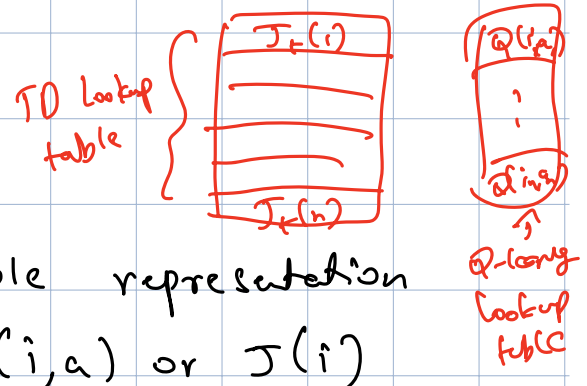


# Reinforcement learning with function approximation

Why approximation?



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

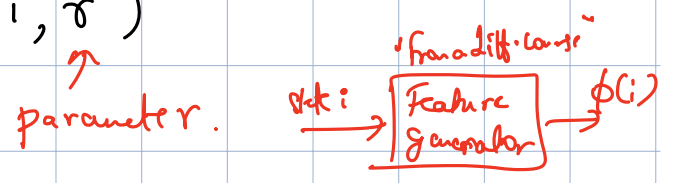
On MDPs with large state spaces, these algorithms may not even be implementable.

e.g. Go:  $10^{170}$  states, Chess  $10^{21}$  states  
 other practical applications have large state spaces.

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

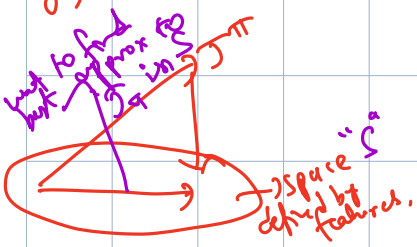
Let's look at approximation in value space i.e.,

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



parameter.

E.g., linear function approximation



$$\tilde{J}(i, \tau) = \phi(i)^T \tau$$

feature vector

parameter

$$\phi(i) \in \mathbb{R}^d$$

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

$|x| \gg d$

So, no lookup table. In place of  $J^\pi(i)$ , we use  $\tilde{J}(i, r)$

Question: What "r" to use 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 & acting 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.  $\|J - J^*\|_\infty = \epsilon$ ,  $\epsilon \geq 0$

If  $\pi$  is a greedy policy based on  $J$ , i.e., using  $(\pi)$  with  $J$  instead of  $J^*$

$$\|J_\pi - J^*\|_\infty \leq \frac{2\epsilon}{1-\alpha}$$

Further, one can choose an  $\epsilon_0$  s.t.  $\forall \epsilon < \epsilon_0$ ,  $\pi$  is an optimal policy.

Pf:

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

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

$\Delta^1$  ineq.  $\rightarrow$

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

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

$T_\pi$  is a  $\alpha$ -contraction in  $\|\cdot\|_\infty$

$$\leq \alpha \|J_\pi - J\|_\infty + \|TJ - J^*\|_\infty$$

$J^* = TJ^*$   
 $T$  is a  $\alpha$ -contraction in  $\|\cdot\|_\infty$

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

$$\leq \alpha \|J_\pi - J^*\|_\infty + 4\|J^* - J\|_\infty + \alpha \|J - J^*\|_\infty$$

$$= 2\|J_\pi - J^*\|_\infty + 2\epsilon$$

$$\|J_\pi - J^*\|_\infty \leq \alpha \|J_\pi - J^*\|_\infty + 2\epsilon$$

So,

$$\|J_\pi - J^*\|_\infty \leq \frac{2\epsilon}{1-\alpha}$$

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

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

Choose  $\epsilon$  s.t.  $\frac{2\epsilon}{1-\alpha} < \delta$  & let  $\bar{\pi}$  be the greedy policy with this  $\epsilon$ .

Then  $\|J_{\bar{\pi}} - J^*\| < \delta \Rightarrow \bar{\pi} = \pi^*$ .



## Approximate policy evaluation using TD-type algorithms

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

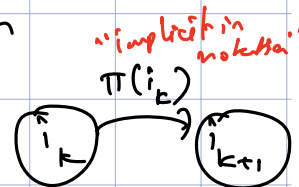
Consider a finite-state MDP.

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

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

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

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



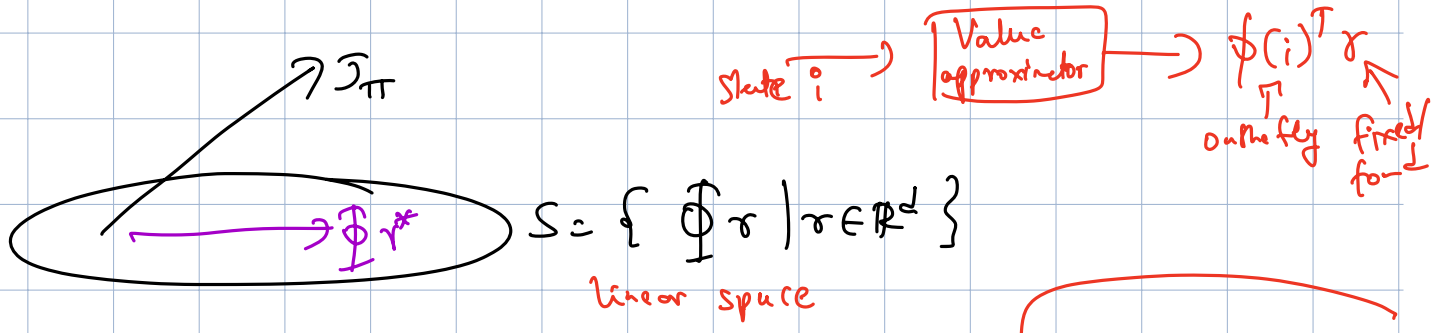
Aim: Estimate  $J_{\pi}(i) = E \left( \sum_{k=0}^{\infty} \alpha^k g(i_k, i_{k+1}) \mid i_0 = i, \pi \right)$   
no  $\pi$  here, it is implicit



# Linear function approximation:

$$J(i, \tau) = \phi(i)^T \tau, \quad i=1 \dots n$$

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



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

"(can't) also work with  $\Phi^T$ "  
 exp. usually work with  $\Phi^T$

$\Phi \rightarrow$  big Phi (matrix)  
 $\phi(i) \rightarrow$  small phi "feature vector"  
 Feature matrix  
 "Tall matrix"  
 $n \times d, n \gg d$

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

$n \times d$  matrix  $\rightarrow$  a tall one.

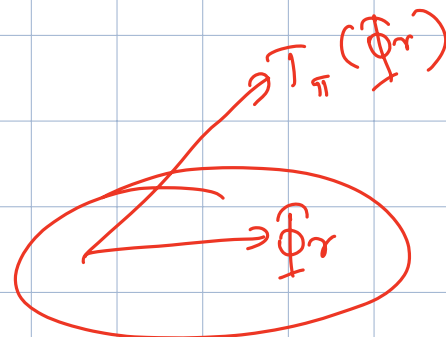
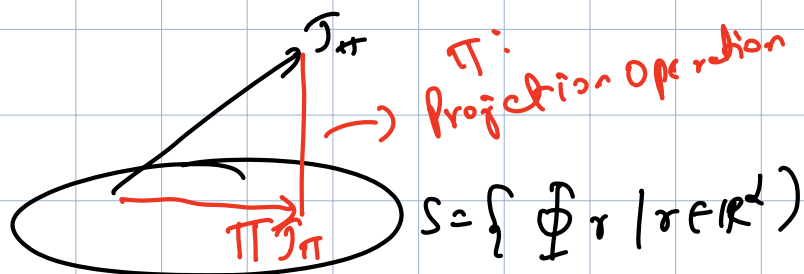
$$J_\tau = [ J(1, \tau) \quad \dots \quad J(n, \tau) ]^T \rightarrow \text{approx. value function}$$

$$J_\tau = \Phi \tau$$

"Want to bound  $\| \Phi \tau - J^* \|^2$ "  
 $\rightarrow$  some norm

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

Lecture-30 (continued after Markov chains review)



Q: In  $T_\pi \approx \Phi r$ , what "r" to choose?

For regular policy evaluation, one solves

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

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

This fixed point relation does not necessarily hold

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

doesn't make sense since  $T_\pi(\Phi r^*)$  need not be the same as  $\Phi r^*$  (unless  $T_\pi = \Phi r^*$ )

Projected fixed point equation

$$\Phi r^* = \Pi T_\pi (\Phi r^*) \quad (2)$$

projection onto linear space S

Projected fixed point equation

$\tilde{J} = \Pi T_\pi \tilde{J}$   
 contracts  $\tilde{J}$  is unique  
 $\Downarrow$   
 r is also unique?  
 where  $\tilde{J} = \Phi r$

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

(can solve (\*) if  $\Pi T_\pi$  is a contraction.  
 & we will show it is the case.

TD with linear function approximation

Assumptions:

(1) The Markov chain underlying policy  $\pi$  is irreducible, positive recurrent  
i.e.,  $\exists$  a stationary distribution  $\{\xi_1, \dots, \xi_n\}$   
for this chain  $[\xi = \xi P]$

(2) The matrix  $\Phi$  has full column rank  
or  $\text{rank}(\Phi) = d$  [Note: we assume  $n \gg d$ ]

Lecture-3)\*

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

for a  $n$ -vector  $J$ , define  $\rightarrow$  Stationary distribution (assumed to exist)

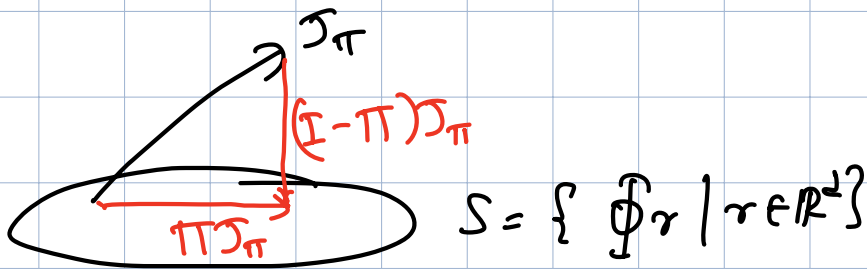
$$\|J\|_{\xi}^2 = \sum_{i=1}^n \xi_i (J(i))^2$$

$\rightarrow$  weighted  $l_2$  norm

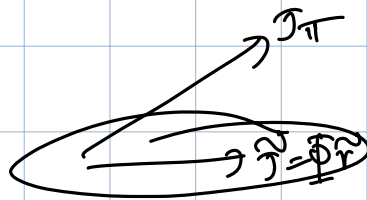
$$D = \begin{bmatrix} \xi_1 & & & \\ & \xi_2 & & \\ & & \ddots & \\ & & & \xi_n \end{bmatrix}$$

$\xi_i > 0, \forall i$

$$\|J\|_{\xi}^2 = J^T D J, \quad \forall J \in \mathbb{R}^n$$



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



$\Pi J_{\pi}$  is the "unique" vector in  $S$  that minimizes  $\|J_{\pi} - J\|_{\xi}$ , over all  $J \in S$ .

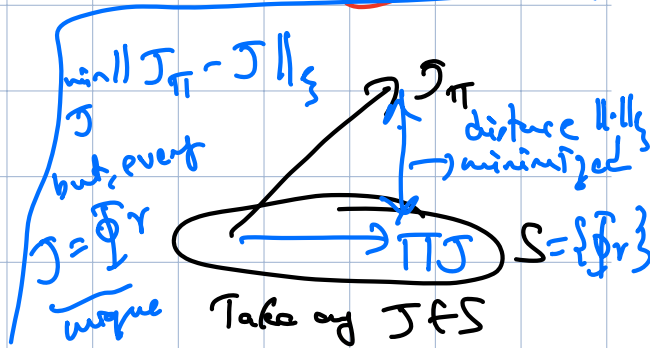
Any  $J \in S$  is of the form  $\Phi r$

Any  $J \in S$  can be uniquely written as some  $\Phi r$  since  $\Phi$  is full column rank

So,  $\tilde{r} = \operatorname{argmin}_{r \in \mathbb{R}^2} \|\tilde{J}_{\pi} - \Phi r\|_{\xi}^2$

and  $\Pi J_{\pi} = \Phi \tilde{r}$

Projection operator



To find  $\vec{r}$  :

Cannot solve  $\Phi \vec{r} = \mathbb{T}(\Phi \vec{r})$

Can be outside

$$\nabla \|\mathbb{T}_\pi - \Phi \vec{r}\|_2^2 = 0$$

So, project  $\mathbb{T}_\pi(\Phi \vec{r})$

ie.,  $\Pi(\mathbb{T}_\pi(\Phi \vec{r}))$

& then solve

$$\Rightarrow \nabla (\mathbb{T}_\pi - \Phi \vec{r})^T D (\mathbb{T}_\pi - \Phi \vec{r}) = 0$$

$$\Phi \vec{r} = \Pi(\mathbb{T}_\pi(\Phi \vec{r}))$$

$$\Rightarrow \Phi^T D (\mathbb{T}_\pi - \Phi \vec{r}) = 0 \quad (*)$$

$$\Rightarrow \Phi^T D \mathbb{T}_\pi - \Phi^T D \Phi \vec{r} = 0$$

$\Phi$  is  $n \times d$

$D$  is  $n \times n$

$\vec{r}$  is  $d \times 1$

$$\Rightarrow \Phi^T D \Phi \vec{r} = \Phi^T D \mathbb{T}_\pi$$

$$\Rightarrow \vec{r} = (\Phi^T D \Phi)^{-1} \Phi^T D \mathbb{T}_\pi$$

Why is this invertible?

$\Phi \rightarrow$  full rank and

diagonal elements of  $D > 0$ .

Now,  $\Pi \mathbb{T}_\pi = \Phi \vec{r}$

$$\Rightarrow \Pi = \Phi (\Phi^T D \Phi)^{-1} \Phi^T D$$

projection operator

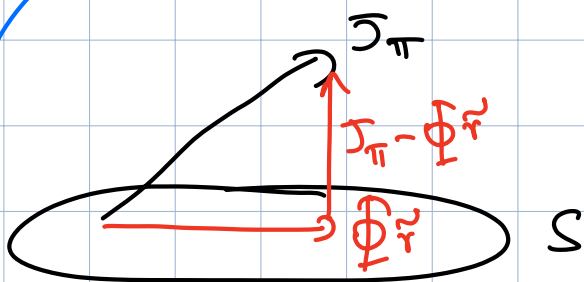
diagonal matrix with stationary distribution values

From (\*),

$$\Phi^T D (\mathcal{J}_\pi - \Phi \tilde{r}) = 0$$

$$\Leftrightarrow \tilde{r}^T \Phi^T D (\mathcal{J}_\pi - \Phi \tilde{r}) = 0$$

$$\Leftrightarrow (\Phi \tilde{r})^T D (\mathcal{J}_\pi - \Phi \tilde{r}) = 0$$



This equation  $\Leftrightarrow$

$$\Phi \tilde{r} \perp (\mathcal{J}_\pi - \Phi \tilde{r})$$

↑  
orthogonal  
in  $\|\cdot\|_\xi$  norm

$$\Leftrightarrow \langle (\Phi \tilde{r}), \mathcal{J}_\pi - \Phi \tilde{r} \rangle_\xi = 0$$

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

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

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

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

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

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

Projected fixed point equation:

$$\Phi_{\pi} z^* = \Pi T_{\pi}(\Phi_{\pi} z^*)$$

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  $\tilde{z} = \Phi_{\pi} z^*$ , we have

$$\tilde{z} = \Pi T_{\pi}(\tilde{z}) \leftarrow \text{Projected eqn}$$

Contrast with regular fixed point equation:

$$z_{\pi} = T_{\pi} z_{\pi}$$

$\Pi$ : Come in because  $\tilde{z} \in S$  "linear space".

Lemma 1:  $\|PJ\|_{\xi} \leq \|J\|_{\xi} \quad \forall J \in \mathbb{R}^n$

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

Stationary distribution vector

PF:  $\|PJ\|_{\xi}^2 = \sum_{i=1}^n \xi_i \left( \sum_{j=1}^n P_{ij} J(j) \right)^2$

Jensen's inequality  
①  $P_{ij}^2 \leq P_{ij}$

$\leq \sum_{i=1}^n \xi_i \sum_{j=1}^n P_{ij} J(j)^2 \quad \text{--- } (**)$

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

this step requires stationarity

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

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

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

So,  $\|PJ\|_{\xi}^2 \leq \|J\|_{\xi}^2$



"Projection is non-expansive"

Lemma 2:  $\|\pi J - \pi J'\|_{\xi} \leq \|J - J'\|_{\xi}, \quad \forall J, J' \in \mathbb{R}^n$

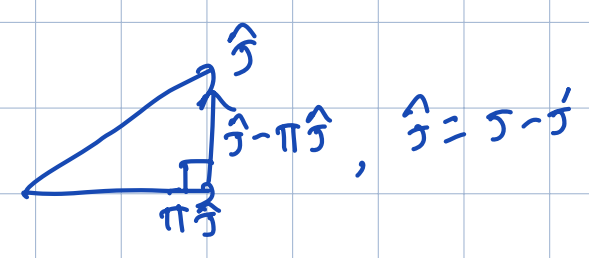
PF: PTO



$$\|\pi \mathcal{J} - \pi \mathcal{J}'\|_{\mathcal{E}}^2 = \|\pi(\mathcal{J} - \mathcal{J}')\|_{\mathcal{E}}^2$$

Pythagoras theorem  
 $\pi \mathcal{J}$  &  $(\mathbb{I} - \pi)\mathcal{J}$   
 are orthogonal

$$\leq \|\pi(\mathcal{J} - \mathcal{J}')\|_{\mathcal{E}}^2 + \|(\mathbb{I} - \pi)(\mathcal{J} - \mathcal{J}')\|_{\mathcal{E}}^2 \quad (*)$$



Note:

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

$$\begin{aligned} & \|\pi(\mathcal{J} - \mathcal{J}')\|_{\mathcal{E}}^2 + \|(\mathbb{I} - \pi)(\mathcal{J} - \mathcal{J}')\|_{\mathcal{E}}^2 \\ &= \|\pi(\mathcal{J} - \mathcal{J}') + (\mathbb{I} - \pi)(\mathcal{J} - \mathcal{J}')\|_{\mathcal{E}}^2 \\ &= \|(\mathcal{J} - \mathcal{J}')\|_{\mathcal{E}}^2 \end{aligned}$$

Hence, from (\*), we obtain

$$\|\pi \mathcal{J} - \pi \mathcal{J}'\|_{\mathcal{E}}^2 \leq \|\mathcal{J} - \mathcal{J}'\|_{\mathcal{E}}^2$$



Main claim:

$T_{\pi}$  and  $\pi T_{\pi}$  are contraction mappings w.r.t.  $\|\cdot\|_{\mathcal{E}}$ , and have modulus  $\alpha$  ( $\rightarrow$  discount factor)

PF:

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

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

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

$$\begin{array}{l} \text{Using Lemma 1} \\ \|PJ\|_{\xi} \leq \|J\|_{\xi} \end{array} \Rightarrow \leq \alpha \|J - J'\|_{\xi} \quad \text{--- (**)}$$

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 w.r.t.  $\|\cdot\|_{\xi}$  as well)

Next

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

$$\begin{array}{l} \Pi \text{ is} \\ \text{non-expansive} \\ \text{(Lemma 2)} \end{array} \Rightarrow \leq \|T_{\pi} J - T_{\pi} J'\|_{\xi}$$

$$\begin{array}{l} \text{From (**)} \end{array} \Rightarrow \leq \alpha \|J - J'\|_{\xi}$$

So,  $\Pi T$  is contractive w.r.t.  $\|\cdot\|_{\xi}$  with modulus  $\alpha$ . ■

Implication:

$r^*$  is unique because (i)  $\Pi T_\pi$  is contractive, (ii)  $\Phi$  is full col. rank

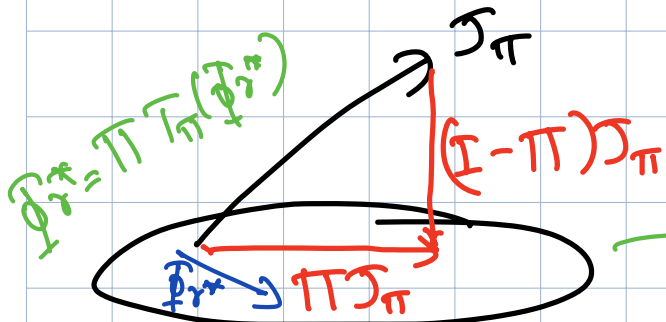
$$\Phi r^* = \Pi T_\pi (\Phi r^*)$$

Cannot solve  $T_\pi r = \Phi r$  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.,  $\Phi r_0 \xrightarrow{\Pi T_\pi} \Phi r_1 \dots \rightarrow$  asymptotically converges to  $\Phi r^*$ )



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

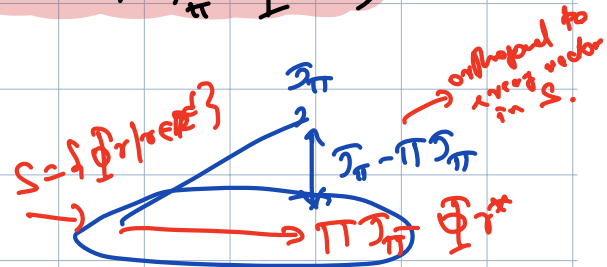
$\Pi T_\pi = \Phi r^*$  (Don't confuse  $r^*$  with  $\Phi r^*$ )

Not-so-main claim!

Recall  $r^*$  is the fixed point of  $\Pi T_\pi$ , i.e.,

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

$$\| T_\pi - \Phi r^* \|_\xi^2$$



because  $(T_\pi - \Pi T_\pi) \perp (\Pi T_\pi - \Phi r^*)$  is orthogonal to S

$$= \| T_\pi - \Pi T_\pi \|_\xi^2 + \| \Pi T_\pi - \Phi r^* \|_\xi^2$$

$$= \| T_\pi - \Pi T_\pi \|_\xi^2 + \| \Pi T_\pi T_\pi - \Pi T_\pi \Phi r^* \|_\xi^2$$

$\Pi T_\pi$  is a  $\lambda$ -contraction

$$\leq \| T_\pi - \Pi T_\pi \|_\xi^2 + \lambda^2 \| T_\pi - \Phi r^* \|_\xi^2$$

Rearranging:

$$\| T_\pi - \Phi r^* \|_\xi^2 \leq \frac{1}{1-\lambda^2} \| T_\pi - \Pi T_\pi \|_\xi^2$$

# Lecture-32

Matrix form of  $\Phi r^* = \Pi T_\pi(\Phi r^*)$  i.e.,  $(\sigma^* = d$

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

$$T_\pi J = g + \alpha P J$$

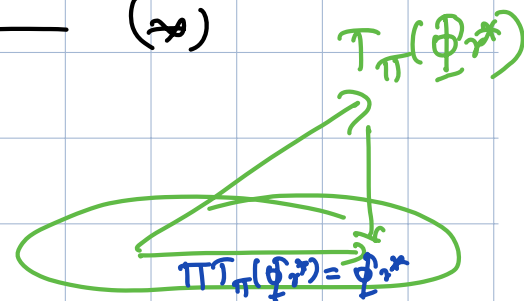
$\Pi J \rightarrow$  linear

$T_\pi J \rightarrow$  linear

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

↓

"linear" system of equations



Want to write (\*) as  $C r^* = d$

$\Pi J$ : minimize distance between  $J$  &  $\tilde{J} \in S$  in  $\|\cdot\|_S$

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

$$\tilde{r} = \arg \min_{\tilde{J} \in S} \|\Phi \tilde{r} - \tilde{J}\|_S^2$$

$$\Pi J = \Phi \tilde{r}$$

arg(\*)

$$= \arg \min_{r \in \mathbb{R}^d} \|\Phi r - (g + \alpha P \Phi r^*)\|_S^2$$

$T_\pi g + P$

$$= \arg \min_{r \in \mathbb{R}^d} (\Phi r - (g + \alpha P \Phi r^*))^T D (\Phi r - (g + \alpha P \Phi r^*))$$

min over "r"

nd-n variable

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

Differentiating the expression being minimized, we obtain

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

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

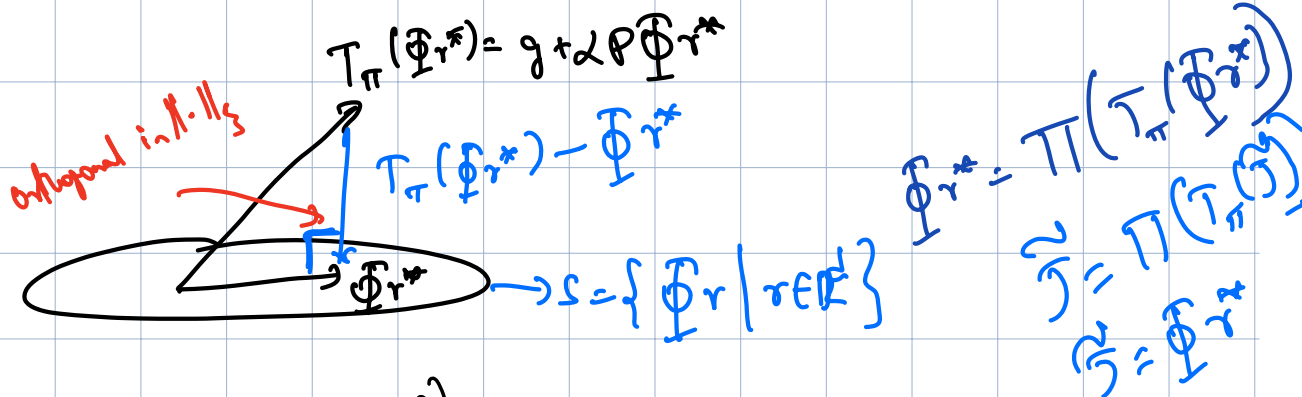
Intuitively: using (\*\*), we have

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

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

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

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



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

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

$$\Phi^T Dg = \Phi^T D(I - \alpha P)\Phi r^* \quad \text{or, equivalently}$$

$$C r^* = d, \text{ where } C = \Phi^T D(I - \alpha P)\Phi, \quad d = \Phi^T Dg$$

$\downarrow$   
 This is the same as  $\Phi r^* = \Pi T_\pi(\Phi r^*)$

Explicit solution!  $r^* = C^{-1} d$

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

$$J_{k+1}(i) = J_k(i) + \beta \left( TD_{\text{input}} \right)$$

$$J_k(i) \approx \phi(i)^T r^*$$

### Projected value iteration:

We know  $\Pi T_\pi$  is a  $\alpha$ -contraction

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

$$\Phi r_0 \xrightarrow{\Pi T_\pi} \Phi r_1 \rightarrow \dots \rightarrow \Phi r^*$$

Converge to

$$\Phi r_{k+1} = \Pi T_\pi(\Phi r_k), \quad k = 0, 1, \dots \quad \text{--- (**)}$$

$\uparrow$   
 Value iteration + projection

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

$$r_{k+1} = \arg \min_{r \in \mathbb{R}^d} \left\| \underbrace{\Phi r}_{\text{variable}} - \underbrace{(g + \alpha P \Phi r_k)}_{\text{fixed}} \right\|_2^2$$

Differentiating,

$$\Phi^T D (\Phi r_{k+1} - (g + \alpha P \Phi r_k)) = 0 \quad \left( \begin{matrix} \times \\ \times \end{matrix} \right)$$

$$r_{k+1} = r_k - (\Phi^T D \Phi)^{-1} (C r_k - d) \quad \left( \begin{matrix} \times \times \\ \times \times \end{matrix} \right)$$

Check this by substituting expressions C, d in  $\left( \begin{matrix} \times \\ \times \\ \times \end{matrix} \right)$ .

↳ the same as  $\Phi r_{k+1} = \Pi_{\Pi}(\Phi r_k)$

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  $\left( \begin{matrix} \times \times \\ \times \times \end{matrix} \right)$

Lecture 34\*

$x_0 \xrightarrow{\pi(x_0)} x_1 \xrightarrow{\pi(x_1)} x_2 \rightarrow \dots$  observe single step costs

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

Recall  $C = \Phi^T D (I - \alpha P) \Phi$ ,  $d = \Phi^T D g$

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

to form estimates of C, d.  $\hat{\Phi} = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix}$ ,  $\hat{\Phi} = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix}$

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

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

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

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

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

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

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

unknown in an RL setting

Using  $\pi$ , 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  $\Phi^T D \Phi$ ,  $\Phi^T D P \Phi$  &  $\Phi^T D g$ .



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

$(i_0, i_1, \dots, i_k)$   
 $\frac{1}{k+1} \phi(i) \phi(i)^T + \frac{1}{k+1} \phi(i) \phi(i)^T$

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

$(i_0, i_1, \dots, i_k)$   
 $\frac{\#N(i)}{k+1} \xrightarrow{k \gg 0} \rho_i$

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

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

Simple-based approx of  $C$  &  $d$

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

$$C_k \approx C \quad (C = \Phi^T D (\mathbb{I} - \alpha P) \Phi) \quad d_k \approx d \quad (= \Phi^T D g)$$

LSTD: Solve

$$C_k r_k = d_k$$

Least-squares temporal difference.

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

Temporal difference.

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

$$TD\text{-term} = \tilde{V}(i_t) - (g(i_t, i_{t+1}) + \alpha \tilde{V}(i_{t+1}))$$

Recall projected system of equations:

$$C = \Phi^T D (\mathbf{I} - \alpha P) \Phi, \quad d = \Phi^T D g \quad (\Leftrightarrow) \quad \Phi_{\sigma^*} = \prod_{\pi} T_{\pi}(\Phi_{\sigma^*})$$

↑  
vec stationary dist<sup>n</sup>  
( $\| \cdot \|_{\infty}$ )

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

Sample from  $\sum_i \phi(i)$ ,  $i=1 \dots n$

$$C_2 := \Phi^T D P \Phi = \sum_{i=1}^n \sum_{j=1}^n \sum_i \phi(i) \phi(j)^T P_{ij}$$

}  $\sum P = \{ \sum_i P_{ij}, i, j=1 \dots n \}$

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

}           

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

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

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

Using  $\hat{\xi}_i, \hat{p}_{i,j}$ , we estimate  $C_1, C_2, d$  as follows:

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

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

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

replace  $\hat{\xi}_i$  &  $\hat{p}_{i,j}$  by  $\xi_i$  &  $p_{i,j}$  would get us to  
 $C_1 = \Phi^T D \Phi$   
 $C_2 = \Phi^T D P \Phi$  &  
 $d = \Phi^T D g$ .

How to approximate  $C_1 \approx d$ ?

$$\hat{A} \theta = \hat{b}$$

$$\hat{A} = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) \phi(i_t)^T$$

$$\hat{b} = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) g(i_t)$$

Least squares regression

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

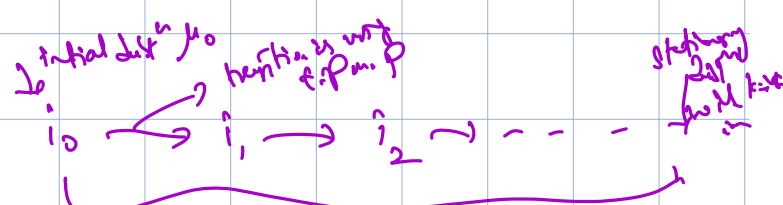
$$C_k = \hat{C}_1 - 2 \hat{C}_2$$

$$d_k = \hat{d}$$

② Solve  $C_k \hat{r}_k = d_k$

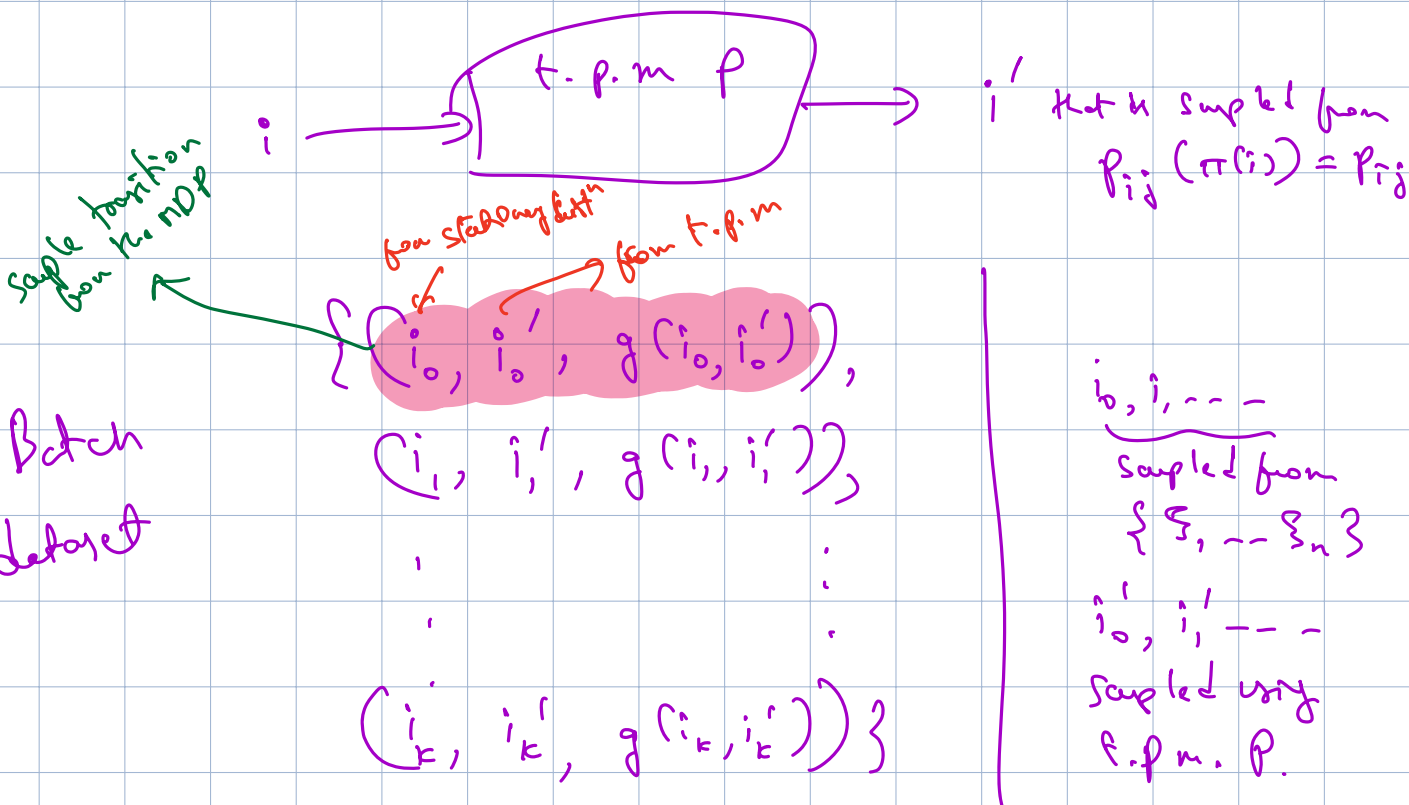
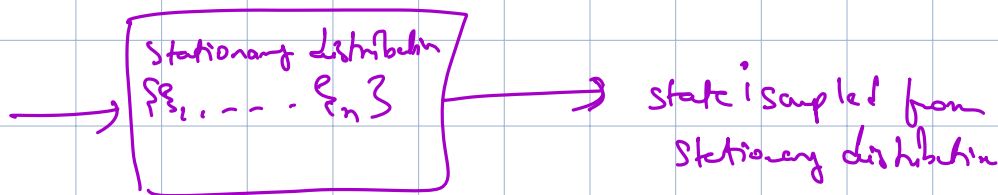
$\hat{r}_k \rightarrow$  LSTD solution  $\rightarrow$  batch algorithm - uses a dataset (not incremental)

Some analysis:-



# Generative model (for analysis)

$$\mu_0 P^n \xrightarrow{\text{diff.}} \xi$$



$$\text{form } \hat{c}_1, \hat{c}_2, \hat{d}$$

$$C_k = \hat{c}_1 - \alpha \hat{c}_2, \quad d_k = \hat{d}$$

$$\text{Solve } C_k \hat{c}_k = d_k$$

LSTD solution

$$\hat{\xi}_i = \frac{\# \text{ visits to } i \text{ in } (i_0, \dots, i_k)}{k+1} \xrightarrow{k \rightarrow \infty} \xi_i$$

Suppose using the  $\xi$

$$\hat{\xi}_i = \frac{\# \text{ visits to } i \text{ in } (i_0, \dots, i_k)}{k+1} + \frac{\# \text{ visits to } i \text{ in } (i_{k+1}, \dots, i_{k+L})}{k+L+1}$$

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

As  $k \rightarrow \infty$ ,  $\hat{\Sigma}_k \rightarrow \Sigma$  w.p.1  
 $\hat{P}_{ij} \rightarrow P_{ij}$  w.p.1  
 (Version of SLLN)

$$\boxed{C_k \xrightarrow{k \rightarrow \infty} C, \quad d_k \rightarrow d} \quad \rightarrow \text{LLN-type result for "LSTD"}$$

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

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

$$\hat{C}_1 = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) \phi(i_t)^T$$

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

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

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)

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

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

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

So,  $C_k$  is not necessarily invertible

As an alternative, solve

$$(C_k + \beta I) \hat{r}_k = d_k$$

$$\Leftrightarrow \hat{r}_k = (C_k + \beta I)^{-1} d_k$$

this is invertible if  $\beta$  is large enough

### Lecture 35\*

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

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

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

$$\text{Since } C_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) (\phi(i_t) - \alpha \phi(i_{t+1}))^\top$$
$$d_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) g(i_t, i_{t+1})$$

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

is a temporal difference term because

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

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

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

So  $(*) \Leftrightarrow$  TD-term.

"LSTD( $\lambda$ ) can be worked out using TD( $\lambda$ ) + LSTD idea"

## Lecture-34\*

TD(0) with linear function approximation

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

Want to solve  $\mathcal{J}_\pi = \bar{T}_\pi \mathcal{J}_\pi$

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

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

Tabular TD(0) or TD(0) in full state representations

$$\mathcal{J}_{t+1}(i) = \mathcal{J}_t(i) + \beta_t \left( g(i, \bar{i}) + \alpha \mathcal{J}_t(\bar{i}) - \mathcal{J}_t(i) \right)$$

$g(i, \bar{i}) + \alpha \mathcal{J}_t(i)$  is a proxy for  $\underbrace{E_{\bar{i}} \left( g(i, \bar{i}) + \alpha \mathcal{J}_\pi(\bar{i}) \right)}_{\text{TD-term}}$

Onto linear function approximation case:

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

Cannot do this

$$r_{t+1} \neq r_t + \beta_t \left( g(i_t, i_{t+1}) + \alpha r_t^T \phi(i_{t+1}) - r_t^T \phi(i_t) \right)$$

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

TD(0) with linear function approximation

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

$$r_{t+1} = r_t + \beta_t \phi(i_t) \left( g(i_t, i_{t+1}) + \alpha r_t^T \phi(i_{t+1}) - r_t^T \phi(i_t) \right)$$

$\underbrace{r_{t+1}}_{\in \mathbb{R}^d} = \underbrace{r_t}_{\in \mathbb{R}^d} + \underbrace{\beta_t}_{\in \mathbb{R}} \underbrace{\phi(i_t)}_{\in \mathbb{R}^d} \underbrace{\left( g(i_t, i_{t+1}) + \alpha r_t^T \phi(i_{t+1}) - r_t^T \phi(i_t) \right)}_{\in \mathbb{R}}$

↖ a line algorithm

Go back to projected fixed point:

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

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

$$d = \Phi^T D g$$

Want to use a sample path  $(i_0, i_1, \dots)$  to find  $r^*$  s.t.  $C r^* = d$ .



Sto-iter-algo:  $r_{t+1} = r_t - \beta_t (Cr_t - d)$  — (\*)

Where would  $r_t$  converge? Ans: wherever  $Cr^* = d$   
 $r_t \rightarrow r^*$

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

$$= \sum_{i,j} \sum_i P_{ij} \phi(i) \left( \phi(i)^T r^* - 2\phi(j)^T r^* - g(i,j) \right)$$

picked from  $\xi$ 
picked from  $P_{ij}$

$$= E \left( \phi(i) \left( \phi(i)^T r^* - 2\phi(j)^T r^* - g(i,j) \right) \right)$$

taken with  $\xi P$  distribution

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

$$E \left( \phi(i) \left( \phi(i)^T r^* - 2\phi(j)^T r^* - g(i,j) \right) \right) = 0$$

$i$  picked from  $\xi = (\xi_1, \dots, \xi_n)$ 
 $j \sim \dots \sim (P_{ij})$

Let  $(i_t, i_{t+1})$  be a sample transition  
 Then, do the following update

$$(*) \rightarrow r_{t+1} = r_t - \beta_t \phi(i_t) \left( \phi(i_t)^T r_t - 2\phi(i_{t+1})^T r_t - g(i_t, i_{t+1}) \right)$$

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

## Remark:

① In the above, we assumed sampling from the stationary distribution. Under this, it is straightforward to invoke "Stor-iter-algo 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 w/ some distribution " $\mu$ "

Let the Markov chain have t.p.m.  $P$ .

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

$\mu P^k \rightarrow \xi$  as  $k \rightarrow \infty$  (assuming irreducible + positive recurrence)

$(i_0, i_1, \dots, \dots)$   
*transient phase*

↑  
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  $C\gamma^* = d \Leftrightarrow \gamma^* = T(\gamma^*)$

Ref: J.M. 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  
DPOC vol II Chapter 6"