

## Chapter 6 Q-Learning

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

$\underbrace{\max_a Q(S_{t+1}, a)}_{\text{greedy-}\epsilon \text{ policy w.p. } 1-\epsilon \text{ max } \epsilon \text{ rand}}$

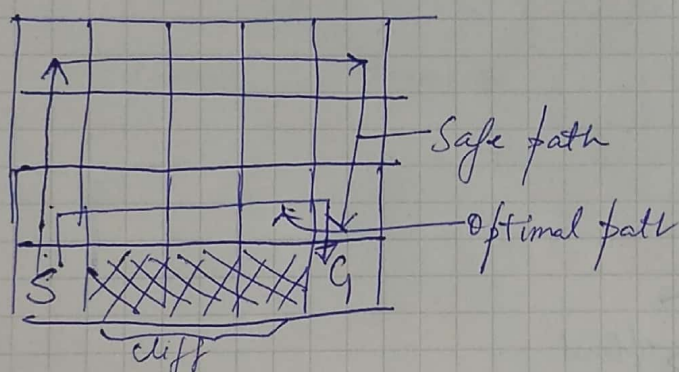
$Q(S_t, A_t)$  = Q-value associated with tuple  $(S_t, A_t)$

will give optimal Q-value.

$$Q(S_t, A_t) \rightarrow Q^*(S_t, A_t)$$

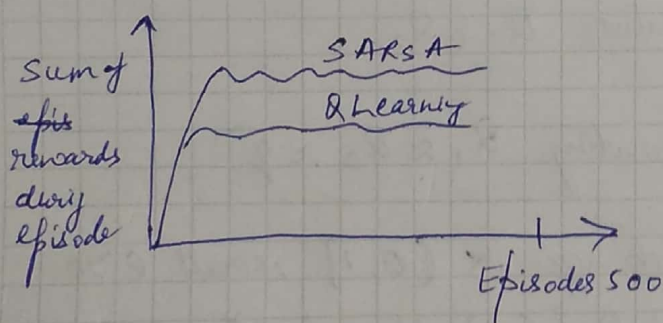
maximising action (optimal action) in state  $S_t$   $\arg \max_a Q^*(S_t, a)$

Cliff Walking:



reward = -1 for any action that takes you to "non-cliff" position.

reward = -100 for any action that takes to the cliff

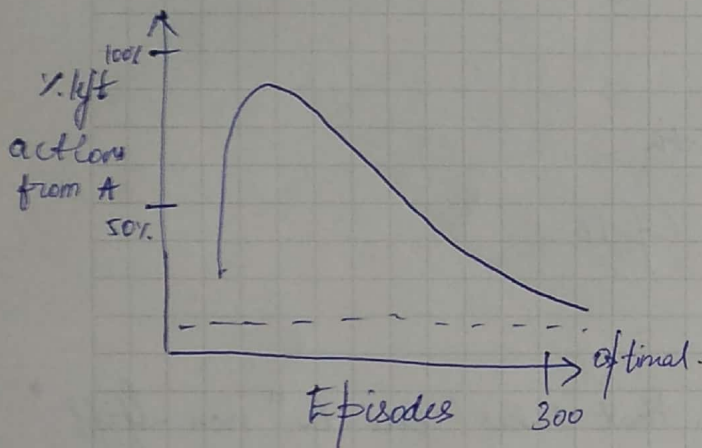
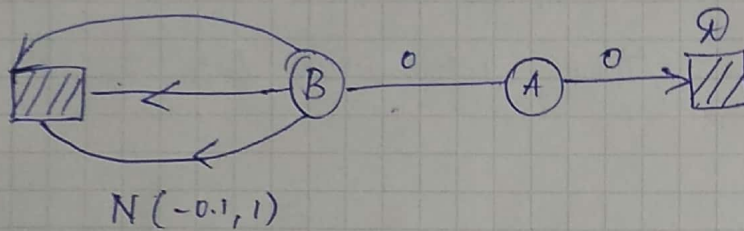


Expected SARSA:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma E_{\pi} [Q(s_{t+1}, A_{t+1}) | s_t] - Q(s_t, a_t)]$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a | s_{t+1}) Q(s_{t+1}, a) - Q(s_t, a_t)]$$

Q-learning suffers from maximization bias:

Double Q-learning (2010-11 algorithm)

Have 2 estimates of Q-value,  $Q_1$  &  $Q_2$ .

Double Q-learning for estimating  $Q_1 \approx Q_2 \approx q^*$

Algorithm parameters: step size  $\alpha \in [0, 1]$ , small  $\epsilon > 0$

Initialize  $Q_1(s, a)$  and  $Q_2(s, a) \forall s \in S^*, a \in A(s)$

st  $Q(\text{terminal}) = 0$



Loop for each episode:

Initialize  $S$

Loop for each episode:

choose action  $A$  in  $S$  using the policy  $\epsilon$ -greedy  
in  $Q_1$  &  $Q_2$

Take action  $A$ , observe  $R, S'$

with prob 0.5,

$$Q_1(S, a) \leftarrow Q_1(S, A) + \alpha (R + \gamma Q_2(S', \arg \max_a Q_2(S', a)) - Q_1(S, A))$$

else,

$$Q_2(S, a) \leftarrow Q_2(S, A) + \alpha (R + \gamma Q_1(S', \arg \max_a Q_1(S', a)) - Q_2(S, A))$$

until  $S$  is terminal.

## Function Approx Based Methods:

Neural N/w based architecture

No of states —  $10^{1000}$  states i.e., large

1. State  $i$  is encoded as  $x = (x_1(i), \dots, x_L(i))$
2.  $x$  is transformed to linearly through a layer as  

$$\sum_{l=1}^L w(k, l) x_l(i), k=1 \dots K$$

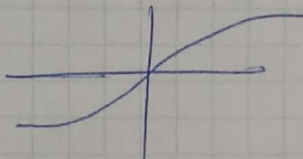
3. Non-linear transformation via a sigmoid  $\sigma(\cdot)$   
 we obtain  $\sigma\left(\sum_{l=1}^L w(k, l) x_l(i)\right)$

$\sigma(\cdot)$  are differentiable &  $n$  satisfying

$$-\infty < \lim_{z \rightarrow -\infty} \sigma(z) < \lim_{z \rightarrow \infty} \sigma(z) < \infty$$

$\sigma(\cdot)$  are non-decreasing

eg:  $\sigma(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

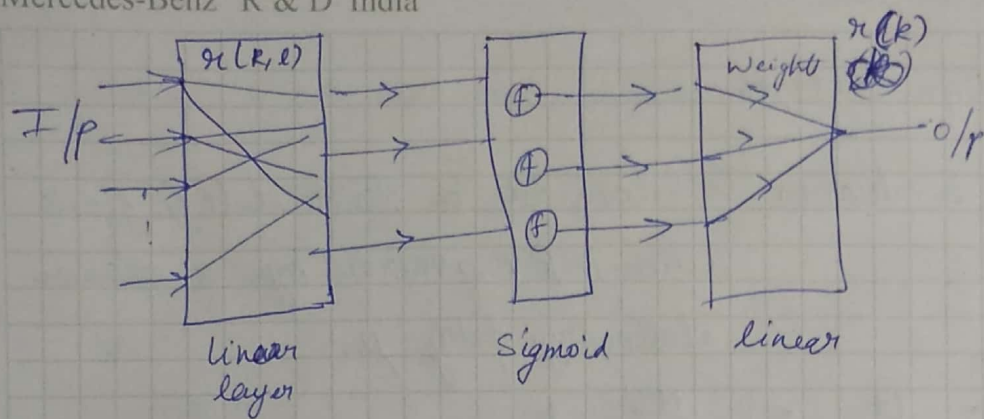


$$\sigma(z) = \frac{1}{1+e^{-z}}$$

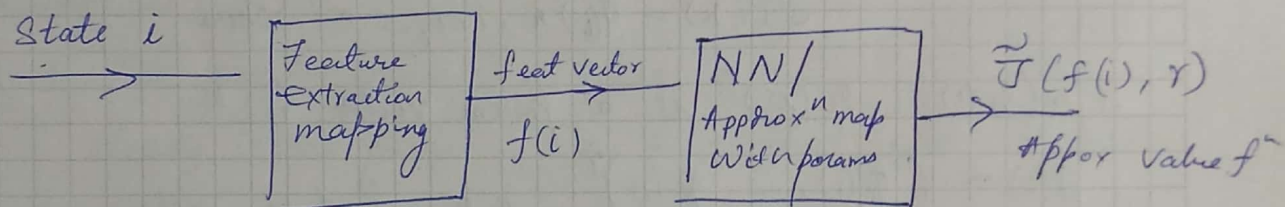
4) Final o/p

$$\tilde{J}(i, r) = \sum_{k=1}^K r(k) \sigma\left(\sum_{l=1}^L w(k, l) x_L(i)\right)$$



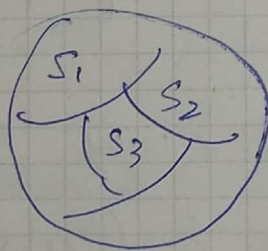


Features: - Capture the most important aspects of a state  
 - When used  $T$  is a complicated  $f^n$  [non linear]



Feature: A mapping  $f_k: S \rightarrow \mathbb{R}$ . Suppose  $f_1, f_2, \dots, f_k$  are  $k$  features

Feature vector  $f(i) \equiv (f_1(i), \dots, f_k(i))$



$$\mu^* = \underset{\pi}{\operatorname{argmin}} \sum_i \left( T(i) - \sum_k \pi(k) \phi(f(i)) \right)^2$$

## Projected Eqn Methods:

**Policy Evaluation:** Consider a finite state & finite action MDP governed by a stationary stationary policy  $\mu$ .

$$J_{\mu}(i) = \lim_{N \rightarrow \infty} E \left[ \sum_{k=0}^{N-1} \gamma^k g(i, k+1) \mid i_0 = i \right] \quad i=1, \dots, n$$

Approximate  $J_{\mu}(i)$  according to

$$\tilde{J}(i, \mathbf{x}) = \phi(i)^T \mathbf{x}, \quad i=1, \dots, n \quad \text{where}$$

$\mathbf{x} = (x(1) \dots x(s))^T$  is a param vector &

$\phi(i) = (\phi_1(i), \dots, \phi_s(i))^T$  is a feature vector associated

with state  $i$  let  $\phi = \begin{bmatrix} \phi(1)^T \\ \phi(2)^T \\ \vdots \\ \phi(n)^T \end{bmatrix}_{n \times s}$

$$\phi = \begin{bmatrix} \phi_1(1) & \dots & \phi_s(1) \\ \phi_1(2) & & \phi_s(2) \\ \vdots & & \vdots \\ \phi_1(n) & \dots & \phi_s(n) \end{bmatrix}_{n \times s}$$

$\phi$  is called **feat matrix**

$$\text{let } \tilde{\mathbf{J}}_{\mathbf{x}} = (\tilde{J}(i, \mathbf{x}), i \in S)^T$$

where  $S = \{1, 2, \dots, n\}$  is the state space

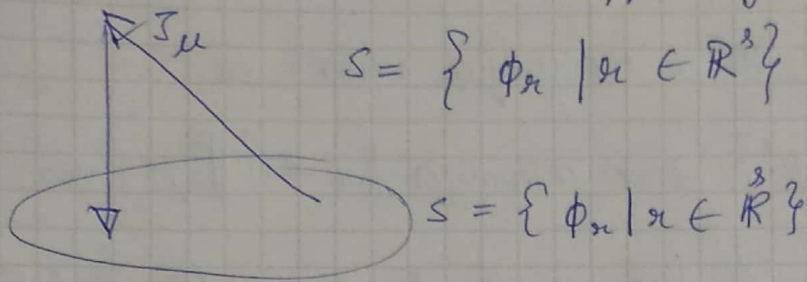
$$\text{Then } \tilde{\mathbf{J}}_{\mathbf{x}} = \phi \mathbf{x}$$

$$\tilde{\mathbf{J}}_{\mathbf{x}} = \begin{pmatrix} \tilde{J}(1, \mathbf{x}) \\ \tilde{J}(2, \mathbf{x}) \\ \vdots \\ \tilde{J}(n, \mathbf{x}) \end{pmatrix} = \begin{pmatrix} \phi(\mathbf{x})_1 \\ \phi(\mathbf{x})_2 \\ \vdots \\ \phi(\mathbf{x})_n \end{pmatrix}$$

eg:  $\phi_1(i) = 1$   
 $\phi_2(i) = i$   
 $\phi_3(i) = i^2$



Goal: Find the best approx of  $J_\mu$  within



Note that

$$\phi_x = \sum_{i=1}^s \phi_i x_i \quad \phi_i = \begin{pmatrix} \phi_i(1) \\ \vdots \\ \phi_i(n) \end{pmatrix}$$

Assumptions:

1. The markov chain  $\{X_n\}$  has steady state prob  $\pi_1, \dots, \pi_n$  that are all positive, i.e.,  $\forall i = 1, \dots, n$ .

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N P(i_k = j \mid i_0 = i) = \pi_j > 0 \quad \forall j = 1, \dots, n$$

The matrix  $\phi$  has rank  $s$  [ $s \leq n$ ]

$$n \approx 10^{1000}, \quad s \approx 10-100$$

Define a weighted Euclidean norm as  $\mathbb{R}^n$  as

$$\|J\|_v = \sqrt{J^T V J} = \sqrt{\sum_{i=1}^n v_i (J(i))^2}$$

Here  $V = \begin{bmatrix} v_1 & & 0 \\ & v_2 & \\ 0 & & v_n \end{bmatrix}$   $J = (J(i), i \in S)^T$

Assume  $v_1, v_2, \dots, v_n > 0$

Let  $\pi$  be the projection operator into  $S$  w.r.t this norm

For any  $J \in \mathbb{R}^n$ ,  $\pi J$  is the unique vector in

$S = \{ \phi x \mid x \in \mathbb{R}^c \}$  that minimizes  $\|J - \hat{J}\|_V^2$  over all  $\hat{J} \in S$ . Since  $\phi$  has rank  $c$ , any vector  $\hat{J} \in S$

is uniquely written as  $\hat{J} = \phi x$  for some  $x \in \mathbb{R}^c$ .

$$\text{Thus } \|J - \hat{J}\|_V^2 = \|J - \phi x\|_V^2 = (J - \phi x)^T V (J - \phi x)$$

Thus,

$$\pi J = \phi x_J \quad \text{where } x_J = \underset{x \in \mathbb{R}^c}{\operatorname{argmin}} \|J - \phi x\|_V^2, J \in \mathbb{R}^n$$

In order to find  $x_J$ ,

$$\nabla_x (\|J - \phi x\|_V^2) = 0$$

$$\text{or } \nabla_x ((J - \phi x)^T V (J - \phi x)) = 0$$

$$\text{or } \phi^T V (J - \phi x_J) = 0$$

$$\text{or } \phi^T V J - \phi^T V \phi x_J = 0$$

$$\Rightarrow \phi^T V \phi x_J = \phi^T V J$$

$$\Rightarrow x_J = (\phi^T V \phi)^{-1} \phi^T V J$$



$$\phi x_J = \pi J \quad \text{Thus,}$$

$$\pi J = \phi (\phi^T V \phi)^{-1} \phi^T V J$$

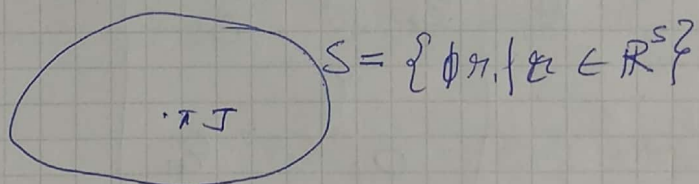
$$\text{Thus } \pi = \phi (\phi^T V \phi)^{-1} \phi^T V$$

Thus  $\pi$  is also called the projection matrix. The optimality condition (\*) can be re-written as,

$$\phi^T (J - \phi x_J) = 0 \quad \text{or}$$

$$J^T V (J - \phi x_J) = 0 \quad \text{where } \bar{J} = \phi x_J \in S$$

The difference / approx error  $(J - \phi x_J)$  is orthogonal to subspace  $S$ .



$$S = \{\phi x_i \mid x_i \in \mathbb{R}^n\}$$

The vectors  $x$  &  $y$  are orthogonal if  $x^T y = 0$

$$\text{or } \sum_{i=1}^n v_i x_i y_i = 0$$

Consider now the mapping,

$T: V \rightarrow V$  where  $V = \text{space of bounded } f^v \text{ from } S \text{ to } \mathbb{R} \text{ defined by}$

$$(TJ)(i) = \sum_{j=1}^n t_{ij} (g(x_j) + x J(j)), \quad i=1 \dots n$$

$$TJ = g + \alpha PJ$$

where  $g = (g_1, \dots, g_n)^T$  with

$$g_i = \sum_{j=1}^n p_{ij} g(j), \quad i=1, \dots, n$$

$$P = [p_{ij}]_{i,j=1}^n$$

Bellman Eqn:

$$J = TJ$$

$$\boxed{\phi r = \pi T(\phi r)}$$

$$T(\phi r) = g + \alpha P(\phi r)$$

↓  
projected Bellman Eqn

$V$  is replaced with

$$Q = \begin{bmatrix} \zeta_1 & & 0 \\ & \zeta_2 & \\ 0 & & \ddots \\ & & & \zeta_n \end{bmatrix}$$

$\zeta_i = \zeta_i$  = steady state prob of markov chain in state  $i$



Lemma 1:  $\|Pz\|_{\pi} \leq \|z\|_{\pi} \quad \forall z \in \mathbb{R}^n$  where

$P$  is the transition probab matrix for the markov chain for which  $\pi = \{ \pi_1, \dots, \pi_n \}$  is the stationary distribution.

Proof  $\|Pz\|_{\pi}^2 = \sum_{i=1}^n \pi_i \left( \sum_{j=1}^n p_{ij} z_j \right)^2 \quad \text{--- (1)}$

according to Jensen's inequality,

$$E[z|i]^2 \leq E[z^2|i]$$

$\therefore$  (1) becomes

$$\begin{aligned} &\leq \sum_{j=1}^n \pi_j \sum_{i=1}^n p_{ij} z_j^2 \\ &= \sum_{j=1}^n \sum_{i=1}^n \pi_i p_{ij} z_j^2 \\ &= \sum_{j=1}^n \pi_j z_j^2 = \|z\|_{\pi}^2 \end{aligned}$$

$$\text{or } \|Pz\|_{\pi} \leq \|z\|_{\pi} \quad \forall z \in \mathbb{R}^n$$

Lemma 2: The projection map  $\pi$  is ~~non expansive~~ <sup>expansive</sup> i.e.,

$$\|\pi J - \pi \bar{J}\|_V \leq \|J - \bar{J}\|_V \quad \forall J, \bar{J} \in \mathbb{R}^n$$

Proof: Note that,

$$\|\pi(J - \bar{J})\|_V^2 \leq \|\pi(J - \bar{J})\|_V^2 + \left\| \frac{(I - \pi)}{I - \pi} (J - \bar{J}) \right\|_V^2$$

Note that  $\pi(J - \bar{J}) \perp ((J - \bar{J}) - \pi(J - \bar{J}))$   
 $\in S$

by Pythagorean theorem, we have

$$\begin{aligned}\|\pi(J-\bar{J})\|_V^2 + \|(I-\pi)(J-\bar{J})\|_V^2 &= \|\pi(J-\bar{J}) + (J-\bar{J}) - \pi(J-\bar{J})\|_V^2 \\ &= \|J-\bar{J}\|_V^2\end{aligned}$$

$$\Rightarrow \|\pi(J-\bar{J})\|_V \leq \|J-\bar{J}\|_V \quad \forall J, \bar{J} \in \mathbb{R}^n$$

Proposition: The mappings  $T$  &  $\pi T$  are contractions of modulus  $\alpha$  with respect to  $\|\cdot\|_{\xi}$  where  $\xi = (\xi_1, \dots, \xi_n)^T$  is the steady state distribution of the Markov chain.

Proof: Recall that,

$$\begin{aligned}TJ &= g + \alpha PJ \\ T\bar{J} &= g + \alpha P\bar{J}\end{aligned}$$

$$\Rightarrow (TJ - T\bar{J}) = \alpha P(J - \bar{J})$$

$$\begin{aligned}\|TJ - T\bar{J}\|_{\xi} &\leq \alpha \|P(J - \bar{J})\|_{\xi} \\ &\leq \alpha \|J - \bar{J}\|_{\xi} \quad \text{by lemma 1}\end{aligned}$$

$T$  is a contraction of modulus  $\alpha$  w.r.t  $\|\cdot\|_{\xi}$

$$\begin{aligned}\text{Now } \|\pi TJ - \pi T\bar{J}\|_{\xi} &\leq \|TJ - T\bar{J}\|_{\xi} \\ &\leq \alpha \|J - \bar{J}\|_{\xi}\end{aligned}$$

$\pi T$  is a contraction of modulus  $\alpha$  w.r.t  $\|\cdot\|_{\xi}$



Proposition 2: Let  $\phi r^*$  be the fixed point of  $\pi J$ . Then,

$$\|J_\mu - \phi r^*\|_\xi \leq \frac{1}{\sqrt{1-\alpha^2}} \|J_\mu - \pi J_\mu\|_\xi$$

Proof:

Note that,

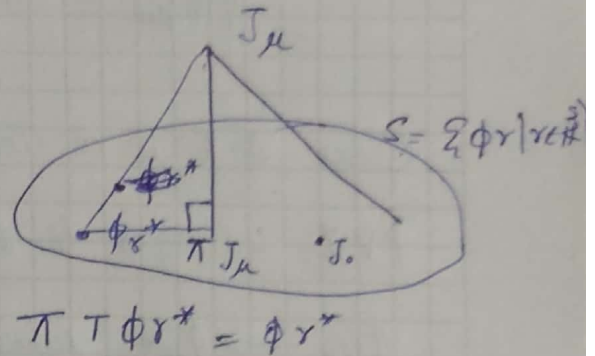
$$\|J_\mu - \phi r^*\|_\xi^2 = \|J_\mu - \pi J_\mu\|_\xi^2 + \|\pi J_\mu - \phi r^*\|_\xi^2 \quad \text{by Pythagoras theorem}$$

$$= \|J_\mu - \pi J_\mu\|_\xi^2 + \|\pi J_\mu - \pi T(\phi r^*)\|_\xi^2$$

$$\leq \|J_\mu - \pi J_\mu\|_\xi^2 + \alpha^2 \|J_\mu - \phi r^*\|_\xi^2$$

$$\Rightarrow (1-\alpha^2) \|J_\mu - \phi r^*\|_\xi^2 \leq \|J_\mu - \pi J_\mu\|_\xi^2$$

$$\text{or } \|J_\mu - \phi r^*\|_\xi \leq \frac{1}{\sqrt{1-\alpha^2}} \|J_\mu - \pi J_\mu\|_\xi$$



The solution to the projected Bellman eqn is the vector  $J = \phi r^*$  where  $r^* = \underset{r \in \mathbb{R}^s}{\operatorname{argmin}} \| \phi r - g + \alpha P \phi r^* \|_\xi^2$

$\phi r^*$  is the projection of  $g + \alpha P \phi r^*$  to the subspace  $S$

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

$$\text{or } \phi^T D ((I - \alpha P) \phi r^* - g) = 0$$

$$\text{or } \underbrace{\phi^T D (I - \alpha P)}_C \phi r^* = \phi^T D g$$

$$C r^* = \phi^T D g \quad \text{or } r^* = C^{-1} \phi^T D g$$

$$D = \begin{bmatrix} \tau_1 & & 0 \\ & \tau_2 & \\ 0 & & \ddots \\ & & & \tau_n \end{bmatrix}$$

Numerical Sol<sup>n</sup> to PBE:

Value iteration:

$$\phi_{x_{k+1}} = \pi T(\phi_{x_k}), \quad k=0, 1, 2, \dots$$

Since  $\pi T$  is a contraction,  $\{\phi_{x_k}\}$  generated by PVI converges to the unique fixed point  $\phi_{x^*}$  of  $\pi T$

$$\phi_{x_k} \rightarrow \phi_{x^*} \text{ where } \phi_{x^*} = \pi T(\phi_{x^*})$$

Note that

$$x_{k+1} = \arg \min_{x \in \mathbb{R}^n} \|\phi_x - (g + \alpha P \phi_{x_k})\|_Z^2$$

$$\text{Consider } \nabla_x (\phi_x - g + \alpha P \phi_{x_k})^T \otimes (\phi_x - g - \alpha P \phi_{x_k})$$

$$= 2 \phi^T \otimes (\phi_{x_{k+1}} - (g + \alpha P \phi_{x_k})) = 0$$

$$= \phi^T \otimes \phi_{x_{k+1}} = \phi_D \otimes (g + \alpha P \phi_{x_k})$$

$$\therefore x_{k+1} = (\phi^T \otimes \phi)^T \phi^T \otimes g + \alpha (\phi^T \otimes \phi)^T (\phi^T \otimes P \phi_{x_k}), \quad k \geq 0$$

$$d = \phi \otimes g.$$

Thus,

$$x_{k+1} = x_k + (\phi^T \otimes \phi)^T d - (\phi^T \otimes \phi)^T \phi^T \otimes P \phi_{x_k} + \alpha (\phi^T \otimes \phi)^T \phi^T \otimes P \phi_{x_k}$$

$$= x_k - (\phi^T \otimes \phi)^T \phi^T \otimes (I - \alpha P) \phi_{x_k} + (\phi^T \otimes \phi)^T d$$

$$= x_k - (\phi^T \otimes \phi)^T (x_k - d) \quad \text{where } c = \phi^T \otimes (I - \alpha P) \phi$$



Proposition 3: The matrix  $C = \Phi^T D (I - \alpha P) \Phi$  is positive definite.

Proof from Pythagorean theorem,

$$\|\Phi y\|_{\mathcal{E}}^2 = \|\Pi \Phi y\|_{\mathcal{E}}^2 + \|(I - \Pi) \Phi y\|_{\mathcal{E}}^2$$

$$\Rightarrow \|\Pi \Phi y\|_{\mathcal{E}}^2 \leq \|\Phi y\|_{\mathcal{E}}^2$$

$$\Rightarrow \|\Pi \Phi y\|_{\mathcal{E}} \leq \|\Phi y\|_{\mathcal{E}} \leq \|y\|_{\mathcal{E}}$$

Note that all vectors from  $\Phi y$  are orthogonal to all vectors of the form  $x - \Pi x$  i.e.,

$$y^T \Phi^T D (x - \Pi x) = 0 \quad \forall y \in \mathbb{R}^S, x \in \mathbb{R}^n$$

Thus  $\forall y \neq 0$ , we have

$$\begin{aligned} y^T C y &= y^T \Phi^T D (I - \alpha P) \Phi y \\ &= y^T \Phi^T D (I - \alpha \Pi P + \alpha (\Pi - I) P) \Phi y \\ &= y^T \Phi^T D (I - \alpha \Pi P) \Phi y + \underbrace{\alpha y^T \Phi^T D (\Pi - I) P \Phi y}_{\textcircled{2}} \end{aligned}$$

$\textcircled{2} = 0$  Since  $\Phi y$  is orthogonal to  $(\Pi - I) P \Phi y$

$$\text{Then } \textcircled{1} = y^T \Phi^T D \Phi y - \alpha y^T \Phi^T D \Pi P \Phi y$$

$$\geq \|\Phi y\|_{\mathcal{E}}^2 - \underbrace{\alpha \|\Phi y\|_{\mathcal{E}} \|\Pi P \Phi y\|_{\mathcal{E}}}_{\text{Cauchy's inequality}}$$

weighted inner product

$$\langle \Phi y, \Pi P \Phi y \rangle_D \leq \|\Phi y\|_{\mathcal{E}} \|\Pi P \Phi y\|_{\mathcal{E}}$$

$$\geq \|\Phi y\|_{\mathcal{E}}^2 - \alpha \|\Phi y\|_{\mathcal{E}}^2 = (1 - \alpha) \|\Phi y\|_{\mathcal{E}}^2 \geq 0 \text{ whenever } \alpha < 1$$

Thus  $C = [\Phi^T D (I - \alpha P) \Phi]$  is a +ve definite matrix

Material: D Bertsekas Optimal Control & Dynamic Program  
 Chap 6  
 [Approx Dynamic Programming]

Chapter 9 : On-policy prediction with Approx

Mean Square Value Error:  $\overline{VE}(w) \triangleq \sum_{s \in S} \mu(s) [V_\pi(s) - \hat{V}(s, w)]^2$

$\mu = (\mu(1), \dots, \mu(n))^T$  [steady state dist. of the markov chain]

$\mu(s) \equiv$  fraction of markov chain spends in state  $s$ .

On-policy distribution for episodic tasks: let  $h(s) = \text{prob that}$

an episode begins in state  $s$ . Let  $\eta(s) = \text{no of times step spent in state 's' on avg}$

Then  $\eta(s) = h(s) + \sum_{\bar{s}} n(\bar{s}) \sum_a \pi(a|\bar{s}) p(s|\bar{s}, a), \forall s \in S$

$$\mu(s) = \frac{\eta(s)}{\sum_{s'} \eta(s')}$$

$\bar{s}$ -avg.