

## ~~Lecture 23~~ Finite horizon Markov decision processes:-

Text books :- ① CS6700 in my website

② Tsitsiklis & Bertsekas "Neuro-dynamic Programming", 1996

③ Sutton & Barto, "Intro. to RL"

④ Bertsekas "DP & OC, Vol I & II".

Example :- Machine replacement/repair

States =  $\{1, 2, \dots, n\}$   
↓                      ↓  
good                  worst

Actions = Do nothing, Repair

On repair: machine goes to "1". Cost = R

Operating cost:  $g(i)$

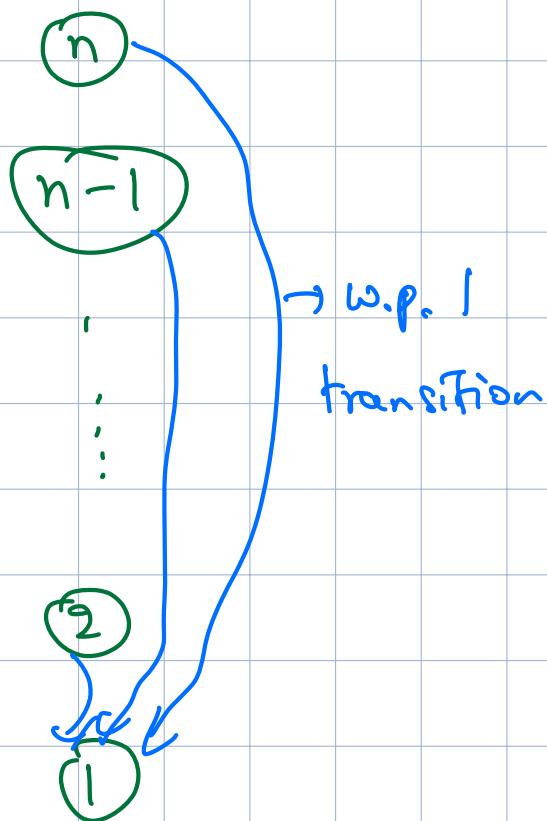
$$g(1) \leq g(2) \leq g(3) \leq \dots \leq g(n)$$

State transitions:

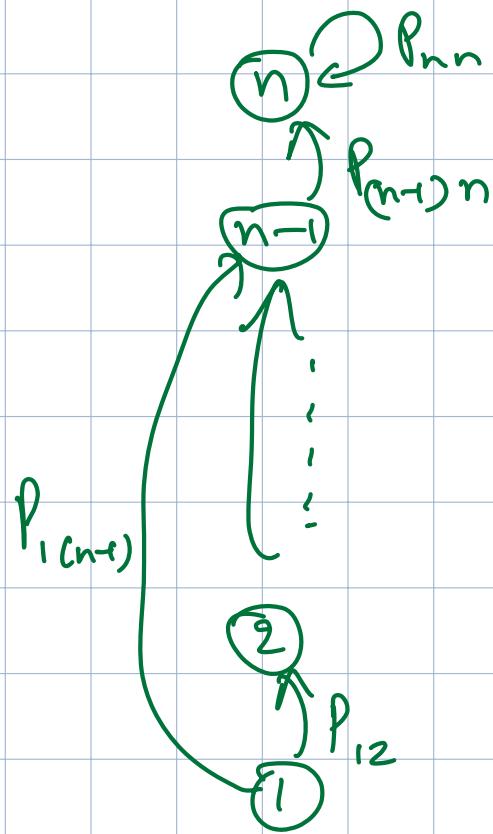
$P_{ij}$ : probability of state transitioning from  $i$  to  $j$

$$P_{ij} = 0 \text{ if } j < i$$

Action: report



Action: Do nothing



Example 2: Chess match

Playing styles (actions): Timid, Bold

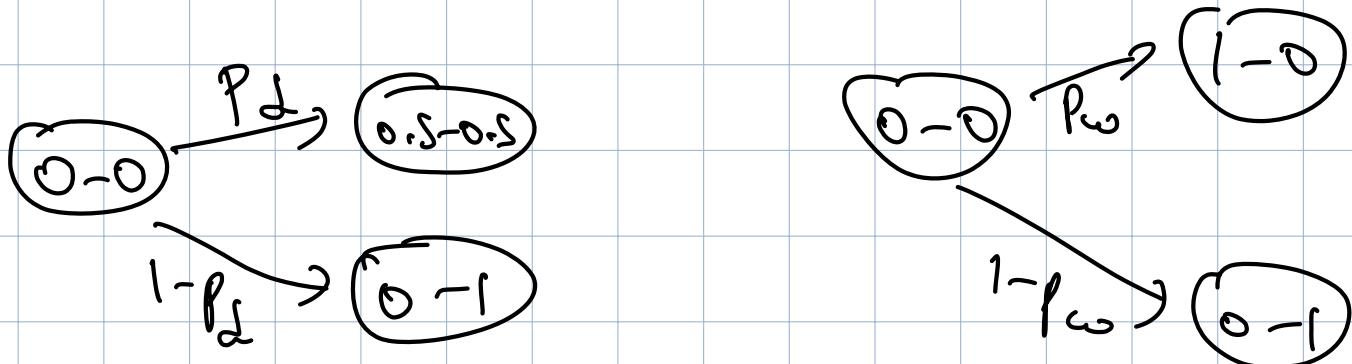
Timid

Bold

draw $P_d$	lose $1 - P_d$
win $P_w$	lose $1 - P_w$

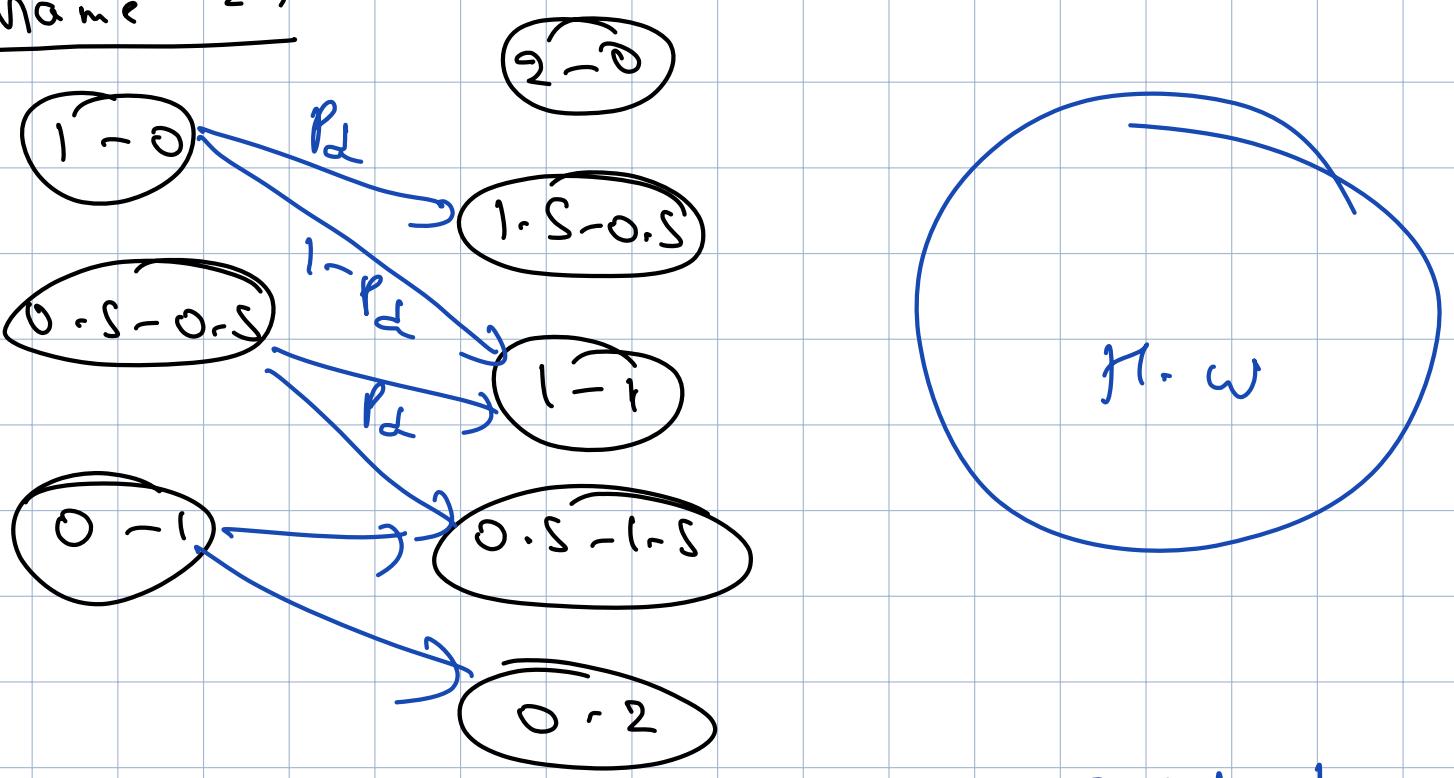
$$P_d > P_w$$

Game 1



Timid

Game 2:



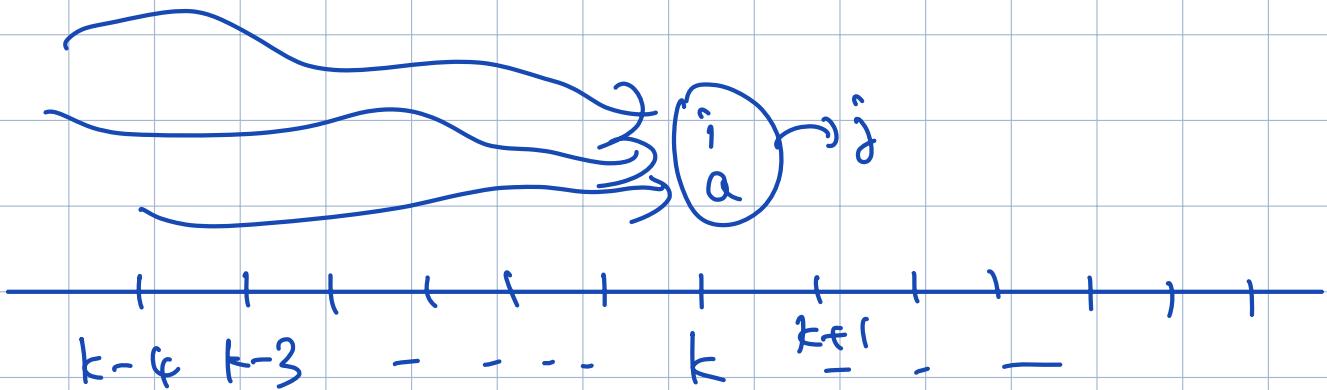
Timid play

Bold play

## Framework

$S$ : State space     $A$ : action space

$$P_{ij}^a(a) = P(S_{k+1} = j \mid S_k = i, a_k = a)$$



~~Single-stage~~ cost:  $g_k(i, a, j)$   
"slot"

Policy  $\pi$ :  $\{\mu_0, \dots, \mu_{N-1}\}$

$N \in \# \text{ of slots}, \mu_i: S \rightarrow A$ .  
 ↪ State space      Action space

Performance metric

$$\mathcal{J}_\pi(s_0) = E \left[ g_N(s_N) + \sum_{k=0}^{N-1} g_k(s_k, \mu_k(s_k), s_{k+1}) \right]$$

↓  
Terminal cost

Goal: Find  $\pi^*$  s.t.

$$J_{\pi^*}(s_0) = \min_{\pi} J_{\pi}(s_0)$$

Optimality principle:

$$\pi^* = (\mu_0^*, \mu_1^*, \dots, \mu_{N-1}^*)$$

A tail sub-problem

$$\min_{\pi^i} E \left( g_N(s_N) + \sum_{k=i}^{N-1} g_k(s_k, \mu_k(s_k), s_{k+1}) \right)$$
$$\pi^i = (\mu_i, \dots, \mu_{N-1})$$

Claim: For this tail sub-problem

$\{\mu_i^*, \dots, \mu_{N-1}^*\}$  is optimal

DP algorithm :

Idea: Solve tail sub-problems of problem  
backwards

Algorithm :

Set  $J_N(S_N) = g_N(S_N)$   $\forall S_N \in \mathcal{S}$ .

For  $k = N-1, \dots, 0$

{

$J_k(S_k) = \min_{a_k \in A} E_{S_{k+1}} \left[ g_k(S_k, a_k, S_{k+1}) + J_{k+1}(S_{k+1}) \right]$ ,  
 $\forall S_k \in \mathcal{S}$ .

}

Apply DP algo to "Machine Example".

States =  $\{1, \dots, n\}$

P.T.O.

Suppose terminal cost is zero

$$J_N(i) = 0$$

$$J_k(i) = \min \left\{ \underbrace{R + g(i) + J_{k+1}(l)}, \text{Repair} \right.$$

$$g(i) + \sum_{j=i}^n p_{ij} J_{k+1}(j) \right\}$$

Do nothing

Chess-example (re-visited)

State = net-score (#wins - #losses)

$$J_k(S_k) = \max \left( \underbrace{p_d J_{k+1}(S_k)}_{\text{Timid}} + (1-p_d) J_{k+1}(S_k-1), \underbrace{p_w J_{k+1}(S_k+1) + (1-p_w) J_{k+1}(S_k-1)}_{\text{Bold}} \right)$$

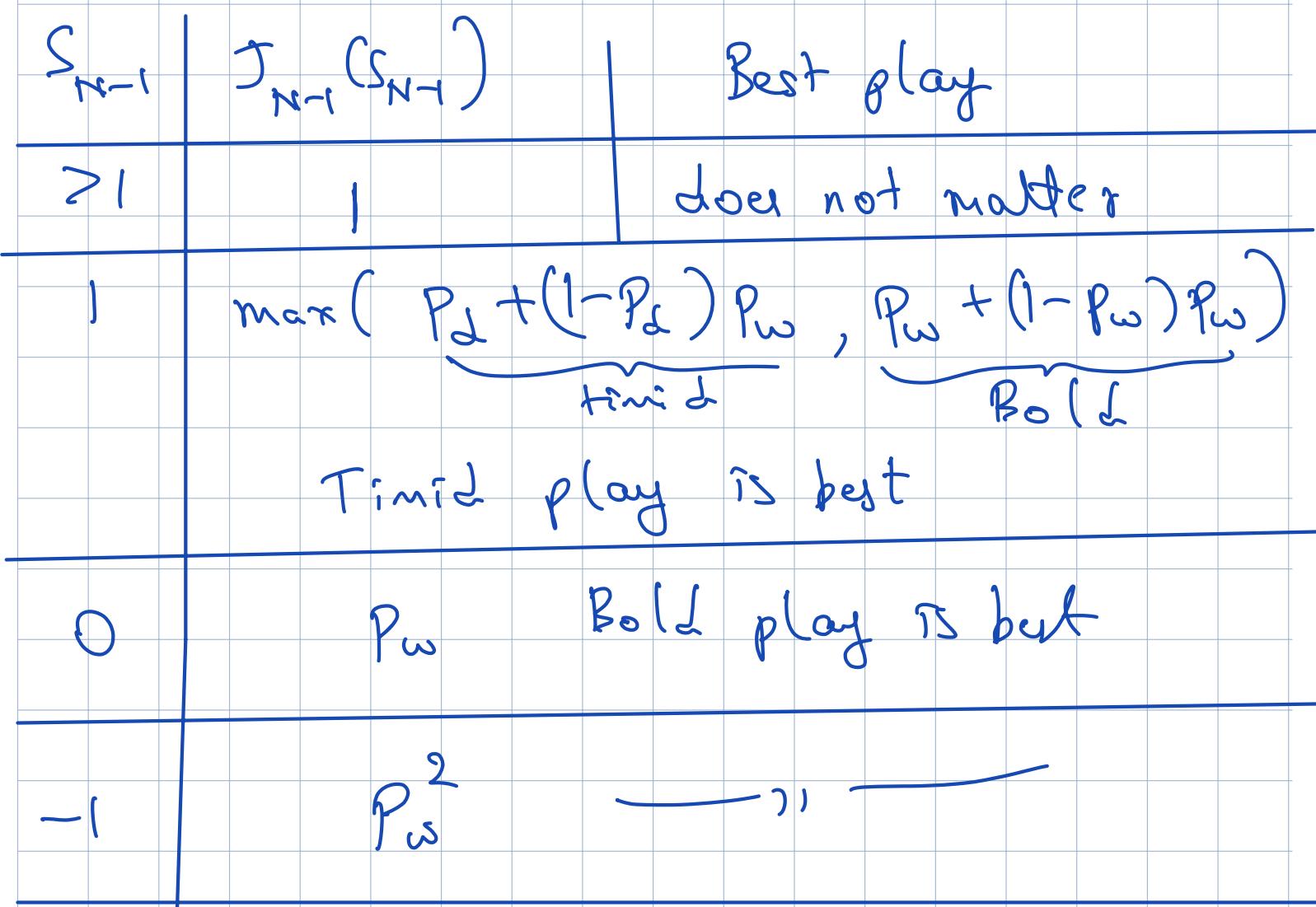
$$p_w J_{k+1}(S_k+1) + (1-p_w) J_{k+1}(S_k-1)$$

Bold (L)

Play bold if  $\frac{P_w}{P_d} \gtrsim \frac{\mathcal{J}_{k+1}(S_k) - \mathcal{J}_{k+1}(S_{k-1})}{\mathcal{J}_{k+1}(S_{k+1}) - \mathcal{J}_{k+1}(S_{k-1})}$

Initial condition:

$$\mathcal{J}_N(S_N) = \begin{cases} 1 & \text{if } S_N > 0 \\ P_w & \text{if } S_N = 0 \\ 0 & \text{if } S_N < 0 \end{cases}$$



$< -1$

0

Leave it

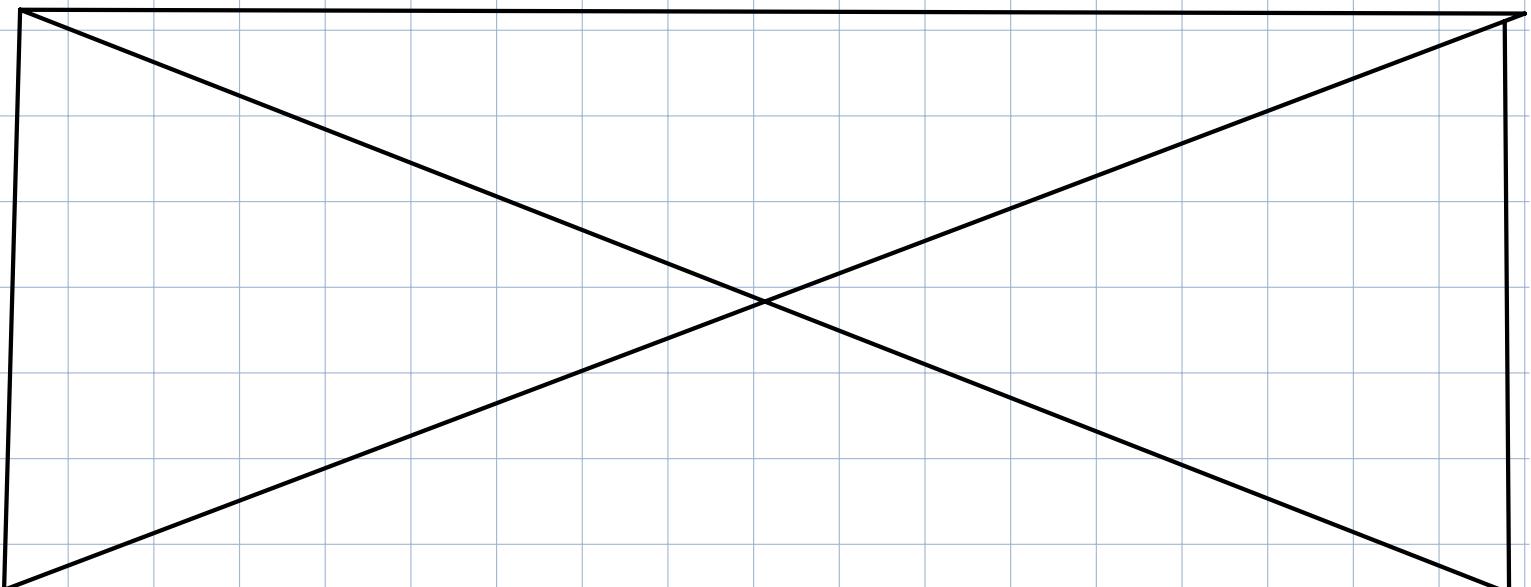
For the 2-game match

$$J_{N-2}(0) = \max \left( P_d P_w + (1-P_d) P_w^2, \underbrace{\dots}_{\text{Timid}} \right)$$

$$P_w (P_d + (1-P_d) P_w) \\ + (1-P_w) P_w^2 \underbrace{\dots}_{\text{Bold}}$$

Check  $\Rightarrow$

$$P_w (P_d + (1-P_d) P_w) + (1-P_w) P_w^2$$



## Lecture - 21

# Infinite horizon discounted cost MDPs:

States:  $\{1, \dots, n\}$

Stationary deterministic policy (SDP)

$$\Pi = \{\mu_0, \mu_1, \dots\} \quad \mu_i: S \rightarrow A$$

Start state  
col

Stationary policy  $\Pi = \{\mu_0, \mu_1, \dots\}$

$$J_\pi(i) = E \left( \sum_{k=0}^{\infty} \alpha^k g(S_k, \pi(S_k), S_{k+1}) \middle| S_0 = i \right)$$

↓  
Discount factor  $0 < \alpha < 1$

Single-slot cost policy

Goal:  $J^\infty(i) = \min_{\pi} J_\pi(i), \forall i \in S.$

Let  $\pi^*$  be an optimal policy

corresponding to  $J^*$

# states

$$J^* = [J^*(1), \dots, J^*(n)]$$

Bellman operator

$$J = (J(1), \dots, J(n))$$

$$TJ = (TJ(1), \dots, TJ(n)), \text{ where}$$

$$(TJ)(i) = \min_{a \in A} \sum_{j=1}^n P_{ij}(a) [g(i, a, j) + \gamma J(j)], \quad \forall i \in S$$

$$(T_\pi J)(i) = \sum_{j=1}^n P_{ij}(\pi(i)) (g(i, \pi(i), j) + \gamma J(j)), \quad \forall i \in S$$

$$\text{Let } P_\pi = \begin{bmatrix} P_{11}(\pi(1)) & & & \\ & \ddots & & \\ & & \ddots & \\ & & & P_{nn}(\pi(n)) \end{bmatrix}$$

$$= \left[ \left[ P_{ij} (\pi(i)) \right] \right]_{i,j=1-n}$$

$$g_\pi = \begin{bmatrix} \sum_{j=1}^n P_{1j} (\pi(1)) g(1, \pi(1), j) \\ \vdots \\ \sum_{j=1}^n P_{nj} (\pi(n)) g(n, \pi(n), j) \end{bmatrix}$$

$$T_\pi J = g_\pi + \alpha P_\pi J$$

Properties of  $T, T_\pi$ :

For any  $J, J'$  s.t.  $J(i) \leq J'(i)$ ,  $\forall i$

- ①  $T^k J(i) \leq T^k J'(i), \forall i, \quad \text{for } k \geq 1$
- ②  $T_\pi^k J(i) \leq T_\pi^k J'(i), \forall i, \quad \text{for } k \geq 1$

Proposition:- (Basis for Value Iteration)

For any (bounded)  $J: S \rightarrow \mathbb{R}$

$$J^*(i) = \lim_{N \rightarrow \infty} (T^N J)(i), \text{ If } S$$

Remark:-

$$\begin{aligned} J_0 &\xrightarrow{T} J_1 = T J_0 & T \xrightarrow{} J_2 = T^2 J_0 \\ && \downarrow \\ && \vdots \quad \text{as } N \rightarrow \infty \\ T^N J_0 &\xrightarrow{\cdot} J^* \end{aligned}$$

(Corollary:  $J_\pi(i) = \lim_{N \rightarrow \infty} (T_\pi^N J)(i), \text{ If } S$ )

for any bounded  $J$ .

# Value iteration for solving a

## discounted MDPs

Fix  $J_0$

Repeatedly apply  $T$

$$J_0 \xrightarrow{T} J_1 \xrightarrow{\dots} \dots \xrightarrow{*} J^*$$


Bellman      optimality      equation

The optimal cost

$$J^*(i) = \min_{\pi_i} J_\pi(i)$$

Satisfies

fixed-point  
equation

$$J^*(i) = T J^*(i)$$

(or)

$$J^*(i) = \min_a \sum_{j=1}^n p_{ij}(a) (g(i, a, j) + \gamma J^*(j))$$

Claim:-  $J^*$  is the unique fixed point

Pf: Suppose not.

$$\exists J', \quad J' = TJ'$$

$$J' = TJ' = T^2J' = \dots = \lim_{N \rightarrow \infty} T^N J' \\ = J^*$$

Corollary: For any stationary  $\pi$ ,

$J_\pi = T_\pi J_\pi$  &  $J_\pi$  is the unique solution.

Illustration of Value Iteration

$$S = \{1, 2\}$$

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

$$P(a) = \begin{bmatrix} P_{11}(a) & P_{12}(a) \\ P_{21}(a) & P_{22}(a) \end{bmatrix} = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} \\ \frac{3}{4} & \frac{1}{4} \end{bmatrix}$$

$$P(b) = \begin{bmatrix} P_{11}(b) & P_{12}(b) \\ P_{21}(b) & P_{22}(b) \end{bmatrix} = \begin{bmatrix} \frac{1}{4} & \frac{3}{4} \\ \frac{1}{4} & \frac{3}{4} \end{bmatrix}$$

Costs:  $g(1, a) = 2, g(1, b) = 0.5,$

$$g(2, a) = 1, g(2, b) = 3$$

$$\lambda = 0.9$$

$$J_0 = (0, 0)$$

Do VI for 2 steps

$$J_1 = (0.5, 1)$$

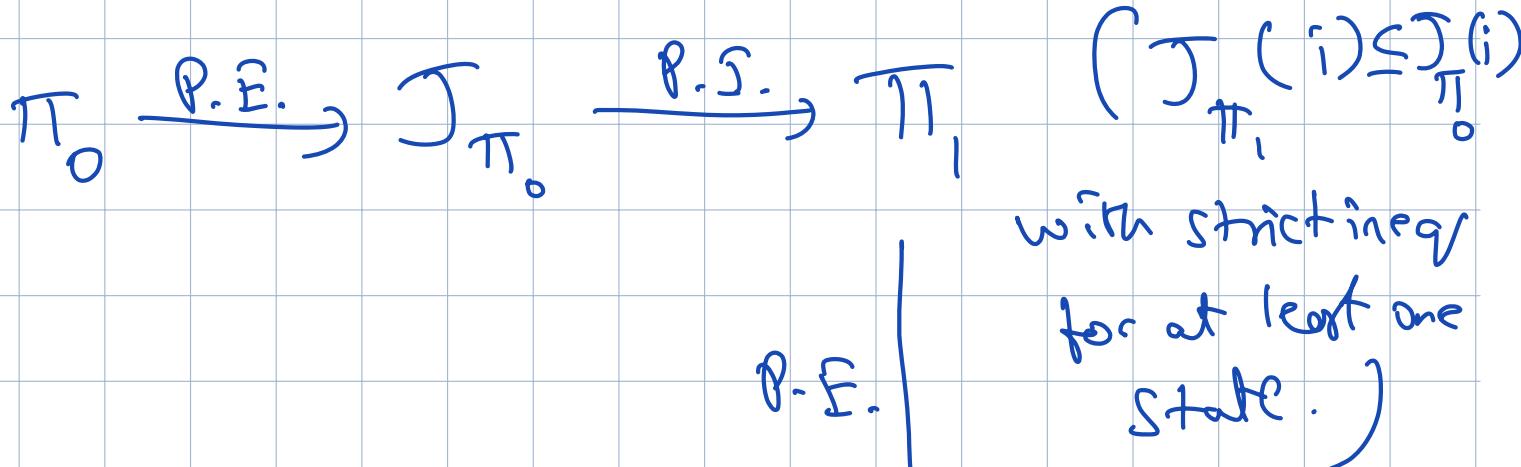
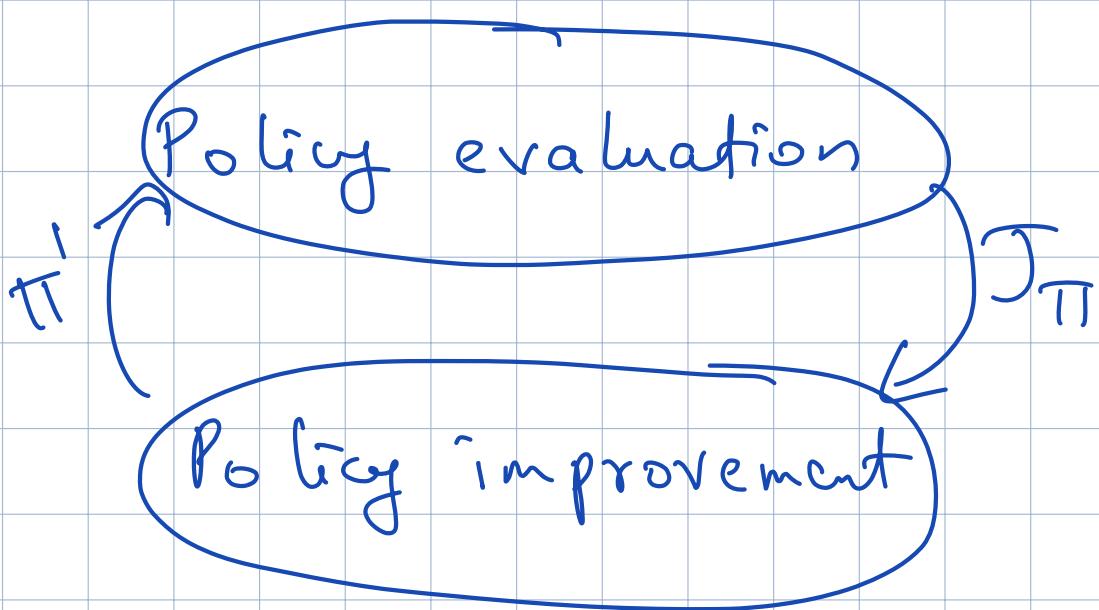
$$TJ(i) = \min \left( g(i, a) + \lambda \sum_{j=1}^2 P_{ij}^a J(j), \right)$$

$$g(i, b) + \gamma \sum_{j=1}^2 p_{ij}(b) J_j^{(t)}$$

$$J_2 = (1.287, 1.562)$$

### Lecture - 22

#### Policy iteration



$$\pi^* \dots \leftarrow \pi_2 \xleftarrow{P.I} J_{\pi_1}$$

Claim :- A policy  $\pi$  is optimal if

and only if  $\pi(i)$  attains the minimum

in the Bellman equation

$$g(i, \pi(i)) + \alpha \sum_{j=1}^n P_{ij}(\pi(i)) J^*(j) \\ = \min_a \left[ g(i, a) + \alpha \sum_{j=1}^n P_{ij}(a) J^*(j) \right]$$

(or)

$$J_\pi J^* = J J^*$$

Proof: Assume  $TJ^* = T_\pi J^*$  for some  $\pi$

We know  $J^* = TJ^*$

So,  $J^* = TJ^* = T_\pi J^*$

$\Rightarrow J^* = J_\pi$  & hence,  $\pi$  is optimal.

Converse:  $\pi$  is optimal

$\Rightarrow J^* = J_\pi$

$\Rightarrow J^* = T_\pi J^*$

$\Rightarrow TJ^* = T_\pi J^*$



Policy iteration algorithm

Step 1: Fix  $\pi_0$

Step 2: Find  $J_{\pi_K}$  corresponding to  $\pi_K$

↓

## Policy evaluation

Step 3:

Obtain  $\pi_{k+1}$  as

Policy improvement

$$T_{\pi_{k+1}} J_{\pi_k} = T J_{\pi_k}$$

If  $\pi_{k+1} = \pi_k$ , stop, else goto Step 2.

---

Claim: Let  $\pi, \pi'$  satisfy

$$T_{\pi'} J_{\pi} = T J_{\pi}$$

Then,  $J_{\pi'}(i) \leq J_{\pi}(i), \forall i$

& if  $\pi$  is not optimal, then

$J_{\pi'}(k) < J_{\pi}(k)$  for at least one state  $k$ .

# Illustration of policy iteration:

$$S = \{1, 2\}, A = \{a, b\}$$

$$P_a = \begin{bmatrix} 3/4 & 1/4 \\ 3/4 & 1/4 \end{bmatrix}$$

$$P_b = \begin{bmatrix} 1/4 & 3/4 \\ 1/4 & 3/4 \end{bmatrix}$$

$$g(1, a) = 2, g(1, b) = 0.5, g(2, a) = 1, g(2, b) = 3$$

$$\alpha = 0.9$$

$$\pi_0(1) = a, \pi_0(2) = b$$

$$\begin{aligned} J_{\pi_0}(1) = & g(1, a) + \alpha P_{11}(a) J_{\pi_0}(1) \\ & + \alpha P_{12}(a) J_{\pi_0}(2) \end{aligned}$$

$$\begin{aligned} J_{\pi_0}(2) = & g(2, b) + \alpha P_{21}(b) J_{\pi_0}(1) \\ & + \alpha P_{22}(b) J_{\pi_0}(2) \end{aligned}$$

$$J_{\pi_0}(1) = 2 + 0.9 \times \frac{3}{4} J_{\pi_0}(1) + 0.9 \times \frac{1}{4} J_{\pi_0}(2)$$

$$J_{\pi_0}(2) = 3 + 0.9 \times \frac{1}{4} J_{\pi_0}(1) + 0.9 \times \frac{3}{4} J_{\pi_0}(2)$$

$$J_{\pi_0}(1) = 24.12, \quad J_{\pi_0}(2) = 25.96$$

Policy Improvement:

$$T_{\pi_1} J_{\pi_0} = T J_{\pi_0}$$

$$(T J_{\pi_0})(1) = \min \text{ of}$$

$$(2 + 0.9 \left( \frac{3}{4} \times 24.12 + \frac{1}{4} \times 25.96 \right))$$

action  $a$

$$0.5 + 0.9 \left( \frac{1}{4} \times 24.12 + \frac{3}{4} \times 25.96 \right)$$

$$= \min (24.12, 23.45) = 23.45 (= \text{Action } b)$$

$$(\mathcal{J}_{\pi_0})(2) = \min \left( \begin{array}{l} 23.12 \\ 25.95 \end{array} \right)$$

$$\pi_1(1) = b, \quad \pi_1(2) = a$$

H.W:- Do one more step of policy iteration

i.e., Find  $\mathcal{J}_{\pi_1}(1), \mathcal{J}_{\pi_1}(2)$

& we that to perform  
policy improvement

You should get

$$\boxed{\pi_2 = \pi_1}$$

Sample path:

$$s_0 \xrightarrow{\pi(s_0)} s_1$$

$$s_1 \sim P_{s_0 s_1}(\pi(s_0))$$

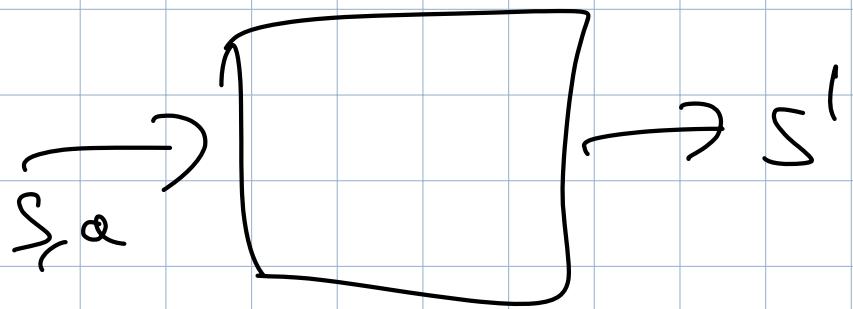
for e.g.  $S = \{1, 2, 3\}$

$$s_0 = 1$$

$$P_{11}(\text{action}) = 0$$

$$P_{12}(\text{action}) = 0.4$$

$$P_{13}(\text{---}) = 0.6$$



$(s_0, a_0, s_1, a_1, s_2, \dots)$

Sample path

Policy evaluation (Prediction)

$(s_0, \pi(s_0), s_1, \pi(s_1), \dots)$

Goal: Find  $\pi^*$

Policy control :-

$(s_0, a_0, s_1, \dots)$

# Reinforcement Learning

Mean estimation:

R.V.  $\mathcal{V}$  with mean  $\mu$  and finite variance.

$\{\mathcal{V}_1, \dots, \mathcal{V}_n\} \leftarrow$  iid samples from the distribution of  $\mathcal{V}$ .

$$\bar{\mathcal{V}}_n = \frac{1}{n} \sum_{i=1}^n \mathcal{V}_i$$

$$\bar{\mathcal{V}}_{n+1} = \frac{1}{n+1} \sum_{i=1}^{n+1} \mathcal{V}_i$$

$$= \frac{n}{n+1} \left( \frac{1}{n} \sum_{i=1}^n \mathcal{V}_i \right) + \frac{1}{n+1} \mathcal{V}_{n+1}$$

$$= \frac{n}{n+1} \bar{\mathcal{V}}_n + \frac{1}{n+1} \mathcal{V}_{n+1}$$

$$\boxed{\bar{\mathcal{V}}_{n+1} = \bar{\mathcal{V}}_n + \frac{1}{n+1} (\mathcal{V}_{n+1} - \bar{\mathcal{V}}_n)}$$

More generally,

$$r_{n+1} = r_n + \eta_n (\vartheta_{n+1} - \vartheta_n)$$

↗ step-size

If  $\sum_{n=1}^{\infty} \eta_n = \infty$ ,  $\sum_{n=1}^{\infty} \eta_n^2 < \infty$ , then

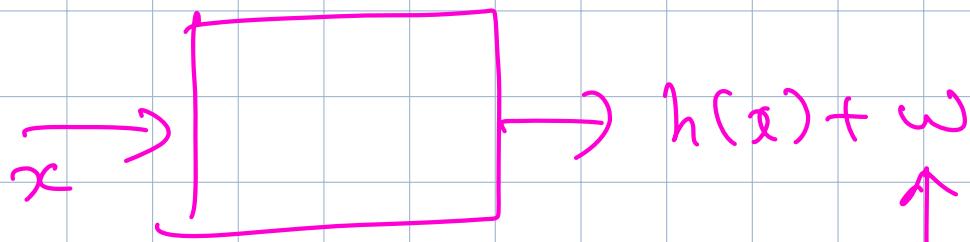
$$r_n \rightarrow \mu \text{ w.p. 1 as } n \rightarrow \infty.$$

Stochastic approximation

(Robbins Monro algorithm)

Suppose you want to solve  $h(x) = 0$ ,

$h$  is Lipschitz



Zero-mean noise  
e.g.  $w \sim N(0, 1)$

$$x_{n+1} = x_n + \eta_n (h(x_n) + w_n)$$

Under some conditions, it can be shown

that  $x_n \rightarrow x^*$  w.p. 1 as  $n \rightarrow \infty$ ,  
where  $h(x^*) = 0$

Two special cases:-

- ①  $\min_x f(x)$  Then  $h = f'$   
or in higher dimensions  
 $h = \nabla f$   
"gradient descent"
- ② Want to solve  $f(x) = x$  } fixed point iteration.  
Set  $h(x) = f(x) - x$

In a MDP context, one usually wants

to solve  $H\gamma = \gamma$

$$H: \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad r \in \mathbb{R}^n$$

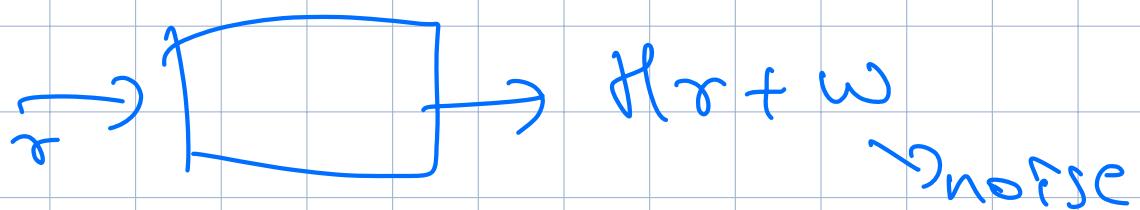
e.g.  $H = T_\pi$

A direct method:

$$x_{n+1} = Hx_n$$

$$(or) \quad x_{n+1} = (1 - \eta_n)x_n + \eta_n(Hx_n)$$

In a RL setting,



$$x_{n+1} = (1 - \eta_n)x_n + \eta_n(Hx_n + w_n)$$

Assuming  $H$  is "well-behaved", & noise  $w_n$  is zero-mean + bounded variance,

One can show that

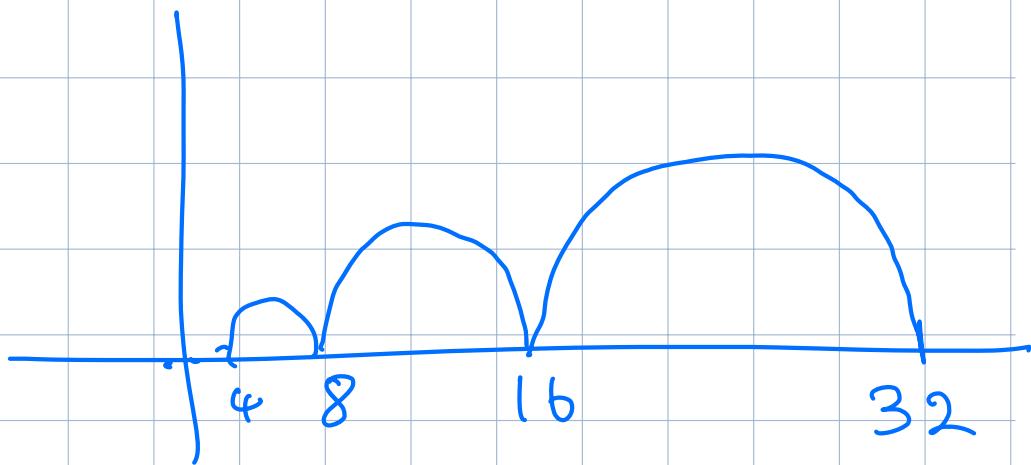
$$x_n \rightarrow x^* \text{ a.s. as } n \rightarrow \infty$$

where  $H_{x^*} = x^*$

---

A brief tour of contraction mappings.

$$f(x) = \frac{x}{2}$$



$$\lim_{n \rightarrow \infty} f^n(x) = 0$$

$$f(0) = 0$$

Def:  $f: \mathbb{R} \rightarrow \mathbb{R}$  is a contraction mapping

with modulus  $\beta \in (0,1)$  if

$$|f(x) - f(y)| \leq \beta |x - y|, \forall x, y \in \mathbb{R}$$

## Contraction properties of $T, T_\pi$

in a discounted MDP setting:

Max-norm: -  $\|\cdot\|_\infty$  on  $\mathbb{R}^n$  is

$$\|\boldsymbol{\gamma}\|_\infty = \max_{i=1 \dots n} |\gamma(i)|$$

Proposition: - For any two bounded

$\boldsymbol{\gamma}, \boldsymbol{\gamma}'$ , we have

$$\|T\boldsymbol{\gamma} - T\boldsymbol{\gamma}'\|_\infty \leq \alpha \|\boldsymbol{\gamma} - \boldsymbol{\gamma}'\|_\infty$$

↑  
Discount factor

$$\|T_\pi \boldsymbol{\gamma} - T_\pi \boldsymbol{\gamma}'\|_\infty \leq \alpha \|\boldsymbol{\gamma} - \boldsymbol{\gamma}'\|_\infty,$$

for any stationary  $\pi$ .

Corollary:

$$\|T^k \tau - T^k \tau' \|_{\infty} \leq \alpha^k \| \tau - \tau' \|_{\infty}$$

Corollary:-

$$\|T^k \tau_0 - \tau^* \|_{\infty} \leq \alpha^k \| \tau_0 - \tau^* \|_{\infty}$$

Lecture-23

Stochastic iterative algorithm

for finding fixed point of a

contraction mapping

Goal: Solve  $H\tau^* = \tau^*$

$H: \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $\tau \in \mathbb{R}^n$

Update rule:

$$r_{t+1}(i) = (1-\eta_t) r_t(i) + \eta_t ((H r_t)(i) + w_t(i))$$

↗  
noise

Assumptions:

$$\mathcal{F}_t = \{r_0(i), \dots, r_t(i), w_0(i), \dots, w_{t-1}(i), \dots\}$$

$T = t - n$

(history)

(A1)  $E(w_t(i) | \mathcal{F}_t) = 0$

$$E(w_t^2(i) | \mathcal{F}_t) \leq A + \beta \|r_t\|^2$$

(A2)  $\sum_{t=1}^{\infty} \eta_t = \infty$ ,  $\sum_{t=1}^{\infty} \eta_t^2 < \infty$ .

(A3)

H is a max-norm

pseudo-contraction

$$\|Hr - r^*\|_\infty \leq \alpha \|r - r^*\|, \forall r$$

Fuel contraction

$$\|Hr - Hr'\|_\infty \leq \alpha \|r - r'\|_\infty$$

$\forall r, r'$

Under A1-A3,

$$r_t \rightarrow r^* \text{ w.p. 1 as } n \rightarrow \infty$$

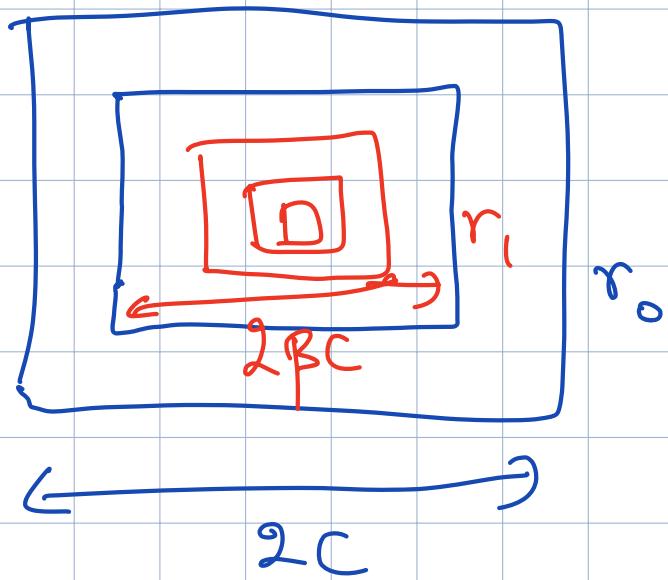
$$\text{where } Hr^* = r^*$$

An intuitive justification:-

$$Ho = 0, \quad r^* = 0$$

$$|r_0(i)| \leq C$$

$$r_{t+1}(i) = (H r_t)(i)$$



$$r_{t+1}(i) = (I - H_t) r_t(i)$$

$$+ \gamma_t (H r_t(i) + w_t(i))$$

Similar behavior, but shrinkage  
is not immediate

Temporal-difference (TD) learning for policy evaluation :-

Fix a policy  $\pi$ .

Single stage wt:  $g(i, i)$

Want to estimate:

$$J_{\pi}(i) = E \left[ \sum_{m=0}^{\infty} \gamma^m g(i_m, i_{m+1}) \mid i_0 = i \right]$$

Bellman equation:-

$$J_{\pi} = T_{\pi} J_{\pi}$$

$$J_{\pi}(i) = E(g(i, \bar{i}) + \gamma J_{\pi}(\bar{i}))$$

TD(0) algorithm:-

$$J_{t+1}(i) = (1 - \eta_t) J_t(i) + \eta_t (g(i, \bar{i}) + \gamma J_t(\bar{i}))$$

$\bar{i} \sim P_{ij}(\pi(i))$

$$\mathcal{T}_{t+\Delta}(i) = \mathcal{T}_t(i) + \eta_t \left( \underbrace{g(i, \bar{i}) + \alpha \mathcal{T}_t(\bar{i}) - \mathcal{T}_t(i)}_{\text{Proxy for}} \right)$$

$$E(g(i, \bar{i}) + \alpha \mathcal{T}_t(\bar{i})) \\ = (\mathcal{T}_\pi \mathcal{T}_t)(i)$$

$$\mathcal{T}_{t+\Delta}(i) = \mathcal{T}_t(i) + \eta_t \left( (\mathcal{T}_\pi \mathcal{T}_t)(i) - \mathcal{T}_t(i) \right. \\ \left. + \underbrace{g(i, \bar{i}) + \alpha \mathcal{T}_t(\bar{i}) - (\mathcal{T}_\pi \mathcal{T}_t)(i)}_{w_t(i)} \right)$$

With result for  $r_t$  above, we can infer

$$\mathcal{T}_t \rightarrow \mathcal{T}_\pi \text{ w.p. 1 as } t \rightarrow \infty$$

$$J_{\pi}(i) = E(g(i, \bar{i}) + \alpha J_{\pi}(\bar{i}))$$

$$= E\left(g(i, \bar{i}) + \underbrace{\alpha g(\bar{i}, \bar{i})}_{\text{next state}} + \alpha^2 J_{\pi}(\bar{\bar{i}})\right)$$

+  
next state  
+  
next-to-next  
state

$T_D(\lambda), \lambda \in [0, 1]$  :

$$\overline{J_{\pi}(i)} = E\left(\sum_{m=0}^l \alpha^m g(i_m, i_{m+1}) + \alpha^{l+1} \overline{J_{\pi}(i_{l+1})}\right)$$

Fix  $\lambda < 1$ , and multiply  $(1-\lambda) \lambda^l$  & sum over  $\lambda$  :

$$J_{\pi}(i) = (1-\lambda) E \left( \sum_{l=0}^{\infty} \lambda^l \left( \sum_{m=0}^{\infty} \lambda^m g(i_m, i_{m+1}) + \lambda^{l+1} J_{\pi}(i_{l+1}) \right) \right)$$

Interchange the two summations & we  $(1-\lambda) \sum_{l=m}^{\infty} \lambda^l = \lambda^m$ , to obtain

$$\begin{aligned} J_{\pi}(i) &= E \left[ (1-\lambda) \sum_{m=0}^{\infty} \lambda^m g(i_m, i_{m+1}) \sum_{l=m}^{\infty} \lambda^l + \sum_{l=0}^{\infty} \lambda^{l+1} J_{\pi}(i_{l+1}) (\lambda^l - \lambda^m) \right] \\ &= E \left( \sum_{m=0}^{\infty} \lambda^m \left( \lambda^m g(i_m, i_{m+1}) + \lambda^{m+1} J_{\pi}(i_{m+1}) - \lambda^m J_{\pi}(i_m) \right) \right. \\ &\quad \left. + J_{\pi}(i) \right) \end{aligned}$$

Letting  $d_m = g(i_m, i_{m+1}) + \lambda J_{\pi}(i_{m+1}) - \lambda^m J_{\pi}(i_m)$

$$J_{\pi}(i) = E \left( \sum_{m=0}^{\infty} (\lambda \lambda)^m d_m \right) + J_{\pi}(i)$$

TD(λ) update rule:-

$$J_{t+1}(i) = J_t(i) + \gamma_t \left( \sum_{m=0}^{\infty} (\alpha \lambda)^m d_m \right)$$

for  $\lambda=0$ ,  $d_0 = g(i, \hat{i}) + \alpha J_\pi(\hat{i}) - J_\pi(i)$   
& we recover TD(0)

### Stochastic Shortest path

$$\mathcal{S} = \{1, \dots, n, T\}$$

$$p_{TT}(a) = 1$$

$$g(T, a, T) = 0$$

$$J_\pi(i) = E \left( \sum_{t=0}^{\infty} g(i_t, \pi(i_t), i_{t+1}) \mid i_0 = i \right)$$

$$(J_\pi(i) = E \left[ g(i, \pi(i), i) + J_\pi(i) \right])$$

$\pi$ : proper if it ensures that terminal state  $T$  is reached with positive probability.

$$\text{TD}(0): J_{t+1}(i) = J_t(i) + \eta_t \underbrace{\left( g(i, \pi(i), i) + J_t(i) - J_t(i) \right)}_{\text{temporal-difference}}$$

for  $i = 1, \dots, n$

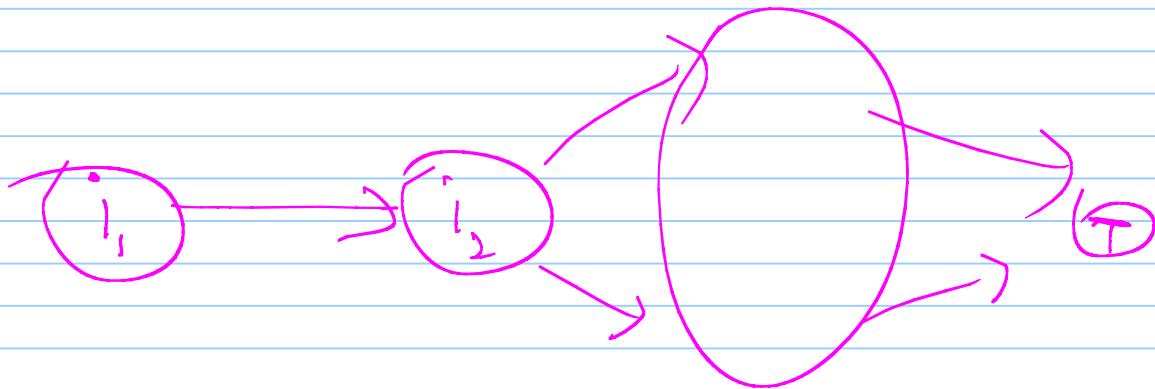
TD(i) :- Simulate episodes of the underlying SSP w.r.t  $\pi$ .

$$(i_0^m, i_1^m, \dots, i_N^m), i_N = T, i_0 = i$$

$$\hat{J}(m) = g(i_0^m, \pi(i_0^m), i_1^m) + \dots + g(i_{N-1}^m, \pi(i_N^m), i_N^m)$$

$$\tilde{J}(i) = \frac{1}{M} \sum_{m=1}^M \hat{J}(m), \quad M = \# \text{ of episodes simulated}$$

TD(0) vs TD(1) :- A toy example



TD(1) :- Start in  $i_1$ , simulate episode, we total cost  $\sum \hat{J}(m)$ ,  
to estimate  $J_{\pi}(i_1)$

TD(0) :- Start in  $i_1$ , simulate one transition and  
we "  $g(i_1, i_2) + J(i_2) - J(i_1)$ " to estimate  $J_{\pi}(i_1)$

## Lecture-24

Q-learning: SSP setting

Let  $J^*$  be the optimal cost-to-go vector

$$Q^*(i, a) = \sum_{j=t+1}^{nT} p_{ij}(a) (g(i, a, j) + J^*(j)), \quad t=1 \dots n$$

Bellman's equation:

$$J^*(i) = \min_a \sum_{j=t+1}^{nT} p_{ij}(a) (g(i, a, j) + J^*(j))$$

\$J^\*(i) = \min\_a Q^\*(i, a)\$

Q-Bellman equation:-

$$Q^*(i, a) = \sum_{j=t+1}^{nT} p_{ij}(a) (g(i, a, j) + \min_b Q^*(j, b))$$

$\overline{Q^*(i, a)} = E(g(i, a, j) + \min_b Q^*(j, b))$

$$Q_{t+1}(i, a) = Q_t(i, a) + \gamma_t \left( g(i, a, \bar{i}) + \min_b Q_t(\bar{i}, b) - Q_t(i, a) \right)$$

Define the operator  $\nu$

$$(\nu Q)(i, a) = \sum_{j=1, \dots, N, \bar{I}} p_{ij}(a) \left( g(i, a, j) + \min_b Q(j, b) \right)$$

Under some conditions,  $\nu Q$  is a contraction mapping.

$$Q_{t+1}(i, a) = Q_t(i, a) + \gamma_t \left( (\nu Q_t)(i, a) + w_t(i, a) \right)$$

$$w_t(i, a) = g(i, a, \bar{i}) + \min_b Q_t(\bar{i}, b) - \sum_{j=1, \dots, N, \bar{I}} p_{ij}(a) \left[ g(i, a, j) + \min_b Q_t(j, b) \right]$$

$Q_t \rightarrow Q^*$  w.p.l at  $t \rightarrow \infty$  under conditions like (A1)-(A3)

## Issue of exploration:-

for Q-learning to converge, all state-action pairs  $(i, a)$  have to be visited infinitely often.

Greedy policy:  $\pi_{t+1}(i) = \min_a Q_t(i, a)$

$\epsilon$ -Greedy policy:-  $\pi_{t+1}(i) = \begin{cases} \text{greedy action w.p. } (1-\epsilon) \\ \text{random action w.p. } \epsilon \end{cases}$

$\epsilon \rightarrow$  small number e.g. 0.05

## Linear function approximation

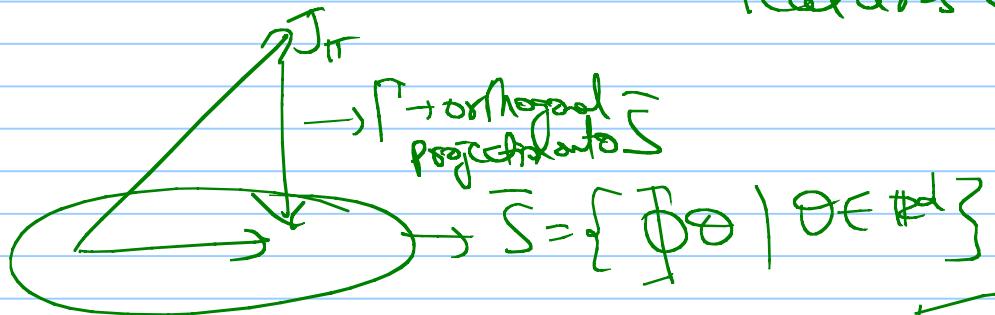
Fix a policy  $\pi$ .

$$J_\pi(i) \approx \phi(i)^T \theta$$

↓

features  $\in \mathbb{R}^d$

$d \ll |\mathcal{S}|$



Cannot solve

Approximate solution

$$\boxed{J_\pi = T_\pi J_\pi}$$

$$\Phi \theta = \Gamma(T_\pi \Phi \theta)$$

for a nice solution  
since  $\Gamma T_\pi$  is a  
contraction.

$$J_\pi \approx \Phi \theta$$

$$\theta_{t+1} = \theta_t + \eta_t (g(i, \pi(i), \bar{i}) + \phi(\bar{i})^T \theta_t - \phi(i)^T \theta_t) \phi(i)$$

TD(0) with linear function approximation

Q-learning with linear function approximation

$$Q(i, a) \approx \phi(i, a)^T \theta$$

$$\begin{aligned} \theta_{t+1} &= \theta_t + \eta_t (g(i, a, \bar{i}) + \min_b \phi(i, b)^T \theta_t - \phi(i, a)^T \theta_t \\ &\quad \times \phi(i, a)) \end{aligned}$$

Greedy action

$$\min_a \theta_t^T \phi(i, a)$$

## Policy gradient methods (PG methods)

The main idea behind PG methods is

the "likelihood ratio" trick (aka score function)

An illustration of this trick in a very simple setting.

Let  $X$  be a <sup>discrete</sup> r.v. with mass function

$$p(\theta, \cdot)$$

↗ parameter

i.e.,  $P(X=x)$  is parameterized by  $\theta$ .

$$\mathcal{J}(\theta) = E f(x)$$

Goal:

$$\min_{\theta} \mathcal{J}(\theta)$$

e.g.  $\theta \in \mathbb{R}^d$

min among a class of  
parameterized r.v.s

Want to find  $\theta$  but using a gradient method

$$\theta_{t+1} = \theta_t - \beta_t \nabla \mathcal{J}(\theta_t) \quad \leftarrow \text{Gradient descent}$$

Need: " $\nabla \mathcal{J}(\theta)$ "

Use "likelihood ratio" trick, which is shown below.

$$\mathcal{J}(\theta) = \sum_x f(x) p(\theta, x) \quad \leftarrow \text{LOTUS}$$

$$\nabla \mathcal{J}(\theta) = \nabla_{\theta} \left( \sum_x f(x) p(\theta, x) \right)$$

need a few conditions  
to justify interchange of  $\sum$  &

$$= \sum_x f(x) \nabla p(\theta, x)$$

mild regularity  
conditions  $\rightarrow$  usually 1st  
one invoke "dominated"

Convergence theorem

$$\nabla J(\theta) = \sum_x f(x) \nabla p(\theta, x)$$

$$\boxed{\nabla J(\theta) = \sum_x \left( f(x) \frac{\nabla p(\theta, x)}{p(\theta, x)} \right) p(\theta, x)}$$

→ RMS is an expectation

$$\nabla J(\theta) = E \left( f \frac{\nabla p}{p} \right) = E(f \nabla \log p)$$

To get an estimate of  $\nabla J(\theta)$ , I can sample from  $\nabla \log p(\theta)$

$$\frac{\nabla p(\theta, x)}{p(\theta, x)} \rightarrow \text{likelihood ratio.}$$

$$\nabla \log p(\theta) = \frac{\nabla p(\theta)}{p(\theta)}$$

Connecting likelihood ratio to RL:

Fix an SSP problem with state space  $\mathcal{X}$ , action space  $\mathcal{A}$ .

$\overline{\Pi}_{\text{det}}$  : class of admissible stationary deterministic policies

$\{ \pi : \pi : \mathcal{X} \rightarrow \mathcal{A} \text{ & it is timeinvariant} \}$

$\xrightarrow{x} \boxed{\text{unparameterized policy}} \rightarrow \pi(x) \rightarrow \text{an action } \pi : \boxed{\begin{matrix} \text{states} \\ \hline \text{actions} \end{matrix}}$

For PGI method, we consider stationary randomized policies, i.e.,

$$\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$$

→ set of all distributions over the actions.  
For simplicity, assume all actions available in all states.

e.g.

$$\pi_{\theta}(x, a) = \frac{\exp(h(\theta, x, a))}{\sum_b \exp(h(\theta, x, b))}$$

*is the probability of choosing action  $a$  in state  $x$*

randomized policy parameterized by " $\theta$ "

$$\pi_{\theta}(x) = [\pi_{\theta}(x, a), \forall a \in \mathcal{A}]$$

distribution over  $\mathcal{A}$ .

$$\pi_{\theta}(x, a) = \frac{\exp(\theta^T \phi(x, a))}{\sum_b \exp(\theta^T \phi(x, b))}$$

$\theta \in \mathbb{R}^d$  vector of size  $d$  (status/ $x$ /action)

Simple example for  $h$ :

$$h(\theta, x, a) = \theta^T \phi(x, a)$$

State-action features

$$\pi_{\theta}(x, a) = \frac{\exp(\theta^T \phi(x, a))}{\sum_b \exp(\theta^T \phi(x, b))}$$

$\theta \in \mathbb{R}^d, \phi(x, \cdot) \in \mathbb{R}^d$

Boltzmann distribution aka Soft-max

Assumption A1: Policy  $\pi_{\theta}$  is a continuously differentiable function of  $\theta$

$\Leftrightarrow \log \pi_{\theta}$  exists.

Note! Every  $\pi_{\theta}$  is identified by its parameter  $\theta \in \mathbb{R}^d$

Goal:

$$\min_{\theta \in \Theta} J_{\pi_{\theta}}(x^0)$$

Find an approximately optimal policy in the class of parameterized policies

& we want to find the best parameter in a class

$$\{ \pi_{\theta} \mid \theta \in \Theta \}$$

e.g.  $\Theta \subset \mathbb{R}^d$

Want to find a  $\theta^* \in \arg \min_{\theta \in \Theta} J_{\pi_{\theta}}(x^0)$

"-" is necessary a convex function of  $\theta$

A formula for policy gradient in an SGP

SSP:  $\mathcal{O}$  is a special -lost-free absorbing state.

Total Cost Objective:

$$J_{\pi_\theta}(x^0) = E \left( \sum_{k=0}^{\infty} g(x_k, a_k, x_{k+1}) \mid x_0 = x^0 \right)$$

$a_k \sim \pi_\theta(x_k, \cdot)$

A trajectory (aka episode) is

$$\gamma = x_0, a_0, x_1, a_1, \dots, x_T, \quad x_T = \mathcal{O}, \text{ the terminal state.}$$

$\uparrow$   
fixed

Probability of seeing a trajectory  $\gamma$  is

$$\pi(x_0, a_0) P(x_1 | x_0, a_0) \cdots P(x_T | x_{T-1}, a_{T-1})$$

$J_{\pi_\theta} \rightarrow$  can be seen as an average of the total cost over episodes.

Let  $D(\gamma)$  denote the total cost r.v.

$$D(\gamma) = \sum_{k=0}^{T-1} g(x_k, a_k)$$

Here  $\tau = (x_0, a_0, \dots, x_{T-1}, a_{T-1}, x_T)$   
in the episode

$$\mathcal{J}_{\pi_\theta}(x^0) = E(D(\tau))$$

↓  
expectation over episodes.

$$\nabla_{\theta} \mathcal{J}_{\pi_\theta}(x^0) = \nabla_{\theta} \sum_{\tau} D(\tau) \text{Prob}(\tau \text{ under } \pi_\theta)$$

$$D(\tau) \text{ is not a function of } \theta \text{ & } \Sigma$$

→ intermediate for backpropagation  
↓  
subtler for finite state spaces.

$$= \sum_{\tau} D(\tau) \nabla_{\theta} \text{Prob}(\tau \text{ under } \pi_\theta)$$

$$= \sum_{\tau} D(\tau) \text{Prob}(\tau \text{ under } \pi_\theta) \nabla \log(\text{Prob}(\tau \text{ under } \pi_\theta))$$

$$= \sum_{\tau} D(\tau) \text{Prob}(\tau \text{ under } \pi_\theta)$$

$$\times \nabla \log \left\{ \pi_\theta(x_0, a_0) P(x_1 | x_0, a_0) \dots \pi_\theta(x_{T-1}, a_{T-1}) P(x_T | x_{T-1}, a_{T-1}) \right\}$$

why??

$$= \sum_{\tau} D(\tau) \text{Prob}(\tau \text{ under } \pi_\theta) \sum_{k=0}^{T-1} \nabla \log \pi_\theta(x_k, a_k)$$

Main claim! PTO

$$\nabla \mathcal{J}_{\pi_\theta}(x^0) = E \left[ D(\tau) \sum_{k=0}^{T-1} \nabla \log \pi_\theta(x_k, a_k) \right]$$

expectation over trajectories.

An unbiased estimate of  $\nabla \mathcal{J}_{\pi_\theta}(x^0)$

given an episode  $(x^0, a_0, \dots, x_T)$

Calculate total cost of this episode, say  $\hat{D}$

Calculate likelihood ratio  $\hat{\psi} = \sum_{k=0}^{T-1} \nabla \log \pi_\theta(x_k, a_k)$

$$\hat{\nabla} \mathcal{J} = \hat{D} \times \hat{\psi}.$$

$$\theta_{t+1} = \theta_t - \alpha(t) \hat{\nabla} \mathcal{J}$$



policy gradient update

"REINFORCE". (Williams, 1992)

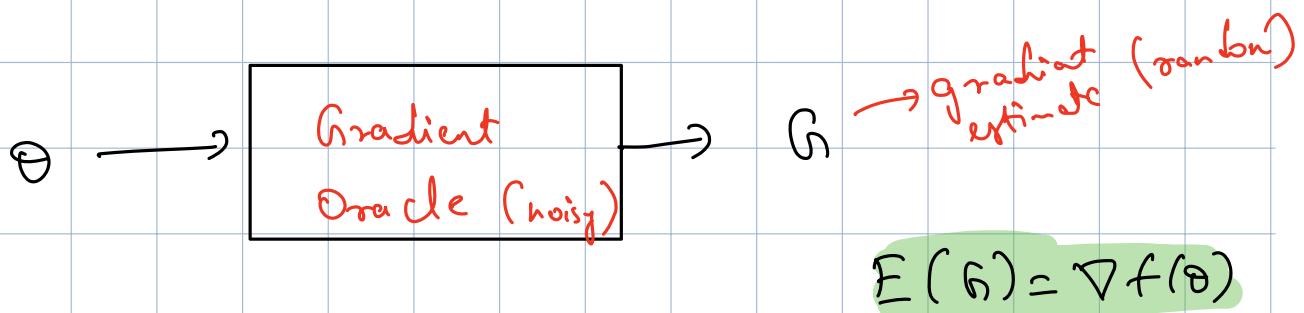
# Lecture - 18

A detour  $\rightarrow$  Non-asymptotic analysis of  
stochastic gradient algorithms

Aim: Solve  $\min_{\theta} f(\theta)$

$\curvearrowright$  smooth

Setting: Unbiased gradient information



SG algorithm:  $\theta_{k+1} = \theta_k - \alpha(k) G(\theta_k, \xi_k)$  ①

Assumption:

(A)  $E_{\xi_k}(G(\theta_k, \xi_k)) = \nabla f(\theta_k), \forall k \geq 1$

$$E_{\xi_k} \|G(\theta_k, \xi_k) - \nabla f(\theta_k)\|^2 \leq \sigma^2$$

for some  $\sigma > 0$ .

(A0)  $f$  is  $L$ -smooth

$$\|\nabla f(\theta_1) - \nabla f(\theta_2)\| \leq L \|\theta_1 - \theta_2\|, \forall \theta_1, \theta_2 \in \mathbb{R}^d$$

Asymptotic Convergence:

(A2)

$$\sum_k \alpha(k) = \infty, \quad \sum_k \alpha(k)^2 < \infty$$

e.g.  $\alpha(k) = \frac{c}{k} \rightarrow$  standard stochastic approximation conditions.

Under (A0) - (A2), the parameter  $\theta_k$  governed by ① converges almost surely to the set  $\{\theta^* \mid \nabla f(\theta^*) = 0\}$

Tool used for proving this claim: Kushner-Clark Lemma

Eq ① ( $\Rightarrow$ )

$$\theta_{k+1} = \theta_k - \alpha(k) (\nabla f(\theta_k) + \xi_k)$$

$$\xi_k = G(\theta_k, \xi_k) - \nabla f(\theta_k)$$

We want to understand the non-asymptotic performance of the Sh algorithm above.

Stationary point, say  $\theta^*$ , has  $\nabla f(\theta^*) = 0$

$\epsilon$ -stationary point:  $\bar{\theta}$  is an  $\epsilon$ -stationary point if  $E \|\nabla f(\bar{\theta})\| \leq \epsilon$ , where the expectation is over the randomness in the SG algorithm.

[Ghadimi-Lan, 2014, SIAM J. Opt, "Stochastic first/zeroth order"]

"Randomized Stochastic gradient" (RSG)

Suppose we run ① for  $N$  iterations.

$\{\theta_1, \theta_2, \theta_3, \dots, \theta_N\} \xrightarrow{\text{(random)}} \text{Set of}$   
iterations visited by SG algorithm.

RSG will pick a random iterate as follows:

$\theta_R$  is picked uniformly at random from  $\{\theta_1, \dots, \theta_N\}$

i.e.,  $P(\theta_R = \theta_i) = \frac{1}{N}$  + i

$$E \theta_R = \frac{1}{N} \sum_{i=1}^N \theta_i \quad \begin{cases} \text{average} \\ \text{iterate.} \end{cases}$$

The non-asymptotic bound for RSG is of the form

$$E \| \nabla f(\theta_R) \|^2 \leq \frac{\text{const}}{\sqrt{N}}$$

under suitable choice of step-size.

FL connection:  $\theta \rightarrow$  policy parameter

Objective $f \rightarrow$	$J_{\pi_\theta}(x^*)$ or $J(\theta)$
	$\downarrow$
	Average cost

$\uparrow$

Value function  
with state  $x^*$   
in a discounted/SSP

Proof of the non-asymptotic bound for RSG:

First  $\rightarrow$  derive a bound for a general step-size

Second  $\rightarrow$  specialize the bound above with a particular step size.

$$\alpha_i \sim \pi_\theta(\cdot | s_i) \quad \xrightarrow{\text{terminal}} \quad \tau = (x_0, a_0, x_1, a_1, \dots, x_T)$$

$$\nabla J_{\pi_\theta}(x^0) = E \left[ D(\tau) \sum_{k=0}^{T-1} \nabla \log \pi_\theta(x_k, a_k) \right]$$

expectation over trajectories  $\tau$

"Each  $\tau$  ends in a terminal state."

Random length trajectories.

but length bounded "w.p. 1"  
is proper".

$\theta$  - parameter

$$J(\theta) = J_{\pi_\theta}(x^0)$$

fixed start state

For RSG analysis/bound to be applicable,  
we need to show

$$(1) \quad \| \nabla J(\theta_1) - \nabla J(\theta_2) \| \leq 2 \| \theta_1 - \theta_2 \|$$

" $J$  is  $L$ -smooth".

(2) Gradient estimate from a sample trajectory  $\tau$ .

$$\hat{h} = D(\tau) \sum_{k=0}^{T-1} \nabla \log \pi_\theta(x_k, a_k)$$

$$E(\hat{h}) = \nabla J(\theta) \leftarrow \begin{matrix} \text{unbiased gradient} \\ \text{estimate.} \end{matrix}$$

To verify ① & ② in an RL context,  
we make the following assumptions:

(B1) State & action spaces are finite

(B2)  $\forall \theta$ ,  $\pi_\theta$  is proper.

$\Leftrightarrow$  "in a finite, say  $M$  steps, there is a positive prob.  
of reaching the terminal state".

(B3)  $\|\nabla \log \pi_\theta(x, a)\| \leq G$

$\|\nabla^2 \log \pi_\theta(x, a)\| \leq H$

$$\left\| \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log \pi_\theta(x, a) \right\| \leq H \quad i, j = 1, 2, \dots, n$$

"Showing  $J$  is  $L$ -smooth or a function of  $\theta$ "

Use this identity: for any twice diff'ble fn  $f$ :

$$\nabla^2 f(\theta) = f(\theta) \left( \nabla^2 \log f(\theta) + \nabla \log f(\theta) \nabla \log f(\theta)^T \right)$$

(Scalar Case:  $f''(\theta) = f(\theta) \left( \frac{d^2}{d\theta^2} (\log f(\theta)) + \left( \frac{d}{d\theta} \log f(\theta) \right)^2 \right)$ )

$$\mathcal{J}(\theta) = \sum_{\tau} \text{Prob}(\tau \text{ under } \pi_\theta) \underbrace{D(\tau)}_{\substack{\text{Total cost} \\ \text{from } \tau.}}$$

$$= \sum_{\tau} p_\theta(\tau) D(\tau)$$

$$\nabla^2 \mathcal{J}(\theta) = \sum_{\tau} p_\theta(\tau) \nabla \log p_\theta(\tau) \nabla \log p_\theta(\tau)^T D(\tau) + \sum_{\tau} p_\theta(\tau) \nabla^2 \log p_\theta(\tau) D(\tau)$$

(\*)

Proper policies (see  $\mathbb{B}_2$ )

Assume: Trajectory length is finite w.p. 1

Let  $T$  be a r.v. denoting the length of a random trajectory  $\tau$ .

Then  $|T| < C_1 < \infty$  w.p. 1.

Fact 2: Finite state-action spaces ( $\mathbb{B}_1$ )

$$\Rightarrow \sup_{x,a} |g(x,a)| < C_2 < \infty.$$

Fact 1 & 2  $\Rightarrow$

$$D(\tau) = \sum_{k=0}^{T-1} g(x_k, a_k)$$

$$|D(\tau)| \leq C_2 C_1$$

Using (B3) ( $\|\nabla \log \pi\| \leq h$ ,  $\|\nabla^2 \log \pi(\cdot | \cdot)\| \leq H$ )

in (\*),

$$\|\nabla^2 J(\theta)\| \leq C_1 C_2 \left( \sum_{\tau} P_\theta(\tau) \|\nabla \log P_\theta(\tau) \nabla \log P_\theta(\tau)^T\| \right.$$

$$\left. + \sum_{\tau} P_\theta(\tau) \|\nabla^2 \log P_\theta(\tau)\| \right)$$

$$\leq C_1 C_2 (H^2 C_1 + H C_1),$$

Since

$$\tau = (x_0, a_0, \dots, x_T)$$

$$P_\theta(\tau) = \Pi_\theta(x_0, a_0) P(x_1 | x_0, a_0) \dots P(x_T | x_{T-1}, a_{T-1})$$

$$\nabla \log P_\theta(\tau) = \sum_{k=0}^{T-1} \nabla \log \Pi_\theta(x_k, a_k) \leq C C_1$$

So,  $\mathcal{T}(\theta)$  is  $L$ -smooth, where

$$L = C_1 C_2 (G^2 C_1 + H C_1)$$

Homework:

$$E \parallel \hat{h} - E \hat{h} \parallel^2$$

$$= E \parallel \hat{h} - \nabla \mathcal{T}(\theta) \parallel^2$$

$$\leq E \parallel \hat{h} \parallel^2$$

$$\hat{h} = D(\tau) \sum_{k=0}^{T-1} \nabla \log \pi_{\theta}(x_k, a_k)$$

$$\leq C_1^2 C_2^2 C_1 E \sum_{k=0}^{T-1} \parallel \nabla \log \pi(x_k, a_k) \parallel^2$$

$$\leq C_1^3 C_2^2 G^2 C_1$$

So, variance of the gradient estimate is bounded.

$\Rightarrow$  RSG analysis is applicable, leading to

$$E \parallel \nabla \mathcal{T}(\theta_R) \parallel^2 \leq \frac{\text{const}}{\sqrt{N}},$$

where  $\theta_k$  chosen unif. at. random from  
 $\{\theta_1, \dots, \theta_N\}$

$$\& \quad \theta_{k+1} = \theta_k - \alpha(k) (\hat{g}_k)$$

↓  
single trajectory estimate  
of  $\nabla J(\theta_k)$