

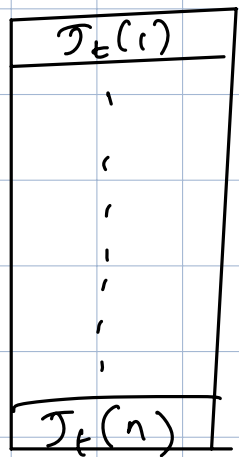
## TD-learning, Q-learning with full state representations

ref: Chapter 5 of NDP book

Full state-representations?

Suppose we want to estimate  $J_{\pi}(i)$ , for a given policy  $\pi$ , & for any  $i \in \mathcal{X}$ .

A full-state algorithm aka **tabular** case, maintains an estimate  $J_t(i)$ ,  $\forall i \in \mathcal{X}$  & updates this estimate incrementally using samples.



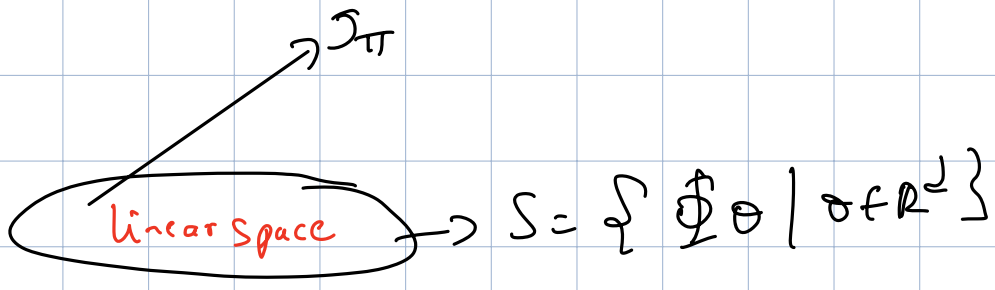
we sample & update each entry

**Problem with this approach:** It doesn't scale. e.g. on large state spaced MDPs it may be computationally infeasible. e.g. consider a MDP with  $10^{30}$  states or game Go has  $10^{170}$  states (rough approx)

On such problem, we require parametric approximation of  $J_{\pi}$ .

e.g.  $J_{\pi}(i) \approx \Theta^T \phi(i)$   $\Theta \in \mathbb{R}^d$   
 $d \leq n$ .

$\nearrow$  parameter  $\nearrow$  feature vector  
 linear function approximation



### Recall: Mean estimation

Want to estimate  $\mu = EX$  & you keep getting samples  $X_1, X_2, \dots$  increment

$$r_{m+1} = r_m + \beta_m (X_{m+1} - r_m)$$

$\nearrow$  estimate of  $\mu$  from  $\{X_1, \dots, X_m\}$     
  $\nearrow$  stepsize e.g.  $\frac{1}{m}$     
  $\nearrow$  sample from distribution of  $X$     
  $\nearrow$  prev. estimate of  $\mu$

### Question of policy evaluation!

Want to estimate!  $J_{\pi}(i) = E \left( \overbrace{\sum_{n=0}^{\infty} g(i_n, i_{n+1})}^{\text{Total cost } r.v.} \mid i_0=i \right)$

Assume a SSP context.

Fixed point equation for  $J_\pi$ :

$$J_\pi = T_\pi J_\pi$$

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

↑  
expectation over next state  $\bar{i}$

The action  $\pi(i)$  is implicit.

or  $E(g(i, \bar{i}) + J_\pi(\bar{i})) - J_\pi(i) = 0$

Idea: sample " $g(i, \bar{i}) + J_\pi(\bar{i})$ " & update incrementally

"Temporal difference"

$$(*) \quad J_{m+1}(i) = J_m(i) + \beta_m (g(i, \bar{i}) + J_m(\bar{i}) - J_m(i))$$

↗ This is the TD(0) algorithm  
↘ temporal difference

sample of the r.v. in expectation highlighted above

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

So, " $g(i, \bar{i}) + J_m(\bar{i})$ " is a proxy for

$$E(g(i, \bar{i}) + J_m(\bar{i})) = (T_\pi J_m)(i)$$

(\*) is equivalent to

$$J_{m+1}(i) = J_m(i) + \beta_m ( (T_\pi J_m)(i) - J_m(i) )$$

$$+ g(i, \bar{i}) + J_m(\bar{i}) - (T_\pi J_m)(i)$$

This is the noise factor

$w_m(i)$  from the general sto. itr. also.

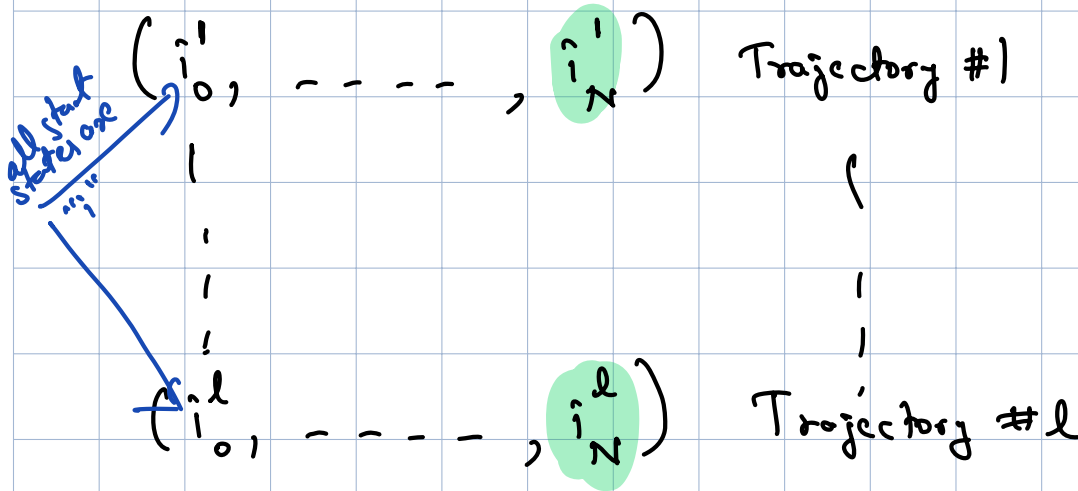
For simplicity, we have removed the dependence on action in single stage cost

↗ This is the update  $H$  in Sto. itr. (2) from part of chapter

# Monte Carlo Policy Evaluation

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

Fix proper policy  $\pi$ . Simulate the SSP to generate  $l$ -trajectories



$N$  is random  
 $i_0^1 \dots i_N^l = T$   
 Actions using  $\pi$ .

Want to estimate:  $J_\pi(i) = E \left( \sum_{n=0}^{\infty} g(i_n, i_{n+1}) \mid i_0 = i \right)$

$\uparrow$   
 This is the r.v. whose expectation we estimate.

Let

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

for  $k=1, \dots, l$

$\hat{J}_k \rightarrow$  total cost in  $k$ th trajectory.

Assume  $i_0^k = i \quad \forall k$ .

$$\hat{J}_l(i) = \frac{1}{l} \sum_{k=1}^l \hat{J}_k(i) \quad \leftarrow \text{Sample average.}$$

# trajectories used for average  $\rightarrow$

$\hat{J}$  can be incrementally computed by

$$J_{m+1}(i) = J_m(i) + \beta_m (\hat{J}_m(i) - J_m(i))$$

with initial condition  $J_0(i) = 0$ .

$\beta_m \rightarrow$  could be a general step size  
 $(\sum \beta_m < \infty, \sum \beta_m^2 < \infty)$

eg.  $\beta_m = \frac{1}{m}$ .

Reusing trajectories:

$(i_0, \dots, i_N) \rightarrow$  a trajectory

$(i_k, \dots, i_N) \rightarrow$  a sub-trajectory

Suppose  
 e.g.  $i_0 = i_k = i$   
 $\mathbb{I}(i, i, \dots, T) \hat{J}(i, i)$   
 $\uparrow$   
 k-thick was  $i$   
 $\mathbb{II}(i, i_{k+1}, \dots, T) \hat{J}(i, k)$   
 $\mathbb{II}$  is a sub-trajectory of  $\mathbb{I}$

$\hat{J}(i, i)$  &  $\hat{J}(i, k)$  are dependent. But, we could still use them both to estimate  $J_{\pi}(i)$

$$J_{\text{est}}(i_k) = J_m(i_k) + \beta_m (g(i_k, i_{k+1}) + \dots + \dots + g(i_{N-1}, i_N) - J_m(i_k))$$

This would estimate  $J_{\pi}$  with start state  $i_k$ .



## Lecture-22\*

Monte Carlo policy evaluation & its relation to  
temporal differences.

Given a MDP trajectory  $(i_k, \dots, i_N)$  simulated  
using policy  $\pi$ .

$$(*) \rightarrow \mathcal{J}_{m+1}(i_k) = \mathcal{J}_m(i_k) + \beta_m \left( \underbrace{g(i_k, i_{k+1}) + \dots}_{\text{One sample of } \sum_{n=0}^{\infty} \gamma^n g(i_{k+n}, i_{k+n+1})} + \dots + g(i_{N-1}, i_N) - \mathcal{J}_m(i_k) \right)$$

[Aside: If  $\beta_m \rightarrow 1$ , then  $\mathcal{J}_m \rightarrow \hat{v}$  a sample mean]

Rewrite (\*) as

$$\begin{aligned} \mathcal{J}_{m+1}(i_k) = \mathcal{J}_m(i_k) + \beta_m & \left[ \underbrace{(g(i_k, i_{k+1}) + \mathcal{J}_m(i_{k+1}) - \mathcal{J}_m(i_k))}_{\text{yellow}} \right. \\ & + \underbrace{(g(i_{k+1}, i_{k+2}) + \mathcal{J}_m(i_{k+2}) - \mathcal{J}_m(i_{k+1}))}_{\text{yellow}} \\ & \vdots \\ & \left. + \underbrace{(g(i_{N-1}, i_N) + \mathcal{J}_m(i_N) - \mathcal{J}_m(i_{N-1}))}_{\text{yellow}} \right] \end{aligned}$$

Note:  $\mathcal{J}_m(i_N) = 0$  (why?  $i_N = T$ )

$TD(0)$  uses only  $d_k$  to update.

$$J_{m \neq l}(i_k) = J_m(i_k) + \beta_m (d_k + d_{k+1} + \dots + d_{N-1}),$$

where

$$d_l = \underbrace{g(i_l, i_{l+1}) + J_m(i_{l+1})}_{\text{Estimate of } J_\pi(i_l) \text{ based on a simulated transition}} - \underbrace{J_m(i_l)}_{\text{current estimate of } J_\pi(i_l)}, \quad l = k, \dots, N-1$$

(\*\*) can be done incrementally as

[Note: going from  $i_l \rightarrow i_{l+1}$  makes  $d_l$  available]

$$J_{m \neq l}(i_k) = J_m(i_k) + \beta_m d_l$$



do this for  $l = k, \dots, N-1$

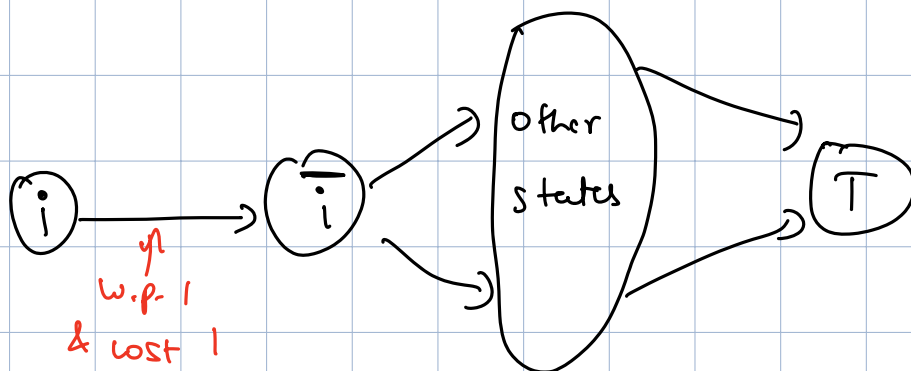
Remark:  $TD(0)$  updates as

$$J_{m \neq l}(i_k) = J_m(i_k) + \beta_m d_k$$

$TD(0)$  does not use all the temporal differences  $(d_k, d_{k+1}, \dots, d_{N-1})$  to update its estimate  $J_m(i_k)$ , & instead relies on a 1-step fixed-point equation (i.e., uses  $d_k$ )



# TD(0) vs. MCPE: A bias-variance tradeoff perspective



Fix some policy  $\pi$ .

From  $i$ , MDP always transitions to  $\bar{i}$  under  $\pi$ ,  
&  $g(i, \bar{i}) = 1$ .

Want to estimate:  $J_{\pi}(i)$

MCPE: Simulate trajectories starting in  $i$  & collect the total cost samples, say  $\{\hat{J}(m)\}_{m=1}^N$

$$\hat{J}(i) = \frac{1}{N} \sum_{m=1}^N \hat{J}(m)$$

Is this an unbiased estimate of  $J_{\pi}(i)$ ? Yes.

**TD(0):** Suppose through some other route (an independent simulation or some other approximate route), we have an estimate  $J(\bar{i})$ .

Then, we can estimate  $J_{\pi}(i)$  by

$$J(i) = J(\bar{i}) + 1$$

MCPE would never use  $J(\bar{i})$  & instead rely on sample trajectories, even at  $\bar{i}$ .

With TD(0), we have a biased estimate  $J(i)$  of  $J_{\pi}(i)$ .

MCPE estimation may suffer from high variance (c.s.p. if  $N$  is small), while a biased estimate  $J(i)$  may do better.

Bottomline: TD(0)  $\rightarrow$  biased estimate

MCPE  $\rightarrow$  unbiased, but possibly high variance.

Is there a middle path?

Yes, TD( $\lambda$ ),  $\lambda \in [0, 1]$ .

TD(0) used the 1-step fixed point equation to arrive at its update rule.

$$J_{\pi}(i) = E_i(g(i, \bar{i}) + J_{\pi}(\bar{i})) \quad (\Rightarrow J_{\pi} = T_{\pi} J_{\pi})$$

Why 1-step? We could also go 2 steps.

2-step fixed point equation

$$T_{\pi}(i) = E_{\bar{i}, \bar{i}} \left( g(i, \bar{i}) + g(\bar{i}, \bar{i}) + T_{\pi}(\bar{i}) \right)$$

$T_{\pi} = T_{\pi}^2 T_{\pi}$

Using the above, we can have the following iterative algorithm

$$T_{m \in \mathbb{C}}(i_k) = T_m(i_k) + \beta_m \left( g(i_k, i_{k+1}) + g(i_{k+1}, i_{k+2}) + T_m(i_{k+2}) - T_m(i_k) \right)$$

Can extend this to  $(l+1)$ -steps, i.e.,

$$T_{\pi}(i) = E_{i_0=i} \left( \sum_{m=0}^l g(i_m, i_{m+1}) + T_{\pi}(i_{l+1}) \right)$$

$(l+1)$ -step fixed point equation.

"Choice of  $l$  is arbitrary".

TD( $\lambda$ ) idea: Form a weighted average of the fixed point equations for different  $l$ .

How do we combine?

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

normalize (since  $\sum \lambda^l = \frac{1}{1-\lambda}$ )      weight for a particular  $l$        $\lambda \in [0, 1]$

$l$ -step fixed point target.



Simplifying the term (B):

$$E \left( \sum_{l=0}^{\infty} (\lambda^l - \lambda^{l+1}) \mathcal{J}_{\pi}(i_{l+1}) \right)$$
$$= E \left[ (1-\lambda) \mathcal{J}_{\pi}(i_1) + (\lambda - \lambda^2) \mathcal{J}_{\pi}(i_2) \right. \\ \left. + (\lambda^2 - \lambda^3) \mathcal{J}_{\pi}(i_3) + \dots \right]$$

$$= E \left[ (\mathcal{J}_{\pi}(i_1) - \mathcal{J}_{\pi}(i)) + \lambda (\mathcal{J}_{\pi}(i_2) - \mathcal{J}_{\pi}(i_1)) \right. \\ \left. + \lambda^2 (\mathcal{J}_{\pi}(i_3) - \mathcal{J}_{\pi}(i_2)) + \dots \right] + \mathcal{J}_{\pi}(i)$$

Added & subtracted terms

$$= E \left( \sum_{m=0}^{\infty} \lambda^m (\mathcal{J}_{\pi}(i_{m+1}) - \mathcal{J}_{\pi}(i_m)) \right) + \mathcal{J}_{\pi}(i)$$

(2)

Combining (1) & (2), we obtain

$$(*) \rightarrow \mathcal{J}_{\pi}(i) = E \left( \sum_{m=0}^{\infty} \lambda^m (g(i_m, i_{m+1}) + \mathcal{J}_{\pi}(i_{m+1}) - \mathcal{J}_{\pi}(i_m)) \right) + \mathcal{J}_{\pi}(i)$$

Recall  $d_m = g(i_m, i_{m+1}) + \mathcal{J}_{\pi}(i_{m+1}) - \mathcal{J}_{\pi}(i_m)$

So,  $(*) \Rightarrow 0 = E \left( \sum_{m=0}^{\infty} \lambda^m d_m \right) - \text{(B)}$

This really is a finite sum with a random # of terms because  $\forall m \geq N, d_m = 0$  since  $i_N = T$ .

Eq (3) is valid because  $E(d_m) = 0 \forall m$ .

How to turn (3) into an iterative update rule?

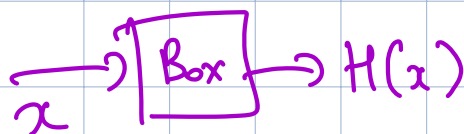
NOTE to typesetters! Change variable  $m$  to something else

$$J_{t+1}(i) = J_t(i) + \beta_t \sum_{m=0}^{\infty} \lambda^m d_m \quad \rightarrow \text{TD } (\lambda) \text{ update iteration.}$$

(4)

A short detour into stochastic update

$$h(x) = \mathbb{E}(H(x))$$



Want to find  $x^*$  s.t.  $h(x^*) = 0$

$$x_{t+1} = x_t + \beta_t (H(x_t))$$

Stochastic root - finding "Robbins-Monro" algorithm (PSI)

$$= x_t + \beta_t (h(x_t) + \underbrace{(H(x_t) - h(x_t))}_{\text{noise}})$$

Under regularity conditions (like Ac-Al before)  $x_t \rightarrow x^*$  a.s. as  $t \rightarrow \infty$

where  $x^*$  satisfies  $h(x^*) = 0$

Remarks:

① If  $\lambda = 1$ , then (2) becomes

$$J_{t+1}(i) = J_t(i) + \beta_t \sum_{m=0}^{\infty} d_m$$

$$J_{t+1}(i) = J_t(i) + \beta_t (d_0 + d_1 + \dots + d_{N-1})$$

↑  
MCPE scheme  $\Leftrightarrow$  TD(1)

↑  
assuming  $i_N = J$

② If  $\lambda = 0$ , then (2) becomes (using  $d^0 = 1$ )

$$J_{t+1}(i) = J_t(i) + \beta_t d_0$$

(Note:  $i_0 = i$ )

$$J_{t+1}(i) = J_t(i) + \beta_t (g(i_0, i_1) + J_t(i_1) - J_t(i_0))$$

↑  
TD(0) update iteration

$i_0, \dots, i_{100}$   
This TD term  $\rightarrow g(i_{100}, i_{101}) + J_t(i_{101}) - J_t(i_{100})$   
should have  
(or horizon  $J_{\pi}(i_0) = J_{\pi}(i)$ )

③ For any  $0 < \lambda < 1$ , the temporal difference  $d_m$  is weighted by  $\lambda^m$ , making a future temporal difference less important while updating the estimate of the current state  $i$ , i.e., the effect of  $d_0$  is more pronounced in TD( $\lambda$ ) update as compared to  $d_m, m > 1$ .

Note:  $\lambda$  weighs temporal differences. Not to be confused with discount ( $\alpha$ ).

P.T.O.

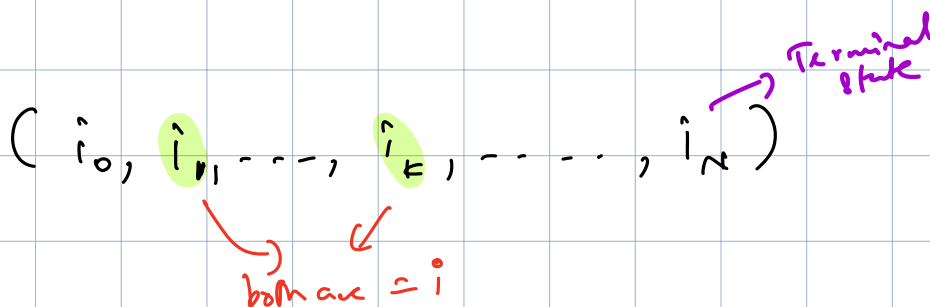
### Lecture-23

TD( $\lambda$ ) variations:

I Every visit vs. first visit

Suppose we have a trajectory  $(i_0, \dots, i_N)$

In this trajectory, a given state, say "i", may be visited more than once.





Every visit:  $(i_1, \dots, i_N)$  &  $(i_k, \dots, i_N)$   
to estimate  $J_{\pi}(i)$

First visit:  $(i_1, \dots, i_N)$  to estimate  $J_{\pi}(i)$

Formally, suppose  $i$  is visited  $M$  times in  $(i_0, \dots, i_N)$ ,  
and  $(m_1, \dots, m_M)$  are the time instants when  
state  $i$  is visited.

Then, TD( $\lambda$ ) would update as follows

$$J_{t+1}(i) = J_t(i) + \beta_t \sum_{j=1}^M \sum_{m=m_j}^{\infty} \lambda^{m-m_j} d_m$$

Why this?

In TD( $\lambda$ ) derivation, if we consider some state  $i_k$   
instead of  $i_0$ , then the fixed point equation becomes

$$J_{\pi}(i_k) = E \left( \sum_{m=k}^{\infty} \lambda^{m-k} d_m \right) + J_{\pi}(i_k)$$

First visit TD( $\lambda$ ) would update as follows:

$$J_{t+1}(i) = J_t(i) + \beta_t \sum_{m=m_1}^{\infty} \lambda^{m-m_1} d_m$$

Question: Are these two variants equivalent?

No.

But, both variants can be shown to converge. (For  $\lambda=1$ , see Section 5.2 of NDP book)  
idea: SLLN

## II Off-line TD( $\lambda$ ) vs online TD( $\lambda$ )

Offline: Simulate entire trajectory  $(i_0, \dots, i_N)$  and update in the end, i.e., after all temporal differences  $d_0, \dots, d_{N-1}$  are available.

Online: Update after each transition, i.e., after a single temporal difference term is available.

Offline TD- $\lambda$ : (Every visit variant)

$$J_{t+1}(i) = J_t(i) + \beta_t \sum_{j=1}^M \sum_{m=m_j}^{\infty} \lambda^{m-m_j} d_m$$

one transition of "d" become available  
 $(i_0, i_1, \dots, i_N)$

Online TD( $\lambda$ ): (Incremental update)

initial estimate

$$O_n(i_0, i_1) : J_1(i_0) = J_0(i_0) + \beta d_0$$

→ using constant step size for simplicity

$$i_0 \xrightarrow{a_0} i_1 \xrightarrow{a_1} i_2$$

On  $(i_1, i_2)$ :  $J_2(i_0) = J_1(i_0) + \beta \lambda d_1$   
 Use  $d_1$  to update estimates of value fn for  $i_0$  &  $i_1$   
 $J_2(i_1) = J_1(i_1) + \beta d_1$

and so on.

In general, on  $(i_k, i_{k+1})$

$$J_{t+1}(i_0) = J_t(i_0) + \beta \lambda^k d_k$$

$$J_{t+1}(i_1) = J_t(i_1) + \beta \lambda^{k-1} d_k$$

⋮

$$J_{t+1}(i_k) = J_t(i_k) + \beta d_k$$

Remark! If a state "i" is visited multiple times, then offline and online TD( $\lambda$ ) result in different estimates.

Example to illustrate the updates of offline & online TD( $\lambda$ )

Suppose we have a trajectory  $(1, 2, 1, T)$   
 TD-terms  $\rightarrow$   $d_0$   $d_1$   $d_2$

Let  $J_0(1)$  &  $J_0(2)$  be initial values. ( $J_0(T) = 0$ )

Offline TD( $\lambda$ ): Denote estimates by  $J_f(1)$  &  $J_f(2)$

State 1's update: (Every visit style)

$(1, 2, 1, T)$  &  $(1, T)$

$$J_f(1) = J_0(1) + \beta (d_0 + \lambda d_1 + \lambda^2 d_2 + d_2)$$

$$= J_0(1) + \beta (g(1, 2) + J_0(2) - J_0(1) + \lambda (g(2, 1) + J_0(1) - J_0(2)) + \lambda^2 (g(1, T) - J_0(1)))$$

from  $(1, 2, 1, T)$

$$+ g(1, T) - J_0(1)$$

from sub-trajectory  $(1, T)$

$$J_f(2) = J_0(2) + \beta (g(2, 1) + J_0(1) - J_0(2) + \lambda (g(1, T) - J_0(1)))$$

from  $(2, 1, T)$

On-line TD( $\lambda$ ) update: (Every-visit style)

On  $(1, 2)$ :

$$J_1(1) = J_0(1) + \beta (g(1, 2) + J_0(2) - J_0(1))$$

$$J_1(2) = J_0(2)$$

process TD-term "do"

On  $(2,1)$ :

$$J_2(1) = J_1(1) + \beta \lambda (g(2,1) + J_1(1) - J_1(2))$$

$$J_2(2) = J_1(2) + \beta (g(2,1) + J_1(1) - J_1(2))$$

On  $(1,T)$ :

$$J_3(1) = J_2(1) + \beta (\lambda^2 (g(1,T) - J_2(1)) + (g(1,T) - J_2(1)))$$

$$J_3(2) = J_2(2) + \beta \lambda (g(1,T) - J_2(1))$$

$(J_3(1), J_3(2)) \rightarrow$  Online TD( $\lambda$ ) estimates.

Compare this with  $(J_f(1), J_f(2))$ :

If we replace  $J_1$  &  $J_2$  by  $J_0$  in online TD update,

then  $(J_3(1), J_3(2)) = (J_f(1), J_f(2))$

$J_1$  &  $J_0$  difference is  $O(\beta)$

$J_2$  &  $J_0$  difference is  $O(\beta^2)$

$J_3$  &  $J_0$  ———  $O(\beta^3)$

or  $O(\beta^{c-1})$

If we take the step-size  $\beta \rightarrow 0$  as we update,  
then offline & online TD update will get close.

## Lecture 26\*

### Convergence of TD: A high-level sketch

Recall 
$$J_{t+1}^{(i)} = J_t^{(i)} + \beta_t (H J_t^{(i)} + w_t^{(i)} - J_t^{(i)}), i=1, \dots, n$$
 (\*)

If (i)  $H$  is a contraction,

(ii)  $E(w_t^{(i)} | \mathcal{F}_t) = 0$ ,  $E(w_t^{(i)2} | \mathcal{F}_t) \leq A + B \|J_t\|^2$

(iii)  $\sum \beta_t = \infty$ ,  $\sum \beta_t^2 < \infty$  ← standard stochastic approximation conditions

then  $J_t \xrightarrow[t \rightarrow \infty]{} J^*$  a.s., where  $HJ^* = J^*$  — (1)

Want to show TD( $\lambda$ ) update is like the Sto-iter-algo (\*)

Condition (iii) requires no effort. e.g.  $\beta_t = \frac{1}{t}$

Condition (i): We will dig deeper into TD( $\lambda$ )'s underlying mapping.

Note: We are doing policy evaluation using TD( $\lambda$ ) for a proper policy.

Side note!

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

$$P = \begin{bmatrix} p_{11} & \dots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \dots & p_{nn} \end{bmatrix}$$

$$P_{ij} = P(X_1=j | X_0=i)$$

$$P(X_2=k | X_0=i) = \sum_j P(X_2=k, X_1=j | X_0=i)$$

↳ split & use Markov property

So, squaring matrix  $P$ , i.e., " $P^2$ " would give two-step transition probabilities

Recall the fixed point equation underlying TD(0):

$$J_{\pi}(i) = (1-\lambda) E \left[ \sum_{l=0}^{\infty} \lambda^l \left( \sum_{m=0}^l g(i_m, i_{m+1}) + J_{\pi}(i_{l+1}) \right) \middle| i_0=i \right] \quad \textcircled{0}$$

normalize (since  $\sum \lambda^l = \frac{1}{1-\lambda}$ )

weight for a particular  $l$

$\lambda \in [0, 1]$

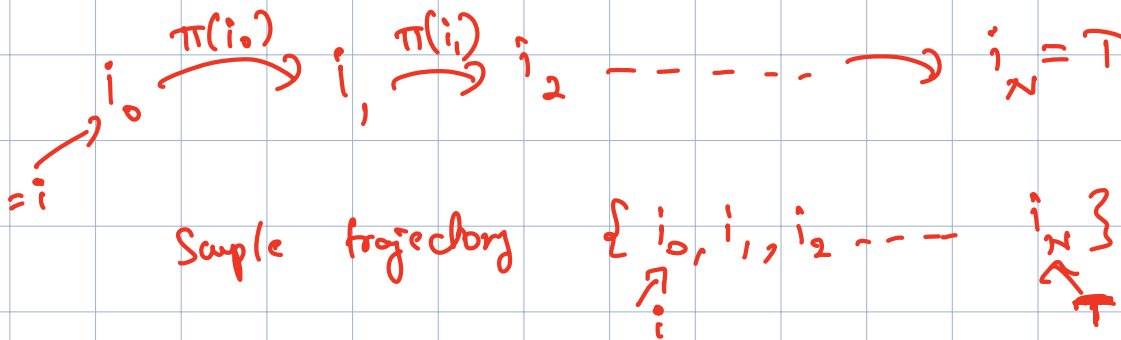
$\Rightarrow$  This is the TD( $\lambda$ ) operator

$$J_{\pi} = G + (1-\lambda) \sum_{l=0}^{\infty} \lambda^l P^{l+1} J_{\pi} \quad \textcircled{1}$$

$E J_{\pi}(i_{l+1})$  is like applying  $P$  ( $l+1$ ) times  
 $G$  is  $(1-\lambda) \sum_{l=0}^{\infty} \lambda^l E(g(-, -))$

transition probability matrix underlying the policy considered.

What is the expectation over in eq  $\textcircled{0}$



Inside expectation is terms like

(i)  $g(i_0, i_1) + J_{\pi}(i_1)$

(ii)  $g(i_0, i_1) + g(i_1, i_2) + J_{\pi}(i_2)$  & so on

each of these terms get a weight i.e.,  $\lambda$  "something"

Eq  $\textcircled{1}$  is like

$$J_{\pi} = H J_{\pi}, \text{ where}$$

$$H = G + (1-\lambda) \sum_{l=0}^{\infty} \lambda^l P^{l+1}$$

$\rightarrow$  TD( $\lambda$ ) fixed point equation

In the all policies proper case, we showed earlier that there exists a positive vector  $\xi$  and some  $\rho \in (0, 1)$  s.t.

$$\sum_{j=1}^n P_{ij}(\pi(i)) \xi(j) \leq \xi(i) - \rho \leq \rho \xi(i),$$

Or, equivalently

$$\|PJ\|_{\xi} \leq \rho \|J\|_{\xi}, \text{ where } \|\cdot\|_{\xi} \text{ is the weighted max-norm}$$

Define  $HJ = G + (1-\lambda) \sum_{l=0}^{\infty} \lambda^l P^{l+1} J$

Want  $\|HJ - HJ'\|_{\xi} \leq \rho \|J - J'\|_{\xi}$  — (xxx)

$\uparrow$   
 $0 < \rho < 1$

$\Rightarrow H$  is a  $\rho$ -contraction

& assuming the noise conditions  $E(w_t(i) | \mathcal{F}_t) = 0$   
 $\& E(w_t^2(i) | \mathcal{F}_t) \leq A + B \|z_t\|^2$   
 are met,

one can infer  $T_D(\lambda)$  iterate converges asymptotically.



If  $(**)$  holds, then the mapping  
 underlying  $TD(\lambda)$  is a contraction &  
 we can claim convergence using ①

$$\|HJ - HJ'\|_{\xi}$$

$$= \left\| \left( \cancel{H} + (1-\lambda) \sum_{l=0}^{\infty} \lambda^l P^{l+1} J \right) \right.$$

$$\left. - \left( \cancel{H} + (1-\lambda) \sum_{l=0}^{\infty} \lambda^l P^{l+1} J' \right) \right\|_{\xi}$$

triangle  
 inequality

$$\leq (1-\lambda) \sum_{l=0}^{\infty} \lambda^l \|P^{l+1}(J - J')\|_{\xi}$$

$$\leq (1-\lambda) \sum_{l=0}^{\infty} \lambda^l e \|J - J'\|_{\xi}$$

$$= (1-\lambda) e \|J - J'\|_{\xi} \sum_{l=0}^{\infty} \lambda^l$$

$$= e \|J - J'\|_{\xi}$$

Since

$$\|P^{l+1}(J - J')\|_{\xi} \leq e^{l+1} \|J - J'\|_{\xi} \leq e \|J - J'\|_{\xi}$$

since  $e \in (0, 1)$

So, we have

$$\|HJ - HJ'\|_{\xi} \leq e \|J - J'\|_{\xi}$$

So, the operator  $T$  underlying TD( $\lambda$ ) update is contractive.

Assuming condition (iii) (leading up to Eq. (1) above) hold, we can claim

$$T_t \rightarrow T^* \text{ a.s. as } t \rightarrow \infty$$

where  $T_t$  is the TD( $\lambda$ ) iterate.

---

## TD( $\lambda$ ) for discounted MDPs

Policy evaluation using TD( $\lambda$ )

### Approach I

Since there is no termination state, to reuse the TD( $\lambda$ ) idea for SSPs,

convert a discounted MDP to its equivalent SSP.

### Approach II

Do something else

Approach I: **Discounted MDP**, fix some policy  $\pi$ .

An MDP trajectory has no end ( $\because$  no termination state)

Add a "termination state" & from each state  
add a prob  $(1-\alpha)$  transition ( $\alpha \in (0,1)$   
discount factor)

**Simulation!** Toss a coin with bias  $\alpha$  in each time instant  
If heads, continue simulation.  
Else, end the trajectory & do TD( $\lambda$ ) update.

$(i_0, i_1, \dots, i_N)$  — episode & trajectory aka sample path  
↑  
termination state.

$N \in \text{r.v. "Geometric"} \quad E(N) = \frac{1}{1-\alpha}$

Use this trajectory to do TD( $\lambda$ ) update.

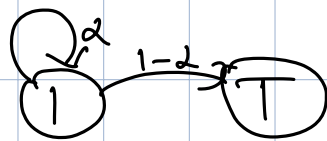
**Drawback of this approach: High variance**

**Example!** Just one state, say 1

Twist: single stage cost is a r.v. with mean 0  
& variance  $\sigma^2$  (note: cost does not depend  
on state).

Approach 1 adds termination state "T"

$$J_{\pi}(1) = 0$$



equivalent SSP

From a trajectory we got single stage costs

Single stage cost samples  $\rightarrow \{g_1, g_2, \dots, g_N\}$  ( $N$  is a r.v.)

$$\hat{J} = \sum_{i=1}^N g_i$$

Variance of this total cost sample

$$= E \left( (g_1 + g_2 + \dots + g_N)^2 \right)$$

Conditional expectation

$$= \sum_{k=1}^{\infty} E \left( (g_1 + g_2 + \dots + g_N)^2 \mid N=k \right) P(N=k)$$

$$= \sum_{k=1}^{\infty} E \left( (g_1 + g_2 + \dots + g_k)^2 \right) P(N=k)$$

$$= \sum_{k=1}^{\infty} k \sigma^2 P(N=k)$$

$$= \sigma^2 \left( \sum_{k=1}^{\infty} k P(N=k) \right)$$

$$= \sigma^2 E N = \frac{\sigma^2}{1-\alpha}$$

$g_i \sim \text{iid}$   
with mean 0  
variance  $\sigma^2$

$$Z = X_1 + \dots + X_k$$

$$X_i \sim N(0, \sigma^2)$$

$\{X_1, \dots, X_k\}$  indep

$$E Z^2 = k \sigma^2$$

Variance of estimate  $\hat{J}$

Note:  $\alpha$  very close to 1 leads to "big" variance.

## Approach 2:

1-Step fixed point equation  $J_{\pi}(i) = E [ g(i, \bar{i}) + \alpha J_{\pi}(\bar{i}) ]$   
2-Step — a —  $J_{\pi}(i) = E [ g(i, \bar{i}) + \alpha g(\bar{i}, \bar{\bar{i}}) + \alpha^2 J_{\pi}(\bar{\bar{i}}) ]$

TD( $\lambda$ ) fixed point equation for discounted case:

$$J_{\pi}(i_0) = (1-\lambda) E \left[ \sum_{l=0}^{\infty} \lambda^l \left( \sum_{m=0}^l g(i_m, i_{m+1}) + \alpha J_{\pi}(i_{l+1}) \right) \right]$$

Repeating all the steps from the SSP TD( $\lambda$ ) derivation, we obtain the following fixed point relation for discounted case:

(\*\*\* $\rightarrow$ )  $J_{\pi}(i_0) = E \left( \sum_{m=0}^{\infty} (\alpha \lambda)^m d_m \right) + J_{\pi}(i_0)$ ,

where  $d_m = g(i_m, i_{m+1}) + \alpha J_{\pi}(i_{m+1}) - J_{\pi}(i_m)$

Can do this update from an intermediate state  $i_k$  in the trajectory.

$$J_{\pi}(i_k) = E \left( \sum_{m=k}^{\infty} (\alpha \lambda)^{m-k} d_m \right) + J_{\pi}(i_k)$$

The TD( $\lambda$ ) update would be

$$J_{t+1}(i_k) = J_t(i_k) + \beta \sum_{m=k}^{\infty} (\alpha \lambda)^{m-k} d_m$$

in comparison to SSP one,  
we have the  $\alpha$  factor here.

Ques: When to end trajectories? To be answered.

### Lecture-25\*

Variance calculation for the single state MDP

Note: We aren't adding the termination state.

Trajectory:  $(i_0, i_1, i_2, \dots)$

Total Cost Sample =  $g_1 + \alpha g_2 + \alpha^2 g_3 + \dots$

$$\text{Variance} = E \left[ \left( \sum_{k=0}^{\infty} \alpha^k g_{k+1} \right)^2 \right]$$

$\{g_k\}$  iid  
mean zero  
variance  $\sigma^2$

$$= \sigma^2 \sum_{k=0}^{\infty} \alpha^{2k} = \frac{\sigma^2}{1-\alpha^2}$$

$$= \frac{\sigma^2}{(1-\alpha)(1+\alpha)}$$

With SSP formulation, we had a variance of  $\frac{\sigma^2}{1-\alpha}$

which is  $>$  than  $\frac{\sigma^2}{1-\alpha^2}$

Approach 2

Offline TD( $\lambda$ ): cut the trajectory at some random time  $\tau$

"Simulate for a finite ( $\tau$ ) # of steps & approximate the discounted cost by

$$g_1 + \gamma g_2 + \dots + \gamma^{\tau-1} g_\tau + \gamma^\tau J(i_\tau)$$

e.g.  $\tau = 200$ .  $\gamma = 0.8$ . Then, after 200 steps, the contribution of costs to  $J_\pi$  is negligible.

Online TD( $\lambda$ )

update value function estimates on every sample transition.

Update rule  $\leftarrow$  very similar to the SSP case.

Just patch in " $\lambda$ ".

Estimate:  $E \left( \sum_{m=0}^{\infty} \gamma^m g_m \right)$

Tweak:  $E \left( \sum_{m=0}^{\tau} \gamma^m g_m \right)$

Take samples of this expectation with a truncated trajectory.

If  $\tau$  is large enough, then intuitively

the contribution of  $E\left(\sum_{m=\tau+1}^{\infty} \alpha^m g_m\right)$  to

the total cost is negligible &

$E\left(\sum_{m=0}^{\tau} \alpha^m g_m\right)$  is a good

enough approximation.

Convergence analysis of TD( $\lambda$ ): skipped.

Check Prop 5.1 of NDP book for details.

---

Why Q-factors in a learning scenario?

Policy evaluation:  $J_{\pi}(i) = E_i(g(i, \pi(i), i) + J_{\pi}(\bar{i}))$

To learn  $J_{\pi}$ , sample from the r.v. inside expectation & do a Sto-iter-algo.

Control:  $J^*(i) = \min_a E(g(i, a, i) + J^*(\bar{i}))$

Can I turn this fixed point relation into a Sto-iter-algo?



# Q-learning

Recall Q-factors:

Def:  $Q^*(i, a) = \sum_j P_{ij}(a) (g(i, a, j) + J^*(j))$

optimal cost starting in  $j$ .

Take action  $a$  in state  $i$  & then follow the optimal policy from state  $j$  onwards.

Bellman equation  $J^* = T J^*$  is equivalent to

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

Combining the two equations leads to the following "Q-Bellman equation".

$$Q^*(i, a) = \sum_j P_{ij}(a) (g(i, a, j) + \min_b Q^*(j, b)) \quad (*)$$

optimal Q-values

Note:  $Q^*$  is the unique solution of  $*$ .

Suppose  $Q$  is a solution of  $*$ , i.e.,

$$Q(i, a) = \sum_j P_{ij}(a) (g(i, a, j) + \min_b Q(j, b)) \quad (1)$$

Then,  $Q = Q^*$ .

This can be seen using the fact that  $J^*$  is the unique solution of  $J^* = T J^*$ .

$Q$  solves (1), so

$\left( \min_a Q(i,a) \right)$  solves the Bellman equation  $J^* = T J^*$

i.e.,  $J^*(i) = \min_a Q(i,a)$

$J^*$  is unique  $\Rightarrow$

$$\min_a Q(i,a) = \min_a Q^*(i,a)$$

Using  $\min_b Q(j,b) = \min_b Q^*(j,b)$  in (1), we

obtain  $Q^* = Q$ .

So,  $Q^*$  is the unique solution.

★ We did not require contraction in this argument.

So, if  $J^* = T J^*$  has a unique solution, then  $Q^* = H Q^*$  also has a unique solution

↑  
is the operator underlying Q-Bellman equation



Value iteration (VI) using Q-factors:

LSPF  
policy gradient

Q-value iteration

$$Q_{t+1}(i, a) = \sum_j P_{ij}(a) (g(i, a, j) + \min_b Q_t(j, b))$$

(This is like  $Q_{t+1} = H Q_t$ , starting with some  $Q_0$ )  
where  $H Q(i, a) = \sum_j P_{ij}(a) (g(i, a, j) + \min_b Q(j, b))$

A variation to VI is

$$Q_{t+1}(i, a) = (1 - \beta_t) Q_t(i, a) + \beta_t \sum_j P_{ij}(a) (g(i, a, j) + \min_b Q_t(j, b))$$

VI requires knowledge of these transition probabilities

Sto-iter-algo version of the above:

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

( $\bar{i} \leftarrow$  sampled from  $P_{ij}(a)$ )

This is the Q-learning algorithm.

Note: The step-size could be iteration dependent i.e.,  $\beta_t$ . Need  $\sum \beta_t = \infty$  &  $\sum \beta_t^2 < \infty$ .

Remark: ① In principle, Q-learning is similar to TD(0).  
Both are based on VI & replace an expectation by its sample.

② There is no straightforward variation of Q-learning that is in the spirit of TD( $\lambda$ ).  
No  $(l+1)$ -step Q-Bellman equation.  
(Think about this!)

Lecture-26\*

## Convergence analysis of Q-learning:

Note: Convergence of Q-learning  $\Rightarrow$  convergence of TD(0)  
(just consider a special MDP with <sup>only</sup> feasible action  $\pi(i)$  in state  $i$ .  
Then Q-BE  $\Leftrightarrow (J^\pi = T^\pi J^\pi)$ )

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

$\uparrow$  sampled for  $R_{ij}(a)$        $\bar{i}$  action in next state  $\bar{i}$

Assumptions:

(A1)  $\sum \beta_t = \infty$ ,  $\sum \beta_t^2 < \infty$   $\leftarrow$  easy to satisfy e.g.  $\beta_t = \frac{c}{t}$

(A2) All policies are proper in the underlying SSP.

Recall from previous chapter:

Suppose Q-learning update in compact notation is

$$Q_{t+1} = (1 - \beta_t) Q_t + \beta_t (H Q_t + w_t)$$

Then, if we show  $(\mathcal{F}_t = \sigma(Q_0, \dots, Q_t, w_0, \dots, w_{t-1}))$

(B1)  $H$  is a contraction

(B2)  $E(w_t | \mathcal{F}_t) = 0$      $E(w_t^2 | \mathcal{F}_t) \leq A + B \|Q_t\|^2$

(B3)  $\sum \beta_t = \infty$ ,     $\sum \beta_t^2 < \infty$

Then,  $Q_t \rightarrow Q^*$  (which is the fixed point of  $H$ )  
a.s. as  $t \rightarrow \infty$ .

Main proof (Theorem comes later) is

Define  $(H Q)(i, a) = \sum_j p_{ij}(a) (g(i, a, j) + \min_b Q(j, b))$   
 $\forall i, a.$

Let  $w_t(i, a) = g(i, a, \bar{j}) + \min_b Q_t(\bar{i}, b) - (H Q_t)(i, a)$

So, Q-learning update iteration is equivalent to

$Q_{t+1}(i, a) = (1 - \beta_t) Q_t(i, a) + \beta_t ((H Q_t)(i, a) + w_t(i, a))$   
↳ (\*\*)

*Annotations:*  
-  $(H Q_t)(i, a)$ : operator without Q-learning update  
-  $w_t(i, a)$ : noise

Verifying conditions on noise:

$$w_t(i, a) = Y_t - E Y_t, \text{ where}$$

$$Y_t = g(i, a, i) + \min_b Q_t(i, b)$$

$$E(w_t(i, a) | \mathcal{F}_t) = 0 \rightarrow w_t \text{ is conditionally (given } Q_t) \text{ zero-mean}$$

$$E(w_t^2(i, a) | \mathcal{F}_t) = E((Y_t - E Y_t)^2 | \mathcal{F}_t)$$

$$\leq E(Y_t^2 | \mathcal{F}_t)$$

$$\begin{aligned} \text{Var}(X) &= E((X - E X)^2) \\ &= E X^2 - (E X)^2 \\ &\leq E X^2 \end{aligned}$$

Assuming single stage cost  $g(\cdot, \cdot, \cdot)$  is bounded, we have

$$E(Y_t^2 | \mathcal{F}_t) = E\left(\left(g(i, a, i) + \min_b Q_t(i, b)\right)^2 | \mathcal{F}_t\right)$$

$$\begin{aligned} (a+b)^2 &\leq 2a^2 + 2b^2 \end{aligned}$$

" $|g(i, a, i)| \leq \text{const}$ "  
 & state-action spaces are both finite

$$\leq K \left(1 + \max_{j, b} Q_t^2(j, b)\right) \quad (*)$$

$$\Rightarrow E(w_t^2(i, a) | \mathcal{F}_t) \leq K \left(1 + \max_{j, b} Q_t^2(j, b)\right)$$

So, we have verified (B2)

(B3) is satisfied if we chose  $\beta_t$  carefully.  
 i.e.,  $\sum \beta_t = \infty, \sum \beta_t^2 < \infty$

Onto (B1):  $\mathcal{H}$  is a contraction.

All policies proper  $\Rightarrow \exists$  a positive vector  $\xi$  & scalar  $c$ ,  
 such that

Fact A

$$\sum_j p_{ij}(a) \xi(j) \leq \epsilon \xi(i)$$

← from SSP Chapter.

Define  $\|Q\|_{\xi} = \max_{i,a} \frac{|Q(i,a)|}{\xi(i)}$

← weighted max-norm for Q-values.

$\xi$  is given by the claim in SSP chptr.

Need to show:  $\|HQ - HQ'\|_{\xi} \leq \epsilon \|Q - Q'\|_{\xi}$   
for some  $\epsilon \in (0,1)$ .

If this holds, then (B1) is satisfied.

Pf of (need to show):

Recall  $(HQ)(i,a) = \sum_j p_{ij}(a) (g(i,a,j) + \min_b Q(j,b))$

$$|(HQ)(i,a) - (HQ')(i,a)| = \left| \sum_j p_{ij}(a) (\min_b Q(j,b) - \min_b Q'(j,b)) \right|$$

De  
ineq.  $\leq \sum_j p_{ij}(a) \left| \min_b Q(j,b) - \min_b Q'(j,b) \right|$

$$\leq \sum_j p_{ij}(a) \max_b |Q(j,b) - Q'(j,b)|$$

$$= \sum_j p_{ij}(a) \left( \max_b \frac{|Q(j,b) - Q'(j,b)|}{\xi(j)} \right) \xi(j)$$

$$\leq \sum_j p_{ij}(a) \left( \max_{j,b} \frac{|Q(j,b) - Q'(j,b)|}{\xi(j)} \right) \xi(j)$$

← "add max over j"

$$\leq \sum_j P_{ij}(a) \|Q - Q'\|_{\xi} \xi(j) = \|Q - Q'\|_{\xi} \sum_j P_{ij}(a) \xi(j)$$

Using Fact A  $\rightarrow$

$$\leq \|Q - Q'\|_{\xi} e \xi(i)$$

So, we get

$$|(HQ)(i,a) - (HQ')(i,a)| \leq \|Q - Q'\|_{\xi} e \xi(i)$$

$\Rightarrow$

$$\frac{|(HQ)(i,a) - (HQ')(i,a)|}{\xi(i)} \leq e \|Q - Q'\|_{\xi}$$

$$\max_{i,a} \frac{|(HQ)(i,a) - (HQ')(i,a)|}{\xi(i)} \leq e \|Q - Q'\|_{\xi}$$

$$\Rightarrow \|HQ - HQ'\|_{\xi} \leq e \|Q - Q'\|_{\xi}$$

$$\Rightarrow H \text{ is a contraction w.r.t } \|\cdot\|_{\xi}$$

Thus, we have

### Theorem (Q-learning convergence)

(A1) All policies proper (A2)  $\sum \beta_c^{2\alpha} < \infty$ ,  $\sum \beta_c^2 < \infty$

(A3) Single stage cost  $g$  is bounded, i.e.,  $\sup_{i,a,j} |g(i,a,j)| \leq M < \infty$

Under (A1) - (A3), the Q-learning algorithm converges a.s. i.e.,

$$Q_t \rightarrow Q^* \text{ a.s. as } t \rightarrow \infty.$$



A variation when (A1) is not satisfied.

Instead, we have "∃ a proper policy" & "improper policies have infinite cost".

Even here  $Q^*$  is the unique solution to the Q-Bellman equation. (This didn't require  $H$  to be contractive. Instead, we only used  $J^* = T J^*$  &  $J^*$  is unique.)

Question: Does Q-learning converge in this case?

(A1') ∃ a proper policy & all improper policies have infinite cost.

Under (A1'),  $H$  is a monotone mapping

$$\text{i.e., } Q \leq Q' \Rightarrow H Q \leq H Q'$$

↓  
"Check this using definition of  $H$ "  
n.w.

Also check  $J_r$

$$H(r - \delta c) \leq H(r - \delta c) \leq H(r + \delta c) \leq H(r + \delta c)$$

$c =$  vector of all ones.

Theorem:  
(Q-learning  
under  
monotonicity)

(A1') + " $\sum \beta_t = \infty, \sum \beta_t^2 < \infty$ " +

"bounded stage cost"

+ "bounded iterate i.e.,  $\sup_{t, i, a} |Q_t(i, a)| < \infty$ "

$\Rightarrow Q_t \rightarrow Q^*$  a.s. as  $t \rightarrow \infty$

where

$\Pi Q^* = Q^*$ .

} using Theorem 2  
from previous  
chapter

A sufficient condition was stated after Theorem 2 in previous chapter.

This condition ensures boundedness of the iterates  $\{Q_t\}$

& is satisfied for Q-learning.

For details, see prop 5.6. of NDP book.

So, the final claim is Q-learning converges under (A1') + step size condition. The boundedness of iterates is implied.

## Lecture - 27

### Q-learning for discounted MDPs:

Q-Bellman equation in a discounted setting:

$$Q^*(i,a) = \sum_j P_{ij}(a) (g(i,a,j) + \alpha \min_b Q^*(j,b))$$

discount  $\swarrow$   $(*)$

VI°:

$$Q_{t+1}(i,a) = \sum_j P_{ij}(a) (g(i,a,j) + \alpha \min_b Q_t(j,b))$$

Q-learning°:

$$Q_{t+1}(i,a) = (1 - \beta_t) Q_t(i,a) + \beta_t (g(i,a,j) + \alpha \min_b Q_t(j,b))$$

$\nearrow$  sampled from  $P_{ij}(a)$

Convergence of Q-learning:

(A1') Assume  $|g(\cdot, \cdot, \cdot)| \leq M < \infty$

(A2')  $\sum \beta_t = \infty$ ,  $\sum \beta_t^2 < \infty$

Then, following the proof in the "all policies are proper" case of SSP-Q-learning, one can infer that

The underlying operator  $HQ$  is a  $\alpha$ -contraction.

$$(HQ)(i, a) = \sum_j p_{ij}(a) (g(i, a, j) + \alpha \min_b Q(j, b))$$

$$\|HQ - HQ'\|_{\infty} \leq \alpha \|Q - Q'\|_{\infty}$$

$\uparrow$   
max-norm  $\|Q\|_{\infty} = \max_{|i|=n} |Q(i)|$

So, Under  $(A1')$ ,  $(A2')$ , we have  
 $Q_t \rightarrow Q^*$  w.p.1 as  $t \rightarrow \infty$

---

Issue of exploration:

MCPE:  $E \left( \sum_{k=0}^{\infty} \alpha^k g(i, \pi(i), j) \right)$

$\downarrow$   
trajectories to do policy evaluation

$\hat{J}(i) \rightarrow$  estimate of the value function/expected wt,

then  $\hat{J}(i) \rightarrow J_{\pi}(i)$  if

you see enough trajectories starting with "i".

$\hat{J}(i)$ : sample average  $\xrightarrow{\text{converge}}$   $J_{\pi}(i)$  if we

See state  $i$  <sup>infinitely</sup> often in the trajectories.

Same logic applies to TD and Q-learning

With TD, for  $J_t(i) \rightarrow J_{\pi}(i)$  as  $t \rightarrow \infty$   
we need to visit  $i$  "i.o." in the trajectories.

But, sample is using a fixed policy  $\pi$ .

Compare with Q-learning:

$$Q_t(i, a) \xrightarrow{t \rightarrow \infty} Q^*(i, a)$$

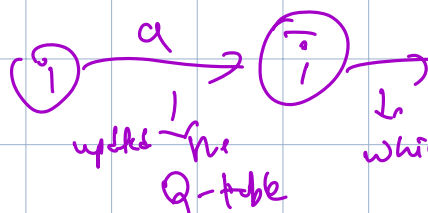
Here, the algorithm is free to choose the action "a" in each state.

The requirement for convergence:

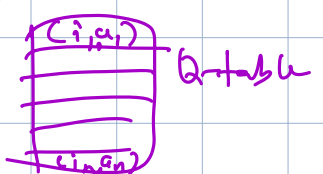
All pairs " $(i, a)$ " are visited frequently.

issue of exploration

Q-learning:



If  $Q_t \approx Q^*$



then select  $a' = \arg \min Q(i, a')$

Why do we need to take the route using Q-factors for finding an optimal policy?

Policy evaluation! we used  $J$  which is a function of state variable.

So, why  $Q(i,a)$  & Q-Bellman equation?

Or, why not use  $J^* = T J^*$  & make a sto-iterative algo to find  $J^*$ ?

Normal Bellman equation:  $J^*(i) = \min_a \sum_j P_{ij}(a) (g(i,a,j) + \alpha J^*(j)) \rightarrow (1)$

Q-Bellman equation:  $Q^*(i,a) = \sum_j P_{ij}(a) (g(i,a,j) + \alpha \min_b Q^*(j,b)) \rightarrow (2)$

General sto-iter-algo: Want to estimate  $\mu = E(X)$

Take samples of  $X$  & do an iterative update.

Eq (2) is in the form  $Q^*(i,a) = E_j (g(i,a,j) + \alpha \min_b Q^*(j,b))$

↳ sample this & do an iterative update.

Eq (1) is of the form  $J^*(i) = \min_a E(\_)$

Since min is outside, an iterative algo is not possible by just replacing the expectation above with a sample.

## How to do exploration:

For Q-learning to converge, we require all state-action pairs to be visited frequently.

Recall Q-learning update:

$$Q_{t+1}(i,a) = (1-\beta_t) Q_t(i,a) + (\gamma(i,a,\bar{i}) + \alpha \min_b Q(i,b))$$

Question! How to choose actions?

I Greedy! In state  $i$ , choose  $\arg \min_a Q_t(i,a)$ , at time instant  $t$ . If  $Q_t \approx Q^*$ , then this choice makes perfect sense. However, if  $Q_t$  isn't close to  $Q^*$ , then "we need to explore".

$Q(i, a_1)$
$Q(i, a_2)$
$Q(i, a_3)$
$\vdots$
$Q(i, a_m)$

Q-table

→ Each entry of this table has to be updated a good # of times for Q-learning to converge

II

 $\epsilon$ -greedy:Fix  $\epsilon > 0$ , usually a small number.At time instant  $t$ , pick the greedy action w.p.  $(1-\epsilon)$ & pick an action unif. at random w.p.  $\epsilon$ .Alternative: Make  $\epsilon$  a function of iteration  $t$ , say  $\epsilon_t$ , &take  $\epsilon_t \rightarrow 0$  as  $t \rightarrow \infty$ , i.e., reduce

exploration as algorithm updates.

III

In state  $i$ , at time instant  $t$ ,action  $a$  is chosen w.p.

$$\xi(i, a) = \frac{\exp(-Q_t(i, a)/\tau)}{\sum_{b \in \mathcal{A}(i)} \exp(-Q_t(i, b)/\tau)}$$

temperature  
parameter  $\tau$ 

$$\xi(i) = [\xi(i, a_1) \dots \xi(i, a_m)]$$

prob. distribution over  
actions in state  $i$ this is the prob. of choosing action  $a$  $\tau \rightarrow$  "temperature"  $\rightarrow$  controls the exploration.Note: If  $\tau$  is very small, the choice is  
"greedy".



Policy evaluation!

$$J_{\pi_k} = T_{\pi_k} J_{\pi_k}$$

↑  
exact

↘  
no noise

MP2:

$$J_{k+1} = T_{\pi_k}^{m_k} J_k$$

↗  
no noise

Compare this with TD(0):

$$J_{t+1}(i) = J_t(i) + \beta_t (T_{\pi_k} J_t(i) + \omega_t(i) - J_t(i))$$

↑  
noise

Running TD(0) for some # of steps, say  $M$ ,

" $J_M$ "  $J_M$  is very close to  $J_{\pi_k}$

Using  $J_M$  for policy improvement,

There is no guarantee of improvement.

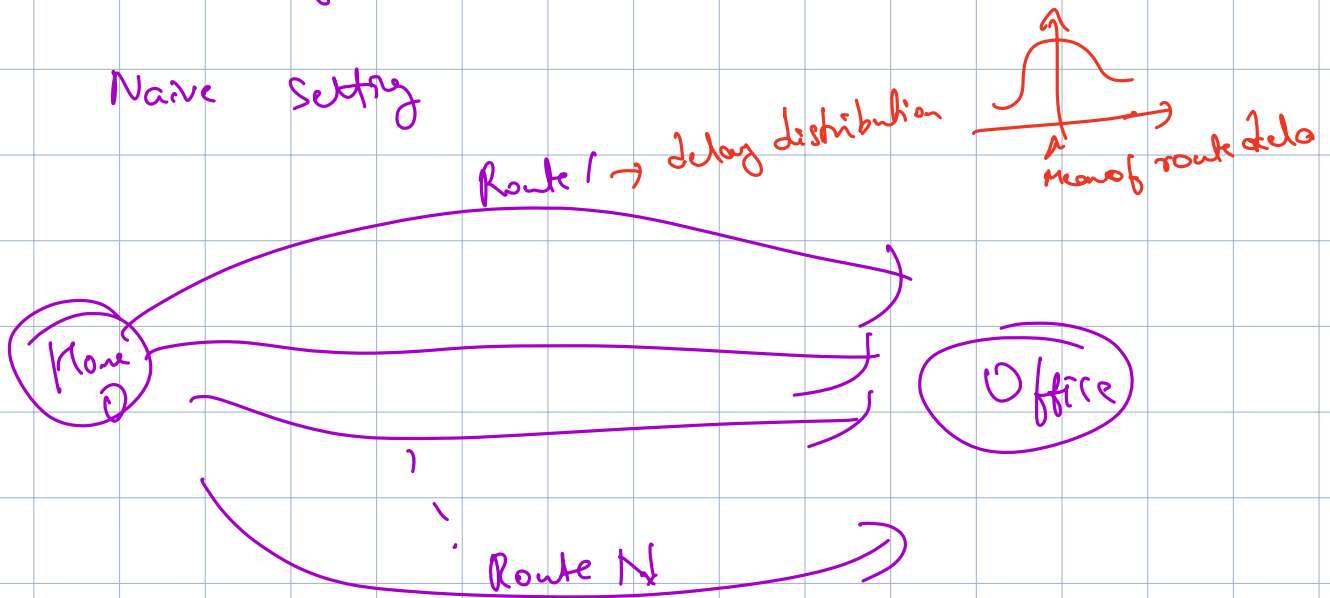
"Konda-Borkar

"Actor-critic algos"

1999, SIAM. J. Control. Opt.

# Bandit angle to exploration:

Naive setting



Routes:  $1, \dots, N$

Means:  $\mu_1, \dots, \mu_N$

$$\mu^* = \min_{i=1 \dots N} \mu_i$$

$$\text{Opt-route } i^* = \arg \min_{i=1 \dots N} \mu_i$$

for day = 1, 2, 3, ...

{

Pick a route  $\text{Route}_{\text{day}}$

Observe a sample delay from distribution

$P_{\text{day}}$

}

On some day → sample average

$\mu_1$   
...  
 $\mu_N$

How to pick route on following day?

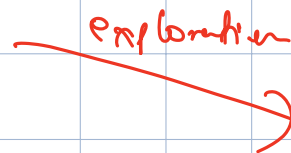
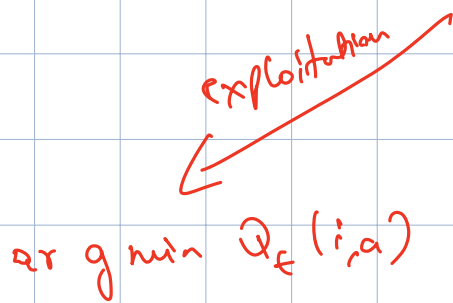


pick the route  
with best sample avg

pick a random route

route to pick  $\rightarrow \arg \max_i \hat{r}_i$

In Q-learning



pick a random action

$\epsilon$ -greedy



w.p.  $(1-\epsilon)$

w.p.  $\epsilon$