

Lecture-28\*  
(continued)

## Reinforcement learning with function approximation

Why approximation!

TD/Q-learning: we look-up table representation  
i.e., need an entry  $Q(i, a)$  or  $\pi(i)$   
for every state  $i$  & action  $a$ .

On MDPs with large state space, these algorithms  
may not even be implementable.

e.g. Go:  $10^{170}$  states.

Other practical applications have large state spaces.

Practical alternative: Approximate  $\pi^\pi$  or  $Q^\pi$ .

Let's look at approximation in value space i.e.,

$$\pi^\pi(i) \approx \tilde{\pi}(i, \gamma)$$

$\uparrow$   
parameter.

E.g., Linear function approximation

$$\tilde{\pi}(i, \gamma) = \phi(i)^T \gamma$$

$\uparrow$   
feature vector

$\phi(i) \in \mathbb{R}^d$ ,  
 $\gamma \in \mathbb{R}^d$ ,  
parameter  $|X| \gg d$

so, no lookup tabu. In place of  $\tilde{J}^*(i)$ , we use  $\tilde{J}(i, r)$

Question: What "r" to we in the approximation?

What features  $\phi(\cdot)$  to employ  $\rightarrow$  out of scope of this course.

$\tilde{J}$   $\rightarrow$  non-linear function of features.

The question of feature selection is orthogonal to choosing the best parameter  $\alpha$  by using these approximations. We focus on the latter.

---

Suppose we use  $\tilde{J}$  in place of  $J^*$  & pick actions using a greedy policy:

$$\pi(i) = \arg \min_a \sum_j p_{ij}(a) (g(i, a, j) + \alpha \tilde{J}(j, r))$$

If  $\tilde{J}$  is close to  $J^*$ , then is  $\pi$  close to  $\pi^*$ ?

Prop 1: (Ref: Chap. 6 of MDP book)

$\alpha \rightarrow$  discount factor Assume finite state & finite action space.

$\|\cdot\|_\infty \rightarrow$  max-norm. We have a discounted MDP.

Suppose we have a vector  $J$  s.t.  $\|J - J^*\|_\infty = \epsilon$ ,  $\epsilon > 0$

If  $\pi$  is a greedy policy based on  $J$ , then

$$\|J_\pi - J^*\| \leq \frac{\alpha \epsilon}{1-\alpha}.$$

Further, one can choose an  $\epsilon_0$  s.t.  $\forall \epsilon < \epsilon_0$ ,

$\pi$  is an optimal policy.

pf:

$$\|J_\pi - J^*\|_\infty = \|T_\pi J_\pi - T_\pi J^*\|_\infty$$

since  $J_\pi$  is the fixed pt of  $T_\pi$

$$D^L \text{ ineq. } \rightarrow \|T_\pi J_\pi - T_\pi J\|_\infty + \|T_\pi J - J^*\|_\infty$$

because  $\pi$  is  
greedy wrt  $J$ ,  
we have  
 $T_\pi J = TJ$

$$T_\pi \text{ is } \alpha \text{-contraction} \quad \|J\|_\infty \rightarrow \leq \alpha \|J_\pi - J\|_\infty + \|TJ - J^*\|_\infty$$

$$J^* \subset J^* \text{ is } \alpha \text{-contraction} \quad \|J\|_\infty \rightarrow \leq \alpha \|J_\pi - J\|_\infty + \alpha \|J - J^*\|_\infty$$

$$\|J_\pi - J\|_\infty \leq \alpha \|J_\pi - J\|_\infty + \alpha \epsilon$$

$$\text{So, } \|J_\pi - J\|_\infty \leq \frac{\alpha \epsilon}{1-\alpha}$$

$$\text{Let } \delta = \min_{\pi'} \|J^{\pi'} - J^*\|_\infty$$

min attained since # of policies finit. So  $\delta > 0$ .

Choose  $\epsilon$  s.t.  $\frac{\alpha}{1-\alpha} < \delta$  & let  $\bar{\pi}$  be the greedy policy with this  $\epsilon$ .

$$\text{Then } \|\pi_{\bar{\pi}} - \pi^*\| < \delta \Rightarrow \bar{\pi} = \pi^*. \blacksquare$$

Approximate policy evaluation  
using TD-type algorithms

Ref: DPOC-vol. II, 4th edition  
Section 6.3

Consider a finite-state MDP.

Fix policy  $\pi \in \omega$  we shall not attach " $\pi$ " to the symbols used & keep the policy implicit.

States = {1--n}

Transition probabilities =  $P_{ij}$  (skipping  $\pi$ )  
in notation here

Fix  $\pi \Rightarrow$  we have a Markov chain at hand.

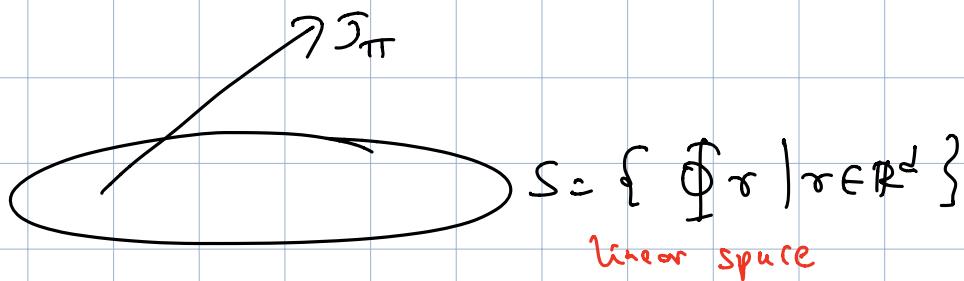
Aim: Estimate  $\pi_{\pi}(i) = E \left( \sum_{k=0}^{\infty} \alpha^k g(i_k, i_{k+1}) \mid i_0=i \right)$   
no  $\pi$  here, it is implicit

in Reaction chanc.

Linear function approximation:

$$\tilde{f}(i, \tau) = \phi(i)^T \tau, \quad i=1 \dots n$$

$\phi(i) \in \mathbb{R}^d$ ,  $\tau \in \mathbb{R}^d$ , "n >> d".



Let  $\Phi = \begin{bmatrix} -\phi(1)^T & - \\ -\phi(2)^T & - \\ \vdots & \\ -\phi(n)^T & - \end{bmatrix}$  Feature matrix

$$\Phi = \begin{bmatrix} \phi_1(1) & \dots & \phi_d(1) \\ \vdots & & \vdots \\ \phi_1(n) & \dots & \phi_d(n) \end{bmatrix}$$

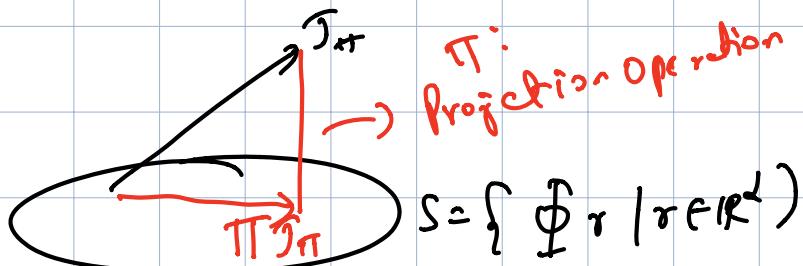
nxd matrix → a tall one.

$$\tilde{f}_\tau = [\tilde{f}(1, \tau) \dots \tilde{f}(n, \tau)]$$

$\tilde{f}_\tau = \Phi \tau$

Aim: Find the best approximation to  $\mathcal{J}_\pi$  in the space  $S = \{ \Phi^\top r \mid r \in \mathbb{R}^d \}$

### Lecture-30 (continued after Markov chains review)



For regular policy evaluation, one solves

$$\mathcal{J}_\pi = T_\pi \mathcal{J}_\pi \leftarrow \text{fixed point relation}$$

$$T_\pi \approx \Phi^\top r$$

doesn't make sense  
since  $T_\pi$  would map  
any  $\Phi^\top r$  outside  $S$ .

$$\Phi^\top r = T_\pi (\Phi^\top r)$$

Projected fixed point equation

$$\Phi^\top r = \underbrace{\Pi_{\mathcal{T}} T_\pi (\Phi^\top r)}_{\text{projection onto linear space } S} - (\infty)$$

Projected fixed point equation

$\Pi$   $\rightarrow$  definition requires stationary distribution of the markov chain underlying policy  $\pi$ .

(we solve  $\pi^*$  if " $\pi \pi^*$ " is a contraction.  
 & we will show it is the case.

### Lecture-31\*

TD with linear function approximation

#### Assumptions:

(C1) The Markov chain underlying policy  $\pi$  is  
 irreducible, positive recurrent

i.e.,  $\exists$  a stationary distribution  $\{\xi_1, \dots, \xi_n\}$   
 for this chain [ $\xi = \xi P$ ]

(C2) The matrix  $\Phi$  has full column rank

or  $\text{rank}(\Phi) = d$  [Note: we assume  
 $n \gg d$ ]

Towards a projected fixed point equation! (Policy  $\pi$  fixed throughout)

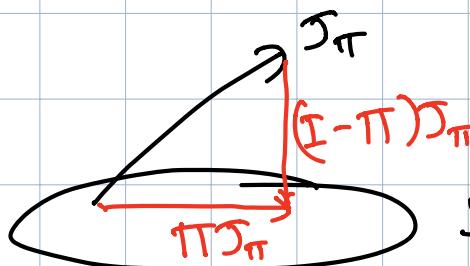
for a  $n$ -vector  $J$ , define  $\rightarrow$  Stationary distribution

$$\|J\|_{\xi}^2 = \sum_{i=1}^n \xi_i (J(i))^2$$

$\rightarrow$  weighted norm

$$D = \begin{bmatrix} \xi_1 & & \\ & \xi_2 & \\ & & \ddots & \\ & & & \xi_n \end{bmatrix}$$

$$\|\mathbf{J}\|_{\xi}^2 = \mathbf{J}^T D \mathbf{J}$$



$$S = \{ \Phi r \mid r \in \mathbb{R}^2 \}$$

$\Pi$ : Projection operator  
 "Projection is orthogonal & performed using  
 $\|\cdot\|_{\xi}$ ".

$\Pi J_{\pi}$  is the "unique" vector in  $S$  that

minimizes  $\|J_{\pi} - \tilde{J}\|_{\xi}$ , over all  $\tilde{J} \in S$ .

Any  $\tilde{J} \in S$  is of the form  $\Phi r$

$$\text{So, } \gamma^* = \underset{r \in \mathbb{R}^2}{\operatorname{argmin}} \|J_{\pi} - \Phi r\|_{\xi}^2$$

$$\text{and } \Pi J_{\pi} = \Phi \gamma^*$$

$\boxed{\begin{array}{l} \text{for } \tilde{J} \in S \\ \text{can be uniquely written as} \\ \text{some } \Phi r \text{ since} \\ \Phi \text{ is full column rank} \end{array}}$

$\boxed{\begin{array}{l} \text{minimize } \|J_{\pi} - \tilde{J}\|_{\xi} \\ \text{but, every } \tilde{J} = \Phi r \\ \text{unique} \end{array}}$

$\boxed{\begin{array}{l} \text{distance } \|J_{\pi} - \tilde{J}\|_{\xi} \\ \text{minimized} \end{array}}$

$\boxed{\begin{array}{l} \text{Take any } \tilde{J} \in S \\ S = \{\Phi r\} \end{array}}$

To find  $r^*$ :

Cannot solve  $\Phi r = T(\Phi r)$

$$\nabla_r \parallel J_\pi - \Phi r^* \|_2^2 = 0$$

$$(=) \quad \nabla_r (J_\pi - \Phi r^*)^\top D (J_\pi - \Phi r^*) = 0$$

$$(=) \quad \boxed{\Phi^\top D (J_\pi - \Phi r^*) = 0} \quad - (x)$$

$$(=) \quad \Phi^\top D J_\pi - \Phi^\top D \Phi r^* = 0$$

$$(=) \quad \Phi^\top D \Phi r^* = \Phi^\top D J_\pi$$

$$(=) \quad r^* = (\Phi^\top D \Phi)^{-1} \Phi^\top D J_\pi$$

Why is this invertible?

$\Phi$  → full rank and

diagonal elements of  $D > 0$ .

$$\text{Now, } \Pi J_\pi = \Phi r^*$$

$$(\uparrow) \quad \Pi = \Phi (\Phi^\top D \Phi)^{-1} \Phi^\top D$$

projection operator

diagonal matrix with  
stationary distribution values

Can be outside

So, project  $T(\Phi r)$   
i.e.,  $\Pi(T(\Phi r))$   
& then solve

$$\Phi r = \Pi(T(\Phi r))$$

$\Phi$  is  $n \times d$

$D$  is  $n \times n$

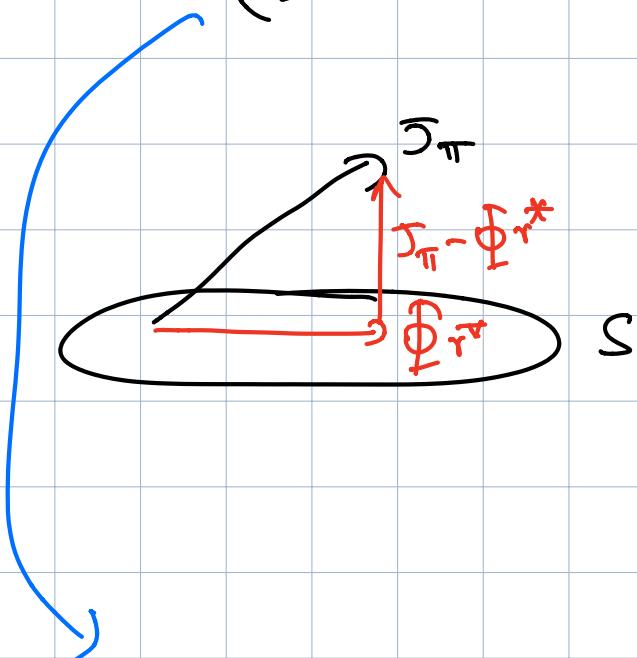
$r$  is  $d \times 1$

From (\*),

$$\Phi^T D (\mathcal{T}_\pi - \Phi_{\gamma^*}) = 0$$

$$(\Rightarrow) \quad \gamma^{*\top} \Phi^T D (\mathcal{T}_\pi - \Phi_{\gamma^*}) = 0$$

$$(\Rightarrow) \quad (\Phi_{\gamma^*})^\top D (\mathcal{T}_\pi - \Phi_{\gamma^*}) = 0$$



This equation ( $\Rightarrow$ )

$$\Phi_{\gamma^*} \perp_{\text{or}} \mathcal{T}_\pi - \Phi_{\gamma^*}$$

orthogonal  
in  $\|\cdot\|_E$  norm

$$(\Rightarrow) \quad \langle (\Phi_{\gamma^*}), \mathcal{T}_\pi - \Phi_{\gamma^*} \rangle = 0$$

$$\langle x, y \rangle = \sum_{i=1}^n x_i y_i$$

Recall, for a bounded  $\mathcal{T} = (\mathcal{T}(1), \dots, \mathcal{T}(n))$

$$(\mathcal{T}_\pi \mathcal{T})(i) = \sum_{j=1}^n p_{ij} (g(i,j) + \alpha \mathcal{T}(j)), \quad \forall i$$

In compact notation :  $\mathcal{T}_\pi \mathcal{T} = g + \alpha P \mathcal{T}$

$$g = (g_1, \dots, g_n) \quad g_i = \sum_j p_{ij} g^{(i,j)}$$

$$P = \begin{bmatrix} p_{11} & & \\ & \ddots & \\ & & p_{nn} \end{bmatrix}$$

Projected fixed point equation:

$$\tilde{\Phi}^{\tau^*} = \Pi \circ T_\Pi(\tilde{\Phi}^{\tau^*})$$

Here  $\Pi T$  is the composition of  $\Pi$  with  $T$ .

IF:  $\Pi T$  is a contraction w.r.t  $\|\cdot\|_S$ , then deriving  
VI or sto-iter-algo variations are straightforward.

" $\Pi T$  is contractive"

Letting  $\tilde{\tau} = \tilde{\Phi}^{\tau^*}$ , we have

$$\tilde{\tau} = \Pi T_\Pi(\tilde{\tau}) \quad \leftarrow \text{Projected eqn}$$

Contrast with regular fixed point equation:

$$\tau_\Pi = T_\Pi \circ \tau$$

$\Pi$ : come in because  $\tilde{\tau} \in S$  "linear space".

Lemma 1:

$$\|P\pi\|_{\xi} \leq \|\pi\|_{\xi} \quad \forall \pi \in \mathbb{R}^n$$

↓  
f.p.m. of M.C.  
underlying policy  $\pi$

↓  
Stationary distribution vector

Pf:

$$\|P\pi\|_{\xi}^2 = \sum_{i=1}^n \xi_i \left( \sum_{j=1}^n p_{ij} \pi(j) \right)^2$$

Jensen's inequality  
 $p_{ij}^2 \leq p_{ij}$

$$\leq \sum_{i=1}^n \xi_i \sum_{j=1}^n p_{ij} \pi(j)^2$$

$$= \sum_{j=1}^n \left( \sum_{i=1}^n \xi_i p_{ij} \right) \pi(j)^2$$

Using  $\xi = \xi P$   
 $\sum_{i=1}^n \xi_i p_{ij} = \xi_j$   
 $\sum_{i=1}^n \xi_i = 1$

$$= \sum_{j=1}^n \xi_j \pi(j)^2$$

$$= \|\pi\|_{\xi}^2$$

$$\text{So, } \|P\pi\|_{\xi}^2 \leq \|\pi\|_{\xi}^2$$

■

"Projection is non-expansive"

Lemma 2:

$$\|\pi\pi\pi - \pi\pi\pi'\|_{\xi} \leq \|\pi\pi\pi' - \pi\pi\pi'\|_{\xi}, \quad \forall \pi, \pi' \in \mathbb{R}^n$$

Pf:

PTO

$$\|\Pi \mathbf{J} - \Pi \mathbf{J}'\|_{\xi}^2 = \|\Pi(\mathbf{J} - \mathbf{J}')\|_{\xi}^2$$

Pythagoras  
theorem  
 $\Pi \mathbf{J}$  &  $(I - \Pi)$   
are orthogonal

$$\leq \|\Pi(\mathbf{J} - \mathbf{J}')\|_{\xi}^2 + \|(I - \Pi)(\mathbf{J} - \mathbf{J}')\|_{\xi}^2$$

(\*)

Note:

$$\underbrace{\Pi(\mathbf{J} - \mathbf{J}')}_{\in S} \perp \underbrace{((\mathbf{J} - \mathbf{J}') - \Pi(\mathbf{J} - \mathbf{J}'))}_{\text{orthogonal to } S}$$

$$\begin{aligned} & \|\Pi(\mathbf{J} - \mathbf{J}')\|_{\xi}^2 + \|(I - \Pi)(\mathbf{J} - \mathbf{J}')\|_{\xi}^2 \\ &= \|\Pi(\mathbf{J} - \mathbf{J}') + (I - \Pi)(\mathbf{J} - \mathbf{J}')\|_{\xi}^2 \\ &= \|(\mathbf{J} - \mathbf{J}')\|_{\xi}^2 \end{aligned}$$

Hence, from (\*), we obtain

$$\|\Pi \mathbf{J} - \Pi \mathbf{J}'\|_{\xi}^2 \leq \|\mathbf{J} - \mathbf{J}'\|_{\xi}^2$$

■

Main claim:

$T_{\Pi}$  and  $\Pi T_{\Pi}$  are contraction mappings w.r.t.  $\|\cdot\|_{\xi}$ , and have modulus of (dis)continuity

Pf:

$$\text{Recall } T_{\pi} \mathbf{j} = g + \alpha P \mathbf{j}$$

For any  $\mathbf{j}, \mathbf{j}' \in \mathbb{R}^n$ ,

$$\|T_{\pi} \mathbf{j} - T_{\pi} \mathbf{j}'\|_{\xi} = \alpha \|P(\mathbf{j} - \mathbf{j}')\|_{\xi}$$

$$\begin{aligned} & \text{Using Lemma 1} \\ & \|P\mathbf{j}\|_{\xi} \leq \|\mathbf{j}\|_k \end{aligned} \quad \leq \alpha \|\mathbf{j} - \mathbf{j}'\|_{\xi} \quad (\star)$$

So,  $T_{\pi}$  is contractive with modulus  $\alpha$ .

(Side note! We showed earlier that  $T_{\pi}$  is a contraction w.r.t. max-norm. Here, we showed  $T_{\pi}$  to be contractive wrt  $\|\cdot\|_{\xi}$  as well)

Next

$$\begin{aligned} & \|\Pi T_{\pi} \mathbf{j} - \Pi T_{\pi} \mathbf{j}'\|_{\xi} \\ &= \|\Pi (T_{\pi} \mathbf{j} - T_{\pi} \mathbf{j}')\|_{\xi} \end{aligned}$$

$$\begin{aligned} & \text{Since } \Pi \text{ is non-expansive (Lemma 1)} \\ & \leq \|T_{\pi} \mathbf{j} - T_{\pi} \mathbf{j}'\|_{\xi} \end{aligned}$$

$$\begin{aligned} & \text{From } (\star) \\ & \leq \alpha \|\mathbf{j} - \mathbf{j}'\|_{\xi} \end{aligned}$$

So,  $\Pi \Pi \Pi$  is contractive wrt  $\|\cdot\|_{\xi}$  with modulus  $\alpha$ . ■

Implication:

$$\Phi^r = \Pi T_\pi (\Phi^r)$$

(cannot solve  $\Pi T_\pi$ )  
in func. approx case.  
So, solve this eqn instead.

Projected equation has a "unique" solution

& we can do value iteration to get the solution

(i.e.,  $\Phi^r_0 \xrightarrow{\Pi T} \Phi^r_1 \rightarrow \dots \rightarrow$  asymptotically converges to  $\Phi^r$ )

### Lecture-32\*

Not-so Main claim!: Recall  $r^*$  is the fixed point of  $\Pi T_\pi$ , i.e.,  
 $\Phi^r = \Pi T_\pi (\Phi^r)$

$$\| J_\pi - \Phi^r \|_2^2$$

$$\begin{aligned} & \text{because } (J_\pi - \Pi J_\pi) \perp \text{es of } (\Pi J_\pi - \Phi^r) \\ & \quad \text{orthogonal to } \\ & \| J_\pi - \Pi J_\pi \|_2^2 + \| \Pi J_\pi - \Phi^r \|_2^2 \\ &= \| J_\pi - \Pi J_\pi \|_2^2 + \| \Pi \Pi J_\pi - \Pi \Phi^r \|_2^2 \\ &\leq \| J_\pi - \Pi J_\pi \|_2^2 + \lambda^2 \| J_\pi - \Phi^r \|_2^2 \end{aligned}$$

Rearranging:  $\| J_\pi - \Phi^r \|_2^2 \leq \frac{1}{1-\lambda^2} \| J_\pi - \Pi J_\pi \|_2^2$

Matrix form of  $\hat{\Phi}r^* = \Pi T(\hat{\Phi}r^*)$

$$\Pi = \hat{\Phi} (\hat{\Phi}^\top D \hat{\Phi})^{-1} \hat{\Phi}^\top D$$

$$T_\Pi J = g + \lambda P J$$

$\Pi J \rightarrow$  linear

$T_\Pi J \rightarrow$  linear

$$\hat{\Phi}r^* = \Pi T_\Pi (\hat{\Phi}r^*) \quad \text{--- (2)}$$

$\downarrow$

"linear" System of equations

Want to write (2) as

$$C r^* = d$$

$\Pi J$ : minimize distance between  $J$  &  $\tilde{J} \in S$  in  $\mathbb{R}^n$

$$r^* = \arg \min_{r \in \mathbb{R}^d} \| \hat{\Phi}r - T_\Pi (\hat{\Phi}r^*) \|_2^2$$

↑  
using (2)

$$T_\Pi r^* = \arg \min_{r \in \mathbb{R}^d} \| \hat{\Phi}r - (g + \lambda P \hat{\Phi}r^*) \|_2^2$$

$$= \arg \min_{r \in \mathbb{R}^d} \frac{1}{2} (\hat{\Phi}r - (g + \lambda P \hat{\Phi}r^*))^\top D (\hat{\Phi}r - (g + \lambda P \hat{\Phi}r^*))$$

$$D = \begin{bmatrix} \xi_1 & 0 \\ \xi_2 & -\xi_n \end{bmatrix}$$

Differentiating the expression being minimized, we obtain

$$\Phi^T D (\Phi r^* - (g + \lambda P \Phi r^*)) = 0 \quad (*)$$

Matrix form of  $\Phi r^* = T T^T (\Phi r^*)$

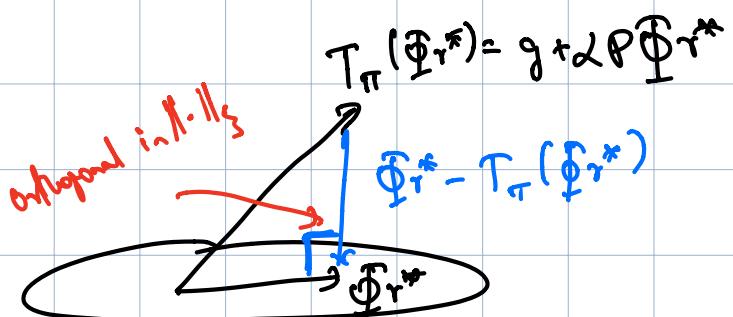
Intuitively: using (\*), we have

$$r^{*T} \Phi^T D (\Phi r^* - (g + \lambda P \Phi r^*)) = 0$$

$$(\Phi r^*)^T D (\Phi r^* - (g + \lambda P \Phi r^*)) = 0$$

$$\langle \Phi r^*, \Phi r^* - (g + \lambda P \Phi r^*) \rangle = 0$$

where  $\langle \cdot, \cdot \rangle$  leads to  $\|\cdot\|_\xi$



$$\Phi^T D (\Phi r^* - (g + \lambda P \Phi r^*)) = 0$$

$$\Phi^T D g = \Phi^T D \Phi r^* - 2 \Phi^T D P \Phi r^* = \Phi^T D (I - 2P) \Phi r^*$$

$$\Phi^T D g = \Phi^T D(I - \lambda P) \Phi r^* \quad \text{Or, equivalently}$$

$C r^* = d$ , where  $C = \Phi^T D(I - \lambda P) \Phi$ ,  $d = \Phi^T D g$

This is the same as

$$\underline{\Phi} r^* = \Pi \Pi_\pi (\underline{\Phi} r^*)$$

Explicit solution!

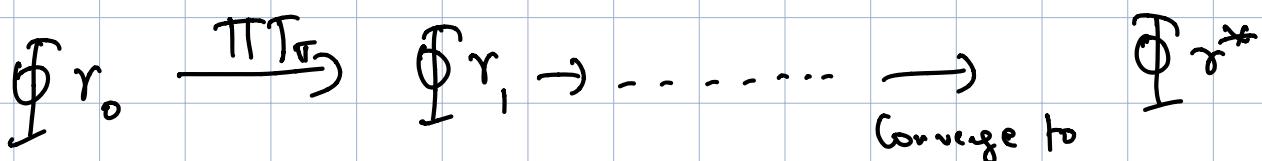
$$r^* = C^{-1} d$$

( $C$  invertible since  
 $\Phi$  full col. rank &  
 $D$  has true diagonals  
 $\Rightarrow (I - \lambda P)$  is invertible)

Projected value iteration:

We know  $\Pi \Pi_\pi$  is a  $\lambda$ -contraction

Start with  $r_0$  & repeatedly apply  $\Pi \Pi_\pi$



$$\underline{\Phi} r_{k+1} = \Pi \Pi_\pi (\underline{\Phi} r_k), \quad k=0, 1, \dots - (\star\star)$$

Value iteration + projection

PVI update ( $\star\star\star$ ) in terms of  $C, d$ :

$$r_{k+1} = \arg \min_{r \in \mathbb{R}^d} \| \hat{\Phi} r - (g + \lambda P \hat{\Phi} r_k) \|^2$$

↑ variable      ↓ fixed

Differentiating,

$$\hat{\Phi}^T D (\hat{\Phi} r_{k+1} - (g + \lambda P \hat{\Phi} r_k)) = 0 \quad \text{--- (1)}$$

$$r_{k+1} = r_k - (\hat{\Phi}^T D \hat{\Phi})^{-1} (C r_k - d) \quad \text{--- (2)}$$

Check this by substituting  $C, d$  in (2).



If the same as  $\hat{\Phi} r_{k+1} = T T_T (\hat{\Phi} r_k)$

**Remark:** For PVI, need knowledge of t.p.m.  $P$  & stationary distribution values (through  $D$ ) to form  $C$  &  $d$ , which are used in update iteration (2).

Solving  $C r^* = d$  using a sample path :

Recall  $C = \hat{\Phi}^T D (I - \lambda P) \hat{\Phi}$ ,  $d = \hat{\Phi}^T D g$

have to estimate:  $\hat{\Phi}^T D \hat{\Phi}$ ,  $\hat{\Phi}^T D P \hat{\Phi}$ ,  $\hat{\Phi}^T D g$

to form estimates of  $C, d$ .

$$\hat{\Phi} = \begin{bmatrix} \phi(1)^T \\ \vdots \\ \phi(n)^T \end{bmatrix} \quad \hat{\Phi}^T = \begin{bmatrix} \phi(1) & \dots & \phi(n) \end{bmatrix}$$

$$\hat{\Phi}^T D \hat{\Phi} = \begin{bmatrix} \phi(1) & \dots & \phi(n) \end{bmatrix} \begin{bmatrix} \xi_1 & 0 \\ 0 & \ddots & \xi_n \end{bmatrix} \begin{bmatrix} \phi(1)^T \\ \vdots \\ \phi(n)^T \end{bmatrix}$$

$$\hat{\Phi}^T D \hat{\Phi} = \sum_{i=1}^n \xi_i \phi(i) \phi(i)^T$$

recall  $P = \begin{bmatrix} p_{11} & \dots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \dots & p_{nn} \end{bmatrix}$   $g = \begin{bmatrix} g_1 \\ \vdots \\ g_n \end{bmatrix}$   $g_i = \sum_{j=1}^n p_{ij} g(i, j)$

$$\hat{\Phi}^T D P \hat{\Phi} = \sum_{i=1}^n \sum_{j=1}^n \xi_i p_{ij} \phi(i) \phi(j)^T$$

$$\hat{\Phi}^T D g = \sum_{i=1}^n \sum_{j=1}^n \xi_i p_{ij} \phi(i) g(i, j)$$

↑ unknown in an RL setting

Using  $\pi_t$ , generate a sample path  $(i_0, i_1, \dots, i_T)$   
 & some  $i_0$

Observe  $g(i_t, i_{t+1})$ ,  $\forall t$

Form sample-based estimate of  $\hat{\Phi}^T D \hat{\Phi}$ ,  $\hat{\Phi}^T D P \hat{\Phi}$  &  $\hat{\Phi}^T D g$ .

$$\Phi^T D \Phi \underset{\text{approximate}}{\approx} \sum_{k=1}^K \sum_{t=0}^K \phi(i_t) \phi(i_t)^T$$

$$\Phi^T D P \Phi \approx \sum_{k=1}^K \sum_{t=0}^K \phi(i_t) \phi(i_{t+1})^T$$

$$\Phi^T D_g \approx \frac{1}{K+1} \sum_{t=0}^K \sum_{\tau=0}^K \phi(i_t) g(i_t, i_{\tau+1})$$

Let  $C_k = \frac{1}{K+1} \sum_{t=0}^K \phi(i_t) (\phi(i_t) - \lambda \phi(i_{t+1}))^T$

$$d_k = \frac{1}{K+1} \sum_{t=0}^K \phi(i_t) g(i_t, i_{t+1})$$

$$C_k \approx C \quad (C = \Phi^T D (I - \lambda P) \Phi) \quad d_k \approx \lambda (=\Phi^T D_g)$$

LSTD:  
↓  
Solve

$$C_k r_k = d_k$$

Least-squares temporal difference.

$$C_k r_k - d_k = \frac{1}{K+1} \sum_{t=0}^K \phi(i_t) \left( \phi(i_t)^T r_k - \lambda \phi(i_{t+1})^T r_k - g(i_t, i_{t+1}) \right)$$

Temporal difference.

$$\phi(i_t)^T r_k \approx J_\pi(i_t), \quad \phi(i_{t+1})^T r_k \approx J_\pi(i_{t+1})$$

$$TD\text{-term} = \tilde{\mathcal{T}}(i_t) - (g(i_t, i_{t+1}) + \gamma \tilde{\mathcal{T}}(i_{t+1}))$$

### Lecture-33

Recall projected system of equations:

$$\begin{aligned} \gamma^* &= 1 \\ C &= \Phi^\top D (\Gamma - \gamma P) \Phi, \quad d = \Phi^\top D g \end{aligned} \quad (\Rightarrow) \quad \Phi \gamma^* = \Gamma \Gamma^\top \pi (\Phi \gamma^*)$$

$$C_1 := \Phi^\top D \Phi = \sum_{i=1}^n \sum_i \phi(i) \phi(i)^\top$$

Sample from  
 $\{x_i, \pi_i, i=1 \dots n\}$

$$C_2 := \Phi^\top D P \Phi = \sum_{i=1}^n \sum_{j=1}^n \sum_i \pi_{ij} \phi(i) \phi(j)^\top \quad \{P = \{\pi_{ij}, j=1 \dots n\}\}$$

$$d := \Phi^\top D g = \sum_{i=1}^n \sum_{j=1}^n \sum_i \pi_{ij} \phi(i) g(i, j)$$

Suppose we obtain a sample path  $\{i_0, i_1, \dots, i_k\}$   
 simulated using policy  $\pi$  & states picked according  
 to the distribution  $\{P\}$   $\Leftrightarrow$  pick an  $i_0$  from  $S$  (state set)  
 & then pick a next state w/  $\pi_{ij}$   
 & repeat.

Empirical frequencies:

$$\hat{\pi}_i = \frac{\sum_{t=0}^k \mathbb{I}(i_t = i)}{k+1}$$

$$\hat{P}_{i,j} = \frac{\sum_{t=0}^T I(i_t=i, i_{t+1}=j)}{k+1}$$

Using  $\hat{\xi}_i$ ,  $\hat{P}_{ij}$ , we estimate  $C_1, C_2, d$  as follows:

$$\hat{C}_1 = \sum_{i=1}^n \hat{\xi}_i \phi(i) \phi(i)^T$$

$$\hat{C}_2 = \sum_{i=1}^n \sum_{j=1}^n \hat{\xi}_i \hat{P}_{ij} \phi(i) \phi(j)^T$$

$$\hat{d} = \sum_{i=1}^n \sum_{j=1}^n \hat{\xi}_i \hat{P}_{ij} \phi(i) g(i, j)$$

How to approximate  $r^* = d$ ?

① Estimate  $C$  by  $C_k$  as follows:

$$C_k = \hat{C}_1 - 2\hat{C}_2$$

$$d_k = \hat{d}$$

② Solve  $C_k \hat{r}_k = d_k$



$\hat{r}_k \rightarrow$  LSTD solution.

Some analysis:-

As the trajectory length  $k$  in  $(i_0, \dots, i_k)$  goes to infinity, do the estimates  $\hat{\xi}_i, \hat{p}_{ij}$  converge?

As  $k \rightarrow \infty$ ,  
 $\hat{\xi}_i \xrightarrow{\text{w.p. 1}} \xi_i$  w.p. 1 (Version of SLLN)  
 $\hat{p}_{ij} \xrightarrow{\text{w.p. 1}} p_{ij}$

$$C_k \xrightarrow{k \rightarrow \infty} C, d_k \rightarrow d$$

$$\text{Hence, } \hat{r}_k \xrightarrow{\text{w.p. 1}} r^* \text{ as } k \rightarrow \infty$$

Remark:  $C_k$  and  $d_k$  can be written alternatively as

$$\hat{C}_1 = \perp \sum_{t=0}^k \phi(i_t) \phi(i_t)^T$$

$$\hat{C}_2 = \perp \sum_{t=0}^k \phi(i_t) \phi(i_{t+1})^T$$

$$\hat{d} = \perp \sum_{t=0}^k \phi(i_t) g(i_t, i_{t+1})$$

Another remark:

LSTD solution!

$$C_k \hat{r}_k = d_k$$

Can we write  $r_k = C_k^{-1} d_k$ ? NO.

Trivial example:  $(i_0, i_1, \dots, i_k)$

Suppose  $i_0 = i_1 = \dots = i_k = i$  (some state)

$$\hat{C}_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) \phi(i_t)^\top$$

$\hat{C}_k$  → with  $d$  columns, but each row is identical

$$\text{So } \text{rank}(\hat{C}_k) < d$$

As an alternative, solve

$$\begin{aligned} (C_k + \beta I) \hat{r}_k &= d_k \\ (\Rightarrow) \hat{r}_k &= \underbrace{(C_k + \beta I)^{-1}}_{\text{this is invertible if } \beta \text{ is large enough}} d_k \end{aligned}$$

Where is "temporal difference" in "Least squares temporal difference (LSTD)"?

$$\text{LSTD: } C_k r_k - d_k = 0$$

$$C_k r_k - d_k = \perp \sum_{t=0}^k \phi(i_t) \left[ \phi(i_t)^\top r_k - d \phi(i_{t+1})^\top r_k - g(i_t, i_{t+1}) \right]$$

$$\begin{aligned} \text{Since } C_k &= \perp \sum_{t=0}^k \phi(i_t) (\phi(i_t) - \lambda \phi(i_{t+1}))^\top, \\ d_k &= \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) g(i_t, i_{t+1}) \end{aligned}$$

the term  $\underbrace{[\phi(i_t)^T r_k - [g(i_t, i_{t+1}) + \alpha \phi(i_{t+1})^T r_k]]}_{(*)}$

is a temporal difference term because

$$\phi(i_t)^T r_k \approx \mathcal{T}_\pi(i_t)$$

$$\phi(i_{t+1})^T r_k \approx \mathcal{T}_\pi(i_{t+1})$$

So  $(*) \approx \mathcal{T}_\pi(i_t) - (g(i_t, i_{t+1}) + \alpha \mathcal{T}_\pi(i_{t+1}))$

So  $(*) \Leftrightarrow \text{TD-term.}$

### Lecture-34\*

TD(0) with linear function approximation

Recall the <sup>story of</sup> TD(0) without function approximation!

Want to solve

$$\mathcal{T}_\pi = \mathbb{E} \mathcal{T}_\pi$$

$$\mathcal{T}_\pi(i) = \underbrace{\mathbb{E}(g(i, \tilde{i}) + \alpha \mathcal{T}_\pi(\tilde{i}))}_{1 \leftarrow}$$

Sto-iter-algo for solving this eqn  
Sampled from RLS.

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

$g(i, \tilde{i}) + \alpha \mathcal{T}_t(i)$  is a proxy for  $\mathbb{E}_{\tilde{i}}(g(i, \tilde{i}) + \alpha \mathcal{T}_\pi(\tilde{i}))$

Onto linear function approximation case:

$$J_{\pi}(i) \approx r^T \phi(i)$$

Cannot do this

$$r_{t+1} \neq r_t + \beta_t (g(i, \tilde{i}) + \alpha r_t^T \phi(\tilde{i}) - r_t^T \phi(i))$$

$r_t \in \mathbb{R}^d \quad \beta_t \in \mathbb{R} \quad \in \mathbb{R}$

TD(0) with linear function approximation

On a transition  $(i_t, i_{t+1})$  in a sample path  $(i_0, \dots)$

$$\tilde{r}_{t+1} = \tilde{r}_t + \beta_t \phi(i_t) (g(i_t, \tilde{i}_{t+1}) + \alpha r_t^T \phi(\tilde{i}_{t+1}) - r_t^T \phi(i_t))$$

$\tilde{r} \in \mathbb{R}^d \quad \beta \in \mathbb{R} \quad \in \mathbb{R}$

Go back to projected fixed point:

$$\hat{\Phi} r^* = \Pi T_{\pi}(\phi r^*)$$

$$(?) C r^* = d, \quad C = \hat{\Phi}^T D (\mathbf{I} - \alpha P) \hat{\Phi}$$

$$d = \hat{\Phi}^T D g$$

Want to use a sample path  $(i_0, i, \dots)$  to find

$$r^* \text{ s.t. } C r^* = d.$$

$$\text{Sto-fir-algo: } r_{t+1} = r_t - \beta_t ((r_t - \lambda)) \rightarrow (\star)$$

Where would  $r_t$  converge? Ans: whenever  $C\gamma^* = \lambda$

$$r_t \rightarrow \gamma^*$$

$$C\gamma^* - \lambda = \hat{\Phi}^T D (I - \lambda P) \hat{\Phi} \gamma^* - \hat{\Phi} D g$$

$$= \sum_{i,j} \underbrace{\xi_i p_{ij}}_{\text{Taken with SP distribution}} \phi(i) (\phi(i)^T \gamma^* - \lambda \phi(j)^T \gamma^* - g(i, j))$$

$$= E \left( \phi(i) (\phi(i)^T \gamma^* - \lambda \phi(j)^T \gamma^* - g(i, j)) \right)$$

$\uparrow$   
Taken with SP distribution

Q: Want to find an  $\gamma^*$  s.t.

$$E \left( \phi(i) (\phi(i)^T \gamma^* - \lambda \phi(j)^T \gamma^* - g(i, j)) \right) = 0$$

$i \uparrow$  picked from  $\mathcal{S} = \{\xi_1, \dots, \xi_n\}$   
 $j \sim \underbrace{(p_{ij})}_{\text{SP}}$

Let  $(i_t, i_{t+1})$  be a sample transition

Then, do the following update

$$(\star) \rightarrow \gamma_{t+1} = r_t - \beta_t \phi(i_t) (\phi(i_t)^T \gamma_t - \lambda \phi(i_{t+1})^T \gamma_t - g(i_t, i_{t+1}))$$

$\uparrow$   
 $TD(0)$  with linear function approximation.

## Remark:

① In the above, we assumed sampling from the stationary distribution. Under this, it is straightforward to invoke "Stoermer's general convergence result" under contraction case.

② TD(0) with linear function approximation would converge even if sampling is not from the stationary distribution.

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

↑

initial state picked using some distribution " $\mu$ "  
Let the Markov chain have t.p.m.  $P$ .

Then, after  $k$  steps  $\rightarrow$  the distribution  $\mu P^k$

$$\mu P^k \rightarrow \xi \quad \text{as } k \rightarrow \infty$$

(assumes irreducible + positive recurrence)

*transient phase*

$$(i_0, i_1, \dots, i_N, i_{N+1}, \dots)$$

↑

after a large #  $N$  of iterations, the Markov chain is in steady state  
i.e., the <sup>state</sup> distribution is  $\xi$ .

4. It can be shown that (x) (TD with LFA)  
would converge to the same fixed point,  
i.e.,  $\gamma^*$  satisfying  $C\gamma^* = 1 \Leftrightarrow \Phi_{\gamma^*} T T^T \Phi_{\gamma^*} = 1$

Rif! J.N.Tsitsiklis & B.V.Roy  
"Analysis of TD with LFA",  
IEEE trans. auto. control, 1997.

(3) Can extend TD(0) to TD( $\lambda$ ) with LFA.

"Read it from NDP book or  
DPDC Vol II Chapter 6 "