

Reinforcement Learning (RL) in Practice

- Reward-based learning
- Successful completion of a task
- Positive feedback
- Machine moves on to the track
(a kind of trajectory tracking)

Usage:

LLMs - Large Language Models

DeepSeek - RLHF

DeepMind - AlphaZero

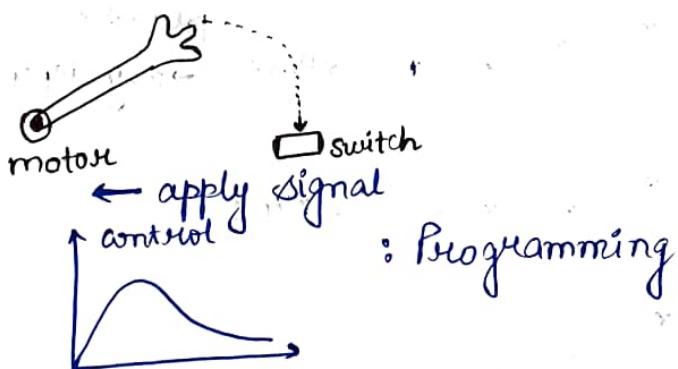
Alpha Star } Read about

Alpha Fold }

Boston Dynamics - robotics application

Think about a robot.

How to make a robot do a "task"?

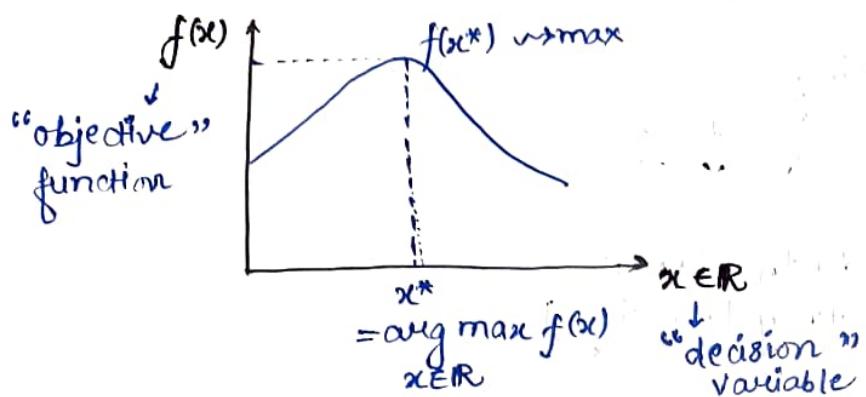


"Method to make this robot learn how to move its arm to achieve a task via reward feedback."

Formal definition:

Achieving an objective is finding maximum of a function
↳ Optimization problem.

Optimization



Eg. Maximize $2x_1 + 3x_2$,

$$x_1 \geq 0, x_2 \geq 0$$

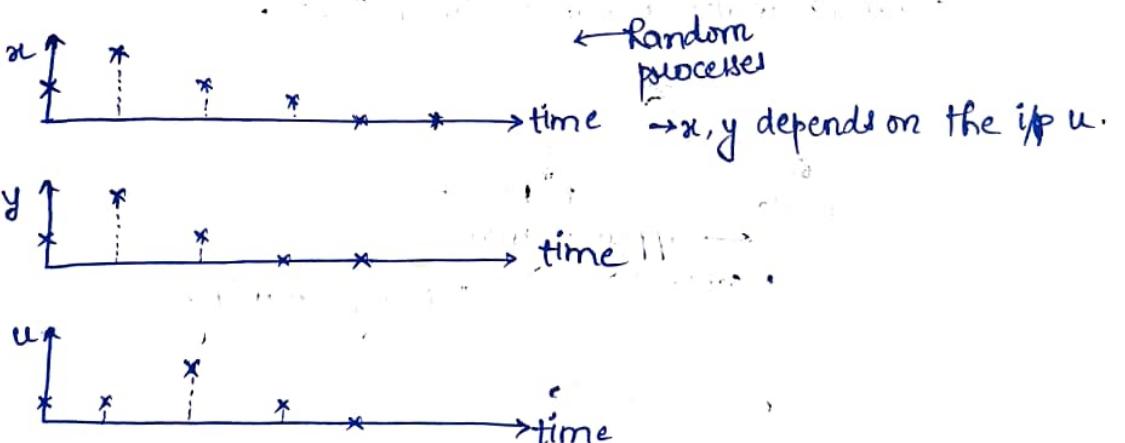
$$x_1 + x_2 \leq 2$$

} A type of
optimization
problem

→ RL is a method to solve sequential optimization problems.

Sequential Optimization (Decision) Problems ⊆ Optimization problems

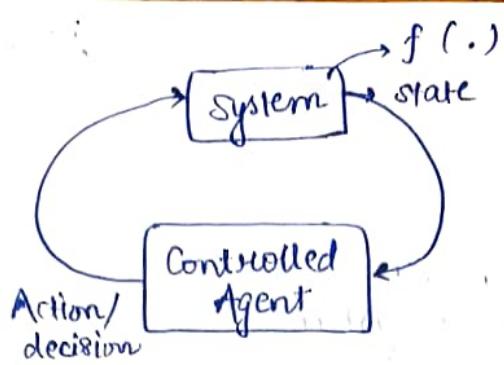
↓
notion of time



$\max f$ (state evolution)
seq. decision

$$\frac{dx}{dt} = g(x, u)$$

Proposition (P) controller : $\frac{du}{dt} = g(x) + K \cdot x$
initial state



Optimal control theory

RL: Mathematical framework of optimizing a system without knowing what the system is.

- ↳ Learns the properties of system.
- ↳ Uses machine learning and deep learning.

Course: 40% theory + 60% Hands-on [From control perspective]

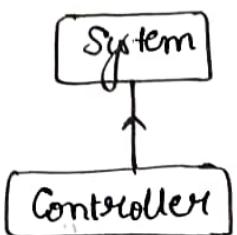
Prerequisites:

- Probability and random processes
- Programming (Python)

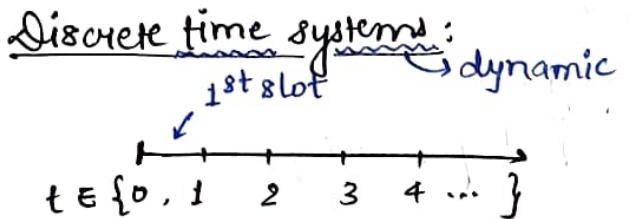
Evaluation:

- No written exams (No endsem exams!)
- Programming lab exam (no internet!)
- class tests
- Project (RL implementation)

Reinforcement Learning



Learn how to control a system.



Dynamical Systems

- ↳ Objects having state which is a function of time.
- discrete-time system: discrete time, t , function.



→ state

↳ vector $s \in \mathbb{R}^d$

$s_t \rightarrow$ notation

$s_t \in S \rightarrow$ state space

↳ fixed

$s_t, s(t)$

↓ discrete ↓ continuous

E.g. @ Grid world

1	2	3	4
5	robot	.	.
.	.	.	.
.	.	.	.
.	.	.	.

$$A = \{n, s, e, w\}$$

$$S = \{1, \dots, 16\}$$

(b)



How does the state change with time?

Control inputs:

At every time t, apply some control input or action, At.

$A_t \in A$ (action space)

Model: Evolution of state (time-variant evolution)

$$S_{t+1} = f(S_t, A_t),$$

f: transition function

↪ can be represented as table

S_t	A_t	S_{t+1}	[Gridworld]
1	e	2	
1	n	1	

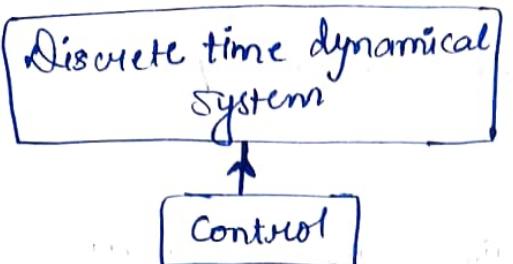
Eg. Bike moving along x-direction.

$$S_t = \begin{pmatrix} x_t \\ v_t \\ a_t \end{pmatrix}, F.m$$

$$S_{t+1} = \begin{pmatrix} x_{t+1} \\ v_{t+1} \\ a_{t+1} \end{pmatrix} = \begin{pmatrix} x_t + v_t + \frac{1}{2} a_t t \\ v_t + a_t t \\ a_t + \frac{F}{m} \end{pmatrix}$$

$$S_{t+1} = \begin{pmatrix} 1 & 1 & 0^{1/2} \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_t \\ v_t \\ a_t \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ F/m \end{pmatrix}$$

$$\therefore S_{t+1} = AS_t + A_t$$

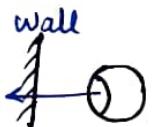
Review

State - $s_t \in S$

Action - $a_t \in A$

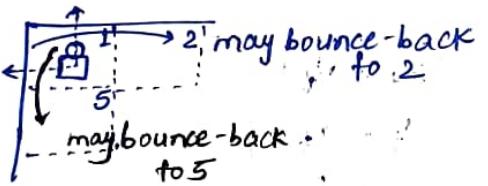
Transitions - $s_{t+1} = f(s_t, a_t)$
 \uparrow
 transition fn

Eg



may not stop in time & bounce back.

Gridworld:



Eg

$$s_{t+1} = \hat{A} \cdot s_t + a_t$$

diff. matrix

$$s_{t+1} = A \cdot s_t + a_t + e_t$$

\hookrightarrow error

\hookrightarrow non-deterministic

In many scenarios,

s_{t+1} is not deterministic, given s_t, a_t .

Model s_{t+1} as a random variable - dependent on s_t & a_t .

Start from a s_0 .

s_1, s_2, \dots etc are all RVs.

For grid world,

suppose $s_0 = 1$, $a_0 = n$.

s_1 has a conditional pmf with

$$\Pr\{s_1 = i \mid s_0 = 1, a_0 = n\} \quad \left[\Pr\{s_1 = i \mid s_0 = 1, a_0 = w\} \text{ is a different pmf.} \right]$$

Transitions:

$$s_{t+1} = F(s_t, a_t) \quad \begin{array}{l} \text{[change of notation]} \\ \text{not deterministic} \end{array}$$

For now, we will assume that S, A are discrete (finite) sets.

$$s_{t+1} = F(s_t, a_t)$$

$$= \Pr\{s_{t+1} = s_{t+1} \mid s_t = s_t, a_t = a_t\} \rightarrow \begin{array}{l} \text{can change} \\ \text{with time or} \\ \text{for simplicity} \\ \text{take as time-} \\ \text{invariant.} \end{array}$$

Policy: \rightarrow Tells how to control the system.

Sequence of control actions (a_0, a_1, a_2, \dots) .

\hookrightarrow A function that looks at the current state and decide what action to take. [A ~~is~~ type of policy — this one is the usual.]

$$\pi(s_t) = a_t$$

\hookrightarrow May not depend on any state, or may depend on past actions.

Eg.

1	2	3	4
5			



Policy e.g.,

s_t	$\pi(s_t) = A_t$	→ deterministic fn of current state
1	e	
2	e	
3	e	
4	s	

Random (stochastic) policies:

$$P \{ A_t = a | s_t = s \}$$

we have two kinds of prob.: Probability of deciding $\xrightarrow{\text{then}}$ Probability of deciding the next stage.
an action

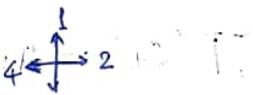
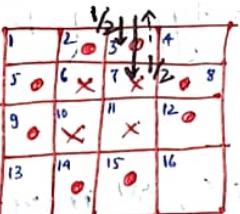
Trajectory: Set of states and actions.

$s_0, A_0, s_1, A_1, s_2, A_2 \rightarrow$ infinite horizon

$\uparrow \quad \uparrow$
 $\pi(s_0) \quad \pi(s_1)$

Horizon = stopping condition
 s_{T-1} (T -horizon) is stop at $(T-1)$.
finite

H.W Simulate a trajectory of a robot moving in the grid-world.



Transition probability model

π -policy

$$P \{ A = d | s_t = s \} = \frac{1}{4}$$

$s_0 = 1$; current state = s_0

for {

current $A =$ sample an action based on $P \{ A = a | \text{current state} \}$

next state = sample from $P \{ s_{t+1} = s_{t+1} | s_t = \text{current state}, A_t = \text{current } A \}$

current state = next state

} store in array to generate trajectory

Recall

Discrete time dynamical system:

State - $s_t, t \in \{0, 1, \dots\}$

State space - S

Action - a_t

Action space - A

Evolutional Transition:

$$s_{t+1} = F(s_t, a_t)$$

$$\Pr\{s_{t+1} = s_{t+1} | s_t = s_t, a_t = a_t\}$$

Policy:

deterministic: $a_t = \pi(s_t)$

stochastic: a_t

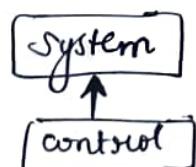
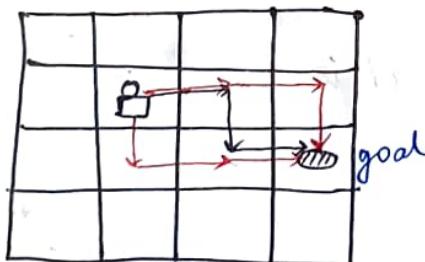
$$\Pr\{a_t = a_t | s_t = s_t\}$$

Trajectory:

$s_0, a_0, s_1, a_1, \dots, s_{T-1}$: finite horizon

$\dots \infty$: infinite horizon

Sample Path: Another term for trajectory.

Rewards and Costs

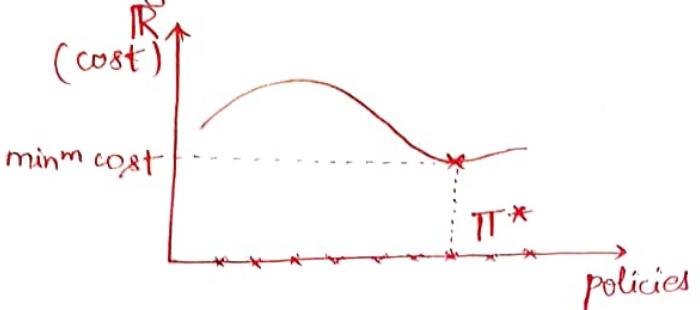
Task/Objective/Goal:

Example: Reach the in minimum no. of steps.

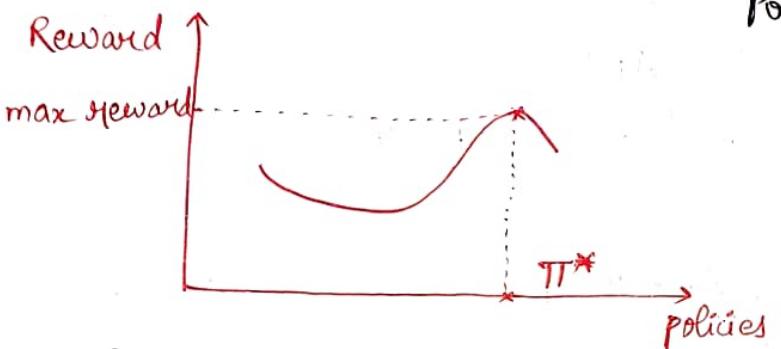
↳ which policy is the best?

↳ Mathematical model?

Ranking the policies (although they're complex mathematical objects):



Policy with lower value is better: cost



Policy with higher value is better: reward

Duality of cost & Reward: Minm cost ~~occurs~~ occurs at max reward.

→ Negative of cost is reward.

→ Usually reward is used in RL.

Stage-wise Rewards:

Example for grid-world:

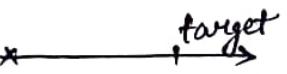
At every time 't', there is a reward

$$R_t = R(s_t, a_t), \dots$$

and for the policy, reward is

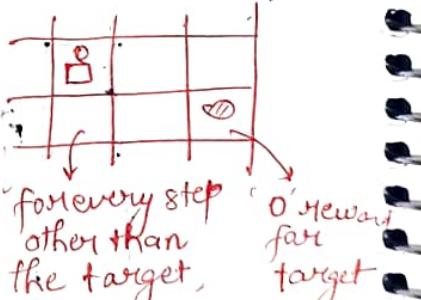
$$\sum_{t=0}^T R_t$$

Eg.



$$S_{t+1} = AS_t + A_t$$

$$R_t = -C_t = -|S_t - S_{\text{target}}|^2$$



Trajectory:

$S_0, A_0, R_0, S_1, A_1, R_1, \dots \rightarrow$ all are random variables
 ↗
 Rewards will be random (in general)
 even when R fn is deterministic.

In general,

$R \rightarrow$ stochastic

$$\Pr\{R_t = r_t | S_t = s_t, A_t = a_t\}$$

$\sum_{t=0}^{\infty} R_t \rightarrow$ scalar but random (sum of RVs)

To rank the policies, we need a single value,
 so take average value of reward (expectation):

$$\mathbb{E} \sum_{t=0}^{\infty} R_t$$

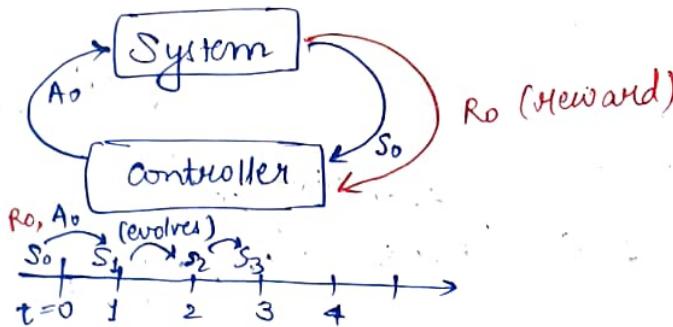
for infinite horizon,

$\sum R_t$ may be infinite, so not possible to rank two policies.

- Discounted Rewards:

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \right], \quad \gamma \in (0, 1).$$

Lec-4 19-08-2025



$$A_t = \pi(s_t)$$

$$\Pr\{A_t = a | s_t = s\}$$

$$\boxed{\pi(a|s)}$$

$$S_{t+1} = F(s_t, A_t)$$

$$\Pr\{S_{t+1} = s_{t+1} | s_t = s, A_t = a\}$$

$$\boxed{P(s_{t+1} | s_t, a_t)} \rightarrow \text{shortcut}$$

$$R_t = R(s_t, A_t)$$

Trajectory: $(s_0, a_0, r_0, s_1, a_1, r_1, \dots)$

$$\begin{pmatrix} s_0 & s_1 & \dots \\ a_0 & a_1 & \dots \\ r_0 & r_1 & \dots \end{pmatrix}$$

Return: Total rewards along a trajectory

$$\sum_{t=0}^{\infty} R_t$$

\downarrow
random

$E\left(\sum_{t=0}^{\infty} R_t\right)$: expected return

Discounted return: For infinite horizon

$$\sum_{t=0}^{\infty} \gamma^t \cdot R_t, \quad \gamma: \text{discount factor } \in (0, 1]$$

Given some policy π , $E\left[\sum_{t=0}^{\infty} \gamma^t \cdot R_t\right]$ can be computed.

Find the policy π with the largest value of the above expected return.

Markov decision Process (MDP)

$$\max_{\pi} E\left[\sum_{t=0}^{\infty} \gamma^t \cdot R_t\right]$$

such that

$$R_t = R(s_t, a_t)$$

$$s_t \in S, \quad a_t \in A$$

$$P = P \left\{ s_{t+1} = s_{t+1} \mid s_t = s_t, a_t = a_t \right\}, \quad \begin{array}{l} \text{transition prob/} \\ \text{frn} \end{array}$$

$$\gamma \in [0, 1]. \quad \begin{array}{l} \text{state} \\ \text{evolution} \end{array}$$

eg:



K : max queue size

Counter (eg, bank counter)
Two persons servicing:
 $\mu_1 < \mu_2 \rightarrow$ more experienced
 $c_1 < c_2 \rightarrow$ payments

In every slot t , a person comes into this system with probability λ .

In every slot, if a person is served, she finishes service with probability μ .

State = num. in queue

S_t :

$$\Pi(S_t) \in \{\mu_1, \mu_2\}$$

$$R(S_t) = \underbrace{(S_t + f_c c_t)}_{\text{To reduce the cost}} \xrightarrow{\text{multiplying factor}} \text{the money paid (additional cost)}$$

$$c_t = \begin{cases} c_1, & \text{if } \Pi(S_t) = \mu_1 \\ c_2, & \text{if } \Pi(S_t) = \mu_2 \end{cases}$$

If $S_t=0$, $A_t \in \{\mu_1, \mu_2\}$, [Special case]

$$S_{t+1} = 0, \text{ w.p. } (1-\lambda), \text{ (no one comes)} \\ = 1, \text{ w.p. } \lambda. \text{ (one person comes)}$$

If $S_t=k$, $A_t=\mu_1$,

$$S_{t+1} = k-1, \text{ w.p. } \underbrace{\mu_1(1-\lambda)}_{\text{no one comes, one served}} \\ \quad \quad \quad = k, \text{ w.p. } (1-\lambda)(1-\mu_1) + \lambda \mu_1.$$

$$= k+1, \text{ w.p. } \underbrace{(1-\mu_1)\lambda}_{\text{no one comes, no one served}} \quad \underbrace{\lambda \mu_1}_{\text{one came, one served}}$$

If $S_t=j$, $A_t=\mu_1$,

$$S_{t+1} = j-1, \text{ w.p. } (1-\lambda)\mu_1$$

$$j, \text{ w.p. } (1-\mu_1)(1-\lambda) + \mu_1 \lambda \\ j+1, \text{ w.p. } \lambda(1-\mu_1).$$

\rightarrow Transition fn.

$$MDP = (S, A, P, R, \gamma)$$

Read ch-1
- MARKOV
 ~ due to the conditional prob. that next state depends upon S_t, A_t
 - Episodic & Continuity
 [Book: Mathematical model of RL]

Markov decision process:

$$(S, A, P, R, \gamma)$$

Eg Machine repair problem

$S = \{G_1, B\}^5$, no. of G_1 (good) components.

$A = \{\text{run}, \text{repair}\}$

$P_\delta = \{S_{t+1} = g' | S_t = g, A_t = \text{run}\}$

$$g' = g - \Delta,$$

$$\Delta \in \{0, \dots, g\}$$

$$P_\delta \{ \Delta = \delta \} = \binom{g}{\delta} p^\delta (1-p)^{g-\delta}$$

$$P_\delta \{ S_{t+1} = g - \delta | S_t = g, A_t = \text{run} \}$$

$$= \binom{g}{\delta} p^\delta (1-p)^{g-\delta}$$

$$P_\delta \{ S_{t+1} = 5 | S_t = g, A_t = \text{repair} \} = 1 \quad [\text{everything is good}]$$

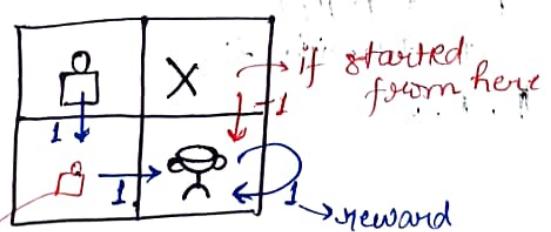
$$\text{Reward } R(S_t, \text{run}) = R$$

$$R(S_t, \text{repair}) = -C$$

$$E \left[\sum_{t=0}^{\infty} \gamma^t R_t \right] \rightarrow \text{to maximise this}$$

State value $\rightarrow E \left[\sum_{t=0}^{\infty} \gamma^t R_t | S_0 = s_0 \right]$: Metric is dependent on the initial state.

Eg:



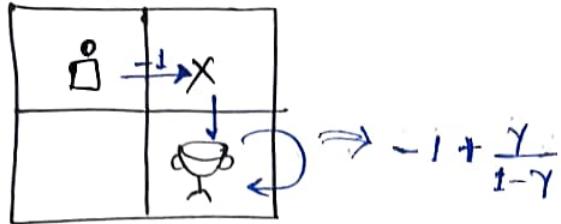
if started from here what is the metric for this policy?

$$\rightarrow \frac{1}{1-\gamma}, \frac{1}{1-\gamma}, -1 + \frac{\gamma}{1-\gamma}$$

if started from \square

if started from X

Eg.



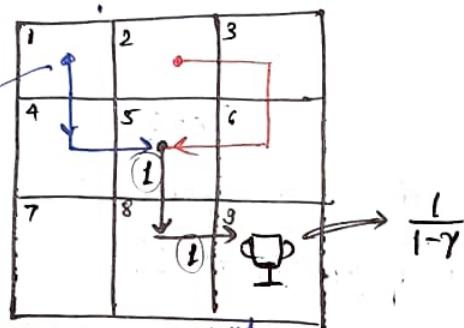
$$-1 + \frac{1}{1-\gamma}$$

$(-1 + \frac{1}{1-\gamma})$ or $(\frac{1}{1-\gamma})$ → which is better?

↓
compare
the policy

Eg.

if started from here



$$\frac{1}{1-\gamma}$$

state → reward
 $(5, 1), (8, 1) \rightarrow (9, 1) \dots$
 $(1, 1), (4, 1), (5, 1), (8, 1) \rightarrow (9, 1) \dots$
 $(2, 1), (3, 1), (6, 1), (5, 1), (8, 1) \rightarrow (9, 1) \dots$

$$\mathbb{E} \left[\sum_{t=t}^{\infty} \gamma^{t-t} R_t \mid S_t=s \right]$$

↓
expected return
from time 't'

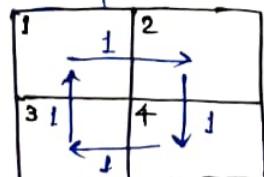
$$= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0=s \right] = V^\pi(s)$$

→ The value of the state does not depend on the time when we started.

→ doesn't depend
on time; depends
on the initial state.

Find the policy that maximizes $V^\pi(s)$.

→ How to compute $V^\pi(s)$?



$$1 + \gamma + \gamma^2 + \dots = \frac{1}{1-\gamma}$$

$$V^\pi(1) = V^\pi(2) = V^\pi(3) = V^\pi(4) = \frac{1}{1-\gamma}$$

$V^\pi(1)$	$= 1 + \gamma V^\pi(2)$
$V^\pi(2)$	$= 1 + \gamma V^\pi(3)$
$V^\pi(3)$	$= 1 + \gamma V^\pi(4)$
$V^\pi(4)$	$= 1 + \gamma V^\pi(1)$

↓
same quantities

solve for these

: Bootstrapping
 ↓
 V^π written in terms of itself.

Read ch-2

Lec-6 [26-08-2025]

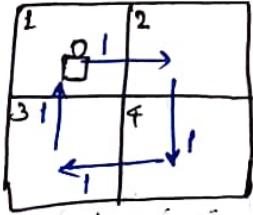
$V^\pi(s)$: expected return (sum of discounted return)
 for the policy π starting from state 's'.

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t | s_0 = s \right]$$

↓
 defined $\forall s \in S$

↳ Also called State value function or state function.

Eg:



$$\begin{bmatrix} 1 & -\gamma & 0 & 0 \\ 0 & 1 & -\gamma & 0 \\ 0 & 0 & 1 & -\gamma \\ -\gamma & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} V^\pi(1) \\ V^\pi(2) \\ V^\pi(3) \\ V^\pi(4) \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

How to express $V^\pi(s)$ (in terms of $V^\pi(s')$, $\forall s \in S$) in general?

E.g.

	1	2	3	4
5	6	7	8	
9	10	11	12	

$\pi(a|s)$: probability that action is 'a', given current state is 's'. \rightarrow policy

$P'(s'|s, a)$: pmf of $s_{t+1} = s'$, state is 's', action is 'a'.

$P(r|s, a)$: distribution of reward. \hookrightarrow Transition fn.

$$V^\pi(s) = \mathbb{E}[R_0 | s_0 = s] + \gamma \cdot \mathbb{E} \left[\underbrace{\sum_{t=1}^{\infty} \gamma^{t-1} \cdot R_t}_{\text{if we had } Z} \mid s_0 = s \right]$$

if we had

$$\mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} \cdot R_t \mid s_1 = s_1 \right] = V^\pi(s_1).$$

$$\mathbb{E}[Z | s_0 = s_0] = \sum_Z \Pr\{Z = z | s_0 = s_0\} \cdot z$$

$$= \sum_{s_1} \sum_Z \Pr\{Z = z, s_1 = s_1 | s_0 = s_0\} \cdot z : \text{Marginal of joint distn.}$$

$$= \sum_{s_1} \sum_Z \Pr\{Z = z | s_1 = s_1, s_0 = s_0\} \cdot \Pr\{s_1 | s_0\} \cdot z$$

$$= \sum_{s_1} \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} \cdot R_t \mid s_0 = s_0, s_1 = s_1 \right] \cdot \Pr\{s_1 | s_0 = s_0\}.$$

$$= \sum_{s_1} V^\pi(s_1) \cdot \Pr\{s_1 = s_1 | s_0 = s_0\}.$$

$$V^\pi(s) = \underbrace{\mathbb{E}[R_0 | s_0 = s]}_{\text{Immediate reward}} + \underbrace{\gamma \cdot \sum_{s_1} P(s_1 = s_1 | s_0 = s_0) \cdot V^\pi(s_1)}_{\text{Average of rewards for next states}}.$$

$$= \sum_a P(a | s_0) \cdot r + \gamma \sum_{s_1} P(s_1 = s_1 | s_0 = s_0) \cdot V^\pi(s_1).$$

$$\Rightarrow V^\pi(s) = \sum_a \underbrace{\pi(a | s_0)}_{\substack{\text{avg. over actions} \\ \downarrow \text{marginalize over actions}}} \cdot \underbrace{\sum_r \Pr(a = a | s_0, a)}_{\substack{\text{avg. over rewards} \\ \downarrow \text{avg. over rewards}}} + \gamma \sum_a \pi(a | s_0) \cdot \sum_{s_1} P(s_1 = s_1 | s_0 = s_0, a = a) \cdot V^\pi(s_1)$$

Bellman's Equation

How to compute $V^\pi(s)$?

Read about Richard Bellman (Wikipedia)

$$H^\pi(s_0) = \sum_a \pi(a|s_0) \cdot \sum_{s_1} P(s_1|s_0, a_0) H$$

$$\sum_a \pi(a|s_0) \cdot \sum_{s_1} \underbrace{P(s_1|s_0, a_0)}_{\text{prob. of action}} \cdot \underbrace{V^\pi(s_1)}_{\text{prob. of next state}}$$

$$= \sum_{s_1} \left(\sum_a \pi(a|s_0) \cdot P(s_1|s_0, a_0) \right) \vdots V^\pi(s_1)$$

$$\left[\equiv \sum a_{s_1} b_{s_1} \rightarrow \text{dot product of 2 quantities} \right]$$

$$= (P^\pi(1|s_0) \quad p^\pi(2|s_0) \dots) \begin{pmatrix} V^\pi(1) \\ V^\pi(2) \\ \vdots \\ V^\pi(s) \end{pmatrix}$$

$$V^\pi(s_0) = H^\pi(s_0) + \gamma (P^\pi(1|s_0) \quad p^\pi(2|s_0) \dots) \begin{pmatrix} V^\pi(1) \\ V^\pi(2) \\ \vdots \\ V^\pi(s) \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} V^\pi(1) \\ V^\pi(2) \\ V^\pi(3) \\ \vdots \\ V^\pi(s) \end{pmatrix} = \begin{pmatrix} H^\pi(1) \\ H^\pi(2) \\ \vdots \\ H^\pi(s) \end{pmatrix} + \gamma \begin{pmatrix} P^\pi(1|1) \quad P^\pi(2|1) \dots \\ P^\pi(1|2) \quad P^\pi(2|2) \dots \\ \vdots \quad \vdots \\ \vdots \quad \vdots \end{pmatrix} \begin{pmatrix} V^\pi(1) \\ V^\pi(2) \\ \vdots \\ V^\pi(s) \end{pmatrix}$$

$$\Rightarrow \bar{V}^\pi = \bar{H}^\pi + \gamma \bar{P}^\pi \bar{V}^\pi$$

$$\Rightarrow \bar{V}^\pi = (I - \gamma \bar{P}^\pi)^{-1} \bar{H}^\pi$$

\downarrow
Matrix inverse exists

\hookrightarrow computationally very difficult

\downarrow
iterative algorithm

Read section 2.7
for invertibility

State value : $v^\pi(s)$

Metric for evaluating / ranking policies?

$v^\pi(\cdot) \rightarrow$ state value function

$\bar{v}^\pi = (v^\pi(1) \dots v^\pi(s)) \rightarrow$ vector

Bellmann's equation:

$$v^\pi(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) v^\pi(s') + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \bar{v}^\pi(s')$$

In metric form,

$$\bar{v}^\pi = \pi^\pi + \gamma \cdot P^\pi \bar{v}^\pi$$

$$\bar{v}^\pi = (I - \gamma P^\pi)^{-1} \cdot \pi^\pi$$

Claim: \bar{v}^π - unique so \exists v^π exists.

Iterative Algorithm for finding \bar{v}^π

Start with $\bar{v}_0 = (v_0(1), v_0(2), \dots, v_0(s))^T$ [all zero]

Compute $\bar{v}_{k+1}^\pi = \bar{v}^\pi + \gamma P^\pi \bar{v}_k^\pi$, with $\bar{v}_0^\pi = \bar{v}_0$ to start with.

Repeat.

$$\lim_{k \rightarrow \infty} \bar{v}_k^\pi = \bar{v}^\pi$$

Stop when a close enough value is achieved.

$$\max_s |v_{k+1}(s) - v_k(s)| < \epsilon v_k(s),$$

then stop these iterations.

$$v_k(s) \approx v^\pi(s).$$

Assignment-①: Prove that $\lim_{k \rightarrow \infty} \bar{v}_k^\pi = \bar{v}^\pi$.

[Read about "contraction maps" (wiki).]

[Apply contraction mapping theorem.]

Assignment-②: Implement Value iteration for policy evaluation in Python.

Ex. ISI, IAI

0.1 0.9
0.5 0.5 } Prob. for 2 states

Reward \rightarrow deterministic

2

3

4

Transition prob.: 0.1, 0.9
 $P(s'|s, a)$

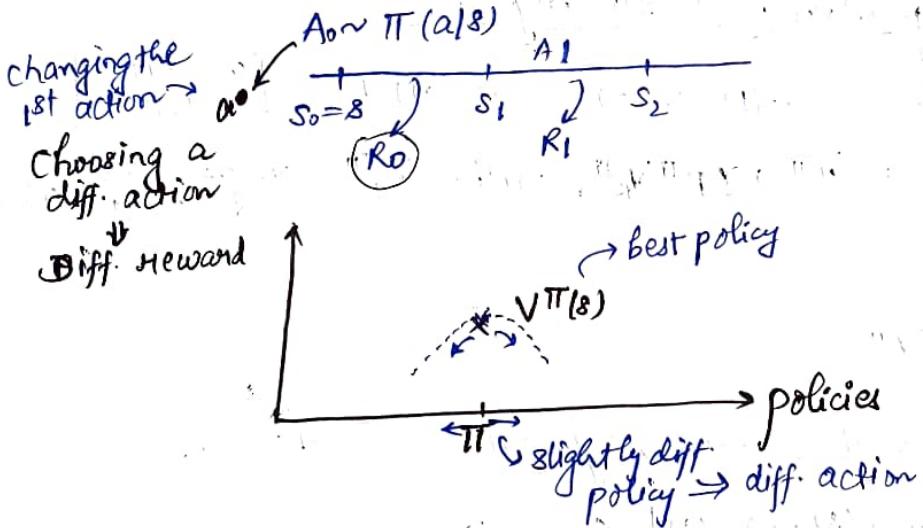
State-action value function:

| # Q-learning

Recall:

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t | S_0 = s \right]$$

\hookrightarrow fn of state given a policy



$$q^\pi(s, a) = \mathbb{E}[R(s, a)] + \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^t R_t | S_0 = s, A_0 = a \right]$$

\hookrightarrow Q-function

\hookrightarrow Diff. value at the 1st instant

\hookrightarrow further values depend on the action taken.

\hookrightarrow further behaviour will depend on the policy:
same from $V^\pi(s)$ & $q^\pi(s, a)$.

$$q^{\pi}(s, a) = \mathbb{E}[R(s, a)] + \gamma \mathbb{E}\left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t | s_0 = s_0, A_0 = a\right]$$

$\sum_{s_1} P(s_1 | s_0, a) \cdot V^{\pi}(s_1)$

$$\Rightarrow q^{\pi}(s, a) = \mathbb{E}[R(s, a)] + \gamma \sum_{s_1} P(s_1 | s_0, a) \cdot V^{\pi}(s_1)$$

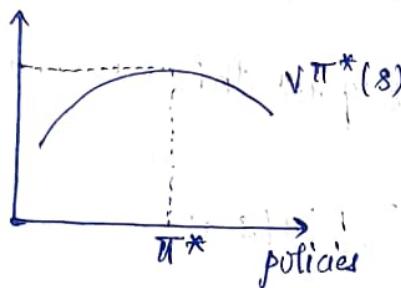
$$V^{\pi}(s) = \sum_a \pi(a|s) \left(\mathbb{E}[R(s, a)] + \gamma \sum_i P(s_i | s, a) \cdot V^{\pi}(s_i) \right)$$

$$\Rightarrow V^{\pi}(s) = \sum_a \pi(a|s) \cdot q^{\pi}(s, a)$$

Read section 2.8

Optimal policies : MDP = (S, A, P, R, γ)

18 01-09-2025



What we have now:

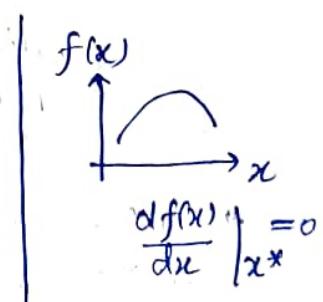
Given a policy π , $V^{\pi}(s)$, $\forall s \in S$,
or, V^{π} and $q^{\pi}(s, a)$ \forall

Defn: A policy π^* is optimal if
 $V^{\pi^*}(s) \geq V^{\pi}(s)$ for every
 $s \in S$ and any policy π .

$V^{\pi^*}(s)$: Optimal state value

Ranking: Condition for policies

π_1 is better than π_2 , if $V^{\pi_1}(s) \geq V^{\pi_2}(s)$, $\forall s \in S$.



$$\pi^* \rightarrow \pi^*(a|s) \xrightarrow[\text{pmf}]{\text{conditional}} \left. \begin{array}{l} \text{Given} \\ \hookrightarrow \text{Actually unknown} \end{array} \right\}$$

We know,

$$V^\pi(s) = \sum_a \pi^*(a|s) \cdot \sum_\mu p(\mu|s, a) \cdot r + \gamma \sum_a \pi^*(a|s) \cdot \sum_{s'} p(s'|s, a) \cdot V^\pi(s')$$

$\xrightarrow{s} \xrightarrow{\pi(a|s)} \xrightarrow{\pi^*}$

For $\pi(a|s)$, only (first) control: choose π . # e.g. $\pi(a|s) = (0.1, 0.8, 0.1)$

$$\max_\pi \left(\sum_a \pi(a|s) \sum_\mu p(\mu|s, a) \cdot r + \gamma \sum_a \pi(a|s) \cdot \sum_{s'} p(s'|s, a) \cdot V^\pi(s') \right) = V^\pi(s)$$

→ maxm/optimal after 1st step

↳ what is this? (interpret this total reward)

$$V^\pi(s) = \max_\pi \left\{ \sum_a \pi(a|s) \sum_\mu p(\mu|s, a) \cdot r + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \cdot V^\pi(s') \right\}$$

Bellman's
Optimality
Equation
(BOE)

↳ Necessary condition and sufficient condition

$\pi^* \rightarrow \max V^\pi(s)$.

$$V^{\pi^*}(s) \geq V^\pi(s)$$

$$\begin{aligned} & \mathcal{H}_0^* + \gamma \mathcal{H}_1^* + \gamma^2 \mathcal{H}_2^* + \dots \geq \\ & \mathcal{H}_0 + \gamma \mathcal{H}_1 + \gamma^2 \mathcal{H}_2 + \dots \stackrel{\text{may not be true}}{\geq} \end{aligned} \quad \begin{array}{l} \text{Digression: Tradeoff b/w} \\ \text{immediate rewards \& future} \\ \text{value} \end{array}$$

only total reward dominates; not individual rewards

$$\rightarrow V(s) = \max_\pi \left\{ \sum_a \pi(a|s) \sum_\mu p(\mu|s, a) \cdot r \right.$$

$$\left. + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \cdot V(s') \right\}$$

then $V(s)$ is $V^{\pi^*}(s)$.

The argmax is $\pi^*(s)$.

↳ jointly maximise for both π^* and $V(s)$.

$$\begin{aligned}
 & \max_{\pi_0, \pi_1, \pi_2} \left\{ \mathbb{E} \sum_{t=0}^{\infty} \gamma^t R_t \mid s_0 = s \right\} \\
 &= \max_{\pi_0, \pi_1, \pi_2} \left\{ \mathbb{E}[R_0 \mid s_0 = s] + \gamma \cdot \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t \mid s_0 = s \right] \right\} \\
 &= \max_{\pi_0, \pi_1, \pi_2} \left\{ \underbrace{\sum_a \pi_0(a|s)}_{\pi_0 \text{ controls this}} \cdot \underbrace{\sum_h p(h|s, a) \cdot h}_{\pi_1, \pi_2, \dots \text{ control this}} \right. \\
 &\quad \left. + \gamma \sum_a \pi_0(a|s) \cdot \sum_{s'} p(s'|s, a) \cdot \left(\underbrace{\mathbb{E} \sum_{t=1}^{\infty} \gamma^{t-1} R_t \mid s_1 = s'}_{\text{original problem (with different time)}} \right) \right\} \\
 &= \max_{\pi_0} \left\{ \sum_a \pi_0(a|s) \sum_h p(h|s, a) \cdot h \right. \\
 &\quad \left. + \gamma \sum_a \pi_0(a|s) \sum_{s'} p(s'|s, a) \cdot \max_{\pi_1, \pi_2} \mathbb{E} \sum_{t=1}^{\infty} \gamma^{t-1} R_t \mid s_1 = s' \right\} \\
 &\hookrightarrow \pi_0 = \pi_1 = \pi_2 \text{ maximizes this.}
 \end{aligned}$$

Read ch. 3 & 4.

Lec-9

Recall: $V\pi^*(s) \geq V\pi(s)$, $\forall s \in S$, and any other π .

To find π^* and $V\pi^*(s)$,

find any solution to the Bellman's Optimality equation.

$$\begin{aligned}
 V(s) = & \underset{\pi}{\operatorname{argmax}} \left\{ \sum_a \pi(a|s) \sum_h p(h|s, a) \cdot h \right. \\
 & \left. + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \cdot V(s') \right\}
 \end{aligned}$$

then a soln,

$$V(s) = V\pi^*(s)$$

$\underset{\pi}{\operatorname{argmax}} \pi = \pi^*(s) \Rightarrow$ Optimal policy

How to solve? + Some important points regarding this eqn:

- Compare with Bellman's equation for a particular policy π .
 - ↪ Here there are two unknowns: V and $\underset{\pi}{\operatorname{argmax}}$.
 - ↪ \max

$$\max_{\pi} \left\{ \sum_a \pi(a|s) \left\{ \sum_u p(u|s,a) u + \gamma \sum_{s'} p(s'|s,a). V(s') \right\} \right\}$$

like $q(s,a)$, $\rightarrow f^n$ of a,s .

$$= \max_{\substack{\pi(a|s) \\ \sum_a \pi(a|s) = 1 \\ \pi(a|s) \geq 0}} \left\{ \sum_a \pi(a|s) \cdot q(s,a) \right\}$$

Eg. $a \in \{1, 2, 3\}$

$q(8, 1) = 2$	[
$q(8, 2) = 10$	
$q(8, 3) = 11 \xrightarrow{\text{give max weight to the max of these}} 11$	

$\pi(3|8) = 1$

$$= \max_a q(s,a)$$

\Rightarrow Optimal policy is deterministic.

$$V(s) = \max_a \left\{ \sum_u p(u|s,a) u + \gamma \sum_{s'} p(s'|s,a) V(s') \right\}$$

argmax of 'a' gives $\pi^*(s)$.

For MDPs (S, A, P, R, γ) , it is enough to look for deterministic policy.

$\sum_u p(u|s,a) u$: expectation of rewards
 $= \bar{R}(s,a)$

$\gamma \sum_{s'} p(s'|s,a) \cdot V(s')$ \rightarrow dot product

$$= \gamma \langle \quad , V(\cdot) \rangle$$

$s \xrightarrow{P(s'|s,a)}$ if the action is a
 matrix of row vectors
 tells prop. of evolving from s to s'

Machine e.g.

$$\begin{pmatrix} 0 & 1 & 2 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 2 & 0 & 0 \end{pmatrix}$$

$0 \xrightarrow{P(0 0)} 1$	$0 \xrightarrow{P(0 1)} 1$	$0 \xrightarrow{P(0 2)} 1$
$0 \xrightarrow{P(1 0)} 0$	$0 \xrightarrow{P(1 1)} 0$	$0 \xrightarrow{P(1 2)} 0$
$0 \xrightarrow{P(2 0)} 0$	$0 \xrightarrow{P(2 1)} 0$	$0 \xrightarrow{P(2 2)} 0$

run

repair

Solve

$$\bar{V} = \max_{\pi(s)} \left\{ \bar{\pi}(., \pi(s)) + \gamma P \bar{V} \right\} \rightarrow \text{Generalization}$$

↳ How to choose the new
depends on $\pi(s)$.

↳ This is a contraction map.

$V(1) = \max_{a_1} \left\{ \bar{\pi}(1, a_1) + \gamma \left[p(1|1, a_1), p(2|1, a_1) \left(\frac{V(1)}{V(2)} \right) \right] \right\}$

$V(2) = \max_{a_2} \left\{ \bar{\pi}(2, a_2) + \gamma \left[p(1|2, a_2), p(2|2, a_2) \left(\frac{V(1)}{V(2)} \right) \right] \right\}$

↳ Compiling [only 2-states, here]

$$\bar{V} = \begin{pmatrix} V(1) \\ V(2) \end{pmatrix} = \max_{\substack{(a_1, a_2) \\ (\pi_1, \pi_2)}} \left\{ \bar{\pi}(., a_1, a_2) + \gamma \begin{pmatrix} p(1|1, a_1) & p(2|1, a_1) \\ p(1|2, a_2) & p(2|2, a_2) \end{pmatrix} \begin{pmatrix} V(1) \\ V(2) \end{pmatrix} \right\}$$

$$\bar{V} = T(\bar{V}) \rightarrow \text{Transformation} \rightarrow \text{Unique solution}$$

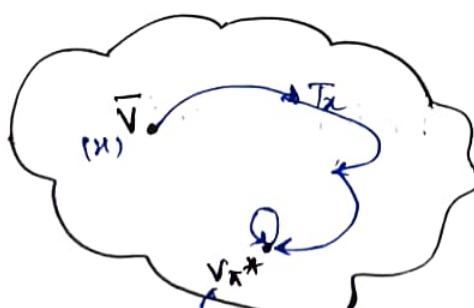
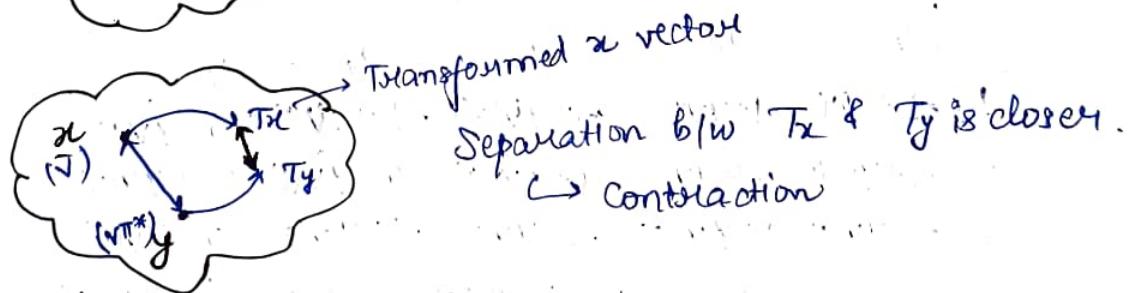
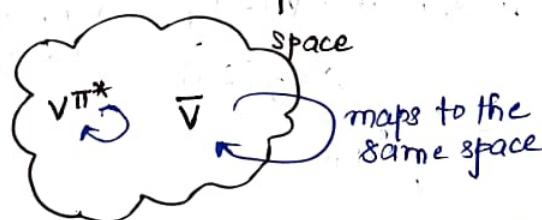
We can compare with policy evaluation ↳ But not unique optimum policy;
just at least one deterministic optimum policy.

$$\bar{V}\pi = \bar{\pi}\pi + \gamma P\bar{V}\pi$$

If we get a solⁿ to this, it is called a fixed point.

T → function / mapping

↳ contraction map



$T \rightarrow$ moves closer to V^{π^*} .

↳ doesn't move on transformation

Procedure for solving: Value Iteration

Start with an arbitrary \bar{V}_0 .

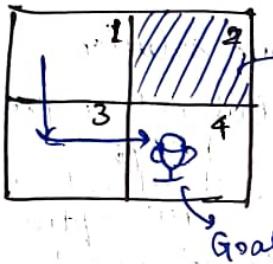
Put \bar{V}_0 in the map to get \bar{V}_1 and so on.

$$\bar{V}_{k+1} = \max_{\pi(s)} \{ \bar{V}_k(s) + \gamma P \cdot \bar{V}_k \}$$

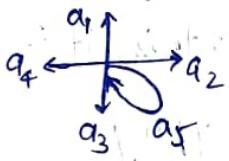
$$\lim_{K \rightarrow \infty} \bar{V}_K = V^*$$

[$\bar{V}_{K+1} \rightarrow$ Not values of some policy]
 ↳ Just iterates

Example:



→ forbidden
 go outside = -1 (bounce back)
 forbidden = -1.
 goal = +1



$$V(s) = \max_{\{a_1, \dots, a_5\}} \{ \bar{\pi}(s, a) + \gamma \bar{V}(s') \}$$

→ no distribution
 ↓
 $s' = f(s, a)$ ↳ deterministic
 ↓
 $q(s|a)$

	a_1	a_2	a_3	a_4	a_5
1	-1 + $\gamma V(0)$	-1 + $\gamma V(2)$	0 + $\gamma V(3)$	-1 + $\gamma V(1)$	0 + $\gamma V(4)$
2	-1 + $\gamma V(2)$	-1 + $\gamma V(2)$	1 + $\gamma V(4)$	0 + $\gamma V(1)$	-1 + $\gamma V(2)$
3	0 + $\gamma V(1)$	1 + $\gamma V(4)$	-1 + $\gamma V(3)$	-1 + $\gamma V(3)$	0 + $\gamma V(3)$
4	-1 + $\gamma V(2)$	-1 + $\gamma V(4)$	-1 + $\gamma V(4)$	0 + $\gamma V(3)$	1 + $\gamma V(4)$

Start with $\bar{V}_0 = (0, 0, 0, 0)^T$

Take row-wise maxm after putting all $V(i) = 0$.

$$\therefore \bar{V}_1 = (0, 1, 1, 1)^T$$

Rewrite the table.

$$\begin{bmatrix} -1 & -1+\gamma & \gamma & -1 & 0 \\ -1+\gamma & -1+\gamma & 1+\gamma & 0 & -1+\gamma \\ 0 & 1+\gamma & -1+\gamma & -1+\gamma & \gamma \\ -1+\gamma & -1+\gamma & -1+\gamma & -1+\gamma & 1+\gamma \end{bmatrix}$$

$$\bar{v}_2 = (\gamma, 1+\gamma, 1+\gamma, 1+\gamma)^T$$

Repeat.

Value Iteration Algorithm:

MDP = (S, A, P, R, γ) , $\epsilon \downarrow$ tolerance
Initialize $\bar{v}_0 (=0)$. \downarrow i/p

$\bar{v}_{\text{new}} = S \in$ (eg.)

while $\|\bar{v}_{\text{new}} - \bar{v}_0\| > \epsilon$

$v_0 = v_{\text{new}}$

for s in S .

$$v_{\text{new}}(s) = \max_a \{ \bar{r}(s, a) + \gamma \sum p(s'|s, a) v_0(s') \}.$$

Lec-10 02-09-2025

Obtaining Optimal Policies

Recall:

If we find a soln to

$$V(s) = \max_a \{ \bar{r}(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s') \}$$

$$\text{then } V^{\pi^*}(s) = V(s)$$

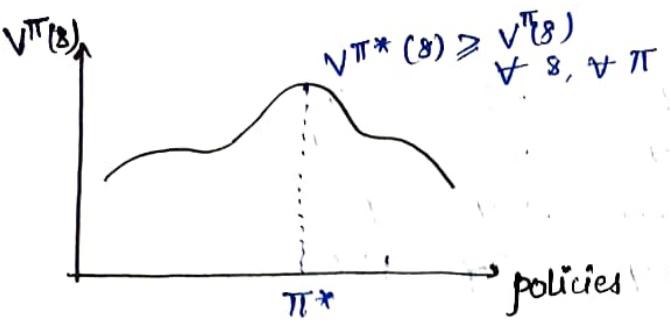
and $\pi^*(s) = \arg\max_a$ for every s .

Value Iteration Procedure: (VI) \rightarrow like $q_K(s, a)$

$$V_{K+1}(s) = \max_a \{ \bar{r}(s, a) + \gamma \sum_{s'} p(s'|s, a) V_K(s') \},$$

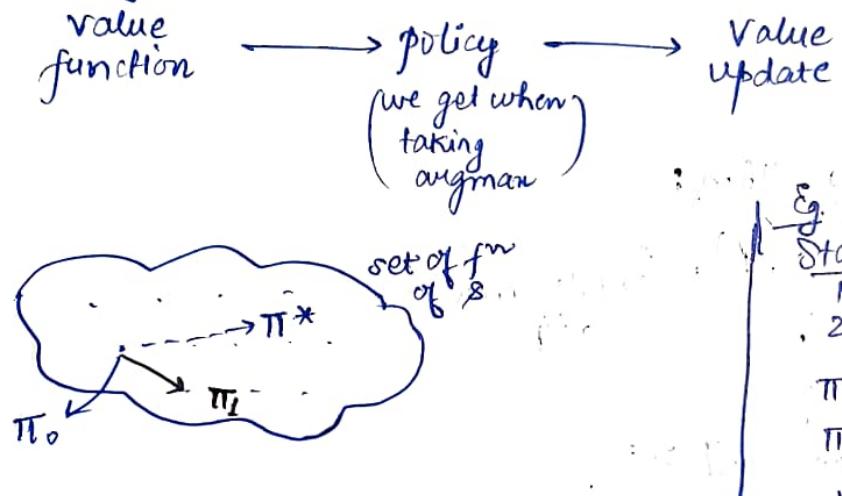
$$\text{with } V_0(s) = 0$$

Terminate if $\|V_{K+1} - V_K\| < \epsilon$.



Lecture notes:
Siva Theja
Maguluru
Georgia Tech

At every step,



Eg.
State. action
1. a_1 ,
2. a_2 ,
 a_3
 $\pi(1) \in \{a_1, a_2, a_3\}$
 $\pi(2) \in \{a_1, a_2, a_3\}$
↳ Total 9

Policy Iteration:

Start with π_0 , $V^{\pi_0}(s)$. ← policy evaluation ①

$$\text{Find } q^{\pi_0}(s, a) = \bar{\pi}(s, a) + \gamma \sum p(s' | s, a) V^{\pi_0}(s')$$

Suppose $\exists a$ such that $q^{\pi_0}(s, a) > V^{\pi_0}(s)$.

Get a time-independent policy π_1

$$\pi_1(s) = \underset{a}{\operatorname{argmax}} q^{\pi_0}(s, a)$$

↳ Policy Improvement Theorem (Ref. text)
②

Policy Iteration Algorithm:

Initialize π_0 .

do {

evaluate π_k to get $V^{\pi_k}(s)$ and $q^{\pi_k}(s, a)$ using policy evaluation.

$$\pi_{k+1}(s) = \underset{a}{\operatorname{argmax}} q^{\pi_k}(s, a).$$

} while $\| V^{\pi_{k+1}} - V^{\pi_k} \| \geq \epsilon$.

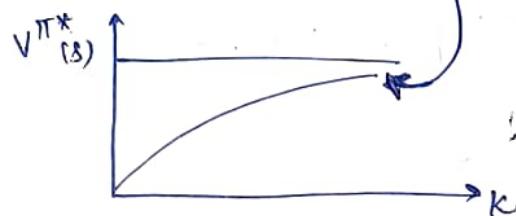
Why should policy iteration converge?

$$\pi_k \rightarrow \pi_{k+1}$$

$$V^{\pi_{k+1}}(s) \geq V^{\pi_k}(s)$$

$$\pi_t(s, a) \leq \pi_{\max}$$

$$V^{\pi^*}(s) \leq \frac{\pi_{\max}}{1-\gamma}$$



→ Increasing sequence

{ diverges X
oscillated X }

→ converges ✓

Greedy policy

(Greedy myopic policy)

↳ Just look at the current reward

Multi-armed Bandits

$$\left[\begin{array}{c} R_0 + \gamma R_1 + \gamma^2 R_2 + \dots \\ \leq \pi_{\max} \leq \pi_{\max} \leq \pi_{\max} \end{array} \right]$$

Policy iteration → faster than value iteration

↳ Requires policy evaluation

↳ Value iteration

→ Truncated policy iteration: Do VI for policy evaluation for j steps.

↳ with just 1 step, it is value iteration

1 step → ∞ step
pure VI pure PI

Review

Markov Decision Processes (MDPs)

↳ Formulate an optimization problem

- Objective: control system using a policy π - $\max_{\pi} V\pi(s) \quad - (S, A, P, R, \gamma)$ - $V(s) = \max_a \{ r(s, a) + \gamma \sum p(s'|s, a) \cdot V(s') \}$

$$\pi^*(s) = \operatorname{argmax}_a \{ \dots \}$$

$$VI: V_{k+1}(s) = \max_a \{ r(s, a) + \sum p(s'|s, a) V_k(s') \}$$

As $k \rightarrow \infty$, $V_k(s) \rightarrow V^{\pi^*}(s)$ - $\pi_k \rightarrow V^{\pi_k^*}(s) \rightarrow q^{\pi_k}(s, a)$

$$\pi_{k+1} = \operatorname{argmax}_a q^{\pi_k}(s, a)$$

As $k \rightarrow \infty$

- Truncated PI.

→ Structural Properties of π^* : Book: Bertsekas

(Dynamic Programming and Optimal Control)

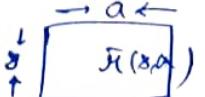
Transition probability:

$$p(s'|s, a) \rightarrow s' \left(\begin{array}{c} \xleftarrow{s} \\ \xrightarrow{a} \\ p(s'|s, a) \end{array} \right)$$

$$s = (\text{e}, \#, !)$$

↓ map to no.

$$s = (0, 1, 2)$$

We need $\sum_{\mu} p(\mu|s, a) \cdot \mu = \bar{\mu}(s, a)$: Average reward

- # Stochastic matrix \rightarrow Row sum to 1.
- # Doubly stochastic matrices \rightarrow Row & columns sum to 1.
- # $\sum p(s'|s,a) \cdot R(s',s,a) = \bar{R}(s,a) \rightarrow$ Reward depending on the next state
↳ More general.
- # Span of matrix = max - min
- # Average Reward: $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} R_t \rightarrow$ Relative Value Iteration method

Lec-12 09-09-2025

MDP Toolbox

↳ Python variant : pymdptoolbox

↳ Eg. Forest management
clinical decision making (Puterman)

MDP Examples:

Ref.: MDPToolbox (Farest)

Puterman (for other examples)

Steps:  finite discrete

- ① MDP = (S, A, P, R, γ).
 - ② Solve using VI, PI (or variants).

Example-1 :

Forest management

Trees \in forest

Every year, decide cut or wait

Age of the forest: 0, 1, 2, ..., A

If we wait, a fire will destroy the forest with prob. P.

$$S = \{0, 1, \dots, A\}$$

$$A = \{c, w\} \rightarrow \text{cut}, \text{wait}$$

$$P^C : \quad A^{+1}$$

$$A+I \begin{pmatrix} 1 & & \\ & 1 & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

$$P^W: \uparrow \quad \leftarrow A+1 \quad \longrightarrow$$

$$A+1 \downarrow$$

α $\begin{pmatrix} 0 & 1 & \dots & a+1 \dots A \\ p & 1-p & 0 & \cdots & 0 \\ p & 0 & 1-p & 0 & \vdots \\ p & 0 & 0 & 1-p & \leftarrow 0 \rightarrow \\ \vdots & \vdots & \vdots & 0 & 1-p \\ p & 0 & \vdots & 0 & 1-p \end{pmatrix}$

Reward:

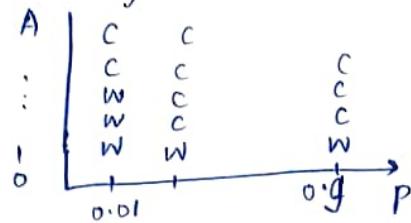
$$H(8, a) \subset W$$

$$A = \begin{pmatrix} 0 & 0 & 0 \\ \vdots & \vdots & \vdots \\ A & 2 & 4 \\ \downarrow & \downarrow & \downarrow \\ \text{Choose} \end{pmatrix}$$

```
# example · forest ()  
from mdptoolbox
```

Assign: To solve this problem.

How optimal policies behave as a function of p and γ ?



Example-2:

Decision making for organ transplants (Liver):

health = { $1 \xrightarrow{\text{best health}} H \xleftarrow{\text{lowest health}} \Delta \xrightarrow{\text{patient is dead}}$ }

organ = { $1 \xleftarrow{\text{best}} L \xleftarrow{\text{worst health}} \Delta \xrightarrow{\text{organ}}$ }

Every day, the health changes.

Every day, the patient gets an organ with one of the health states or not get an organ with prob. ϕ .
decisions = {reject, accept}

State space: includes both patient and organ healths.

$$\{1 \dots H\} \cup \{\Delta\} \times \{1 \dots L\} \cup \{\emptyset\}$$

e.g. $s_0 = (h, \emptyset)$
 $s_0 = (h, l)$

Reject action:

health transition prob.: $P(h'|h)$

dying prob.: $\delta(h)$

prob. of getting organ: $P(l|h)$

prob. of not getting organ: $\phi(h)$

Reward:

$$r(s, a, j) \xrightarrow{\Delta(\text{died})} 0$$

If you accept: $R(h, l)$ depends on patient & organ healths

Uncertainty
Equivalence

Reinforcement Learning:

Recall

$$\text{MDP} = (S, A, P, R, \gamma)$$

Obtain a policy π to

$$\max_{\pi} V^{\pi}(s), \forall s \in S$$

Methods: VI, PT and variants \rightarrow Model Based \sim Known μ and P

$$\max_{\pi} \left\{ \mu(s,a) + \gamma \sum_{s'} p(s'|s,a) V(s') \right\}$$

Bellman operator acts on $V(s)$.

Model-free algorithm

Objective of finding

$$\pi^* = \operatorname{argmax}_{\pi} V^{\pi}(s)$$

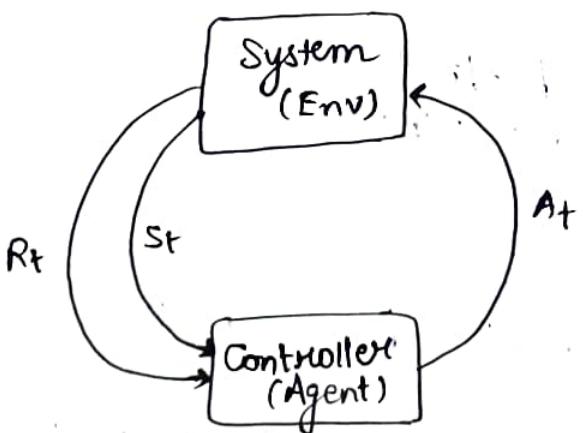
without knowledge of
 $\mu(s,a)$ and $p(s'|s,a)$
 $\forall s, a$ $\forall s', a$

There is still a system
 (no change in system)
 (i.e., $V^{\pi}(s)$) but
 controller does not
 know it.

Still, S and A are assumed to be known.

π_0 : initial policy

$$\operatorname{argmax}_a q^{\pi_0}(s, a) \rightarrow \pi_1$$



Generate $(S_0, A_0, R_0, S_1, A_1, R_1, \dots)$

Recall

- MDP Model : (S, A, P, R, γ)

Bellman Operation:

$$\text{argmax}_{\alpha} \left\{ q(s, a) + \gamma \sum p(s'|s, a) V(s') \right\}$$

Model-based methods: VI, PI

Use Bellman operation

Uses knowledge of π , $p(s'|s, a)$.

RL \Rightarrow MDP model

(Model-free) but controller does not know $[s_t, p]$
Still wants to max $V^\pi(s)$.

π_0
 \downarrow policy evaluation \Rightarrow How to do this
model-free?

$$q^{\pi_0}(s, a)$$

\downarrow policy improvement

$$\pi_1(s) = \text{argmax}_a q^{\pi_0}(s, a)$$

$$V^{\pi_1}(s) \geq V^{\pi_0}(s)$$

$$q^{\pi}(s, a) = \mathbb{E} \left\{ \begin{array}{l} R(s, a) + \pi(s_1) \\ \downarrow \\ \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \gamma^3 R(s_3, a_3) + \dots \end{array} \right\}$$

Using π

$$\{ \cdot \} \downarrow G_0 \rightarrow \text{Random return}$$

(total discounted reward)

$$= \mathbb{E}[G_0 | s_0 = s, a_0 = a]$$

In general,

$$x, \mathbb{E}X = \int f_X(x) x dx$$

How to estimate $\mathbb{E}X$ from samples of the R.V. X .

Eg. $X = \{1, 2, 3\}$

$$P_X = \{0.2, 0.4, 0.4\}$$

x_1, x_2, \dots, x_N

$$\frac{1}{N} \sum_{i=1}^N x_i \approx \mathbb{E}X$$

↓
R.V. ↓
constant if $N \rightarrow \infty$, otherwise R.V.

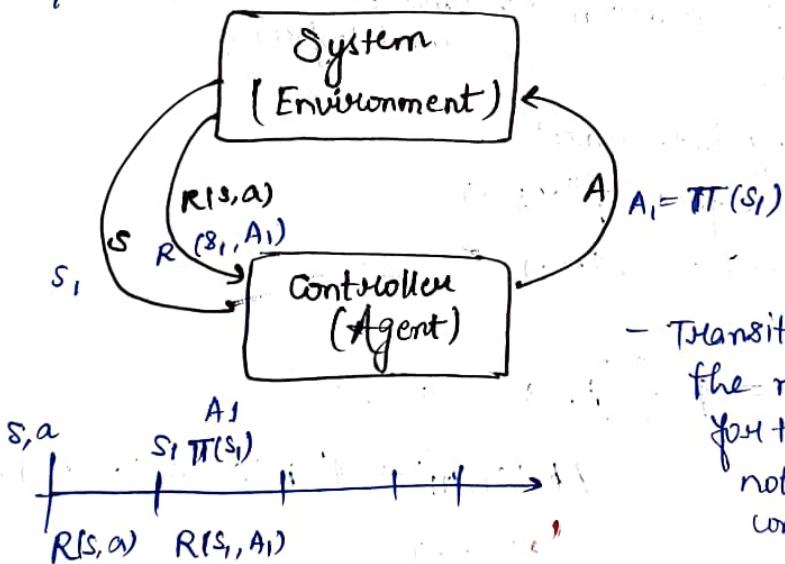
Monte Carlo (MC) Estimate

$$\lim_{N \rightarrow \infty} \sum_{i=1}^N x_i = \mathbb{E}X : \text{Law of Large Numbers}$$

How can we estimate $q^\pi(s, a)$

$$= \mathbb{E}[G_0 \mid s_0 = s; a_0 = a]$$

Samples of the return G_0 :



- Transition happens a/c to the model that is assumed for the system, which is not known to the controller.

$s, a, R_0, s_1, a_1, R_1, s_2, a_2, R_2, \dots$: Trajectory

↳ Truncate : Episode

Samples of G_0 :

$$G_0^1 = R_0' + \gamma R_1' + \gamma R_2' + \dots \quad | \text{ episode}$$

$$G_0^2 = R_0^2 + \gamma R_1^2 + \gamma R_2^2 + \dots \quad |$$

N episodes:

$$q^\pi(s, a) \approx \frac{1}{N} \sum_{i=1}^N G_0^i$$

- Interact with the system for total $|S| \times |A| \times N$ episodes.

- Estimate the $q^\pi(s, a)$ for every s and a .

It should be possible to start the system in any s :

Exploring starts

↪ First Visit Monte Carlo Algorithm : Starting in state 's'.
for policy evaluation

↪ Sample - inefficient.

To measure the quality of estimate :

• Variance
• Bias → How far the mean of estimate is from the actual value.

(Get zero bias and low variance.)

$$R_2 + \gamma R_3 + \gamma^2 R_4 + \dots$$

$$s, a, s_1, a_1, \boxed{s_2, a_2, \dots}$$

$$q^\pi(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a]$$

Every Visit Monte Carlo : Every time we visit a the state 's' and take action 'a', the expectation we get is same if we had started with state 's' and action 'a'.

Once we have estimated $q^{\pi_k}(s, a)$

$$\pi_{k+1} = \arg \max q^{\pi_k}(s, a).$$

Review

Model-free methods:

$$\text{MDP} = (S, A, P, R, \gamma)$$

$$\pi^* = \text{argmax } v^\pi(s)$$

- Motivated by PI

$$\pi_0 \rightarrow q^{\pi_0}(s, a)$$

$$\pi_1 \leftarrow \text{improve } \pi_1(s) = \text{argmax}_a q^{\pi_0}(s, a)$$

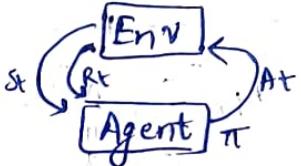
- Last class

Estimate $q^\pi(s, a)$ - form a sample of the return

- what is return?

$$R_0, R_1, R_2$$

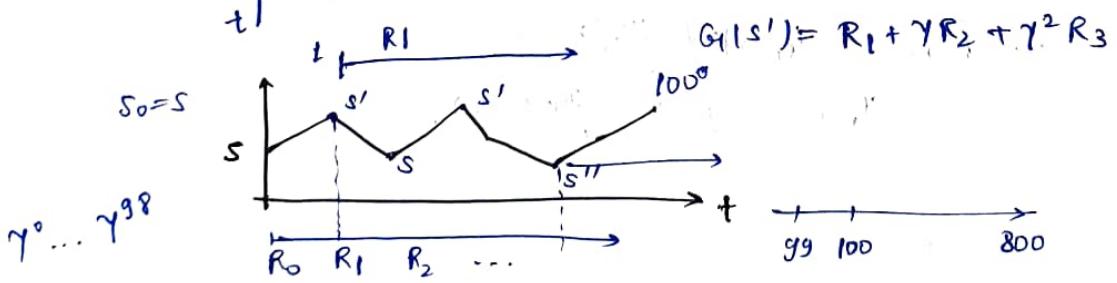
$$G_0 = R_0 + \gamma R_1 + \gamma^2 R_2$$



- Multiple episodes

$$G_0^{(1)}, G_0^{(2)}, \dots$$

$$q^\pi(s, a) = \mathbb{E} G_0 \underset{\substack{\text{Monte} \\ \text{Carlo}}}{\approx} \frac{1}{N} \sum_{i=1}^N G_0^{(i)} \rightarrow \text{First visit Estimate}$$



$$G_0 = R_0 + \gamma R_1 + \gamma^2 R_2 + \dots$$

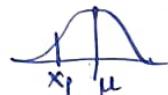
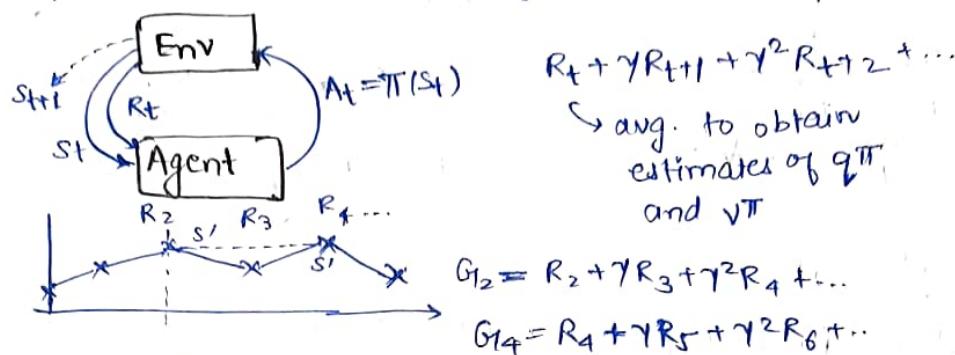
$$\bar{x}_N = \frac{1}{N} \left(\sum_{i=1}^N x_i \right)$$

$$\bar{x}_{N+1} = \frac{1}{N+1} \sum_{i=1}^{N+1} x_i$$

$$x_{N+1}, \frac{\bar{x}_N N + x_{N+1}}{N+1}$$

FVMC for policy evaluation - implementation

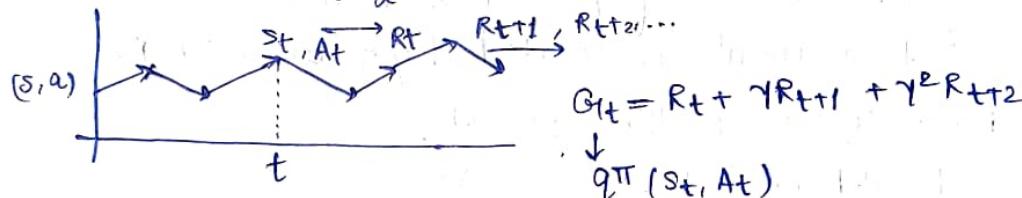
$q^{\pi}(s, a)$ or $v^{\pi}(s) \rightarrow$ agent does not know (P, R)



$$x_1, x_2, \dots : \frac{1}{N}, \sum x_i \rightarrow \mu, \frac{\sigma}{\sqrt{N}}$$

$$x_1, x_2, \dots : \sum x_i \rightarrow x_1$$

$$q^{\pi}(s, a) \Rightarrow \underset{a}{\operatorname{argmax}} q^{\pi}(s, a) = \pi_{\text{new}}(s)$$



$$\begin{aligned} \pi ? \\ \epsilon' &\xrightarrow{\frac{1}{|A|}} \frac{1}{|A|} \\ &\xrightarrow{\frac{1}{|A|}} \frac{1}{|A|} \\ &\xrightarrow{\frac{1}{|A|}} \frac{1}{|A|} \\ 1-\epsilon' &\xrightarrow{\frac{1}{|A|}} \frac{1}{|A|} \\ &\xrightarrow{\frac{1}{|A|}} \frac{1}{|A|} \\ &\xrightarrow{\frac{1}{|A|}} \operatorname{argmax} q^{\pi}(s, a) \end{aligned}$$

$$\pi(a|s) = \epsilon' / |A|$$

$\forall a \neq \operatorname{argmax}$

$a = \operatorname{argmax}$

$$\pi(a|s) = \epsilon' / |A| + (1-\epsilon')$$

How to estimate $q^{\pi}(s, a) \forall s, \forall a$ in the FVMC code (H/W)

$$\pi = [0 \ 0 \ 0 \ 1 \ 0 \ 1] \quad q^{\pi}(s, \neq 0)$$

$$s_t = (0, 1, 1, 2, \dots, 0)$$

$$a_t = (0, 0, 0, 0, \dots, 0)$$

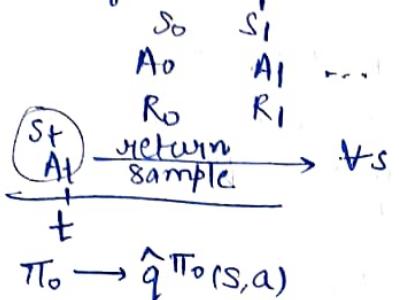
$$\pi_0 \leftarrow q^{\pi_0}(s, a)$$

$$\pi_1(s) = \operatorname{argmax} q^{\pi_0}(s, a)$$

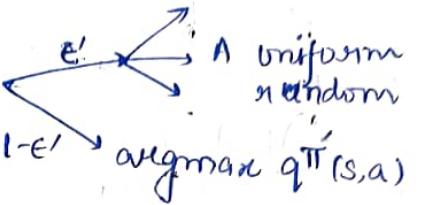
$$q^{\pi_1}(s, \neq \pi_1(s))$$

Recall:EVMC based control

In the first step,



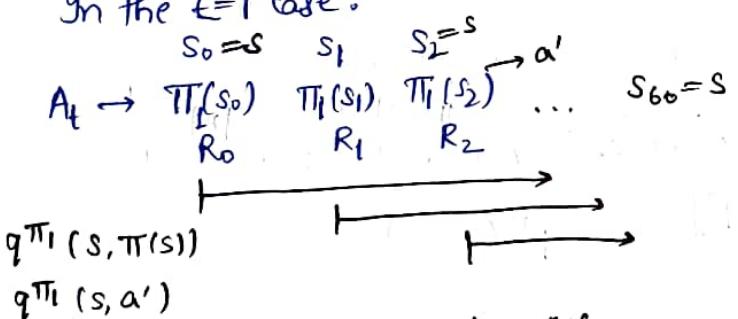
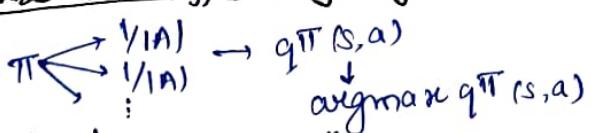
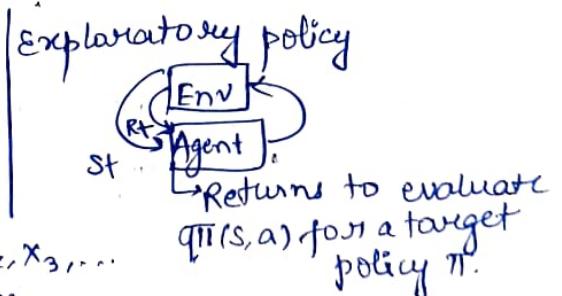
At from a randomised policy



What will be the policy in the 2nd episode if randomization is not there?

$$\pi_1(s) = \operatorname{argmax} \hat{q}^{\pi_0}(s, a) \quad \text{why is this a good idea?}$$

$$\pi_1(s) = \operatorname{argmax} q^{\pi_0}(s, a) \quad [\text{PI variant}]$$

Work with estimate of $q^{\pi_0}(s, a)$ and will get policy improvement (motivated by truncated PI)We have $\pi_1(s)$ - is this optimal?How should we improve $\pi_1(s)$?Get $q^{\pi_1}(s, a)$ or its estimate.If no randomization, then $\pi_1(s)$ is deterministic.In the $\epsilon=1$ case:On-policy and off-policy algorithmsImportance sampling:Monte-Carlo estimate: $X \rightarrow x_1, x_2, x_3, \dots$
 $EX = \frac{1}{N} \sum_{i=1}^N x_i$ 

Suppose we get samples of another random variable Y:

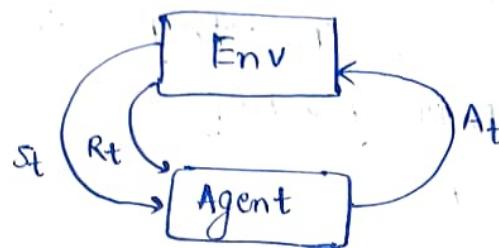
$$Y_1, Y_2, \dots \rightarrow \frac{1}{N} \sum_{i=1}^N Y_i \approx EX$$

$$\begin{aligned} EX &= \sum p_x(x) \cdot x \\ \frac{1}{N} \sum_{i=1}^N y_i &\leftarrow EY = \sum p_y(y) \cdot y = \sum p_y(y) \cdot \frac{p_x(x)}{p_x(x)} \cdot x = \sum p_x(x) \cdot \frac{p_y(y)}{p_x(x)} \cdot x = \sum p_x(x) \cdot f(x) = Ef(x) \end{aligned}$$

Temporal Difference Method

chap-7

Recall



s_0, s_1, s_2, \dots

r_0, r_1, r_2, \dots

a_0, a_1, a_2, \dots

Value function for a policy π :

$$A_t = \pi(s_t) \text{ or } \sim \pi(a|s_t)$$

$$\begin{aligned} V^\pi(s) &= E(G_{t_0}) \\ &\approx \frac{1}{N} \sum_{i=1}^N G_{t_0}^i \end{aligned} \quad \left. \right\} \text{MC}$$

- Extended to $q^\pi(s, a)$.

Another way of estimating $V^\pi(s)$:

$$\begin{aligned} V^\pi(s) &= h(s, a) + \gamma \sum p(s'|s, a) \cdot V^\pi(s') \\ &= E \left[R_0 + \gamma V^\pi(s_1) \right] \end{aligned}$$

\nwarrow samples

$$R_0 + \gamma V^\pi(s_1)$$

is like a sample of G_{t_0} .

$$R_t + \gamma V^\pi(s_{t+1})$$

is like a sample of G_{t+1} .

So, in TD, we use

$$R_t + \gamma \hat{V}^\pi(s_{t+1}).$$

$$\hat{V}_{k+1}^{\pi}(s_t) = \hat{V}_k^{\pi}(s_t) + \alpha_k(R_t + \gamma \hat{V}_k^{\pi}(s_{t+1}) - \hat{V}_k(s_t))$$

can be time-dependent

$$\hat{V}_{t+1}^{\pi}(s_t) = \hat{V}_t^{\pi}(s_t) + \alpha(R_t + \gamma \hat{V}_t^{\pi}(s_{t+1}) - \hat{V}_t(s_t))$$

... "Bootstrap Estimate"

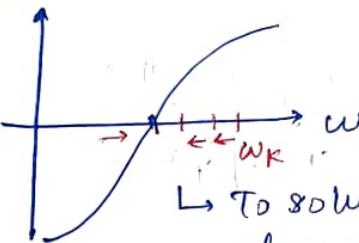
$R_t + \gamma \hat{V}_t^{\pi}(s_{t+1}) - \hat{V}_t(s_t)$: TD (Temporal Difference) error

$R_t + \gamma \hat{V}_t^{\pi}(s_{t+1})$: TD Target

We are trying to get $\hat{V}_t^{\pi} \rightarrow V^{\pi}$.

Digression (chap-6)

$$g(w) = 0$$

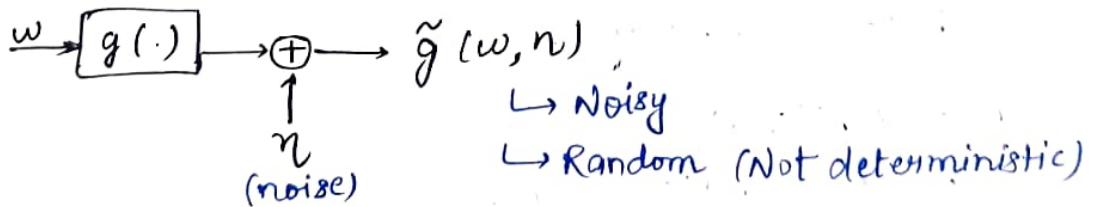


↳ To solve $g(w) = 0$, take two values, check sign and reduce the interval.

$$\left[\begin{array}{l} \text{To solve } \frac{d J(w)}{d w} = 0 \text{ : Optimization problem} \\ w_{k+1} = w_k - \alpha_k \frac{d J(w_k)}{d w} \end{array} \right]$$

$$w_{k+1} = w_k - \alpha_k g(w_k)$$

Take small α_k : $0 < \alpha_k < 1$



SGD: Stochastic Gradient Descent

$$w_{k+1} = w_k - \alpha \tilde{g}(w_k, n_k) : \text{Robbin-Murow Algorithm}$$

↳ Used to solve $g(w) = 0$ when the observation is noisy.

$$V^\pi(s) = \mathbb{E}(s, \pi(s)) + \gamma \sum p(s'|s, a) \cdot V^\pi(s')$$

$$\Rightarrow \vec{v} = \mathcal{H}(s, \pi(s)) + \gamma \sum p(s'|s, a) \cdot V^\pi(s') \perp V^\pi(s)$$

$$\Rightarrow \vec{v} = \vec{\pi} + M\vec{V}^\pi \quad [\text{In Matrix form}]$$

$$\vec{v} = \boxed{g(\vec{v}^\pi)} \rightarrow \begin{array}{l} \text{we don't know this} \\ \hookrightarrow \begin{bmatrix} \text{Random reward} \\ \text{Random/Sampled } M \end{bmatrix} \end{array}$$

RM Algorithm:

$$\bar{V}_{k+1} = \bar{V}_k - \alpha_k \tilde{g}(\bar{v}, \eta_k) \xrightarrow{\text{similar to in TD}} (R_t + \gamma \hat{V}_k(s_{t+1}) - \hat{y}_k(s_t))$$

Conditions: $\sum \alpha_k = \infty$
 $\sum \alpha_k^2 < \infty$

How can we estimate $q^\pi(s, a)$?

- How to do this for Monte-Carlo or FVMCP?

s_0, s_1, \dots
 R_0, R_1, \dots $\xrightarrow{\text{first reward a/c}}$ first reward a/c
 $a \quad (R_0 a + \gamma R_1 + \dots)$
choose all actions for 1st stage (not a/c to policy)
 g_0 - return is still used.

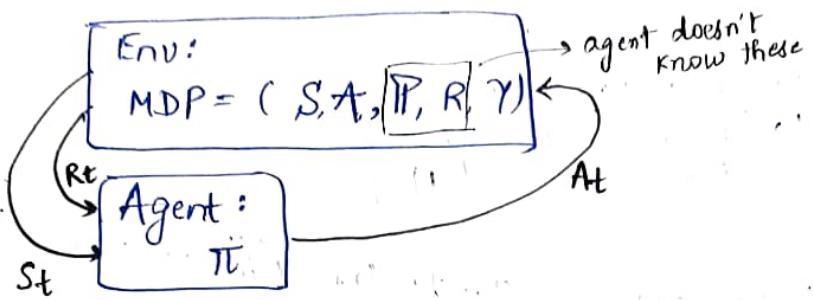
$$s \rightarrow a \quad R_0 + \gamma \hat{q}^\pi(s_{t+1}, \underbrace{A_{t+1}}_{\text{a/c to policy}})$$

$$\hat{q}^\pi(s_{t+1}, A_{t+1}) = \gamma \hat{V}^\pi(s_{t+1})$$

$$\hat{V}^\pi(s) = \sum_a \frac{1}{|A|} \hat{q}(s, a)$$

Temporal difference Algorithm

Recall:



$$v^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) v^\pi(s')$$

$$v^*(s) = \max_a \{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v^*(s') \}$$

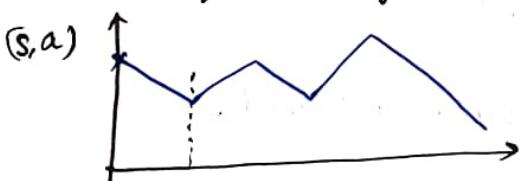
value iteration

Policy improvement and iteration

$$q^\pi(s, a) \rightarrow \arg\max_a q^\pi(s, a)$$

(if we know the model (MDP))

Estimate $q^\pi(s, a)$ from samples.



$$G_0 = R_0 + \gamma R_1 + \gamma^2 R_2 + \dots$$

$$= R_0 + \gamma \underbrace{(R_1 + \gamma R_2 + \dots)}_{\text{Replace by some other estimate}}$$

↓
 Replace by some other estimate
 (current estimate), $\hat{v}(s_1)$

ID Algorithm for policy evaluation

$$\pi, \quad v^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid s_0 = s \right] \quad [\text{Here, } A_t = \pi(s_t)]$$

↓
 Known estimate for every state

$$= \mathbb{E}[G_0 \mid s_0 = s] \approx \frac{1}{N} \sum_i G_{t_0}^{(i)}$$

$$\hat{v}_N(s) = \frac{1}{N} \sum_i G_{t_0}^{(i)}$$

$$\hat{v}_{N+1}(s) = \hat{v}_N(s) + \alpha_N (G_{t_0}^{(N+1)} - \hat{v}_N(s))$$

$$R_0 + \gamma \hat{v}_N(s_1) \quad [\text{Replace/Approx. } G_{t_0} \text{ by}]$$

an estimate

$$\hat{V}_{N+1}(s) = \hat{V}_N(s) - \alpha_N (\hat{V}_N(s) - (R_0 + \gamma \hat{V}_N(s_1))) \dots \text{TD update equation}$$

$\alpha_N = \frac{1}{N+1}$ or something else.

$$\sum \alpha_N = \infty$$

$$\sum \alpha_N^2 < \infty$$

For any state,

$$v^\pi(s) \approx \frac{1}{N} \sum_i G_{t,i}$$

$$G_{t,i} = R_t + \gamma (R_{t+1} + \gamma R_{t+2} + \dots)$$

$$\therefore \hat{V}_{N+1}(s_t) = \hat{V}_N(s_t) - \alpha_N (\hat{V}_N(s_t) - (R_0 + \gamma \hat{V}_N(s_{t+1})))$$

Disadvantage: We only know random R_0 , for G_0 , we estimate $\gamma(R_1 + \gamma R_2 + \dots)$ which may not be correct.

↳ Biased.

$\pi \rightarrow v^\pi$ or its estimate \hat{v}^π

For controls, we need Q^π .

TD form $Q^\pi(s, a)$:

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid s_0 = s, A_0 = a \right] \\ &\quad \forall t > 0, A_t = \pi(s_t) \\ &= \mathbb{E} [G_{t,0} \mid s_0 = s, A_0 = a] \end{aligned}$$

$$Q^\pi(s, a) = h(s, a) + \gamma \sum p(s'|s, a) v^\pi(s')$$

$$(s, a) \longrightarrow \begin{array}{c} \textcircled{R_0^{(1)}} \\ \textcircled{R_0^{(2)}} \\ \textcircled{R_0^{(3)}} \end{array} + \begin{array}{c} \textcircled{\hat{V}(s_1^{(1)})} \\ \textcircled{\hat{V}(s_1^{(2)})} \\ \textcircled{\hat{V}(s_1^{(3)})} \end{array} \xrightarrow{\text{ss}} \begin{array}{c} \downarrow \\ h(s, a) \\ + \gamma \sum p(s'|s, a) v^\pi(s') \end{array} \rightarrow Q^\pi(s, a)$$

[we have to look every episode]

Estimate $Q^\pi(s; a)$ [something similar to \hat{v}_N]:

$$Q^\pi(s, a) = \mathbb{E} [G_0 \mid s_0 = s, A_0 = a]$$

$$\hat{Q}_N(s, a) \approx \frac{1}{N} \sum_i G_0^{(i)}$$

$$\hat{Q}_{N+1}(s, a) = \hat{Q}_N(s, a) + \alpha_N (G_0^{(N+1)} - \hat{Q}_N(s, a))$$

$$G_0^{(N+1)} \approx R_0 + \gamma \hat{Q}_N(s_1, A_1)$$

$\mathbb{E}(R_1 + \gamma R_2 + \gamma R_3 + \dots)$ gives this some Q value?

$$\downarrow \\ Q^\pi(s_1, A_1)$$

we do not know this!

$$\hat{Q}_N(s_1, A_1)$$

$$\therefore \hat{Q}_{N+1}(s, a) = \hat{Q}_N(s, a) + \alpha_N (\hat{Q}_N(s, a) - (R_0 + \gamma \hat{Q}_N(s_1, A_1)))$$

$$\hat{Q}_{N+1}(s_t, a_t) = \hat{Q}_N(s_t, a_t) + \alpha_N (\hat{Q}_N(s_t, a_t) - (R_t + \gamma \hat{Q}_N(s_{t+1}, a_{t+1})))$$

- s, a are the state and action at time 't'.

$$[s = s_t, a = A_t]$$

[N is taken as t .]

... Q-function TD update equation

$\underset{t}{\overset{T}{\sum}} \underset{A}{\overset{A}{\sum}} R \underset{t+t+1}{\overset{t+1}{\sum}} \underset{A}{\overset{A}{\sum}}$

Temporal Difference (TD) algorithms for RL:

$$\pi_0 \rightarrow q^{\pi_0} (s, a) \xrightarrow{\text{argmax}} \pi_1 \rightarrow \dots \quad \pi^* = \underset{\pi}{\text{argmax}} V^{\pi}(s)$$

MDP Model = (S, A, P, R, V)

$$(s, a) : G_0^{\pi_0} = R_0 + \gamma R_1 + \gamma^2 R_2 + \dots$$

$$\begin{aligned} q^{\pi_0}(s, a) &= \mathbb{E}[G_0^{\pi_0} | s_0 = s, A_0 = a] \\ &\approx \underbrace{\frac{1}{N} \sum_{j=1}^N}_{\text{---}} G_0^j \end{aligned}$$

$$\begin{aligned} \hat{q}_{N+1}^*(s, a) &= \hat{q}_N(s, a) + \alpha (\underbrace{G_1^{N+1}}_{R + \gamma \hat{q}_N(s, a)} - \hat{q}_N(s, a)) \\ &\downarrow \\ &R + \gamma \hat{q}_N(s, a) \end{aligned}$$

$$\begin{aligned} \hat{q}_{N+1}^*(s, a) &= \hat{q}_N(s, a) + \alpha (G_t^{N+1} - \hat{q}_N(s, a)) \quad [N \text{ is actually } t] \\ &\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ &s_t \quad A_t \quad s_t \quad A_t \quad t \quad s_t \quad A_t \end{aligned}$$

$$\textcircled{1} \quad \hat{q}_{t+1}^*(s_t, A_t) = \hat{q}_t(s_t, A_t) + \alpha (R_t + \gamma \hat{q}_t^*(s_{t+1}, A_{t+1}) - \hat{q}_t(s_t, A_t))$$

We need $s_t, A_t, R_t, s_{t+1}, A_{t+1}$: SRSA

At t , update (s_t, A_t) .

② Others remain the same.

$$\hat{q}_{t+1}(s', a') = \hat{q}_t(s', a'),$$

$$(s', a') \neq (s_t, A_t)$$

H-w: We get a $V^{\pi}(s) \approx$ MDP soln.

SARSA Variants

- n-step SARSA:

$$R_t + \gamma \hat{q}_t(s_{t+1}, a_{t+1}) \longrightarrow R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^n R_{t+n} + \gamma^{n+1} \hat{q}(s_{t+n}, a_{t+n})$$

- Expected SARSA:

$$R_t + \gamma \hat{q}_t(s_{t+1}, a_{t+1}) \longrightarrow R_t + \gamma \mathbb{E}_{\bar{a}}[\hat{q}(s_{t+1}, \bar{a}_{t+1})]$$

↓ can be calculated ↓ we know the distribution

Q-Learning Algorithm

Recall:

$$\pi_0 \rightarrow \hat{q}_{\pi_0}(s, a) \rightarrow \text{argmax} \rightarrow \pi_1$$

$$\hat{q}_{t+1}(s_t, a_t) = \hat{q}_t(s_t, a_t) + \alpha_t \left\{ (R_t + \gamma \hat{q}_t(s_{t+1}, a_{t+1})) - \hat{q}(s_t, a_t) \right\}$$

Main idea in the code:

- Use a randomized policy
 - uniform at random $a \in A$
 - $\text{argmax}_a \hat{q}_t(s, a)$
- Estimate \hat{q}_t .
- do policy improvement!
- SRSA
 - + variants - n-step expected SRSA

We also saw

$$\hat{q}_{t+1}(s_t, a_t) = \hat{q}_t(s_t, a_t) + \alpha_t \left\{ R_t + \gamma \max_a \hat{q}_t(s_{t+1}, a) - \hat{q}_t(s_t, a_t) \right\}$$

observation:

$\hat{q} \approx \text{MDP result}$

$\text{argmax}_a \hat{q}(s, a) = \text{opt. policy}$

In Q-learning:

Directly solve the Bellman's Optimality Equation.

$$V^{\pi^*} = \max_a \{ \pi(s, a) + \gamma \sum_{s'} p(s' | s, a) V^{\pi^*}(s') \}$$

$$q^{\pi^*}(s, a) = \pi(s, a) + \gamma \sum_{s'} p(s' | s, a) \cdot \max_{a'} q^*(s', a')$$

$$V^{\pi^*}(s) = \max_a q^{\pi^*}(s, a)$$

If we know find a $q(\cdot)$ s.t.,

$$q(s, a) = \pi(s, a) + \gamma \sum_{s'} p(s' | s, a) \cdot \max_{a'} q(s', a')$$

To understand Q-learning,

let's start with the model-based procedure for q .

Start with $q_0(s, a)$. $\mathbb{E} R(s, a)$

$$q_{k+1}(s, a) = \boxed{R(s, a)} + \gamma \left[\sum_{a'} p(s'|s, a) \max_{a'} q_k(s', a') \right]$$

$$\mathbb{E} \left[\max_{a'} q_k(s', a') \right]$$

random variable
with prob.
 $p(s'|s, a)$

What we would have seen is

$$V_{k+1}(s, a) = \max_a \left\{ R(s, a) + \gamma \sum_{s'} p(s'|s, a) V_k(s') \right\}$$

$$V_0(s) = \max_a q_0(s, a)$$

$$q_1(s, a) = R(s, a) + \gamma \sum_{s'} p(s'|s, a) V_0(s')$$

$$V_1(s) = \max_a q_1(s, a)$$

$$q_{k+1}(s, a) = \mathbb{E} R(s, a) + \gamma \mathbb{E}_{s, a} \max_{a'} q_k(s, a')$$

$$\approx \frac{1}{N} \sum_{i=1}^N \left(R^{(i)}(s, a) + \gamma \max_{a'} q_k(s^{(i)}, a') \right)$$

arbitrary for initialized
with some value.

What if we used an incremental update for this?

- Incremental on 'i'.

$$q_{k+1,i}(s, a) = q_{k+1,i-1}(s, a) + \alpha_{k+1,i} (R^{(i)}(s, a) + \gamma \max_{a'} q_{k+1,i-1}(s^{(i)}, a') - q_{k+1,i-1}(s, a))$$

At every time step t (combination of k & i , $(k+1, i)$), update q , each associated with an. (s_t, A_t)

$$q_{t+1}(s_t, A_t) = q_t(s_t, A_t) + \alpha_t (R_t + \gamma \max_{a'} q_t(s_t, a) - q_t(s_t, A_t))$$

- Unlike SRSAs, we need to estimate

$$\begin{aligned} R^{(i)}(s, a) \\ q_k(s^{(i)}, a') \end{aligned} \quad \left. \begin{array}{l} \text{Neither depend on the} \\ \text{policy} \end{array} \right\}$$

Type of policies for Q-learning:

- Behaviour policy : Used to run the system.
↳ generates these $R^{(i)}$ and $s^{(i)}$ samples from the unknown R , P .

- Target policy.

Q -learning = off policy algorithm
 $S\pi\text{SA}$ = on-policy algorithm. } online learning

- Q -learning can be used on datasets, called offline RL.

Ch-7

Lec-24

27-10-2025

Value function Approx. Methods

| Ch-8 of the book

Recall -

$$MDP = (S, A, P, R, \gamma)$$

$\boxed{V^\pi(s)}$ or $\boxed{q^\pi(s, a)}$ ^{finite discrete} → tables / arrays

Other RL problems :

These assumptions do not hold.

Motivating problems:

① Machine Repair Problem

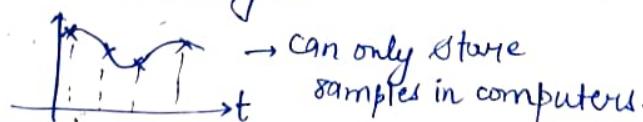
Temp. of a particular component
 $\in S \subseteq [t_{\min}, t_{\max}]$

② Lunar Lander

Real valued (x, y, z)

Orientation (θ, ϕ)

like continuous signal



Additionally, control variable

$$u \in [u_{\min}, u_{\max}]$$

↳ Continuous action space.

© Uncountable/discrete

e.g. queue problem \rightarrow [Sennott]

$$S = \{0, 1, 2, 3, \dots\}$$

state - no. of people in the queue.

Action space:

$$\text{service rate } \in [\mu_{\min}, \mu_{\max}]$$

Generalization:

$S, A \rightarrow$ continuous infinite

$P: f_{S'|S,a}(s') \rightarrow$ pdf
(PMF)

$F_{S'|S,a}(s') \rightarrow$ CDF

$R: r_{S'|S,a}(x)$

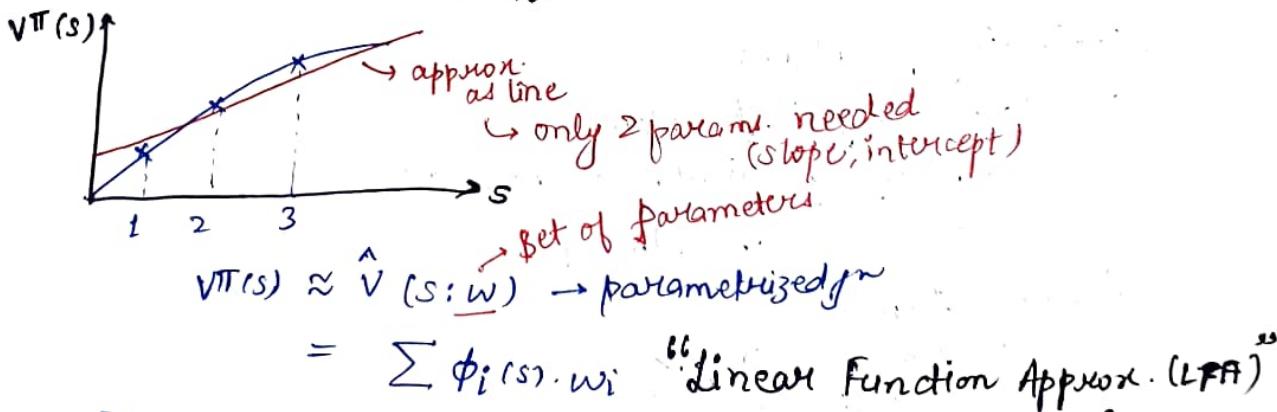
$$\int s \cdot r_{S'|S,a}(s) ds = \bar{r}(s, a) \rightarrow \text{average reward}$$

$$\pi: \pi_{als} = N(\alpha | \mu_s, \sigma_s^2)$$

| Book: Bertsekas

- For continuous problems, Bellman optimality equation may not hold. This needs to be verified.

How to obtain $v^\pi(s)$, $s \in S$ (General Space)



Suppose we approximate $v^\pi(s)$ using a 'line'. }

$$w_0 + w_1 s$$

$$\phi_0(s) = 1$$

$$\phi_1(s) = s$$

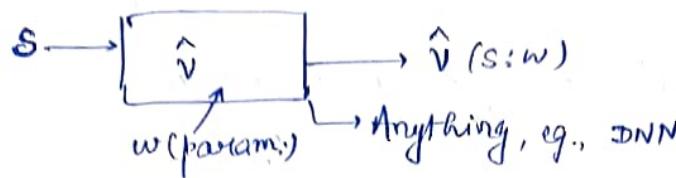
$$\text{If } w_0 + w_1 s + w_2 s^2$$

$$\phi_0(s) = 1$$

$$\phi_1(s) = s, \phi_2(s) = s^2$$

"linear in param.
(may not be linear)

[Linear Regression]



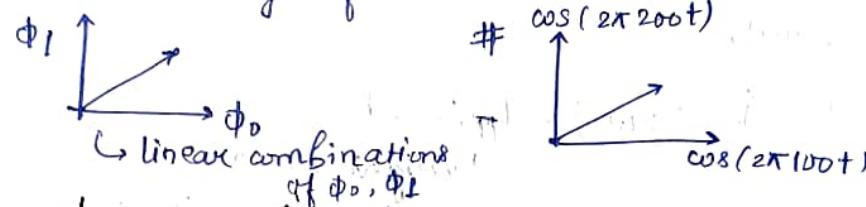
$$V^\pi(s) \approx \hat{v}(s:w)$$

$$= \phi(s)^T \cdot w \quad [\text{dot-product}]$$

$$w = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \end{pmatrix}, \quad \phi(s) = \begin{pmatrix} \phi_0(s) \\ \phi_1(s) \\ \vdots \end{pmatrix}$$

Question to which ϕ_i -s to use?

- Depends on problem-to-problem \rightarrow Guess (may not be possible)
- Use DNN to get features.



Look into simple form of the problem

Suppose ϕ_i -s are given/known — how to find \vec{w} ?

$$\min_w \int (V^\pi(s) - \hat{v}(s:w))^2 ds \rightarrow \text{total error}$$

\downarrow

fⁿ of s

Single value

V^π(s) x
 ϕ_p x
 ϕ_0

Suppose s is finite and discrete.

$$\sum_{j=1}^N (V^\pi(s_j) - \hat{v}(s_j:w))^2 \leftarrow \text{find } \vec{w} \text{ that minimizes this}$$

To get exact approx., use impulse basis functions.

$$\phi_0(s) = \delta(s_0)$$

[finite, discrete]

$$\phi_1(s) = \delta(s_1)$$

$$\phi_2(s) = \delta(s_2)$$

$$w_0 = V^\pi(s_0)$$

$$w_1 = V^\pi(s_1)$$

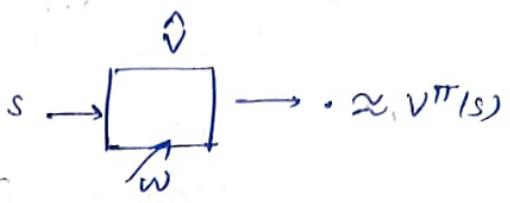
$$w_2 = V^\pi(s_2)$$

Value function approximation \rightarrow General s and A spaces.

Recall

$$V^\pi(s) \approx \hat{v}(s; w)$$

$$= \sum_{i=0}^{M-1} \phi_i(s) \cdot w_i$$



w such that $\int (V^\pi(s) - \hat{v}(s; w))^2 ds \leftarrow$ not proper!

$$\sum_{i=1}^N (V^\pi(s_i) - \hat{v}(s_i; w))^2 = J(w) : \text{objective function.}$$

Find w to $\min J(w)$.

Objective $J(w)$:

$$\int (V^\pi(s) - \hat{v}(s; w))^2 d\pi(s) ds$$

weight factor

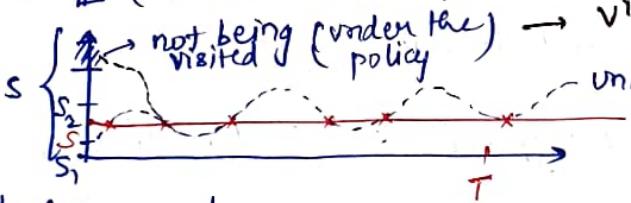
↳ decays with 's'

↳ to make the integral converge (finite)

Write it as expectation which can be estimated as average.

$d\pi(s)$ can be thought of as some pdf (of some state variable).
varying in s

$$= E(V^\pi(s) - \hat{v}(s; w))$$

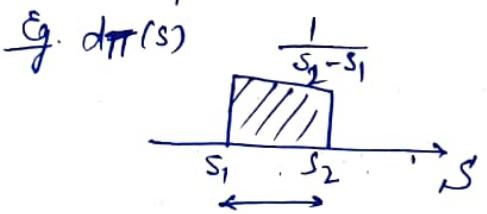


$V^\pi(s)$ is defined for the policy

under π something like
visited once

$$\hat{v}(s; w) \approx V^\pi(s)$$

↳ approx. may not be good for the one being visited once or few times.



It means the approx. state is good for (s_1, s_2) .

↳ Give high weight to the one being visited again & again.

s → being visited more-and-more.

↳ Fraction of time ↑ (for discrete spaces) \Rightarrow e.g. for $t=1000$, it is being visited 100 times

↳ being visited ↗ distribution ↑

↳ Pdf ↑

↳ frac. = $\frac{t}{10}$.

Define $d\pi(s)$ as the fraction of time a state is visited.

Suppose [Genie] $\xrightarrow{\text{gives}} v^\pi(s)$

$$\frac{1}{N} \sum_{t=1}^N (v^\pi(s_t) - \hat{v}(s_t; w))^2 = \mathbb{E} (v^\pi(s) - \hat{v}(s; w))^2$$

at time 't'

$d\pi$: dist. of the state for the policy

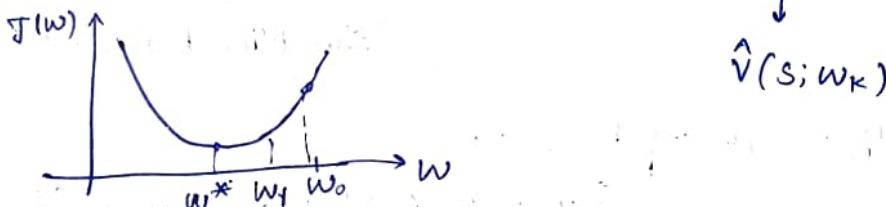
↳ stationary probability distribution. [Read!]

→ States get picked w.r.t $d\pi(s)$.

$$J(w) = \mathbb{E}_{d\pi} [v^\pi(s) - \hat{v}(s; w)]^2 \rightarrow \text{minimize}$$

using gradient descent.

Start with parameter $w_0 \rightarrow w_1 \dots w_k \xrightarrow{\downarrow} w_{k+1}$



$$w_{k+1} = w_k - \alpha_k \cdot (\nabla_w J(w_k))$$

Gradient,

$$\nabla_w J(w) = \nabla_w \mathbb{E} (v^\pi(s) - \hat{v}(s; w))^2$$

$$= \begin{pmatrix} \frac{\partial}{\partial w_0} \mathbb{E} (v^\pi(s) - \hat{v}(s; w))^2 \\ \frac{\partial}{\partial w_1} \mathbb{E} (v^\pi(s) - \hat{v}(s; w))^2 \\ \vdots \end{pmatrix}$$

$$= \mathbb{E} \frac{\partial}{\partial w_0} (v^\pi(s) - \hat{v}(s; w))^2$$

$$= -2 \mathbb{E} [(v^\pi(s) - \hat{v}(s; w)) \frac{\partial \hat{v}(s; w)}{\partial w_0}]$$

$$\text{vector} = \begin{bmatrix} -2 \mathbb{E} (\underbrace{v^\pi(s) - \hat{v}(s; w)}_{\text{scalar}}) \frac{\partial \hat{v}(s; w)}{\partial w_0} \\ -2 \mathbb{E} (\underbrace{v^\pi(s) - \hat{v}(s; w)}_{\text{scalar}}) \frac{\partial \hat{v}(s; w)}{\partial w_1} \\ \vdots \end{bmatrix} \xrightarrow{\text{vector}}$$

$$= -2 \underbrace{\mathbb{E} (v^\pi(s) - \hat{v}(s; w))}_{\text{scalar}} \cdot \underbrace{\nabla_w \hat{v}(s; w)}_{\text{vector (fn of vector w)}}$$

$$w_{k+1} = w_k + \alpha_k \cdot \mathbb{E} \left[(\nabla \pi(s) - \hat{v}(s; w)) \cdot \nabla_w \hat{v}(s; w) \right]$$

↓ ↓
unknown → cannot sum!

["2' absorbed in α_k]

Variant of G.D: Stochastic Gradient Descent (SGD).

- Problem to minimize objective $f^n : \mathbb{E}_x f(x; \theta) \rightarrow f^n$ of θ

$$\theta_{k+1} = \theta_k - \alpha_k \nabla_\theta \mathbb{E}_x f(x; \theta)$$

$\theta_{k+1} = \theta_k - \alpha_k \nabla_\theta f(x_k; \theta) \rightarrow$ expectation replaced by the value.

↳ expectation not needed,
only needs samples.

[Batch method vs. sample method]

[~~Minibatch~~ Mini-batch → take few samples & average]

Change the GD step to SGD:

$$w_{k+1} = w_k + \alpha_k (\nabla \pi(s_k) \cdot \hat{v}(s_k; w)) \nabla_w \hat{v}(s_k; w)$$

\downarrow k-th step

RecallGeneral state and action spaces, S, A Approximation $\rightarrow V^\pi(s)$

$$\hat{V}(s; w) \approx V^\pi(s)$$

$$\hat{V}(s; w) = \sum w_i \phi_i(s) \quad \text{Basis fn}$$

Objective fn:

$$J(w) = \int_S (\hat{V}(s; w) - V^\pi(s))^2 d\pi(s) ds. \quad \begin{matrix} \rightarrow \text{pdf = stationary} \\ \text{distribution} \end{matrix}$$

$$= \mathbb{E} (\hat{V}(s; w) - V^\pi(s))^2$$

Policy evaluation \rightarrow Choose w so that $J(w)$ is minimized. $\min_w J(w)$ using G.D.

$$w_{k+1} = w_k - \alpha_k \cdot \nabla_w J(w)$$

$$= w_k - \tilde{\alpha}_k \mathbb{E} [(\hat{V}(s; w_k) - V^\pi(s)) \cdot \nabla_w \hat{V}(s; w_k)]$$

\hookrightarrow G.D with expectation $(\because \tilde{\alpha} = 2\alpha_k)$

\hookrightarrow we can use stochastic G.D

$$w_{k+1} = w_k + \alpha_k (V^\pi(s_k) - \hat{V}(s_k; w_k)) \cdot \nabla_w \hat{V}(s_k; w_k)$$

\hookrightarrow for doing update, requires $V^\pi \Rightarrow$ Not implementable.

\hookrightarrow If w_k is known $\Rightarrow \hat{V}$ can be evaluated, but V^π is unknown.

Estimate V^π . Approximate from return.

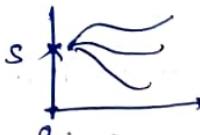
$$V^\pi(s_k) = \mathbb{E} G_t$$

$$V^\pi(s_t) = \mathbb{E} G_t$$

$$w_{k+1} = w_k + \alpha_k (\square \leftarrow \text{estimate} - \hat{V}(s; w_k)) \nabla \hat{V}(s; w_k)$$

Monte Carlo:

$$V^\pi(s) \approx \frac{1}{N} \sum_i G_{t,i}(s)$$

for every state s :

$$s_t \xrightarrow{G_t} V^\pi(s_t) \quad t \xrightarrow{R_t + \gamma \hat{V}(s_{t+1}; w_t)}$$

$$MC: w_{t+1} = w_t + \alpha_t (\alpha_t - \hat{V}(s_t; w_t)) \nabla \hat{V}(s_t; w_t)$$

\hookrightarrow implementable.

$$TD: w_{t+1} = w_t + \alpha_t (r_t + \gamma \hat{V}(s_{t+1}; w_t) - \hat{V}(s_t; w_t)).$$

$$\nabla \hat{V}(s_t; w_t)$$

OpenAI's gym.

$$\hat{V}(s, w) = \sum w_i \phi_i$$

$$\nabla_w (\hat{V}(s, w)) = \begin{pmatrix} \phi(0) \\ \phi(1) \\ \vdots \end{pmatrix}$$

SARSA:



Lec-27

04-11-2025

Recall

$$V^\pi(s) \rightarrow \hat{V}(s; w)$$

$$q^\pi(s, a) \rightarrow \hat{q}(s, a; w)$$

linear fn approx.

$$\hat{V}(s) = \sum w_i \phi_i(s)$$

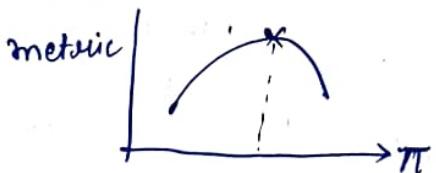
SGD update rules:

$$w_{t+1} = w_t + \alpha_t (TD/MC - \hat{V}(s_t; w_t)) \nabla_w \hat{V}(s_t; w_t)$$

$$= w_t + \alpha_t (TD/MC - \hat{q}(s_t, A_t, w_t)) \nabla_w \hat{q}(s_t, A_t, w_t)$$

$\hat{q}^\pi \approx \hat{q} \rightarrow$ policy improvement.

Policy gradient: (θ, q)

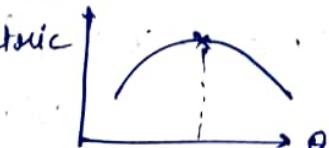


How do we implement policies?

$$\pi(s) = \sum \theta_i \cdot \psi_i(s)$$

parameters different basis fn

Optimize some metric with θ .



- Not always deterministic; we also do exploration.
- Represent the policy as distribution.

$$\pi(a|s) = \mathcal{N}(a : \mu(s, \theta), \sigma^2(s, \theta)) \rightarrow \text{Randomized policy}$$

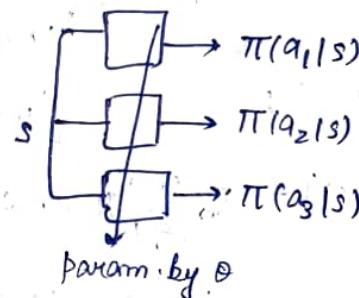
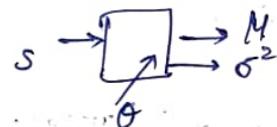
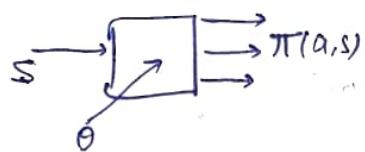
↳ mean & variance
for each state

↳ Used when A is continuous.

When A is discrete,

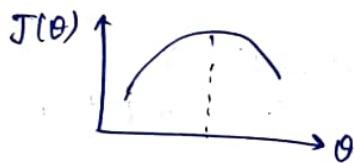
$$\pi(a|s) = \frac{e^{h(a; s, \theta)}}{\sum_{a'} e^{h(a'; s, \theta)}} \stackrel{\text{some fn of } h}{\rightarrow} \text{softmax function}$$

Implementations:



Optimize the policy with the parameter ' θ '.

Metric $\rightarrow ?$



$$\max_{\theta} J(\theta)$$

$$\rightarrow \text{argmax} = \theta^*$$

$$\pi(a|s; \theta^*)$$

What is $J(\theta)$?

- $V^\pi(s_0) \rightarrow \text{argmax}_{\pi} V^\pi(s_0)$

- $\mathbb{E}_d V^\pi(s) \rightarrow \text{denote as } \bar{V}^\pi \checkmark$

↳ average ~~with~~ with some distribution 'd'

- $\mathbb{E}_d r(s, \pi(s)) \rightarrow \text{like average reward}$ to choose initial state.

We will use

$E_d V^{\pi}(s)$ - denote $V^{\pi_\theta} \rightarrow J(\theta)$,
i.e., parameterized by θ , $\pi_\theta(s, a)$

How to maximize?

↪ Gradient Ascent

$$\bar{\theta}_{k+1} = \bar{\theta}_k + \alpha_k \nabla_{\theta} J(\theta)$$

↳ vector

Gradient $\nabla_{\theta} J(\theta) \in$

$$\nabla_{\theta} \bar{V}^{\pi} = \nabla_{\theta} \left(\sum d(s) V^{\pi_\theta}(s) \right)$$

↳ pdf

$$\nabla_{\theta} (V^{\pi_\theta}(s)) : \quad [\theta \rightarrow \text{vector}]$$

$$\begin{aligned} \nabla_{\theta} (V^{\pi_\theta}(s)) &= \nabla_{\theta} \left(\sum_a \pi(a|s; \theta) \cdot q^{\pi_\theta}(s, a) \right) \\ &= \sum_a [\nabla_{\theta} \pi(a|s; \theta) \cdot q^{\pi_\theta}(s, a)] + \sum_a \pi(a|s, a) \cdot \nabla_{\theta} q^{\pi_\theta}(s, a) \end{aligned}$$

A special case: finite horizon with only 1 time-step

$$q^{\pi_\theta}(s, a) = \mu(s, a)$$

$$\Rightarrow \nabla_{\theta} (V^{\pi_\theta}(s)) = \nabla_{\theta} \pi(a|s; \theta) \cdot \mu(s, a) \quad (\because \text{second term} \rightarrow \text{zero})$$

$$\begin{aligned} \nabla_{\theta} q^{\pi_\theta}(s, a) &= \nabla_{\theta} (\mu(s, a)) + \gamma \sum_{s'} p(s'|s, a) V^{\pi_\theta}(s') \\ &= \gamma \sum_{s'} p(s'|s, a) \cdot \nabla_{\theta} V^{\pi_\theta}(s') \end{aligned}$$

∴ we obtain

$$\begin{aligned} \nabla_{\theta} (V^{\pi_\theta}(s)) &= \sum_a \nabla_{\theta} \pi(a|s; \theta) \cdot q^{\pi_\theta}(s, a) \\ &\quad + \gamma \sum_a \pi(a|s; \theta) \cdot \sum_{s'} p(s'|s, a) \cdot \nabla_{\theta} V^{\pi_\theta}(s'). \end{aligned}$$

↪ vector eqn for a single 's'

Take dimension (θ) = 1, for simplicity.

$$\frac{d}{d\theta} V^{\pi_\theta}(s) = \sum_a \frac{d}{d\theta} \pi(a|s; \theta) \cdot q^{\pi_\theta}(s, a) + \gamma \left(\sum_a \pi(a|s; \theta) \cdot \sum_{s'} p(s'|s, a) \cdot \frac{d}{d\theta} V^{\pi_\theta}(s') \right)$$

↓
same thing on both sides

$$\Rightarrow \begin{pmatrix} \frac{d}{d\theta} V^{\pi_\theta}(\cdot) \\ u(s) \end{pmatrix} = \begin{pmatrix} u(s) \\ \downarrow \end{pmatrix} + \gamma P^\pi \begin{pmatrix} \frac{d}{d\theta} V^{\pi_\theta}(\cdot) \\ p^{\pi_\theta}(s'|s) \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} \frac{d}{d\theta} V^{\pi_\theta}(\cdot) \\ u(s) \end{pmatrix} = \begin{pmatrix} I_{|S|} - \gamma P^\pi \\ u(s) \end{pmatrix} = \begin{pmatrix} u(s) \\ \downarrow \end{pmatrix}$$

$$\Rightarrow (I_{|S|} - \gamma P^\pi) \begin{pmatrix} \frac{d}{d\theta} V^{\pi_\theta}(\cdot) \\ u(s) \end{pmatrix} = \begin{pmatrix} u(s) \\ \downarrow \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} \frac{d}{d\theta} V^{\pi_\theta}(\cdot) \\ u(s) \end{pmatrix}_{|S| \times |S|} = (I_{|S|} - \gamma P^\pi)^{-1} \begin{pmatrix} u(s) \\ \downarrow \end{pmatrix} \quad \rightarrow \text{independent of } q^{\pi_\theta}$$

something like

$$\nabla_\theta \pi(a|s, \theta) \cdot q^{\pi_\theta}(s, a).$$

Lec-28

Now,

$$\begin{pmatrix} \frac{d}{d\theta} V^{\pi_\theta}(\cdot) \\ u(s) \end{pmatrix} = \underbrace{(I - \gamma P^\pi)^{-1}}_{I + \gamma P^\pi + \gamma^2 (P^\pi)^2 + \gamma^3 (P^\pi)^3 + \dots} \begin{pmatrix} \sum_a \pi(a|s; b) \cdot q^{\pi_\theta}(s, a) \\ s \boxed{s'} \downarrow \end{pmatrix}$$

for a particular ' s' ,

$$\begin{aligned} \frac{d}{d\theta} V^{\pi_\theta}(s) &= \sum_a \nabla_\theta \pi(a|s, \theta) \cdot q^{\pi_\theta}(s, a) + \\ &\quad \gamma \sum_{s'} p^\pi(s'|s) \cdot \sum_a \nabla_\theta \pi(a|s', a) \cdot q^{\pi_\theta}(s', a) + \\ &\quad \gamma^2 \sum_{s'} (P^\pi(s'|s))^2 \sum_a \nabla_\theta \pi(a|s', \theta) \cdot q^{\pi_\theta}(s', a) + \dots \end{aligned}$$

$$= \sum_a \nabla_\theta \pi(a|s, \theta) \cdot q^{\pi_\theta}(s, a) +$$

$$\sum_a \sum_{s'} (\gamma P^\pi(s'|s) + \gamma^2 (P^\pi(s'|s))^2 + \dots) \cdot \nabla_\theta \pi(a|s', \theta) q^{\pi_\theta}(s', a)$$

About value function approximation:

$$d\pi = d\pi \cdot p\pi^2 \quad \sum_s P_n(s) \cdot P_n(s'|s) = p_n(s') \\ \downarrow \text{stationary distribution} \quad d\pi = p\pi$$

Represent it by f , instead.

$$\sum_s f(s) \frac{d}{d\theta} V^{\pi_\theta}(s) = \sum_s f(s) \left(\sum_a \bar{V}_\theta \cdot \pi(a|s; \theta) \cdot q^{\pi_\theta}(s, a) \right. \\ \left. + \sum_{a' s'} (\gamma P^\pi(s'|s) + \gamma^2 (P^\pi(s'|s))^2 + \dots) \cdot \nabla_\theta \pi(a|s'; \theta) \cdot q^{\pi_\theta}(s', a) \right)$$

$$\sum_a \sum_{s' s} f(s) (\gamma P^\pi(s'|s) + \gamma^2 (P^\pi(s'|s))^2 + \dots) \cdot \nabla_\theta \pi(a|s'; \theta) \cdot q^{\pi_\theta}(s', a)$$

$$= \sum_a \sum_{s'} f(s') \cdot \frac{\gamma}{1-\gamma} \cdot \nabla_\theta \pi(a|s'; \theta) \cdot q^{\pi_\theta}(s', a) \\ = \sum_a \sum_s \underbrace{\frac{f(s)}{1-\gamma}}_{\text{Not a pdf; doesn't sum to 1;}} \cdot \boxed{\nabla_\theta \pi(a|s; \theta)} \cdot q^{\pi_\theta}(s, a) \rightarrow \text{No } \nabla_\theta q^{\pi_\theta}$$

$$\therefore \nabla_\theta J(\theta) = \sum_a \sum_s \underbrace{\frac{f(s)}{1-\gamma}}_{\frac{1}{1-\gamma} \rightarrow \text{constant} \rightarrow \text{can be taken with } \alpha} \cdot \nabla_\theta \pi(a|s; \theta) \cdot q^{\pi_\theta}(s, a) \rightarrow \text{No gradient of } q^{\pi_\theta}; \text{ only of } \pi.$$

log-trick:

$$\nabla_\theta \ln \pi(a|s; \theta) = \frac{\nabla_\theta \pi(a|s; \theta)}{\pi(a|s; \theta)}$$

$$\Rightarrow \pi(a|s; \theta) \cdot \nabla_\theta \ln \pi(a|s; \theta) = \nabla_\theta \pi(a|s; \theta)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{s \sim p(s)} \cdot \mathbb{E}_{A \sim \pi(a|s; \theta)} (\nabla_{\theta} \ln \pi(A|s; \theta)) \cdot q^{\pi_{\theta}}(s, A)$$

\uparrow
 $\circled{S_t, A_t} \rightarrow$ at a time 't' are samples from these distributions.

$$\theta_{t+1} = \theta_t + \alpha_t (\nabla_{\theta} \ln \pi(A_t|s_t; \theta) \cdot q^{\pi_{\theta}}(s_t, A_t)) \rightarrow \text{Update equation}$$

↳ Not implementable as q is not known.

Estimate ' q ' by MC or TD.

Monte-Carlo estimate: Replace by G_t .

$$\theta_{t+1} = \theta_t + \alpha_t (\nabla_{\theta} \ln \pi(A_t|s_t; \theta)) \cdot G_t \dots \text{REINFORCE algorithm}$$

How to compute $\nabla_{\theta} \ln \pi(a|s; \theta)$ for a discrete A .

$$\pi(a|s; \theta) = \frac{e^{\beta h(a|s; \theta)}}{\sum_{a'} e^{\beta h(a'|s; \theta)}} \quad (\because \text{softmax})$$

$$\sum_{a'} e^{\beta h(a'|s; \theta)}$$

to control overflow

if h is large ($e^h \rightarrow$ much larger)

↳ (Temperature)

$$\Rightarrow \ln \pi(a|s; \theta) = \beta \cdot h(a|s; \theta) - \ln \sum_{a'} e^{\beta h(a'|s; \theta)}$$

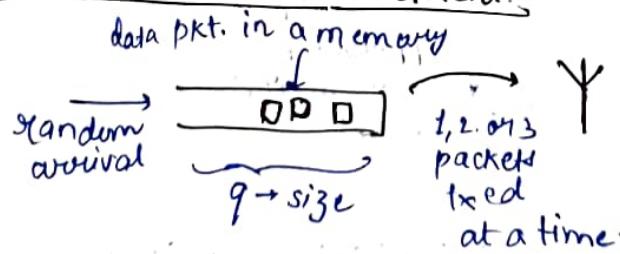
$$\Rightarrow \nabla_{\theta} \ln \pi(a|s; \theta) = \beta \boxed{\nabla_{\theta} h(a|s; \theta)} - \frac{\sum_{a'} (e^{\beta h(a'|s; \theta)} \cdot \boxed{\nabla_{\theta} h(a'|s; \theta)})}{\sum_{a'} e^{\beta h(a'|s; \theta)}}$$

$$h(a|s; \theta) = \sum_i \theta_i \psi_i(a, s)$$

$$\Rightarrow \nabla_{\theta} h(a|s; \theta) = (\psi_i(a, s))$$

↳ vector

Queue Control Environment:



Actor-Critic AlgorithmsRecall:

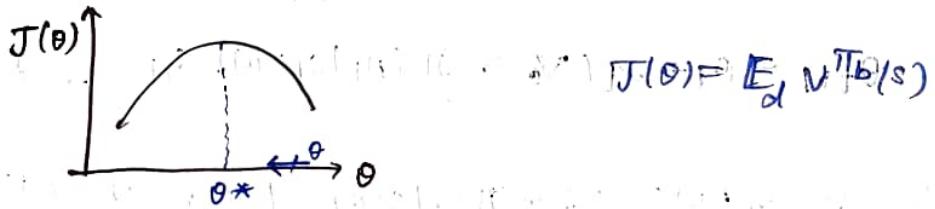
Policies: Parameterised functions $\pi(s)$, $\pi(a|s)$
 $\pi(s; \theta)$ $\pi(a|s; \theta)$

$$\pi(a|s; \theta) = \frac{e^{h(s, a; \theta)}}{\sum_{a'} e^{h(s, a'; \theta)}} \quad (\text{softmax})$$

$\hookrightarrow a \in \text{discrete}$

$$h(s, a; \theta) = \sum_i \phi_i \psi_i(s; a)$$

$$\pi(a|s; \theta) = N(a | \mu(s; \theta), \sigma^2(s)) \rightarrow \text{Gaussian} \quad (\text{continuous})$$



$$\text{Update: } \theta_{k+1} = \theta_k + \alpha_k \nabla_\theta J(\theta) \Big|_{\theta=\theta_k}$$

Gradient of $J(\theta)$,

$$\nabla_\theta J(\theta) = E_d (\nabla_\theta V^{\pi_\theta}(s))$$

$$\nabla_\theta J(\theta) = E_s E_{A \sim \pi(a|s; \theta)} (\nabla_\theta \ln \pi(A|s; \theta) \cdot q^{\pi_\theta}(s, A))$$

↓ example

$$\nabla_\theta \ln \pi(A_t|s_t; \theta) \cdot q^{\pi_\theta}(s_t, A_t)$$

↓

$$\text{SGD: } \theta_{t+1} = \theta_t + \alpha_t \nabla_\theta \ln \pi(A_t|s_t; \theta) \cdot q^{\pi_\theta}(s_t, A_t)$$

Unknown

Reinforce:Replace by Get → MC-ReturnApproximate $q^{\pi_\theta}(s_t, A_t) \approx \hat{q}(s_t, A_t)$ by TD method:

↓ approx. by a
parametrized q

$$\hat{q}(s, a, w)$$

SGD update:

$$\theta_{t+1} = \theta_t + \alpha_t (\nabla_\theta \ln \pi(A_t|s_t; \theta)) \cdot \hat{q}(s_t, A_t; w_t)$$

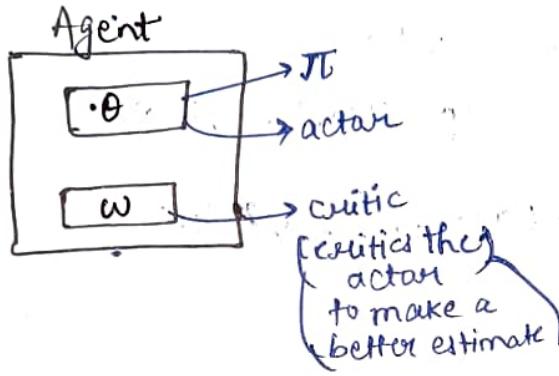
 \hookrightarrow Actor update

But we don't know \hat{q}_t . So update ' w ' also side-by-side to make correct approx. to \hat{q} .

$$w_{t+1} = w_t + \alpha_t^w \cdot [R_t + \gamma \cdot \hat{q}(s_{t+1}, a_{t+1}; w_t) - \hat{q}(s_t, a_t; w_t)].$$

$\nabla_w \hat{q}(s_t, a_t; w_t)$

↳ critic-update ! $\hat{q}(s_t, a_t)$



↳ "Actor-critic algorithm"

(Q-actor critic) [QAC]

[Actor-update]

↳ Requires correct q-estimate

QAC with a small modification:

A₂C - Advantage Action Critic

$$\mathbb{E}_{A \sim \pi(a|s; \theta)} (\nabla_\theta \ln \pi(a|s; \theta)) \cdot (q^\pi(s, a; w) - b(s))$$

choose as $v^\pi(s)$

some fn

= same gradient as before

$$q^\pi(s, a; w) \approx R_t + \gamma \hat{q}(s_{t+1}, a_{t+1}; w_t)$$

$$v^\pi(s) \approx \text{avg}(R_t) + \gamma \cdot \text{avg}(\hat{q}(s_{t+1}, a_{t+1}; w_t))$$

$$\mathbb{E}_{A \sim \pi(a|s; \theta)} (\nabla_\theta \ln \pi(a|s; \theta)) \cdot (\hat{q}(s, a; w) - v^\pi(s; w))$$

↑ cancels.

ideally $\hat{q}(s, a) - v^\pi(s)$: Advantage function

small variance

$A^\pi(s, a)$, ↳ A₂C uses this fn instead of q

$$\therefore \mathbb{E}_{A \sim \pi(a|s; \theta)} (\nabla_\theta \ln \pi(a|s; \theta)) (A^\pi(s, a))$$

$$= \mathbb{E}_{A \sim \pi(a|s; \theta)} \nabla_\theta \ln (\pi(a|s; \theta)) \cdot q^\pi(s, a)$$

$$\mathbb{E}_{\pi^{\theta}} \ln \pi(a_t | s_t; \theta) - V^{\pi}(s_t, a_t)$$

$$= \mathbb{E} V^{\pi}(s_t) \nabla_{\theta} \left(\sum_a \pi(a_t | s_t; \theta) \right) = 0$$

$$\theta_{t+1} = \theta_t + \alpha_{\theta_t} \cdot \nabla_{\theta} \ln \pi(A_t | s_t; \theta_t) \cdot \hat{A}(s_t, A_t; \theta_t)$$

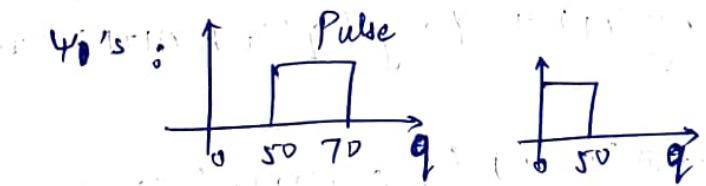
$$\left. \begin{array}{l} A^{\pi}(s, a) = q^{\pi}(s, a) - V^{\pi}(s) \\ \quad = R(s, a) + \gamma \mathbb{E} V^{\pi}(s_t) - V^{\pi}(s) \\ \text{Estimate using } R_t + \gamma \hat{V}^{\pi}(s_{t+1}) - \hat{V}^{\pi}(s_t) \end{array} \right\}$$

Final update:

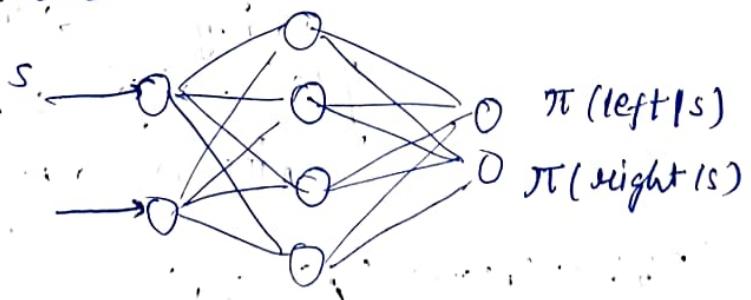
$$\left. \begin{array}{l} \theta_{t+1} = \theta_t + \alpha_{\theta_t} \cdot \nabla_{\theta} \ln \pi(\theta_t | s_t; \theta_t) \\ \quad (\text{Actor}) \qquad \qquad \qquad (R_t + \gamma \hat{V}(s_{t+1}; w_t) - \hat{V}(s_t; w_t)) \\ w_{t+1} = w_t + \alpha_w \cdot (\hat{V}(s_{t+1}; w_t) - \hat{V}(s_t; w_t)) \\ \quad (\text{Critic}) \qquad \qquad \qquad \nabla_w \hat{V}(s_t; w_t) \end{array} \right\}$$

Implementations:

Basis fns (Ψ):



Minimal RL



Model-Based RL methods:

Up to now, model force

$$\rightarrow x_{t+1} = \boxed{A}x_t + \boxed{B}a_t + N_t$$

↓ ↓
Unknown

Let's estimate A and B.

$$\rightarrow R_t = -(x_t^T Q x_t + u_t^T R u_t) \quad \leftarrow \text{LQR}$$

$$a_t = K x_t$$

DYNA algorithm:

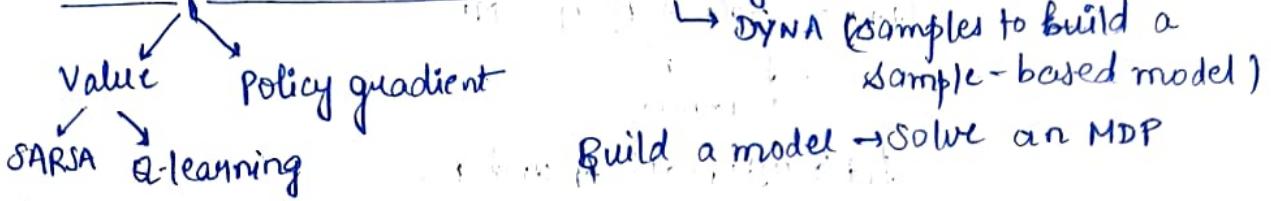
↳ Uses experience samples

↳ we see (s_t, a_t, r_t, s_{t+1}) .

- Actual model: $\pi(s, a), P$

- Sample-based model: Train the system/weights on the model (if sufficient memory is present)

S	A	R	S ⁺
0	1	5	2
1	0	-2	1

RecallModel-free and Model-based RL algorithms

function approximations: $\hat{V}(s; \theta), \hat{q}(s, a; \theta), \pi(a|s; \theta)$

Linear function approx.: $\sum w_i \phi_i(s)$ or $\sum w_i \phi_i(s, a)$

$$\sum \theta_i \psi_i(s)$$

(Basis function: Appropriate functions - very difficult!)

REINFORCE: $\pi(a|s, \theta)$

Neural Network:

$$n: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$$(w_1x_1 + w_2x_2 + \dots + w_nx_n + b)$$

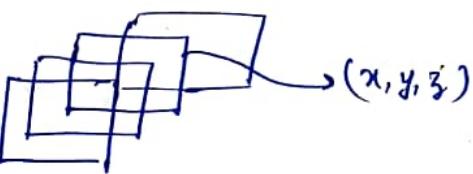
↓
sigmoid, ReLU

→ Deep RL Handbook - Maxim Laiyan → GitHub

(Chap-3)

Tensors:

L
multi-dimensional array

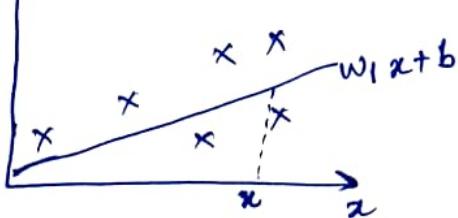


| Scalar, vector, matrices

$$n: w_1x_1 + w_2x_2 + w_3x_3 + b$$

$$n: w_1x + b$$

$$\begin{pmatrix} 2 \\ 3 \end{pmatrix}$$

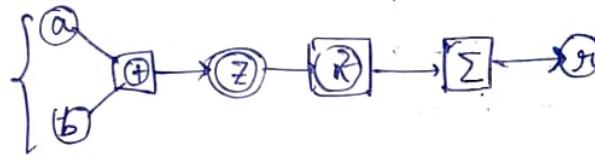


$$\left. \begin{aligned} w_{k+1} &= w_k - \alpha_k \left(\frac{\partial L}{\partial w} \Big|_{w_k} \right) \\ b_{k+1} &= b_k - \alpha_k \left(\frac{\partial L}{\partial b} \Big|_{b_k} \right) \end{aligned} \right\} \quad \nabla \quad \overline{\theta}_{k+1} = \overline{\theta}_k - \alpha_k \nabla L$$

n: $w_1 x + b$ $\nabla = \begin{pmatrix} x \\ 1 \end{pmatrix}$
 $(w_1 + \Delta w_1) x + b$

$$f = f_2(f_1(x)) \quad \left[\frac{df}{dx} = \underbrace{\frac{df}{df_2}}_{df_2} \cdot \underbrace{\frac{df_2}{df_1}}_{df_1} \cdot \underbrace{\frac{df_1}{dx}}_{dx} \right]$$

→ Micrograd implementation: [YouTube : Karpathy]



Minimal RL Codes

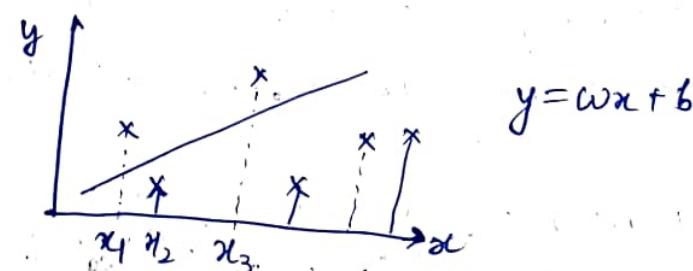
Recall:

Lec-31 | 11-11-2025

Tensors: multi-dimensional array

- functions
- Autograd — optimization

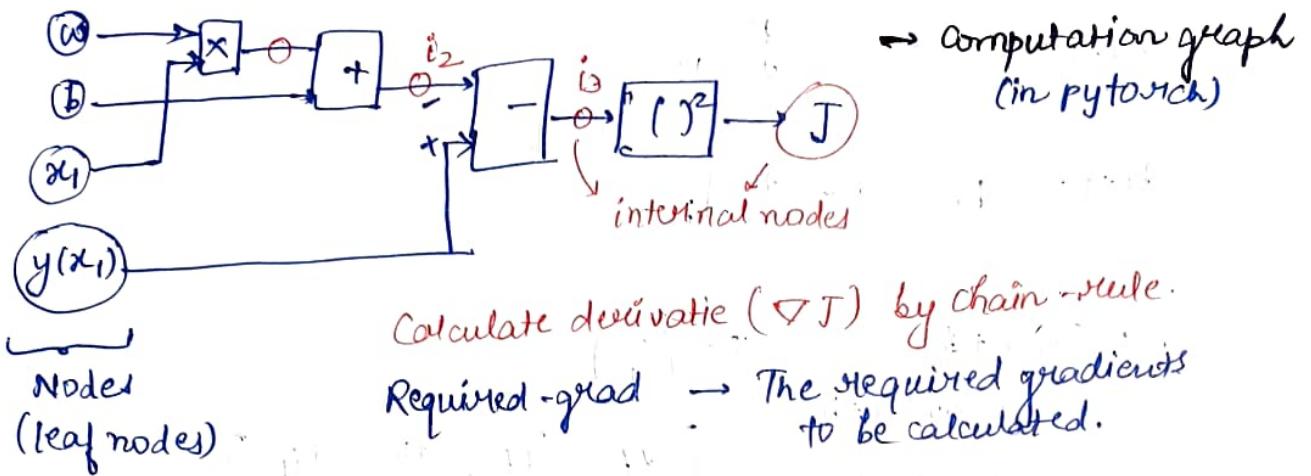
$\nabla f(x)$



$$\sum_{i=1}^N (y(x_i) - (wx_i + b))^2 = J(w, b).$$

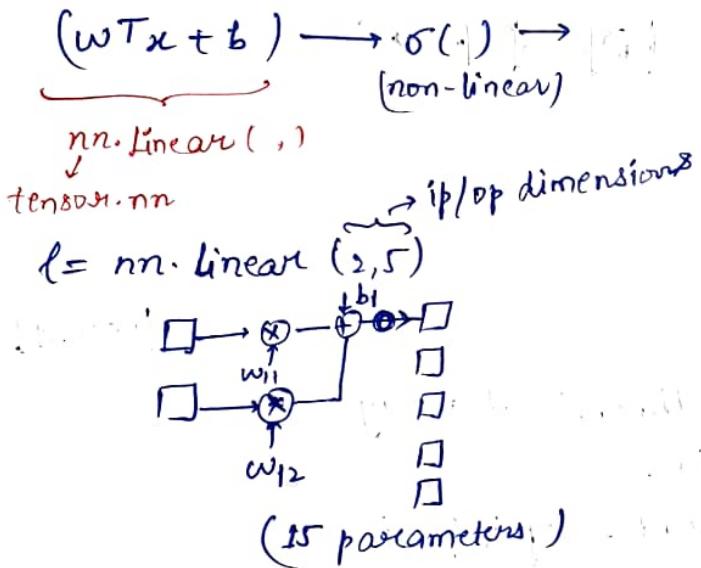
$$(w, b)_{k+1} = (w, b)_k - \alpha_k \cdot \nabla J(w, b)$$

w, b Two variables



$$\frac{\partial J}{\partial b} = (2i_3)(-1)(1) = -2i_3 \\ = -2(y(x_1) - (w(x_1) + b))$$

Neural Networks: function approximators $\hookrightarrow \text{backward}()$

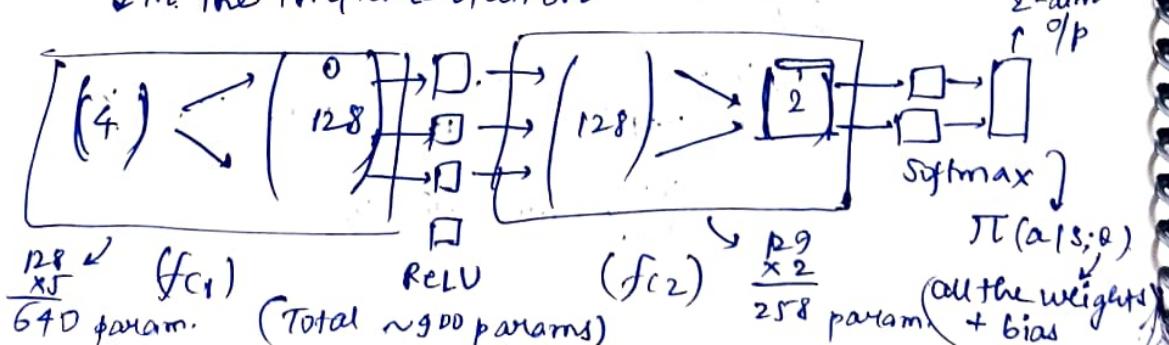


- Zero the gradient → `zero_grad()`
- Cascade the blocks → `nn.Sequential(, ,)`

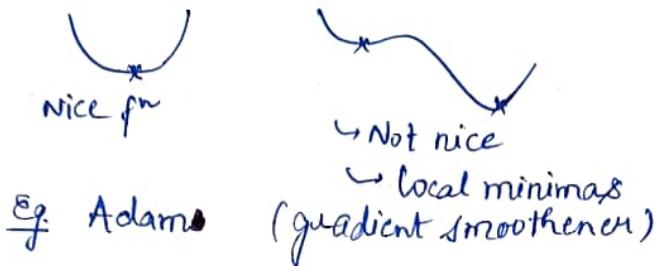
Recall REINFORCE algorithm: policy gradient theorem

$$\pi(a|s; \theta) \xrightarrow{\text{Is } \cancel{\pi(s)} \text{?}} \nabla \ln(\pi(A|s, \theta)) \cdot q(s, a) \xrightarrow{G_t} \nabla_\theta \ln \pi(A_t | s_t; \theta_t) \cdot G_t$$

in the implementation



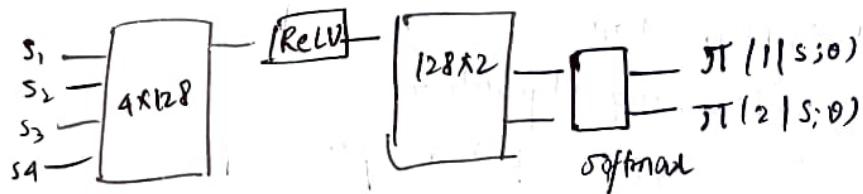
Optimizers: Helps in local minimas



Lec-32 [12-11-2025]

REINFORCE

$$\pi(a|s; \theta)$$



$$J(\theta) = \mathbb{E} V^\pi(s)$$

$$\nabla J(\theta) = \mathbb{E}_s \mathbb{E}_{A \sim \pi(\cdot | s, \theta)} (\nabla_\theta \pi(A | s; \theta)) q^\pi(s_t, A_t)$$

$$\nabla J(\theta) = \nabla_\theta \pi(A_t | s_t, \theta) G_t$$

$$R_1, R_2, \dots, R_T$$

$$\pi(A_0 | s_0), \dots, \pi(A_T | s_T; \theta)$$

Deep QN (Q-Learning):

Q-learning: Directly learn optimal q
(Tabular) by solving Bellman optimality equation.

- Chapter-8 → sec-8.4
- Book: Lajan

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) + \alpha_t (R_t + \gamma \max_a q(s_{t+1}, a) - q(s_t, a_t))$$

Continuous: Parameterized by weight 'w'.

$$q(s, a, \underline{w})$$

$$\text{Update } w: w_{t+1} = w_t + \alpha_t (R_t + \gamma \max_a q(s_{t+1}, a; w_t) - q(s_t, a_t; w_t)) \nabla_w q(s_t, a_t; w)$$

Minimize the error: $(R_t + \gamma \max_a q^{target}(s_{t+1}, a; \underline{w}) - q(s_t, a_t; \underline{w}))^2$

- Error is also parameterised by 'w'; which of which it is no longer $\propto q$.
- Introduce one more parameterization, q_1, q_2, \dots, q_n , q_{target} , q_t , and q_{target} should not or slowly change.

Experience Replay.

[Lec-33]

[17-11-2025]

DQN Implementation

↳ Q-learning algorithm

↳ finding $q(s; a)$ for an optimal policy.

In SARSA, at t , we are at s_t, a_t ,

$R_t + \gamma \hat{q}(s_{t+1}, a_{t+1}) \rightarrow$ update the current
 \uparrow estimate of s_t, a_t
 for a particular policy

Q-learning target:

$$R_t + \gamma \max_a \hat{q}(s_{t+1}, a)$$

$$\hat{q}_{t+1}(s_t, a_t) = \hat{q}_t(s_t, a_t) + \alpha_t (R_t + \gamma \max_a \hat{q}(s_{t+1}, a) - \hat{q}(s_t, a_t)) \dots \text{Tabular update}$$

$\hat{q}(s_t, a_t; w_t)$
 ↳ update w_t so that q-estimate
 and q-target error is minimized.

Error is measured:

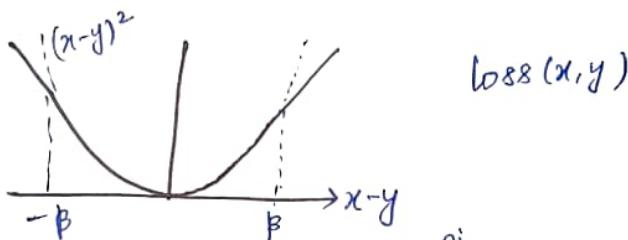
$$(R_t + \gamma \max_a \hat{q}(s_{t+1}, a; w_t) - \hat{q}^{\text{target}}(s_t, a_t; w_t))^2$$

$$w_{t+1} = w_t + \alpha_t / (R_t + \gamma \max_a \hat{q}(s_{t+1}, a; w_t) - \hat{q}^{\text{target}}(s_t, a_t; w_t)) \times \nabla_w \hat{q}(s_t, a_t; w_t)$$

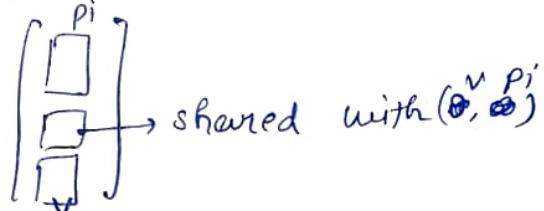
↑ update this

Recall

- ① $V(s; w) \xrightarrow{\text{NN}} \text{change w s.t. } V(s; w) \approx \boxed{\text{target}}$
 (2 NNs)
- ② Experience Replay buffer.
 fixed target

Smooth L₁ loss :

Parameter vector:

computation graph :