

RL Notes

Day: _____

Date: _____

AI

to mimic the behavior of intelligent human beings
(simulation of human intelligence in machines)

Rule based
(Deduction)

rules of logic
defined by
humans

Example based
(Induction)

Machine learning

to learn patterns

and behaviors from
data or examples

Subcategories

- Supervised learning
- Unsupervised learning
- Reinforced learning

In reinforcement learning, reinforcement refers to the feedback (rewards or penalties) that an agent receives as it interacts with its environment, guiding it to improve its decisions or behavior.

(Independent)

An agent in AI is an autonomous entity that perceives its environment, makes decisions based on that perception, and takes actions to achieve specific goals or objectives.

Day: _____

Date: _____

Goal: _____

Rewards

A reward is a numerical signal that provides feedback to an agent about the success of its actions in achieving a desired goal.

Goal: _____

It is maximizing the total reward (cumulative reward) over time.

Rewards provide (immediate feedback) about whether the actions are beneficial (moving it closer to the goal) or detrimental (moving it away from the goal).

policy: optimal behavior strategy that defines how an agent acts in an environment.

Difference between supervised and reinforcement

Train on a labeled dataset

reducing errors

Learn to play by trial and error

maximizing its rewards

delayed feedback

Sequential Decision Making

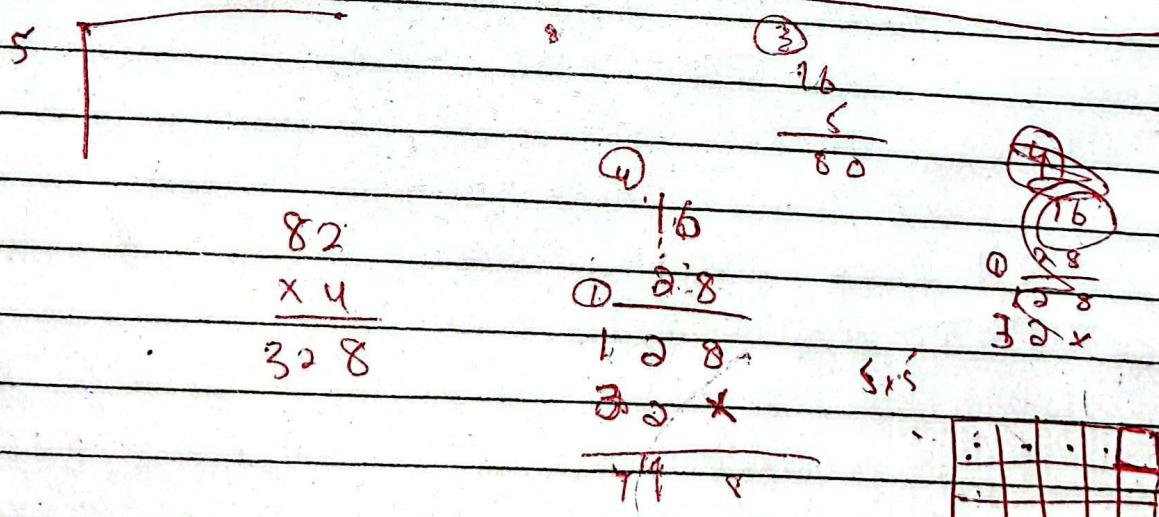
- Goal: select actions to maximize total future reward
 - Actions may have long term consequences
 - Rewards may be delayed

~~Agent at each stop~~

- Execute action A_t
 - Receive observation O_t
 - Received scalar reward R_t

Environment.

- Receives action A_t
 - Emits observation O_{t+1}
 - Emits scalar reward R_{t+1}



$$\Rightarrow \left(\frac{5}{5} \text{ row } + \frac{3}{3} \right) \times \left(\frac{5}{5} + \frac{1}{1} \right)$$

Day:

Date:

States: $(4 + \underline{\hspace{2cm}})$

5x5 grid - 25 states

25 passenger states ($4 \text{ loc} + \text{inside taxi}$)

4 destination points

$$(3,1) \quad 25 \times 5 \times 4 = 500$$

state 328

(5 + row)

$$(\text{row} \times 5) + \text{column}$$

$$(5 + \text{pass_locn}) \times 4 +$$

destination

Taxi at (3,1) P: Y D: R

Pass loc: $(R:0, G:1, Y:2, B:3, \text{Intaxi:4})$

$$\text{Dest loc: } (3 \times 5 + 1) \times (5 + 2) \times 4 + 0$$

$$(1,0,G:1,Y:2, B:3) \quad 15 + 1 = 16 + 0$$

State No = $((\text{Taxi_row} \times 5 + \text{taxi_column}) \times 5 + \text{passenger index}) \times 4 + \text{Destination}$

$$= ((3 \times 5 + 1) \times 5 + 2) \times 4 + 0$$

$$= (16 \times 5 + 2) \times 4$$

$$= 82 \times 4$$

82

328

$$4 \overline{) 328} \quad \underline{328} \quad 0$$

Destination: 0

$$5 \overline{) 82} \quad \text{Passenger} \\ 80 \quad 60 : 2$$

$$5 \overline{) 16} \quad 5 \overline{) 13} \\ 10 \quad 10 \\ (3,1)$$

Date: _____

Policy defines how the agent selects actions in each state to maximize cumulative rewards

Solution Methods in RL

tabular
solution
MAB
MDP

Approximate
solution

Multi-Armed Bandits

- The simplest RL problem
- One state
- Actions don't have delayed effects

The term "bandit" comes from the analogy of a slot machine, also known as a one-armed bandit, because it has a lever (a_m) and "steals" money from players over time.

True value of an action is the expected value of the rewards the agent can get for that action.

$$q_*(a) = E[R_t | A_t = a], \forall a \in \{1, \dots\}$$

The multi-armed bandit problem involves choosing between multiple actions (like slot machine arms) to maximize reward over time. It's named after slot machines ("one-armed bandits") imagine choosing which of many slot machines to play.

Date: _____

Estimated value

$$Q(a_i) \approx q^+(a_i)$$

Action value methods

The action value method estimates the expected reward of each action based on past experiences

types

Sample Average Method

$$Q_n = \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$$

Incremental average method

$$Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n]$$

$$Q_{n+1} = Q_n + \frac{1}{n} (R_n - Q_n) \quad \text{New Estimate} = \text{Old Estimate} + \text{Step Size} \cdot (\text{Target} - \text{Old Estimate})$$

Exploration vs exploitation dilemma

It involves deciding when to explore new actions to gather more information about their potential rewards (exploration) versus when to choose the action that currently has the highest estimated value (exploitation)

$\epsilon \rightarrow$ exploration ϵ value balance

$1 - \epsilon \rightarrow$ exploitation, θ and β

We choose ~~act~~ random action with certain probability ϵ and ~~on~~ θ and β on $1 - \epsilon$

Day: _____

Date: _____

10 armed Bandit Problem

Stationary

Non Stationary

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

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

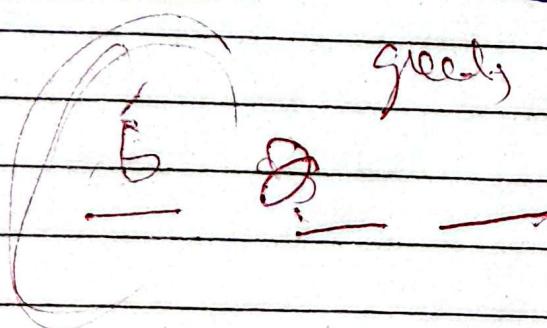
$k=10$

10 armed tested

In the 10 armed test bed the two values $q^*(a)$ for each of the ten actions are drawn from a normal distribution with mean 0 and variance 1.

For the epsilon-greedy algorithm, different epsilon values control the probability of exploring random actions.

An initial optimistic value can be assigned to encourage exploration.



Day: _____

Date: _____

for $k = 10$ (10 actions)

and for $actions = 1000$
as

for each action

Create 10 actions

for each action we max run (and
get reward value 1000 times)

No repeat the whole episode

200 times

in actions

get rewards true value from N(0)

Run Exp

Initialization

how many times to repeat ($n = 2000$)

how many steps per for (≈ 1000)

Create array to store rewards

For each n (200 times),

Create 10 actions with random true values

For each time step (1000 times),

choose action

If random < epsilon,

Select random

else select action with highest estimate val

Get reward,

from normal distribution

Update estimates.

new estimate = old estimate + $(reward - \frac{old\ estimate}{old\ estimate})$

true value

Day: _____

Date: _____

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

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

$$Q_{n+1} = \frac{1}{n} \left[\sum_{i=1}^{n-1} R_i + R_n \right]$$

$$Q_{n+1} = \frac{1}{n} \left[\frac{n-1}{n-1} \sum_{i=1}^{n-1} R_i + R_n \right]$$

$$\Rightarrow \frac{1}{n} [(n-1) Q_n + R_n]$$

$$\Rightarrow \frac{1}{n} [n Q_n - Q_n + R_n]$$

Incremental

$$Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n]$$

Sample Avg Incremental

Step	Action taken	Reward		
1	A1	2	2	$\frac{2+1(2-0)}{2}$
2	A2	5	5	5
3	A1	4	3	3
4	A3	1	1	1
5	A1	3	3	3
6	A2	6	$5 \cdot 5$	$5 \cdot 5$
7	A3	4	$2 \cdot 5$	$2 \cdot 5$
8	A2	3	$4 \cdot 67$	$4 \cdot 67$

Day: _____

Date: _____

$$Q_{n+1} = Q_n + \alpha [R_n - Q_n]$$

$$= Q_n + \alpha R_n - \alpha Q_n$$

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

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

$$Q_{n+1} = \alpha R_n + (1-\alpha)(\alpha R_{n-1} + (1-\alpha) Q_{n-1})$$

$$= \alpha R_n + \alpha(1-\alpha)R_{n-1} + (1-\alpha)^2 Q_{n-1}$$

$$= \alpha R_n + \alpha(1-\alpha)R_{n-1} + \alpha(1-\alpha)Q_{n-2} +$$

$$(1-\alpha)^3 Q_{n-3}$$

$$= \alpha R_n + \alpha(1-\alpha)R_{n-1} + \dots$$

$$(1-\alpha)^n Q_1$$

$$= (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha(1-\alpha)^{n-i} Q_i$$

Date: _____

Optimistic Initial Values

We set initial values other than 0

This optimism encourages action-value methods to explore.

UCB

Upper Confidence Bound Action Selection

Select actions according to

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

c controls the degree of exploration.

If $N_t(a) = 0$ then a is considered

To be a maximizing action.

$N_t(a)$ denotes the no of times that action a has been selected prior to time t .

each

time an action is selected its uncertainty reduced.

Day: _____

Date: _____

Optimistic initial values involve setting the initial estimates of action values to high, high values instead of 0.

Encourages exploration
helps balance exploration and exploitation
without requiring an explicit exploration parameter,

UCB

$$A_t = \operatorname{argmax}_a \left(Q(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right)$$

$Q(a)$: Current estimate of the action value.

c : Exploration parameter

$\ln t$: Natural log of the total number of steps taken so far

$N_t(a)$: Number of times action a has been chosen up to time t)

Actions that are less explored (small $N_t(a)$) have a higher uncertainty term, encouraging the agent to try them.

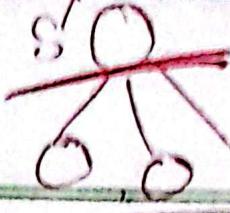
MAB	
1) Only one state; no state transition.	Involves multiple states; the state changes based on actions.
Action influences future.	Action influences future states and long term rewards
2) Actions provide immediate rewards only, no effect on future rewards.	
3) No delayed effects, rewards are immediate	Actions may have delayed effects.
4) Maximizes immediate total rewards from repeated actions.	Maximize cumulative rewards over a sequence of steps.
5) Simple problem, no state transitions or policies	More complex due to state transitions and policies
B)	Focuses on the exploration and exploitation trade off to find the best action
	Requires learning an optimal policy for selecting actions in each state,

$$V(s, a)$$

$$P(s'|s, a)$$

$$\sum_r \sum_s p(s'|s, a) \sum_r P(r|s, a)$$

Day:



Date:

MDP

$$p(s'|s,a) = \sum_r p(s',r|s,a)$$

$$r(s,a) = \sum_r \sum_{s'} p(s',r|s,a)$$

Expected
value
of reward

$$p(s',r|s,a) = p(s'|s,a) \cdot p(r|s',s,a)$$

$$r(s,a,s') = \sum_r r \frac{p(s',r|s,a)}{p(s'|s,a)}$$

$$p(s'|s,a) = \sum_r p(s',r|s,a)$$

$$r(s,a) = \sum_r r \sum_{s'} p(s',r|s,a)$$

$$r(s,a,s') = \sum_r r \frac{p(s',r|s,a)}{p(s'|s,a)}$$

$$r(s,a,s) = \sum_r r \frac{p(s',r|s,a)}{p(s'|s,a)}$$

Day: _____

Date: _____

$$S = \{s_0, s_1, s_2\}$$

$$A = \{a, b\}$$

s_0 ,
a

$$p(s_1, r=2 | s_0, a) = 0.4$$

$$p(s_0, r=0 | s_0, a) = 0.6$$

b

$$p(s_1, r=1 | s_0, b) = 0.5$$

$$p(s_2, r=3 | s_0, b) = 0.5$$

 s_1

a:

$$p(s_0, r=1 | s_1, a) = 1$$

b:

$$p(s_2, r=2 | s_1, b) = 1.0$$

s_2
zero

$$p(s' | s, a) = \sum_r p(s', r | s, a)$$

$$p(s_0 | s_0, a_1) = \sum_r p(s_0, r | s_0, a) \\ = 1$$

$$p(s_1 | s_0, a)$$

Day: _____

Date: _____

 $s_0 : E$ $s_1 : P$ $s_2 : F$

a: Load

b: Drive

 $s_0 \rightarrow a$ $s_1 \rightarrow s_0, a$

$$(s_1, r=2 | s_0, a) = 0.4$$

$$(s_2, r=3 | s_0, a) = 0.3$$

$$(s_0, r=0 | s_0, a) = 0.3$$

$$(s_1, r=1 | s_0, b) = 0.5$$

$$(s_2, r=5 | s_0, b) = 0.2$$

$$(s_0, r=0 | s_0, b) = 0.3$$

$$(s_2, r=4 | s_1, a) = 0.6$$

$$(s_0, r=1 | s_1, a) = 0.3$$

$$(s_1, r=0 | s_1, a) = 0.1$$

$$(s_0, r=2 | s_1, b) = 0.2$$

$$(s_2, r=3 | s_1, b) = 0.7$$

$$(s_1, r=0 | s_1, b) = 0.1$$

Day:

Set $S = \{\text{high, low}\}$

Date:

 $A = \{\text{search, wait, recharge}\}$ $a = A_{\text{high}} (\text{search, wait})$ $A_{\text{low}} = \{\text{search, wait, recharge}\}$ recharged \rightarrow low rewardhigh \rightarrow low reward & low R_{low}

$$P(s'|s, a) = \sum_r P(s', r|s, a)$$

$$P(s'|s, a) = \sum_{r,s'} P(s', r|s, a)$$

pt

$$r(s, a) = \sum_{r,s'} r P(s', r|s, a)$$

S

R

$$r(s, a, s') = \sum_r r P(s', r|s, a)$$

s'	a	s'	$P(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
high	search	low	$1-\alpha$	r_{search}
high/low	wait	low	β	r_{wait}
low	search	high	$1-\beta$	-3
high	wait	high	1	r_{wait}
high	wait	low	0	-
low	wait	high	0	-
low	wait	low	1	r_{wait}
low	recharge	high	1	0
low	recharge	low	0	0

Day: _____

Date: _____

Lecture 7

A₁, A₂, A₃

exploit

0.8

exploit

Optimistic initial values

Setting up

A₁, A₂, A₃

Initial

5, 5, 5

b₁ = 1, s' = 2, s = 5

Values

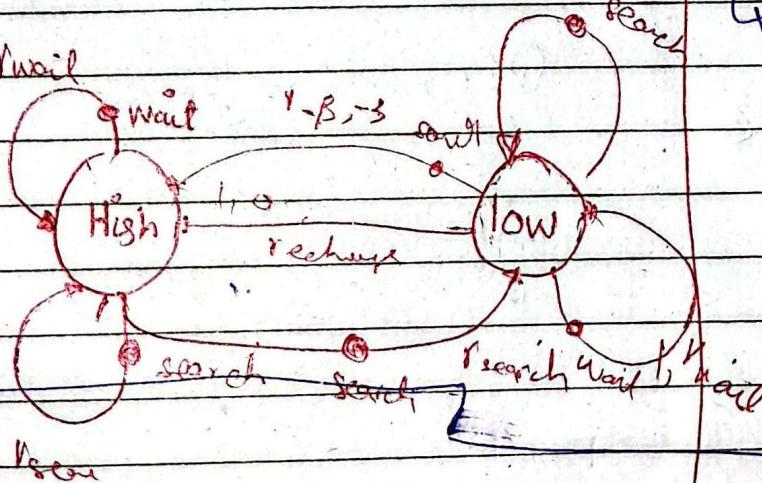
t₂Other than
can encourage
explorationV_{start}

$$Q_1 = Q_0 + \frac{1}{n} [R_0 \cdot Q_0]$$

$$= 5 + \frac{1}{3} [2 - 5]$$

$$= 5 - 3$$

2



Day: _____

Date: _____

RL reward

$$p(s' | s_1, a) = \sum_r p(s'_r | s_1, a).$$

$$\sum_r p(s_0, r | s_1, a) = 0.3 \quad r(s_1, a) = \sum_r r \cdot \sum_{s'} p(s'_r | s_1)$$

$$p(s_0, r | s_1, b) = 0.2$$

$$a \quad p(s_0 | s_1, a) = 0.3$$

$$p(s_1 | s_0, a) = 0.4$$

$$p(s_0 | s_1, b) = 0.2$$

$$p(s_1 | s_0, b) = 0.5$$

$$a \quad p(s_0 | s_0, a) = 0.3$$

$$p(s_1 | s_0, a) = 0.1$$

$$p(s_0 | s_0, b) = 0.3$$

$$p(s_1 | s_0, b) = 0.1$$

$$a \quad p(s_0 | s_1, a) = 0$$

$$p(s_1 | s_0, a) = 0$$

$$p(s_0 | s_2, b) = 0$$

$$p(s_1 | s_2, b) = 0$$

$$p(s_1 | s_0, a) = 0.3$$

(state)

$$p(s_1 | s_0, b) = 0.2$$

$$p(s_1 | s_1, a) = 0.6$$

$$p(s_1 | s_1, b) = 0.7$$

$$p(s_1 | s_2, a) = 0$$

$$p(s_1 | s_2, b) = 0$$

Day: _____

Date: _____

Episodic:

$$G_T = \sum_{t'=t+1}^T R_{t'}'$$

Continuous

$$G_T = \sum_{t'=t+1}^{+\infty} R_{t'}'$$

$$\sum_{t'=t+1}^T R_{t'}$$

$$G_T = \sum_{k=t+1}^{+\infty} r^k$$

$$G_T = \sum_{k=t+1}^{+\infty} r^k R_{t+1+k}$$

$$0 < r < 1$$

$$\sum_{t'=t+1}^{\infty} R_{t'}$$

$$G_T = R_{t+1} + r G_{t+1}$$

$$\sum_{k=t+1}^T r^{k-t-1} R_k$$

$$G_T = \sum_{k=t+1}^T r^{k-t-1} R_k$$

$$k=t+1$$

Unstable if T is ∞ and $r=1$.

$$V_T$$

$$\pi(a|s)$$

$$\pi(a)$$

$$V_\pi(s) = \mathbb{E}_\pi(G_{T+1}|s_t=s)$$

$$= \mathbb{E}_\pi[G_{T+1} | \sum_{k=0}^T r^k R_{t+k+1} | s_t=s]$$

$$= \pi(a|s) \leq \sum_{s'} P(s'|s, a) \cdot [r + \gamma]$$

Day: _____

Date: _____

Policies

$$\pi(a|s)$$

Value function

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t, s]$$

$$V_\pi(s) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t, s \right]$$

Action value func

$$Q_\pi(s, a) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t, s]$$

$$= \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} | S_{t+1}, s]$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_\pi [G_{t+1} | S_{t+1}, s']]$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s' | s, a) [r + \gamma V_\pi(s')]$$

Day: _____

Date: _____

Deterministic
Policy

A policy map
each state to
a given action

$$\pi(s) = a$$

Stochastic
Policy

A policy assign
probabilities to each
action in each state

$$\pi(a|s)$$

$$\sum_{a \in A(s)} \pi(a|s) = 1$$



Policies can be only
depend on current state

$$G_t = \sum_{k=0}^{\infty} r_{t+k+1}$$

state value is basically a future
reward that an agent expect to -
receive starting from a particular state

state value
rewards
action values

$$V(s) = E[G_t | S_t = s]$$

$$q(\gamma, a) = E_{\pi}[G_t | S_{t+1}, A_{t+1}, a]$$

$$V(s) = \bar{\pi}(a) \sum_s \sum_a p(s'|a|s, a)$$

$$[i.e. V(s')]$$

Day: _____

Date: _____

$$\gamma = 0.9$$

State-value functions represent the

expected return from a given state
under a specific policy.

Action-value functions represent the expected

return from a given state after taking
a specific action, later following a specific
policy.

Bellman Equation

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | s_t = s]$$

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} (G_t | s_t = s, a_t = a)$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$V_{\pi}(s) = \mathbb{E}_{\pi} (R_{t+1} + \gamma G_{t+1} | s_t = s)$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s'|r|s, a) [r + \gamma \mathbb{E}_{\pi}(G_{t+1} | s_{t+1} = s')]$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s'|r|s, a)$$

$$[\gamma + \gamma V_{\pi}(s')]$$

$$V_{\pi}(s) = \sum_a \sum_{s'} p(s'|s, a) [r + \gamma \mathbb{E}_{\pi}(G_{t+1} | s_{t+1} = s')]$$

$$q_{\pi}(s, a) = \sum_{s'} \sum_r p(s'|r|s, a)$$

$$[G_{t+1} | s_{t+1} = s, a_t = a]$$

$$= \sum_a \sum_{s'} \sum_r p(s'|r|s, a) [r + \gamma \mathbb{E}_{\pi}(G_{t+1} | s_{t+1} = s')]$$

$$= \sum_a \sum_{s'} \sum_r p(s'|r|s, a)$$

$$[r + \gamma \mathbb{E}_{\pi}(G_{t+1} | s_{t+1} = s')]$$

$$= \sum_a \sum_{s'} \sum_r p(s'|r|s, a)$$

$$[\gamma + \gamma \sum_a \pi(a|s) \mathbb{E}_{\pi}(G_{t+1} | s_{t+1} = s)]$$

Day: _____

Date: _____

In ϵ -greedy action selection, for the case of two actions and $\epsilon = 0.5$, what is the probability that the greedy action is selected?

$$\epsilon = 0.5$$

$$\text{Greedy action} = 1 - \epsilon + \frac{\epsilon}{\text{no of actions}}$$

$$= 1 - 0.5 + \frac{0.5}{2}$$

$$= 1 - 0.5 + 0.25$$

$$= 0.5 + 0.25$$

$$= 0.75$$

Day: _____

Date: _____

 $k = 4$ actionsAction $\rightarrow 1, 2, 3, 4$

$$Q_1(a) = 0$$

time step	1	2	3	4	
A1	1				$Q(1) = 1$
A2		1			$Q(2) = 1$
A3		2			$Q(2) = 1.5$
A4		2			
A5			0		

	1	2	3	4
random	0	0	0	0
random	1			
random + sue		2		
greedy		2		
random			0	

Day: _____

Date: _____

Flappy bird Game

States

- bird's vertical position (high, medium, low)
- bird's velocity (rising, falling)
- Distance to the next pipe

Actions

Flap (move up)

Do nothing (let the bird fall)

Goal

pass through many pipes as possible
without crashing

episodes end when bird crashes into the pipe
or ground

Rewards

- +10 for successfully passing through the pipe
- 0 for staying alive without passing a pipe
- 5 for crashing into a pipe

Space rover

States

- grid coordinates
- destination
- starting point
- hazards
- collected sample

Actions

Move (north, east, west, south)

- Collect sample
- Recharge
- Avoid hazards

Rewards

- +15 for collecting sample
- 10 for encountering hazard
- 5 low energy without recharging
- 0 for exploring between health & sample

Day: _____

Date: _____

Car Refueling Trip

Starts

- home
- Driving to gas station
- Waiting in line
- Refueling car
- Returning home

- Choose nearest station
- Wait or skip long line
- Refuel
- Return home

Rewards

+10. Refueling successful

-5. Running low on fuel before
reaching the station

-3. Wait long in line

s	a	s'	r	p(s', r s, a)
high	wait	high	wait	1 - 0
✓ high	wait	low	0	-
high	search	high	search	0
high	search	low	search	1 - 0
high	recharge	high	1	-
high	recharge	low	0	-
low	wait	high	0	-
low	wait	low	wait	1
low	search	high	-3	0
low	search	low	search	1 - 0
low	recharge	high	1 0	0
low	recharge	low	0	-

Day: _____

Date: _____

Car Refueling Trip

states

- home
- Driving to gas station
Waiting in line
- Refueling car
- Returning home

- Choose nearest station

- Wait or skip long line.
- Refuel
- Return home

Rewards

+10. Refueling successful

-5. Runis low on fuel before
reaching the station

-3. Waiting long in line

s	a	s'	r	pts, rts(a)
high	wait	high	wait	1 0
✓ high	wait	low	0	-
high	search	high	search	0
high	search	low	search	1 - 0
trish	recharge	high	1	-
high	recharge	low	0	-
low	wait	high	0	-
low	wait	low	wait	1
low	search	high	-3	3
low	search	low	search	1 - 3
low	recharge	high	1 0	0
low	recharge	low	0	-

Day: _____

Date: _____

$$\gamma = 0.5$$

$$R_1 = 1, R_2 = 2, R_3 = 6, R_4 = 3, R_5 = 2$$

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

$$T_5 = 5$$

$$G_T = G_{T+1}$$

$$G_T = R_{T+1} + \gamma G_T$$

$$\text{Final } G_0, G_1, \dots, G_5.$$

$$G_4 = R_5 + \gamma G_5$$

$$G_4 = 2 + (0.5)(0) = 2$$

$$G_3 = R_4 + \gamma G_4$$

$$= 3 + (0.5)(2) = 4$$

$$G_2 = R_3 + \gamma G_3$$

$$= 6 + (0.5)(4) = 8$$

$$G_1 = R_2 + \gamma G_2$$

$$= 2 + (0.5)(8) = 6$$

$$G_0 = R_1 + \gamma G_1$$

$$= 1 + (0.5)(6) = 2$$

$$\gamma = 0.9$$

$R_1 = 2$ followed by an infinite sequence of T_5

Find G_1 & G_0

$$G_1 = R_1 + \sum_{k=0}^{\infty} R_k \cdot \gamma^k$$

$$\Rightarrow R_k \cdot \frac{1}{1-\gamma}$$

$$= 7 \cdot \frac{1}{1-0.9} = 70$$

$$G_0 = 2 + (0.9)70 - 2 + (0.9)(70) = 63$$

Day: _____

Date: _____

$$\gamma = 0.8$$

$$R_1 = 5$$

$$+ 0$$

$$G_0 = \gamma G_1 + R_{t+1}$$

$$G_1 = R_t \frac{1}{1-\gamma}$$

$$5 \quad \frac{1}{0.2}$$

$$50$$

$$Cost \ 50 + 10$$