

MAY '23						
	JUNE '23					
Week	22	23	24	25	26	27
Sunday	1	8	15	22	29	
Monday	2	9	16	23	30	
Tuesday	3	10	17	24	31	
Wednesday	4	11	18	25		
Thursday	5	12	19	26		
Friday	6	13	20	27		
Saturday	7	14	21	28		

Week	22	23	24	25	26
Monday	5	12	19	26	
Tuesday	6	13	20	27	
Wednesday	7	14	21	28	
Thursday	1	8	15	22	29
Friday	2	9	16	23	30
Saturday	3	10	17	24	
Sunday	4	11	18	25	

Thursday
June
2023

01

22nd Week | 152nd Day

Reinforcement Learning

points

RL is basically a learning model that learns from its environment. Unlike supervised ML where a dataset for training is already given, RL models basically learns from a trial & error method.

Elements of RL

Agent - the thing performing all the tasks

Environment - the surrounding ^{area} e.g. the game. ~~area~~ of the agent i.e. the place where agent will perform task (like battlefield of any game)

Policy - agent's way of behaving at a given time

Reward signal - defines goal of a RL prob. At each step, the agent receives certain reward which we have to maximize.

Value function - total amt. of reward the agent expects to accumulate over future.

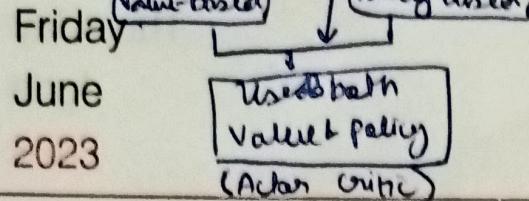
Model (not needed in model-free learning) - basically used to train the model on a dummy-environment before running it on the actual environment.

State - information used to determine what happens next. (In markov state, state is sufficient for predict future, history is not needed)

Q

Rewards are short-term indicators but valueⁿ is the cumulative sum of rewards hence, long-term indication

02



Model-based

Week	18	19	20	21	22
Monday	1	8	15	22	29
Tuesday	2	9	16	23	30
Wednesday	3	10	17	24	31
Thursday	4	11	18	25	
Friday	5	12	19	26	
Saturday	6	13	20	27	
Sunday	7	14	21	28	

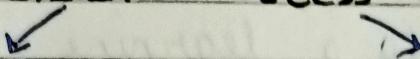
Model free

Week	22	23	24	25	26
Monday	5	12	19	26	
Tuesday	6	13	20	27	
Wednesday	7	14	21	28	
Thursday	1	8	15	22	29
Friday	2	9	16	23	30
Saturday	3	10	17	24	25
Sunday	4	11	18	25	

153rd Day | 22nd Week

Appointments

Markov Decision Process (MDP)



finite MDP
(fini states)

Infinite MDP (infinite states)

⇒ In this, agent directly observes the environment state

$$O_t = S^{\alpha} = \begin{matrix} \uparrow \\ \text{agent state} \end{matrix} \quad \begin{matrix} \uparrow \\ \text{environment state} \end{matrix}$$

$$S_t = \underbrace{\dots}_{\text{agent state}} \quad \underbrace{\dots}_{\text{environment state}}$$

A state is markovian if and only if state is sufficient statistic of future

$$\text{i.e. } P[S_{t+1}|S_t] = P[S_{t+1}|S_1, S_2, \dots, S_t]$$

State transition matrix - defines the probability of going into another state.

$P_{ss'}$ → prob. of transition probabilities from all states s to all successor states s'

$$P = \begin{bmatrix} P_{11} & & & P_{1n} \\ & \ddots & & \vdots \\ & & P_{nn} & P_{nn} \end{bmatrix}$$

Ex

	JULY '23				
Week	26	27	28	29	30
Monday	31	3	10	17	24
Tuesday	4	11	18	25	
Wednesday	5	12	19	26	
Thursday	6	13	20	27	
Friday	7	14	21	28	
Saturday	1	8	15	22	29
Sunday	2	9	16	23	30

	AUGUST '23				
Week	31	32	33	34	35
Monday	7	14	21	28	
Tuesday	1	8	15	22	29
Wednesday	2	9	16	23	30
Thursday	3	10	17	24	31
Friday	4	11	18	25	
Saturday	5	12	19	26	
Sunday	6	13	20	27	

Saturday

June

2023

03

22nd Week | 154th Day

Appointments

Markov Chain \Rightarrow need just 2 parameters

9.00

 $S \rightarrow f_{t+1}$ set of states

10.00

 $P \rightarrow$ state transition prob matrix

11.00

Markov Reward Process \rightarrow 4 params S, P, R, γ $R \rightarrow$ reward function ; $\gamma \rightarrow$ discount factor

12.00

Return $\Rightarrow G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

1.00

 \Rightarrow represents future rewards.

2.00

if $\gamma \rightarrow 0$: Myopic $\{$ ~~since that means we don't care about future rewards~~ $\}$

3.00

if $\gamma \rightarrow 1$: far-sighted we care consider future rewards

4.00

Value function.

$$V(S) = E [R_{t+1} + \gamma V(S_{t+1}) | S_t = s]$$

5.00

$$= R_s + \gamma \sum_{s' \in S} P_{ss'} V(s')$$

6.00

\downarrow

reward
at current
state

discount value of
the value function of all future
states with their probability of
state transition

7.00

Sunday 04

$$\begin{bmatrix} V(1) \\ \vdots \\ V(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} V(1) \\ \vdots \\ V(n) \end{bmatrix}$$

05

Monday
June
2023

	MAY '23					JUNE '23				
	MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED
MAY 1	1	8	15	22	29	5	12	19	26	2
MAY 2	2	9	16	23	30	6	13	20	27	3
MAY 3	3	10	17	24	31	7	14	21	28	4
MAY 4	4	11	18	25						
MAY 5	5	12	19	26						
MAY 6	6	13	20	27						
MAY 7	7	14	21	28						

2023 Date / 2023 Week

$$\text{d. } v = R + \gamma P v$$

$$\Rightarrow v = (I - \gamma P)^{-1} R$$

But computation complex is $O(n^3)$!

So we need other method for large MRPs.

MDP \rightarrow needs S, A, P, R, γ

A \rightarrow set of actions

policy $\pi \rightarrow$ set of actions over given states.

$$\pi(a|s) = P[A_t = a | S_t = s]$$

$$P_{s,s'} = \sum_{a \in A} \pi(a|s) P_{s,a}$$

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

$$\text{Value } \hat{v}^n(s_0)$$

$$\Rightarrow v_n = E_n [R_{t+1} + \gamma V_n(S_{t+1}) | S_t = s]$$

Action-value $\hat{q}_n^n(s, a)$:

$$\hat{q}_n^n(s, a) = E_n [R_{t+1} + \gamma q_n^n(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

	JULY '23					AUGUST '23					
Week	26	27	28	29	30	Week	31	32	33	34	35
Monday	31	3	10	17	24	Monday	7	14	21	28	
Tuesday	4	11	18	25		Tuesday	1	8	15	22	29
Wednesday	5	12	19	26		Wednesday	2	9	16	23	30
Thursday	6	13	20	27		Thursday	3	10	17	24	31
Friday	7	14	21	28		Friday	4	11	18	25	
Saturday	1	8	15	22	29	Saturday	5	12	19	26	
Sunday	2	9	16	23	30	Sunday	6	13	20	27	

Tuesday

June

2023

06

23rd Week | 157th Day

Appointments

Optimal Value f^n 9.00 → optimal state-value $f^n : V^*(s) = \max_n V_n(s)$ 10.00 → Optimal action-value $f^n : q^*(s, a) = \max_n q_n(s, a)$

11.00 basically acting greedily → to get best performance

12.00
1.00 Optimal policy is the policy π at when, $q^*(s, a) > V^*(s)$.2.00
3.00 $V^*(s) = \max_a q^*(s, a) \Rightarrow$ at particular action 'a' the value f^n is the4.00 {Value of being in a particular state} $\underbrace{\text{under optimal policy}}$ is determined by $\max_a q^*(s, a)$
(selects the action that maximizes expected return) value f^n .5.00 So., $q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_a q^*(s')$ 6.00
7.00 $V^*(s) = \max_a R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s')$ For optimal policy, the maximum π for a set of actions

25

$$\pi^* = \arg\max[\pi] \{ q^*(s, a) \}$$

07

Wednesday
June
2023

188th Day | 23rd Week

Appointments

9:00

10:00

11:00

12:00

1:00

2:00

3:00

4:00

5:00

6:00

7:00

8:00

9:00

MAY '23						JUNE '23					
Week	18	19	20	21	22	Week	22	23	24	25	26
Monday	1	8	15	22	29	Monday	8	15	22	29	
Tuesday	2	9	16	23	30	Tuesday	9	16	23	30	
Wednesday	3	10	17	24	31	Wednesday	10	17	24	31	
Thursday	4	11	18	25		Thursday	11	18	25	29	
Friday	5	12	19	26		Friday	12	19	26	29	
Saturday	6	13	20	27		Saturday	13	20	27	24	
Sunday	7	14	21	28		Sunday	14	21	28	25	

Dynamic Programming \Rightarrow Optimising a program.
 Dynamic becomes the subtask problem because solution can be cached and reused.

\Rightarrow Hence DP is easier for MDP

MDP \Rightarrow requires two kinds of solⁿ

Prediction

\hookrightarrow predicting value $f^n V_n$

control

\hookrightarrow giving optimal policy. π_n

Iterative Policy Evaluation

What is Policy Evaluation?

\Rightarrow process of determining state-value f^n or action-value f^a for a given policy in a MDP in order to estimate how good that policy is.

How to do it?

iteratively updates the value function until it converges to the true value f^n for the policy. In case of iterative policy update we update using Bellman Eqⁿ.

25

JULY '23						
Week	26	27	28	29	30	
Monday	31	3	10	17	24	
Tuesday	4	11	18	25		
Wednesday	5	12	19	26		
Thursday	6	13	20	27		
Friday	7	14	21	28		
Saturday	1	8	15	22	29	
Sunday	2	9	16	23	30	

AUGUST '23						
Week	31	32	33	34	35	
Monday	7	14	21	28		
Tuesday	1	8	15	22	29	
Wednesday	2	9	16	23	30	
Thursday	3	10	17	24	31	
Friday	4	11	18	25		
Saturday	5	12	19	26		
Sunday	6	13	20	27		

Thursday

June

2023

08

23rd Week | 159th Day

Appointments

- 9:00 We can backup the value V^* using synchronous backups
 rather than asynchronous backups.

10:00

Synchronous backup

11:00

$$V_{K+1}(s) = \sum_{a \in A} \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right)$$

12:00

here s' is the successor state of s
 or

2:00

$$V^{K+1} = R^\pi + \gamma P^\pi V^K$$

3:00

So basically if we determine our policy by acting greedily wrt V_K after a number of iterations we will converge to some policy π which is the optimum policy. This is called policy improvement. Policy Evaluation & Policy Improvement go together in RL.

6:00

So, for a given policy π

- Evaluate the policy

$$V_\pi(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \dots | s_t = s]$$

- Improve the policy by acting greedily.

In this process of Policy Evaluation we always get π^* .

09

Friday
June
2023

160th Day | 23rd Week

Appointments

	MAY '23					JUNE '23					
Week	18	19	20	21	22	Week	22	23	24	25	26
Monday	1	8	15	22	29	Monday	5	12	19	26	
Tuesday	2	9	16	23	30	Tuesday	6	13	20	27	
Wednesday	3	10	17	24	31	Wednesday	7	14	21	28	
Thursday	4	11	18	25		Thursday	1	8	15	22	29
Friday	5	12	19	26		Friday	2	9	16	23	30
Saturday	6	13	20	27		Saturday	3	10	17	24	
Sunday	7	14	21	28		Sunday	4	11	18	25	

Value Iteration

Goal is to find the optimal value function which represents the maximum expected return in a MDP.

So we start with an initial estimate of the value function then improve it by iterative updates.

It can be understood as a special case of policy iteration, where we stop after one sweep and update our policy.

Can be written as,

$$v_{k+1}(s) = \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma v_k(s')]$$

The value iteration terminates when it converges to a value v^* .

So, it effectively combines in each step one sweep of policy evaluation & one sweep of



JULY '23						
Week	26	27	28	29	30	
Monday	31	3	10	17	24	
Tuesday		4	11	18	25	
Wednesday	5	12	19	26		
Thursday	6	13	20	27		
Friday	7	14	21	28		
Saturday	1	8	15	22	29	
Sunday	2	9	16	23	30	

AUGUST '23						
Week	31	32	33	34	35	
Monday		7	14	21	28	
Tuesday		1	8	15	22	29
Wednesday	2	9	16	23	30	
Thursday	3	10	17	24	31	
Friday		4	11	18	25	
Saturday		5	12	19	26	
Sunday		6	13	20	27	

Saturday

June

2023

10

23rd Week | 161st Day

Appointments

9:00 policy improvement.

10:00

\Rightarrow Complexity per iteration $\Rightarrow O(mn^2)$ for $V_a(s)$
 ↴ marking
 n states

~~Synchronous~~
~~Dynamic~~
~~Progressive?~~
 1:00 \Rightarrow Complexity $O(m^2n^2)$ per iteration for $q_a(s, a)$

2:00 Asynchronous Dynamic Programming

3 ways \Rightarrow ① In-place dynamic programming

② Prioritised sweeping

③ Real-time dynamic programming

① In-place : stores only 1 copy of value f^n

$$v(s) \leftarrow \max_{a \in A} (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v(s'))$$

② Prioritised sweeping: Backup the state over largest bellman error \Rightarrow ie $\max_{a \in A} (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v(s') - v(s))$

③ Real-time Dynamic Programming :

→ use agent's experience to guide the selection of states

→ After each time step, backup in stat. s_t

$$v(s_t) \leftarrow \max_{a \in A} (R_{s_t}^a + \gamma \sum_{s' \in S} P_{s_t s'}^a v(s'))$$

12

Monday
June
2023

163rd Day | 24th Week

	MAY '23					JUNE '23					
Week	18	19	20	21	22	Week	22	23	24	25	26
Monday	1	8	15	22	29	Monday	5	12	19	26	
Tuesday	2	9	16	23	30	Tuesday	6	13	20	27	
Wednesday	3	10	17	24	31	Wednesday	7	14	21	28	
Thursday	4	11	18	25		Thursday	1	8	15	22	29
Friday	5	12	19	26		Friday	2	9	16	23	30
Saturday	6	13	20	27		Saturday	3	10	17	24	
Sunday	7	14	21	28		Sunday	4	11	18	25	

Appointments

9:00

Model-free Prediction

10:00

- ↳ Two ways:
- ① Monte-Carlo Learning
 - ② Temporal Difference Learning

11:00

~~MC~~

12:00

Monte-Carlo Learning

1:00

→ learn directly from episodes of experience.

Q: What does episode of experience mean?

Ans: So we basically run our agent in the environment to gain experience. ~~From~~

↳ Each run will be an episode. From the episodes of the run experience we will have a set of states with their values, reward etc. We use that in MC learning.

Here the value is the mean return.

Why?

In each to get proper estimate of return it is better to take average of the return at that state in each episode.

	JULY '23				
Week	26	27	28	29	30
Monday	31	3	10	17	24
Tuesday	4	11	18	25	
Wednesday	5	12	19	26	
Thursday	6	13	20	27	
Friday	7	14	21	28	
Saturday	1	8	15	22	29
Sunday	2	9	16	23	30

	AUGUST '23				
Week	31	32	33	34	35
Monday	7	14	21	28	
Tuesday	1	8	15	22	29
Wednesday	2	9	16	23	30
Thursday	3	10	17	24	31
Friday	4	11	18	25	
Saturday	5	12	19	26	
Sunday	6	13	20	27	

Tuesday
June
2023

13

24th Week | 16th Day

Appointments

There are two types of MC-Policy Evaluation.

9:00

① First-visit
(commonly used)

② Every visit

In an episode we may visit a particular state more than once, in that case if we consider only the state visited first it is first-visit other every-visit.

1:00

~~F~~ first-Visit Monte-Carlo Policy Evaluation

2:00

For a policy π we have:

3:00 $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_T, A_T, R_T$
 ↓
 4:00 one episode

5:00 for each episode :

① generate an episode following π : $S_0, A_0, R_1, \dots, S_T, A_T, R_T$

6:00 ② Initialize $G_{t+1}^{(s)}$ $G_{t+1}^{(s)} = 0$

③ for each step of episode, $t = T-1, T-2, \dots, 0$

$G_t^{(s)} = G_{t+1}^{(s)} + R_{t+1}$ $\{ \text{Return} \}$

④ increment counter $N(s) \leftarrow N(s) + 1$

⑤ Value : $V(s) = \frac{1}{N(s)} \sum G_t^{(s)}$

AS

$N(s) \rightarrow \infty$; $V(s) \rightarrow V_\pi(s)$

14

Wednesday
June
2023

165th Day | 24th Week

	MAY '23						JUNE '23					
Week	15	19	20	21	22	Week	22	23	24	25	26	
Monday	1	8	15	22	29	Monday	5	12	19	26		
Tuesday	2	9	16	23	30	Tuesday	6	13	20	27		
Wednesday	3	10	17	24	31	Wednesday	7	14	21	28		
Thursday	4	11	18	25		Thursday	1	8	15	22	29	
Friday	5	12	19	26		Friday	2	9	16	23	30	
Saturday	6	13	20	27		Saturday	3	10	17	24		
Sunday	7	14	21	28		Sunday	4	11	18	25		

Appointments

9:00 Incremental Monte-Carlo Learning10:00 if we update after each ~~episode~~ step in an episode11:00 $S_1, A_1, R_2, \dots, S_T$

12:00 $N(S_t) \leftarrow N(S_t) + 1$

1:00 $V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t))$

we knew the future goal return

2:00 So we care only about the recent episode & not the old episodes

3:00

Hence better than normal MC learning.

4:00

In every visit MC learn

5:00

~~$V(S_t) \leftarrow V(S_t) + \Delta t$~~

6:00

7:00

Incremental Mean

$$\mu_K = \frac{1}{K} \sum_{j=1}^K x_i = \underbrace{\mu_{K-1} + \frac{1}{K} (x_K - \mu_{K-1})}_{\text{increment mean}}$$

JULY '23							AUGUST '23						
SUN	26	27	28	29	30	MON	31	1	2	3	4	5	6
Mond	31	1	2	3	4	Tuesday	7	8	9	10	11	12	13
Tuesday	1	2	3	4	5	Wednesday	14	15	16	17	18	19	20
Wednesday	2	3	4	5	6	Thursday	21	22	23	24	25	26	27
Thursday	3	4	5	6	7	Friday	28	29	30	31	1	2	3
Friday	4	5	6	7	8	Saturday	10	11	12	13	14	15	16
Saturday	5	6	7	8	9	Sunday	17	18	19	20	21	22	23
Sunday	6	7	8	9	10								

Thursday
June
2023

15

24th Week | 166th Day

Appointments

Temporal-Difference Learning

⇒ Learning from incomplete episodes by bootstrapping.

11:00 Simplified \Rightarrow TD - learn algorithm: TD(0)

$$V(S_t) \leftarrow V(S_t) + \alpha \underbrace{(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))}_{= \hat{V}(S_t)}$$

estimated return instead of actual return.

$$\hat{G}_t (\text{TD error}) = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

5:00 Major advantage of TD \Rightarrow can work in non-terminating environment

7:00 TD has low variance & some bias (MC \rightarrow high var, ~~no bias~~)

⇒ TD is more sensitive to initial value.

Bias / Variance Trade-off

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T \Rightarrow \text{unbiased estn of } V(S_t)$$

~~TD target $\Rightarrow R_{t+1} + \gamma V(S_{t+1})$ \Rightarrow biased~~

$G_t \Rightarrow$ depends on many actions, TD target \Rightarrow depends on one action

16

Friday
June
2023n-step return
↳ calculating the return for "n" steps

MAY '23

$$G_t^{(n)} = R_{t+1} + \sum_{k=1}^{n-1} \gamma^k R_{t+k+1} + \gamma^n V(S_{t+n})$$

Week	18	19	20	21	22
Monday	1	8	15	22	29
Tuesday	2	9	16	23	30
Wednesday	3	10	17	24	31
Thursday	4	11	18	25	
Friday	5	12	19	26	
Saturday	6	13	20	27	
Sunday	7	14	21	28	

Week	22	23	24	25	26
Monday	5	12	19	26	
Tuesday	6	13	20	27	
Wednesday	7	14	21	28	
Thursday	1	8	15	22	29
Friday	2	9	16	23	30
Saturday	3	10	17	24	
Sunday	4	11	18	25	

167th Day | 24th Week

Two diff time steps can be combined by averaging
Appointments for combining n-steps, $G_t^{(\lambda)} = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$

9:00 TD(λ)

10:00 TD(λ)

⇒ we used λ-return.

11:00

formular TD(λ)

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^\lambda - V(S_t))$$

1:00 ⇒ basically gives less weight to

2:00 ⇒ doesn't consider whole episode

3:00 Value of λ determines how much weight we want to give to our past values

4:00

It is a midway between TD(0) & MC learning

5:00

if $\lambda \rightarrow 1$ ⇒ more weight to past states & more long-term

6:00

if $\lambda \rightarrow 0$ ⇒ more bias to recent state visits & more in TD(0)

7:00

↳ Bootstrapping - using estimates of future values

Sampling → collects data or experience from the env.

→ Sum of off-the update is identical for forward-win & backward-win TD(λ)

JULY '23						
Week	26	27	28	29	30	
Monday	31	3	10	17	24	
Tuesday	4	11	18	25		
Wednesday	5	12	19	26		
Thursday	6	13	20	27		
Friday	7	14	21	28		
Saturday	1	8	15	22	29	
Sunday	2	9	16	23	30	

AUGUST '23						
Week	31	32	33	34	35	
Monday	7	14	21	28		
Tuesday	1	8	15	22	29	
Wednesday	2	9	16	23	30	
Thursday	3	10	17	24	31	
Friday	4	11	18	25		
Saturday	5	12	19	26		
Sunday	6	13	20	27		

Saturday
June
2023

17

24th Week | 168th Day

Appointments

9:00 Forward View TD(λ)

10:00 → Uses action complete episode
to update value from + down - early

11:00

→ off-the update

12:00

1:00 Backward View TD(λ)

2:00 → ~~no~~ Updates every step for
incomplete sequences.

3:00 → ~~online~~ update

→ Keep an eligibility trace for every states

4:00

$$\Rightarrow V(s) \leftarrow V(s) + \alpha S_t E_t(s)$$

5:00

↑
eligibility trace

6:00

when $\lambda = 0$, only current state's update

$$\text{so } i_t(s) = 1$$

7:00

& TD(λ) becomes equal to TD(0) update

Sunday 18

Eligibility trace:.. maintain a balance b/w recency & frequency-based heuristics.
in simple words, what cause same action?
the most recent action or the most frequent one.
Eligible traces ~~however~~ needs both.

19

Monday
June
2023Q-Learning

(or off-policy TD control)

	MAY 23					JUNE 23					
Week	18	19	20	21	22	Week	22	23	24	25	26
Monday	1	8	15	22	29	Monday	5	12	19	26	
Tuesday	2	9	16	23	30	Tuesday	6	13	20	27	
Wednesday	3	10	17	24	31	Wednesday	7	14	21	28	
Thursday	4	11	18	25		Thursday	1	8	15	22	
Friday	5	12	19	26		Friday	2	9	16	23	
Saturday	6	13	20	27		Saturday	3	10	17	24	
Sunday	7	14	21	28		Sunday	4	11	18	25	

Off

Q-values

$Q^*(s, a) = \text{expected return starting in } s, \text{ taking action } a \text{ as action option}$

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, A_t)]$$

Here learned action value Q directly approximates Q^* , independent of the policy being followed.

so. $Q(s, a) \rightarrow$ expected cumulative reward
we iteratively update until it converges to Q^* .

& we choose actions randomly

we choose random actions greedily but not absolutely greedily we choose random actions with prob. ϵ

JULY '23						
Week	26	27	28	29	30	
Monday	31	3	10	17	24	
Tuesday	4	11	18	25		
Wednesday	5	12	19	26		
Thursday	6	13	20	27		
Friday	7	14	21	28		
Saturday	1	8	15	22	29	
Sunday	2	9	16	23	30	

AUGUST '23						
Week	31	32	33	34	35	
Monday	7	14	21	28		
Tuesday	1	8	15	22	29	
Wednesday	2	9	16	23	30	
Thursday	3	10	17	24	31	
Friday	4	11	18	25		
Saturday	5	12	19	26		
Sunday	6	13	20	27		

Tuesday
June
2023

20

25th Week | 171st Day

Appointments

9.00

Advantage of Q-learning \Rightarrow converges to optimal policy even if we act sub-optimally

10.00

11.00

Decay rate shear or such

$$\sum_{t=0}^{\infty} \alpha_t Q_t(s, a) = \infty \text{ but } \sum_{t=0}^{\infty} \alpha_t^2 Q_t^2(s, a) < \infty$$

1.00

2.00

\Rightarrow Not possible to visit all states (in case of continuous env.)
So we use Approximate Q-learning

3.00

so instead of ~~hours~~ hours a set of Q-functions

4.00

will learn them

5.00

Parameterize Q-function

a which is

6.00

① can be linear f" of feature

$$Q_\theta(s, a) = \theta_0 \cdot j_0(s, a) + \theta_1 \cdot j_1(s, a) + \dots + \theta_n \cdot j_n(s, a)$$

7.00

② neural net

$$\theta_{k+1} \leftarrow \theta_k - \alpha \nabla_{\theta} \int \frac{1}{2} (Q_\theta(s, a) - R(s, a))^2 d\pi_{\theta}(s, a)$$

✓

21

Wednesday
June
2023

172nd Day | 29th Week

Appointment

9:00

So what happens in Approx. Q-Learning

10:00

Instead of storing all the state-action pairs, we use a parameterized f^* like nn to approximate the Q-values.

11:00

we update Θ until we converge to Θ^* such that there is minimum discrepancy between approximate Q-values & true Q-values

12:00

We generally use nn to approx. function

1:00

so we need to approximate Θ^* & nn is a good way to do that

2:00

So what exactly happen in DQN?

3:00

→ We have an experience replay buffer to store agent's experience

4:00

→ We use CNN to estimate Q-values
NN takes States as input and give Q-values as output for each action

5:00

→ We act greedily (or ϵ -greedily) to select action

6:00

→ action to perform, next state is observed & reward is received for that action.

22

	MAY 23					JULY 23					
	18	19	20	21	22		22	23	24	25	26
Monday	1	8	15	22	29		5	12	19	26	3
Tuesday	2	9	16	23	30		6	13	20	27	4
Wednesday	3	10	17	24	31		7	14	21	28	5
Thursday	4	11	18	25			8	15	22	29	6
Friday	5	12	19	26			9	16	23	30	7
Saturday	6	13	20	27			10	17	24	31	8
Sunday	7	14	21	28			11	18	25	30	9

	JULY '23					AUGUST '23									
Week	28	29	30	31	32	33	34	35	Week	31	32	33	34	35	
MONDAY	31	3	10	17	24	Monday	7	14	21	28					
TUESDAY	4	11	18	25		Tuesday	1	8	15	22	29				
WEDNESDAY	5	12	19	26		Wednesday	2	9	16	23	30				
THURSDAY	6	13	20	27		Thursday	3	10	17	24	31				
FRIDAY	7	14	21	28		Friday	4	11	18	25					
SATURDAY	1	8	15	22	29	Saturday	5	12	19	26					
SUNDAY	2	9	16	23	30	Sunday	6	13	20	27					

Thursday
June
2023

22

25th Week | 173rd Day

Appointments

⇒ DQN will sample a batch of expenses from replay buffer

9:00

⇒ for each expense target Q-value is calculated using separate target netw to address moving target issue

10:00

⇒ the params of main network are update using gradient-descent to minimize the estimate Q-value & target Q-value

11:00

Q. What is target Q-value?

Desired Q-values that the network wants to approximate. They are estimated using a target network (copy of the main network). It generates target Q-values by taking the maximum Q-value over all possible actions in the next state acc. to target network.

Since these targets are computed from Q-values they constantly change which cause moving target problem.

By decoupling target values from rapidly changing Q-value we can resolve moving target issue.

25