

State-action-reward

$$R: S \rightarrow R \text{ (Previously)}$$

(new state maps
to reward)

Now:-

$$R: S \times A \rightarrow R \text{ (Now state } \rightarrow \text{action maps } \rightarrow \text{to reward)}$$

$$R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2)$$

\downarrow

↳ Different actions might have different rewards

→ Bellman's equation: (Updated)

$$V^*(s) = \max_a [R(s, a) + \gamma \sum_{s'} P_{sa}(s') V^*(s')]$$

↓
Current
Reward

Future
Reward

Max
over the
whole equation

- * Remember reward depends on current reward & action.
 - * Value iteration can still be used
 - * $\pi^*(s) = \underset{a}{\operatorname{argmax}} R(s, a) + \gamma \sum_{s'} P_{sa}(s') V^*(s')$
- ↓
argmax
on the
whole equation

④ Finite horizon MDP:-

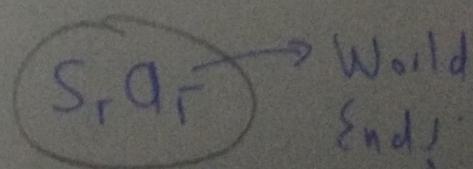
$$\rightarrow (S, A, \{P_{sa}\}, T, R)$$

→ We are going to replace T with T .

$$\rightarrow S_0, a_0$$

→ MDP would run for a definite period of time:-

$$S_0, a_0, S_1, a_1, \dots$$



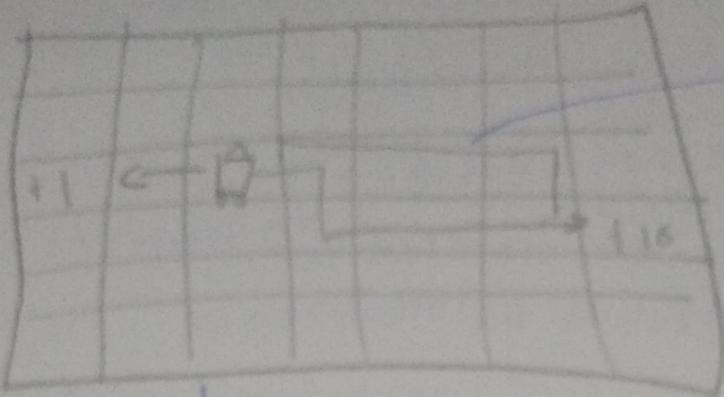
$$[R(S_0, a_0) + R(S_1, a_1) + \dots + R(a_r, S_r)]$$

→ Finite Reward :-

$$\rightarrow E[R(S_0, a_0) + R(S_1, a_1) + \dots + R(a_r, S_r)]$$

→ Goal is to maximize this expected payoff.

(*)



→ If time
allows.

↳ Sometime reward depends on time. If you have two to three time steps left then you should rush towards +1 but if you more than that then you should rush towards +10.

→ Stationary Policy does not change over time.
↳ $\pi^*(s)$ depends on time.
↳ So $\pi^*(s)$ should be written as $\pi^*(s)$ (non-stationary policy/ changes over time).

③ Non-Stationary state transitions:

$$\textcircled{1} \quad S_{t+1} \sim P_{S_t, a_t}^{(t)} \quad \textcircled{1}$$

↳ Non-stationary state-transition Probability

$$\textcircled{2} \quad R^{(t)}(s, a) \quad \textcircled{2}$$

↳ Reward depends on time.

④ When do you need finite horizon
or non-stationary transitions:-

in Changing dynamics

- Weather forecast (Different environment)
- Industrial automation
-

→ How you would solve finite horizon MDP.

$$V_t^*(s) = \mathbb{E}[R(s_t, a_t) + R(s_{t+1}, a_{t+1}) + \dots + R(s_T, a_T) \mid \pi^*, s_0 = s]$$

↑
Expected total payoff starting s at time t execute π^* .

→ Value iteration: (Dynamic Programming)
in this case evolves:-

$$\rightarrow V_t^*(s) = \max_a R(s, a) + \sum_{s'} P_{sa}(s') V_{t+1}^*(s')$$

$$\rightarrow \pi^*(s) = \arg \max_a (\dots)$$

→ So let start with $V_T^*(s)$ the
 T is last time step.

$$V_T^*(s) = \max_a R(s, a)$$

↳ No time step after this

→ The essential idea is that we first compute $V_T^*(s)$ then $V_{T-1}^*(s)$ then V_{T-2}^* and so on till V_0^* and then compute $\pi_t^*(s)$ for each step. Ea station for each step.

→ ↳ How Value iteration for finite time horizon changes.

→ Linear Quadratic Regression (LQR):

④ $(S, A, \{P_{sa}\}, T, R) \rightarrow$ Finite horizon MDP

④ $S \in \mathbb{R}^n, A \in \mathbb{R}^d$

④ $P_{sa}: S_{t+1} = AS_t + Bar + W_t \rightarrow$ Next state evolves as function of p_{t+1} , s_t and a_t .

$R^T \in \mathbb{R}^{n \times n}$ $R \in \mathbb{R}^{n \times d}$ $W \sim (0, \Sigma_w)$

$$\text{④ } R(s, a) = -(s^T Us + a^T Va)$$

↓ $U \in \mathbb{R}^{n \times n}, V \in \mathbb{R}^{d \times d}$

Another Assumption $U, V \geq 0, s^T Us \geq 0, a^T Va \geq 0$

(Q17)

Want $s \approx 0$ (We want this to be stable)

then Choose $U = I$, $V = I$

$$\text{then } R(s, a) = -(\|s\|^2 + \|q\|^2)$$

\hookrightarrow Reward penalizing large movement from stable state.

\rightarrow Just an example.

* Where to get A, B ? \rightarrow Method 1 - Learning A, B

$$S_0^{(1)} \xrightarrow{q_0} S_1^{(1)} \rightarrow \dots S_T$$

←
Collecting
Data

$$S_0^{(m)} \xrightarrow{q_0} S_1^{(m)} \rightarrow \dots S_T^{(m)}$$

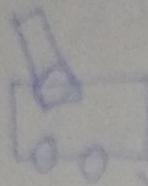
$$\rightarrow S_{t+1} = AS_t + Ba_t$$

$$\rightarrow \min_{A, B} \sum_{i=1}^m \sum_{t=0}^T \| (AS_t + Ba_t) - S_{t+1}^{(i)} \|^2$$

\hookrightarrow You will be able to find
A and B.

(218)

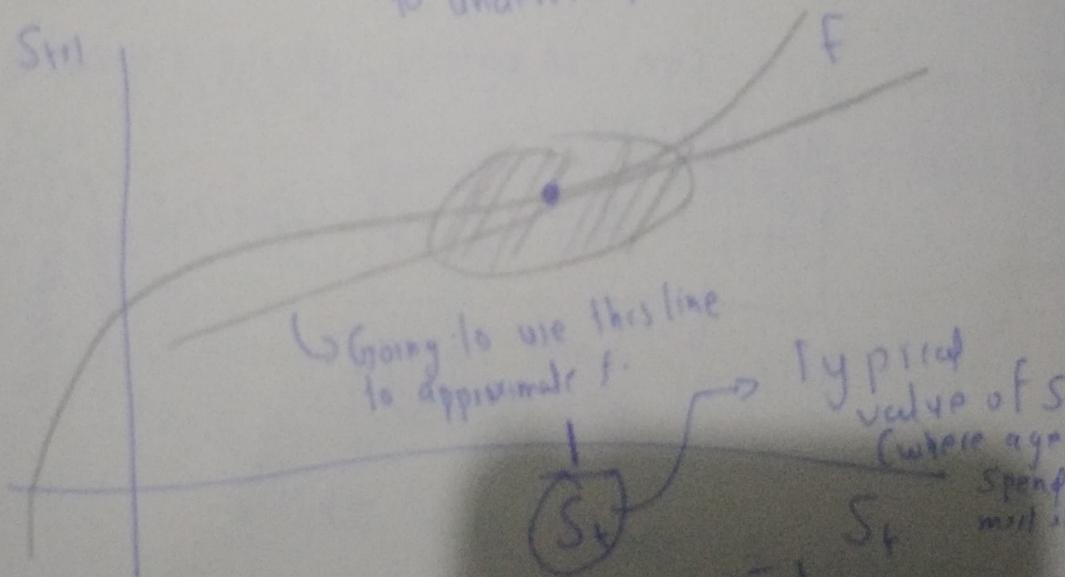
④ → Method 2 (Find A,B) — Linearize a non-linear model



$S_{t+1} = f(S_t, a_t)$ → Physics law + simulator

$$\begin{pmatrix} X_{t+1} \\ \dot{X}_{t+1} \\ \Theta_{t+1} \\ \dot{\Theta}_t \end{pmatrix} = f\left(\begin{pmatrix} X_t \\ \dot{X}_t \\ \Theta_t \\ \dot{\Theta}_t \end{pmatrix}, a_t\right)$$

→ $S_{t+1} = f(S_t)$ (Ignoring action for now)
to understand linearization



$$④ S_{t+1} \approx f'(\bar{S}_t) (S_t - \bar{S}_t) + f(\bar{S}_t) \quad (\text{Approximating } f)$$

\bar{S}_t is a constant

(219)

① And if $S_{t+1} = f(S_t, a_t)$ then linearized

equation would be: (Linearize f around (\bar{S}_t, \bar{a}_t))

$$S_{t+1} \approx f'(\bar{S}_t, \bar{a}_t) \cdot (S_t - \bar{S}_t) + f(\bar{S}_t, \bar{a}_t)$$

$$\begin{aligned} S_{t+1} &= f(\bar{S}_t, \bar{a}_t) + \left(\nabla_{S_t} f(\bar{S}_t, \bar{a}_t) \right)^T (S_t - \bar{S}_t) \\ &\quad + \left(\nabla_{a_t} f(\bar{S}_t, \bar{a}_t) \right)^T (a_t - \bar{a}_t) \end{aligned}$$

↳ Approximation of original function

(Expression S_{t+1} as

linear function $\bar{S}_t + \bar{a}_t$)

→ You can express above equation as follow

$$S_{t+1} = A S_t + B a_t$$

↳ Same concept as method 1.

where $S_t \begin{pmatrix} 1 \\ x \\ a \\ 0 \end{pmatrix}$ → Need to add a intercept term -

→ The above methods show

you model finite horizon MDP
as linear dynamical model

220

→ Main Problem

$$\textcircled{1} \quad S_{t+1} = AS_t + Ba_t + W_t$$

and if

$$\textcircled{2} \quad R(S, a) = -(S^T U_S + a^T V_a)$$

Total Pay off:-

$$\textcircled{3} \quad R(S_0, a_0) + R(S_1, a_1) + \dots +$$

$$R(S_T, a_T)$$

→ Dynamic Programming algorithm
to solve for V^* :

(Final time step) $V_T^*(s) = \max_{a_T} R(S_T, a_T)$

$\leftarrow = \max_{a_T} -S_T^T U_S T - a_T^T V_a -$

Choose action to maximize

$V_T^*(s) - \text{Value}$ of last time $= -S_T^T U_S T$

$\pi_T^*(S_T) = 0$ (choose action 0)

→ Suppose

$$\leftarrow V_{t+1}^*(S_{t+1}) = S_{t+1}^\top \Phi_{t+1} S_{t+1} + \Psi_{t+1}$$

Given
MDP
is linear
System

where

$$\Phi_{t+1} \in \mathbb{R}^{n \times n}, \Psi_{t+1} \in \mathbb{R}$$

and reward
function has
the form
we mentioned
 $V_{t+1}^*(S_{t+1})$ would
be quadratic

Assuming
continuous
state.

$$\rightarrow V_t^*(S_t) = \max_{a_t} R(S_t, a_t)$$

$$\rightarrow \text{Now putting in equations:-} + E_{S_{t+1} \sim P_{S_t a_t}} [V_{t+1}^*(S_{t+1})]$$

$$\rightarrow V_t^*(S_t) = \max_{a_t} -S_t^\top U S_t - Q_t + V_{a_t}$$

$$+ E_{S_{t+1} \sim P_{S_t a_t}} [C]$$

Continue on
next page (or clarify)

(222)

④

$$V_t^*(S_t) = \max_{a_t} -S_t^T A S_t - a_t^T V a_t$$

$$+ E_{S_{t+1}} \sim N(A S_t + B a_t, \Sigma_w) \quad \dots$$

Remember
from part
linear model

$$\dots S_{t+1}^T \Phi_{t+1} S_{t+1} + \gamma_{t+1}$$

$$V_{t+1}^*(S_{t+1})$$

→ This simplifies to a big quadratic
function of a_t

Take derivative wrt a_t
Set to zero and solve for a_t
then you end with following a_t formula

$$a_t = (B^T \Phi_{t+1} B - V)^{-1} B^T \Phi_{t+1} A \cdot S_t$$

therefore

$$\pi_t^*(S_t) = L_t \cdot S_t$$

where

\leftarrow
image
noise system.

(224)

→ Initialize $\Phi_T = -U$, $\gamma_T = 0$

Recursively calculate

Φ_t, γ_t using Φ_{t+1}, γ_{t+1} (for $t = T, T-1, \dots, 0$)
(Equations defined above)

Calculate a_t

$$\pi^*(s_t) = L_t s_t$$

→ Curly to solve a finite horizon
MDP (But MDP has
to be approximated to
linear dynamical system).