

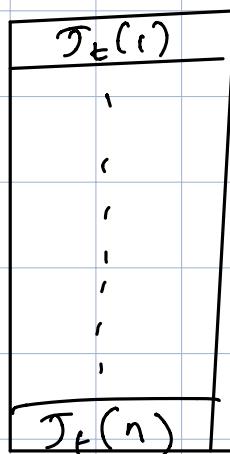
TD-learning, Q-learning with full state representations

ref: Chapter 5 of ND P book

Full state-representation?

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.



$\xrightarrow{\text{we sample & update each entry}}$

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)

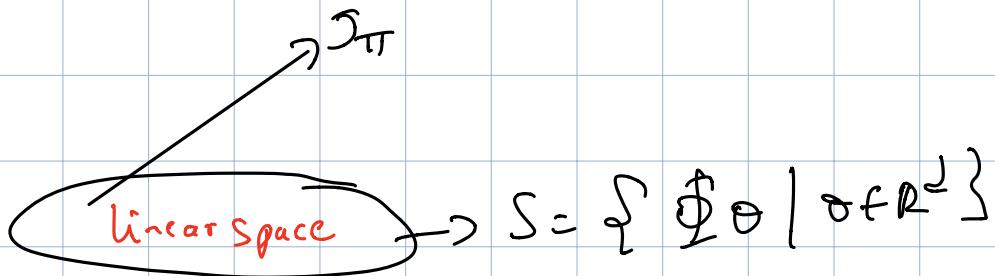
On such problem, we require parametric approximation of  $J_{\pi}$ .

e.g.  $J_{\pi}(i) \approx \theta^T \phi(i)$

$\theta \in \mathbb{R}^d$   
 $d \ll n$ .

↑  
parameter  
feature vector

linear function approximation



Recall: Mean estimation

Want to estimate  $\mu = E X$

$$\hat{\mu}_{m+1} = \hat{\mu}_m + \beta_m (x_{m+1} - \hat{\mu}_m)$$

↑  
stepsize  
e.g.  $\frac{1}{m}$

↑  
sample from  
distribution  
of  $x$

↑  
prev.  
estimate.  
of  $\mu$

estimate  
of  $\mu$

Question of policy evaluation!

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

Assume a SSP context.

Fixed point equation for  $\mathcal{J}_\pi$ :

$$\mathcal{J}_\pi = T_\pi \mathcal{J}_\pi$$

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

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

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

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

$$(*) \quad \mathcal{J}_{m+1}(i) = \mathcal{J}_m(i) + \beta_m (g(i, \bar{i}) + \mathcal{J}_m(\bar{i}) - \mathcal{J}_m(i))$$

This is the  
 $TD(0)$  algorithm  
↑ rapid difference

sample of the r.v. in expectation  
highlighted above

NOTE:  $\bar{i} \sim p_{i,j}(\pi(i))$

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

$$E(g(i, \bar{i}) + \mathcal{J}_m(\bar{i})) = (T_\pi \mathcal{J}_m)(i)$$

(\*) is equivalent to

$$\mathcal{J}_{m+1}(i) = \mathcal{J}_m(i) + \beta_m ((T_\pi \mathcal{J}_m)(i) - \mathcal{J}_m(i))$$

$$+ g(i, \bar{i}) + \mathcal{J}_m(\bar{i}) - (T_\pi \mathcal{J}_m)(i))$$

This is the noise factor

$w_m(i)$  from the general sto. it's.  
also.

## Monte Carlo Policy Evaluation

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

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

$$(i_0^1, \dots, i_N^1) \text{ Trajectory \#1}$$

$$\vdots \qquad \vdots$$

$$(i_0^l, \dots, i_N^l) \text{ Trajectory \#l}$$

$N$  is random

$$i_N^1 = \dots = i_N^l = T$$

Actions w.r.t  $\pi$ .

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

This is the r.v. whose expectation we estimate.

Let

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

for  $k=1, \dots, l$

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

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

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

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

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

with initial condition  $\mathcal{J}_0(i) = 0$ .

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

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

Reusing trajectories:

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

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

$$\mathcal{J}_{m+1}(i_k) = \mathcal{J}_m(i_k) + \beta_m (g(i_k, i_{k+1}) + \dots + g(i_{N-1}, i_N) - \mathcal{J}_m(i_k))$$



This would estimate  $\mathcal{J}_m$  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  $\pi$ .

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

[Aside: If  $\beta_m = 1$ , then  $J_m \rightarrow \text{a sample mean}$ ]

Rewrite (\*) as

$$\begin{aligned} J_{m+1}(i_k) &= J_m(i_k) + \beta_m \left[ (g(i_k, i_{k+1}) + J_m(i_{k+1}) - J_m(i_k)) \right. \\ &\quad + (g(i_{k+1}, i_{k+2}) + J_m(i_{k+2}) - J_m(i_{k+1})) \\ &\quad \cdot \\ &\quad \cdot \\ &\quad \cdot \\ &\quad \left. + (g(i_{N-1}, i_N) + J_m(i_N) - J_m(i_{N-1})) \right] \end{aligned}$$

Note:  $J_m(i_N) = 0$  (why?  $i_N = T$ )

$$\mathcal{T}_{\text{mei}}(i_k) = \mathcal{T}_m(i_k) + \beta_m (d_k + d_{k+1} + \dots + d_{N-1}),$$

(\*)

where

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

(\*) Can be done incrementally as

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

$$\mathcal{T}_{\text{mei}}(i_k) = \mathcal{T}_m(i_k) + \beta_m d_k$$



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

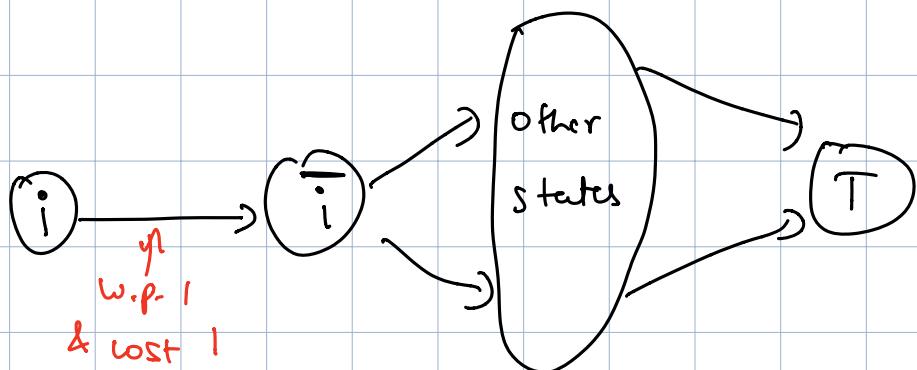
Remark: TD(0) updates as

$$\mathcal{T}_{\text{mei}}(i_k) = \mathcal{T}_m(i_k) + \beta_m d_k$$

$(d_{k+1}, \dots, d_N)$

TD(0) does not use all the temporal differences  $\underline{\text{to}}$  update its estimate  $\mathcal{T}_m(i_k)$ , & instead relies on a (-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  $\{\hat{J}(m)\}_{m=1}^N$

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

$\downarrow$   
 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 (e.g. 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.

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

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

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

Using the above, we can have the following iterative algorithm

$$\begin{aligned} \mathcal{T}_{m+1}(i_k) &= \mathcal{T}_m(i_k) + \beta_m(g(i_k, i_{k+1}) + g(i_{k+1}, i_{k+2}) \\ &\quad + \mathcal{T}_m(i_{k+2}) - \mathcal{T}_m(i_k)) \end{aligned}$$

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

$$\mathcal{T}_\pi(i) = E \left( \sum_{m=0}^l g(i_m, i_{m+1}) + \mathcal{T}_\pi(i_{l+1}) \mid i_0 = i \right)$$

$(l+1)$ -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?

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

$l$ -step fixed point target.

$\lambda$   $\in [0, 1]$

normalize (since  $\sum l \lambda^l = \frac{1}{(1-\lambda)^2}$ )      weight for a particular  $l$

We use the above fixed point equation to arrive at the  $TD(\lambda)$  update rule on Wednesday.

### Lecture-23\*

The fixed point equation that serves as a basis for  $TD(\lambda)$ :

$\lambda$ -step fixed point target.

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

↑                          ↓

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

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

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

(A)                          (B)



Simplifying the term (A):

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

Since  $(1-\lambda) \sum_{l=m}^{\infty} \lambda^l = \lambda^m$ .

Simplifying the term (B):

$$E \left( \sum_{\ell=0}^{\infty} (\ell - \lambda^{\ell+1}) J_{\pi}(i_{\ell+1}) \right)$$

$$= E \left[ (1-\lambda) J_{\pi}(i_1) + (\lambda - \lambda^2) J_{\pi}(i_2) + (\lambda^2 - \lambda^3) J_{\pi}(i_3) + \dots \right]$$

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

$$= E \left( \sum_{m=0}^{\infty} \lambda^m (J_{\pi}(i_{m+1}) - J_{\pi}(i_m)) \right) + J_{\pi}(i) \quad (2)$$

Combining ① & ②, we obtain

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

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

So, (xx)  $\Leftrightarrow$   $J_{\pi}(i) = E \left( \sum_{m=0}^{\infty} \lambda^m d_m \right) + J_{\pi}(i) \quad -B$

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

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

How to turn ③ into an iterative update rule?

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

$$\mathcal{T}_{t+1}(i) = \mathcal{T}_t(i) + \beta_t \sum_{m=0}^{\infty} \lambda^m d_m \quad \xrightarrow{\text{update iteration.}} \text{TD}(i)$$

(4)

Remark:

① If  $\lambda = 1$ , then ④ becomes

$$\mathcal{T}_{t+1}(i) = \mathcal{T}_t(i) + \beta_t \sum_{m=0}^{\infty} d_m$$

$$\mathcal{T}_{t+1}(i) = \mathcal{T}_t(i) + \beta_t (d_0 + d_1 + \dots + d_{n-1})$$



MCPE scheme  $\Rightarrow \text{TD}(i)$

② If  $(\lambda = 0)$ , then ④ becomes (w/  $d_0 = 1$ )

$$\mathcal{T}_{t+1}(i) = \mathcal{T}_t(i) + \beta_t d_0$$

$$\mathcal{T}_{t+1}(i) = \mathcal{T}_t(i) + \beta_t (g(i_0, i_1) + \mathcal{T}_t(i_1) - \mathcal{T}_t(i_0))$$

Note:  
 $i_0 = i$

## TD(0) update iteration

Note to typesetters: See if you can make  $i$  in  $\hat{J}_\pi(i) | \hat{J}_t(i)$  as  $i_0$ .

- ③ For any  $0 < \lambda < 1$ , the temporal difference  $\delta_m$  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  $\delta_0$  is more pronounced in  $\text{TD}(\lambda)$  update as compared to  $\delta_m, m > 1$ .

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

## 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.

$(i_0, i_1, \dots, i_k, \dots, i_n)$

both are  $= i$

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

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

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

$$\mathcal{T}_{t+1}(i) = \mathcal{T}_t(i) + \beta_t \sum_{j=1}^M \sum_{m=m_j}^{\infty} \lambda^{m-m_j} d_m$$

Why this?

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

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

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

$$\mathcal{T}_{t+1}(i) = \mathcal{T}_t(i) + \beta_t \sum_{m=m_1}^{\infty} \lambda^{m-m_1} d_m$$

Question: Are these two variants equivalent?

No.

But, both variants can be shown to converge. (See Section 5.2)  
of ND 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- $\lambda$ :

$$\mathcal{T}_{t+1}(i) = \mathcal{T}_t(i) + \beta_t \sum_{j=1}^M \sum_{m=m_j}^{\infty} \lambda^{m-m_j} d_m$$

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

$$O_n(i_0, i_1): \quad \mathcal{T}_1(i_0) = \mathcal{T}_0(i_0) + \beta d_0$$

using constant step size for simplicity

$$\text{on } (i_1, i_2) : \quad \mathcal{T}_2(i_0) = \mathcal{T}_1(i_0) + \beta \lambda d_1,$$

$$\mathcal{T}_2(i_1) = \mathcal{T}_1(i_1) + \beta d_1,$$

and so on.

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

$$\mathcal{T}_{t+1}(i_0) = \mathcal{T}_t(i_0) + \beta \lambda^k d_k$$

$$\mathcal{T}_{t+1}(i_1) = \mathcal{T}_t(i_1) + \beta \lambda^{k-1} d_k$$

$$\vdots$$

$$\mathcal{T}_{t+1}(i_k) = \mathcal{T}_t(i_k) + \beta d_k$$

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

### Lecture-24\*

Example to illustrate the updates of offline & online  $\text{TD}(\lambda)$ .

Suppose we have a trajectory  $(1, 2, 1, T)$

Let  $\mathcal{T}_0(1)$  &  $\mathcal{T}_0(2)$  be initial values. ( $\mathcal{T}_0(T)=0$ )

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

State 1's update: (Every visit style)

$$(1, 2, 1, T) \text{ & } (1, T)$$

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

$$= J_0(1) + \beta \left( (g(1, 2) + J_0(2) - J_0(1)) \right.$$

$$\quad \quad \quad + \lambda (g(2, 1) + J_0(1) - J_0(2)) \right.$$

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

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

} from  
(1, 2, 1, T)

} from sub-trajectory  
(1, T)

$$J_f(2) = J_0(2) + \beta \left( (g(2, 1) + J_0(1) - J_0(2)) \right.$$

$$\quad \quad \quad + \lambda (g(1, T) - J_0(1)) \right)$$

} from  
(2, 1, T)

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

$$\text{on } (1, 2)^\circ$$

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

$$J_f(2) = J_0(2)$$

On  $(2, 1)$ :

$$\hat{\tau}_2(1) = \tau_1(1) + \beta \lambda (g(2,1) + \tau_1(1) - \tau_1(2))$$

$$\hat{\tau}_2(2) = \tau_1(2) + \beta (g(2,1) + \tau_1(1) - \tau_1(2))$$

On  $(1, T)$ :

$$\hat{\tau}_3(1) = \hat{\tau}_2(1) + \beta (\lambda^2 (g(1,T) - \hat{\tau}_2(1)) + (g(1,T) - \hat{\tau}_2(1)))$$

$$\hat{\tau}_3(2) = \hat{\tau}_2(2) + \beta \lambda (g(1,T) - \hat{\tau}_2(1))$$

$(\hat{\tau}_3(1), \hat{\tau}_3(2)) \rightarrow$  Online TD( $\lambda$ ) estimates.

Compare this with  $(\tau_f(1), \tau_f(2))$ :

If we replace  $\tau_1$  &  $\tau_2$  by  $\tau_0$  in online TD update,

$$\text{Then } (\hat{\tau}_3(1), \hat{\tau}_3(2)) = (\tau_f(1), \tau_f(2))$$

$\tau_1$  &  $\tau_0$  difference is  $O(\beta)$

$\tau_2$  &  $\tau_0$  difference is  $O(\beta^2)$

$\tau_3$  &  $\tau_0$  —————  $O(\beta^3)$

or  $O(\beta^1)$

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

## Convergence of TD: A high-level sketch

Recall

$$J_{t+1} = J_t + \beta_t (H J_t + w_t - J_t) \quad (\star)$$

If (i)  $H$  is a contraction,

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

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

then  $J_t \rightarrow J^*$ , where  $H J^* = J^* - \textcircled{1}$

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

Condition (ii) 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.

Recall the fixed point equation underlying TD(λ):  
 l-step fixed point target.

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

↑  
normalize (since  $\sum \lambda^l = \frac{1}{1-\lambda}$ )      weight for a particular  $l$

$\lambda \in [0, 1]$

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

↑  
transition probability matrix underlying the policy considered.

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

For a proper policy, we showed earlier that there exists a positive vector  $\xi$  and some  $\epsilon \in (0, 1)$  s.t.

$$\sum_{j=1}^n p_{ij}(\pi(i)) \xi(j) \leq \xi(i) - \epsilon \leq e \xi(i),$$

Or, equivalently

$$\|P\mathcal{J}\|_{\xi} \leq e \|\mathcal{J}\|_{\xi}, \text{ where } \|\cdot\|_{\xi}$$

is the weighted max-norm.

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

Want  $\|H\mathcal{J} - H\mathcal{J}'\|_{\xi} \leq \epsilon \|\mathcal{J} - \mathcal{J}'\|_{\xi}$ . — (xx)

$\uparrow$   
 $0 < \epsilon < 1$

If (\*\*) holds, then the mapping

underlying  $T\mathcal{D}(\lambda)$  is a contraction &

we can claim convergence using ①

$$\|H\mathcal{J} - H\mathcal{J}'\|_{\xi}$$

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

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

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

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

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

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

So, we have

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

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

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

$$\pi_t \rightarrow \pi^* \text{ a.s. as } t \rightarrow \infty$$

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

## Lecture-25

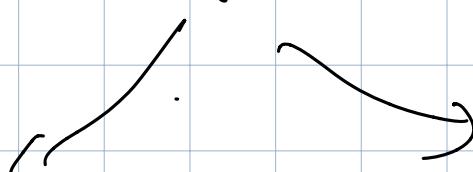
### TD( $\lambda$ ) for discounted MDPs

Policy evaluation using  $TD(\lambda)$

Approach I

Since there is no termination state,

to receive the  $TD(\lambda)$  idea for SSPs,



Approach II

Do something else

Convert a discounted MDP to

its equivalent SSP.

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

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

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

**Simulation:** Toss a coin with bias  $\lambda$  in each time instant

If heads, continue simulation.

Else, end the trajectory & do  $TD(\lambda)$  update.

$\rightarrow (i_0, i_1, \dots, i_N)$   
 $\nwarrow$   
termination stat.

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

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

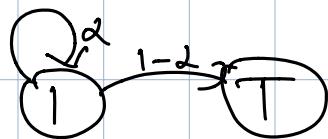
Draw back of this approach: high variance

**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"

$$\mathcal{I}_{\pi}(1) = 0$$



From a trajectory we got single stage costs

$$g_1, g_2, \dots, g_N \quad (\text{N is c.r.v.})$$

Variance of this total cost sample

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

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

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

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

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

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

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

## Approach 2:

1-step fixed point equation  $J_{\pi}(i) = E[g(i, \tilde{i}) + \alpha J_{\pi}(\tilde{i})]$

2-step  $\rightarrow$   $J_{\pi}(i) = E[g(i, \tilde{i}) + \alpha g(\tilde{i}, \tilde{\tilde{i}}) + \alpha^2 J_{\pi}(\tilde{\tilde{i}})]$

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

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

Repeating all the steps from the SSP TD( $\lambda$ ) derivation

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

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

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

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

In comparison to SSP case,  
we have the  $\lambda$  factor here.

Ques: When to end trajectories? To be answered.

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 =  $g_1 + \lambda g_2 + \lambda^2 g_3 + \dots$

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

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

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

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

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

Approach 2 → Offline TD( $\lambda$ ) : cut the trajectory at some random time  $\tau$   
 "Simulate for a fixed ( $\tau$ ) # 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.

On line TD( $\lambda$ )

update value function estimated on every sample transition.

Update rule ← very similar to the SSP Case.

Just patch in " $\lambda$ ".

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

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

Take samples of this expectation  
with a truncated trajectory.

If  $\gamma$  is large enough, then intuitively,

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

the total cost is negligible &

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

enough approximation.

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

Check Prop 5.1 of NDF book for details.

## Lecture-26\*

### Q-learning

Recall Q-factors:

$$Q^*(i, a) = \sum_j p_{ij}(a) (g(i, a, j) + \overbrace{\tau^*(j)}^{\text{optimal cost starting in } j})$$

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

Bellman equation  $\tau^* = T\tau^*$  is equivalent to

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

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

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

(★)

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

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

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

Then,  $Q = Q^*$ .

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

$Q$  solves (1), so

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

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

$J^*$  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

in the operator underlying Q-Bellman equation

Value iteration(VI) using Q-factors:

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

(This is like  $Q_{t+1} = H(Q_t)$ , starting with some  $Q_0$ )

A variation to VI is

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

VI requires knowledge of those transition probabilities

Sto-iter-algo version of the above:

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

( $\bar{i}$  → sampled from  $p_{ij}(a)$ )

This is the Q-learning algorithm.

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

Remarks: 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  $(\lambda+1)$ -step Q-Bellman equation.  
 (Think about this!)

## Convergence analysis of Q-learning:

Note: Convergence of Q-learning  $\Rightarrow$  convergence of TD(0)

(just consider a special MDP  
 with only  
 flexible action  $\pi(i)$   
 in state  $i$ .

$$\text{Then } Q\text{-BE} \Leftrightarrow (J^\pi = T^\pi J^\pi)$$

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

$\uparrow$   
 Sampled for  
 $p_{i,j}(a)$

$b$   $\rightarrow$  Action in next  
 state  $\bar{i}$

Assumptions:

$$(A_1) \quad \sum \beta_t = \infty, \quad \sum \beta_t^2 < \infty$$

(A<sub>2</sub>) All policies are proper

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

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

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

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

---

Main proof (Theorem comes later):

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

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

So, Q-learning update iteration is equivalent to

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

L (\*\*)

Verifying conditions on noise:

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

$$Y = g(i, a, \bar{i}) + \min_b Q_t(\bar{i}, b)$$

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

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

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

$$\begin{aligned} E(Y^2 | \mathcal{F}_t) &= E\left(\left(g(i, \bar{i}) + \min_b Q_t(\bar{i}, b)\right)^2 | \mathcal{F}_t\right) \\ &\leq K \left(1 + \max_{j,b} Q_t^2(j, b)\right) \end{aligned}$$

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

So, we have verified (B2)

(B3) is satisfied if we chose  $\beta_t$  carefully.

Onto (B1):  $H$  is a contraction.

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

Fact A

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

← from SSP  
Chapter.

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

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

If this holds, then (B1) is satisfied.

Pf of (need to show):

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

$$|(HQ)(i,a) - (HQ')(i,a)|$$

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

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

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

Take max over  $j$  as well.

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

Fact A  $\Rightarrow$

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

So, we get

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

$\Rightarrow$

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

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

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

$\Rightarrow$   $H$  is a contraction wrt  $\|\cdot\|_{\xi}$

Thus, we have

Theorem ( $Q$ -learning convergence)

(A1) All policies proper (A2)  $\sum \beta_t^2 < \infty$ ,  $\sum \beta_t < \infty$

(A3) Single stage cost  $g$  is bounded.

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 "There is a proper policy" & "Improper policies have infinite cost".

Even here  $Q^*$  is the unique solution to the Q-Bellman equation.

Question: Does  $\theta$ -learning converge in this case?

(A1')  $\exists$  a proper policy & all improper policies have infinite cost.

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

i.e.,  $Q \leq Q' \Rightarrow HQ \leq HQ'$

"Check this using definition of  $H$ ".  
H.W.

Also check  $J_r$

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

$c$  = vector of all ones.

Theorem's  
 (Q-learning  
 convergence  
 under  
 monotonicity)

(A1') + "  $\sum \beta_t = \infty$ ,  $\sum \beta_t^2 < \infty$ " +  
 " bounded stage cost"  
 + " bounded iterate i.e.,  $\sup_{t,i,a} |Q_t(i,a)| < \infty$

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

where  $HQ^* = Q^*$ .

} very 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:

$$Q^*(i, a) = \sum_j p_{ij}(a) \left( g(i, a, j) + \gamma \min_b Q^*(j, b) \right)$$

discount (A)

VI:

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

Q-learning:

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

sampled from  $p_{i\tilde{j}}(a)$

Convergence of Q-learning:

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

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

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

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

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

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

↑  
max norm

So, Under  $(A1'), (A2')$ , we have

$$Q_t \rightarrow Q^* \text{ w.p.1 as } t \rightarrow \infty$$

Issue of exploration:

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



trajectories to do policy evaluation

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

then  $\hat{f}(i) \rightarrow T_\pi(i)$  if

you see enough trajectories starting with "i".

$\hat{f}(i)$ : simple average  $\xrightarrow{\text{converge}}$   $T_\pi(i)$  if we

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

Same logic applies to TD and Q-learning

With TD, for  $\pi_t(i) \rightarrow \pi_{\pi}(i)$  as  $t \rightarrow \infty$

we need to visit  $i$  "i.0" in the trajectories.

But, sampling is w.r.t a fixed policy  $\pi$ .

Compare with Q-learning:

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

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

## Lecture-28\*

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  $J^*$ ?

$$\text{Normal Bellman equation: } J^*(i) = \min_a \sum_j p_{i,j}(a) (q(i,a,j) + \gamma J^*(j)) \rightarrow \textcircled{1}$$

$$\text{Q-Bellman equation: } Q^*(i,a) = \sum_j p_{i,j}(a) (q(i,a,j) + \gamma \min_b Q^*(j,b)) \rightarrow \textcircled{2}$$

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

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

$$\text{Eq } \textcircled{2} \text{ is in the form } Q^*(i,a) = E_j \underbrace{(q(i,a,j) + \gamma \min_b Q^*(j,b))}_b$$

or sample this

& do an iterative update.

Eq  $\textcircled{1}$  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:

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

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".

$Q(1, a_1)$
$\vdots$
$Q(1, a_m)$
$\hline$
$Q(2, a_1)$
$\vdots$
$Q(2, a_m)$
$\hline$
$Q(I, a_1)$
$\vdots$
$Q(I, a_m)$

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

Q-table

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.

$$\frac{\exp(-Q_t(i,a)/\tau)}{\sum_b \exp(-Q_t(i,b)/\tau)}$$

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".