 paper

2007 Epoch Greedy

2010 Personalized News

2011 Comp / Stat opt algo

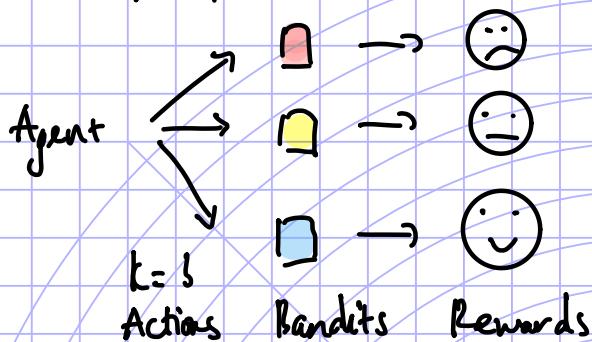
2014 Better Algs

2015 Created first version of "Decision Service"

2015 First RL Service Product "Azure Cognitive Services Personalizer"

→ Summary

- The K-armed Bandit



- Sample-Average Method

$$Q_t(a) = \frac{\sum_{i=1}^{t-1} R_i}{t-1}$$

estimate  $q_*(a)$

- Incremental update rule

$$Q_{t+1}(a) = Q_t(a) + \frac{1}{N_t(a)} (R_t - Q_t(a))$$

New estimate ← old Es + StepSize [target - OldEs]

- Action-Values

$$q_*(a) \doteq E[R_t | A_t = a] \quad \forall a \in \{1, \dots, K\} \rightarrow \text{unknown to the agent}$$

$\downarrow$                    $\downarrow$   
Expected      when taking  
reward received      action a

How do choose when to explore,  
when to exploit?

- Exploration vs. Exploitation

↓  
improve knowledge  
for long-term benefit

↓  
exploit knowledge  
for short-term benefit

- Optimistic Initial Values

↳ OIV >  $q_*(a)$ , the agent will  
systematically explore the actions.

The optimism fades with time and the  
agent eventually stops exploring.

- Epsilon-Greedy Action Selection

$$A_t \leftarrow \begin{cases} \underset{a}{\operatorname{arg\ max}}(Q_t(a)), \text{ prob } 1-\epsilon \\ \text{uniform}(\{a_1, \dots, a_K\}), \text{ prob } \epsilon \end{cases}$$

- UCB - Action Selection

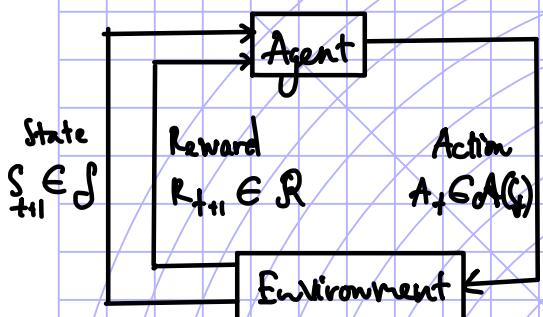
$$A_t \doteq \operatorname{argmax} \left[ Q_t(a) + C \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

↑                  ↑  
Explorit      Explore

## Module 3

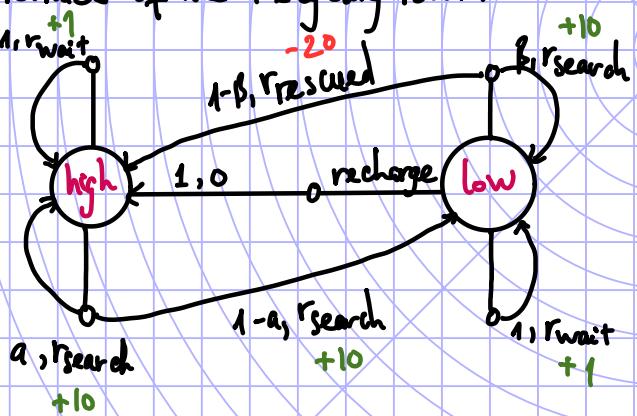
## finite Markov Decision Processes

- Markov Decision Processes



- Examples of MDPs

↳ Dynamics of the Recycling Robot:



→ MDP formalism is abstract and flexible

↳ Task: The goal of the robot is to pick-and-place objects

State : Latest readings of joint angles and velocities

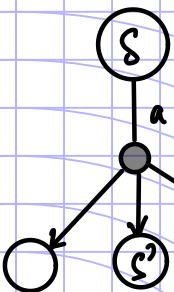
Action : the amount of voltage of applied to each motor

Reward : +100 when an object is successfully placed

- 1 for each unit of energy consumed

- How the dynamics of an MDP are defined

$p(s', r | s, a)$ : prob distribution



→ Markov properties :

$$p: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$$

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1$$

$\forall s \in S, a \in A(s)$

⇒ Sequential - markov problems

$G_5 = 0$

$G_4 =$

## Goal of RL

- How rewards relate to the goal of an agent

Goal of an Agent: Formal definition

$$\text{return } G_t = R_{t+1} + R_{t+2} + \dots$$

- ⇒ Maximize the expected return

$$E[G_t] = E[R_{t+1} + R_{t+2} + \dots + R_t]$$

## The Reward Hypothesis

Whence behavior?

→ Programming (GOFAI)

- Coding
- Human-in-the-loop

→ Examples (LfD)

- Mirror Reward
- Inverse RL

→ Optimization (RL)

- Evolutionary optimization
- Meta RL

→ Identify episodic tasks

↳ In episodic tasks, the agent-env interaction breaks up into episodes

↳ Each episode begins independently of how the previous one ended

↳ At termination, the agent is reset to a start state

## Origin of the hypothesis

Research program:

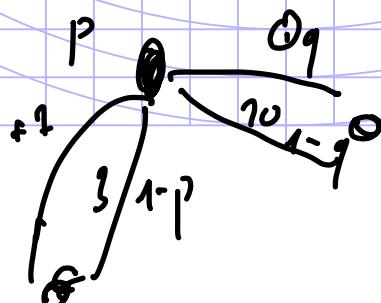
→ Identify where reward signal comes from

→ Develop algos that search the space of behaviors to maximize reward signals

## Goals as Rewards

→ 1 for goal, 0 otherwise  
goal-reward

→ -1 for not goal, 0 once goal reached  
action-penalty



$$p + 3(1-p) = \log(A-q)$$

$$-2p + 3 = \log - \log q$$

$$\gamma = 0, 0.5,$$

$$\frac{1}{4} \cdot -21 + \frac{1}{4} \cdot (-1 + 12) = u$$

## Continuing tasks

- Episodic vs. continuing tasks

- Interaction over episodes
- Terminal state
- Independent

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

- Gives on continually

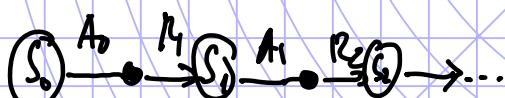
- No terminal state

$$G_t = R_{t+1} + R_{t+2} + \dots \\ = \text{?}$$

$$\Rightarrow G_t = R_{t+1} + \gamma G_{t+1}$$

e.g. games  $\rightarrow$  Episodic tasks

Schedules  $\rightarrow$  Continuing tasks



$\hookrightarrow$  Discounting value  $\gamma$

$\hookrightarrow$  first  $\Rightarrow$  formulate a MDP.

$$G_0 = R_1 + \gamma G_1$$

- Formulate returns using  $\gamma$  by counting
  - Discount the rewards in the future by  $\gamma$
  - $\hookrightarrow \gamma < 1$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{K-1} R_{t+K} + \dots$$

$$= \sum_{k=0}^{\infty} \gamma^K R_{t+k+1} \Rightarrow \text{finite}$$

$$\leq \sum_{k=0}^{\infty} \gamma^K R_{\max} = R_{\max} \sum_{k=0}^{\infty} \gamma^k = R_{\max} \cdot \frac{1}{1-\gamma}$$

- How returns at successive time steps are related to each other

$\gamma = 0 \Rightarrow$  Short-sighted agent!

$$G_t = R_{t+1} \quad (\text{immediate reward})$$

$\gamma \rightarrow 1 \Rightarrow$  Far-sighted agent!

$\hookrightarrow$  takes future rewards into account more strongly