

LSTD, LSPE and TD( $\gamma$ ) Methods

LSTD (Least Square Temporal Difference) method

$$r_k = C_k^{-1} d_k, \quad k=0, 1, 2, \dots$$

- $r_k$  can be obtained by solving  $C_k r_k = d_k$  using a batch of  $k+1$  samples and (30).
- Alternatively, LSTD can be written as

$$C_k r_k - d_k = \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) g_{k,t} = 0 \quad \text{--- (33)}$$

where the scalar  $g_{k,t}$  is the so-called temporal difference (TD for short), associated with  $r_k$  and transition  $(i_t, i_{t+1})$ , given by

$$g_{k,t} = \phi(i_t)' r_k - \alpha \phi(i_{t+1})' r_k - g(i_t, i_{t+1}). \quad \text{--- (34)}$$

- TD may be viewed as a sample of a residual term arising in the projected Bellman equation, i.e.,

$$\tilde{J}(i_t; r_k) = g_{k,t} + g(i_t, i_{t+1}) + \alpha \tilde{J}(i_{t+1}; r_k),$$

associated with  $r_k$  and transition  $(i_t, i_{t+1})$ .

LSPE (Least Square Policy Evaluation) method

$$r_{k+1} = r_k - \gamma G_k (C_k r_k - d_k), \quad k=0, 1, 2, \dots \quad \text{--- (35)}$$

where  $\gamma$  is the stepsize and  $G_k$  is a scaling matrix that converges to some  $G$  such that spectral radius of  $I - \gamma G C$  is smaller than 1.

- One possibility is to choose  $\gamma=1$  and then

$$G_k = \left( \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) \phi(i_t)' \right)^{-1}$$

where one can show that  $G_k \rightarrow G$  as  $C_k \rightarrow C$  and  $d_k \rightarrow d$  due to the

Law of large numbers argument.

- Alternatively, LSPE can be written as

$$r_{k+1} = r_k - \frac{\sigma}{k+1} \sum_{t=0}^k \beta(c_{it}) g_{k,t} \quad (36)$$

### TD(0) method

- May be viewed as a single sample approximation of the simulation-based (general) PVI iteration method with  $G = I$

$$\begin{aligned} r_{k+1} &= r_k - \sigma (C_k r_k - d_k) \\ &= r_k - \frac{\sigma}{k+1} \sum_{t=0}^k \beta(c_{it}) g_{k,t} \quad (37) \end{aligned}$$

$$\approx r_k - \frac{\sigma}{k+1} \beta(c_{ik}) g_{k,k} \quad \text{with } \sigma_k = \frac{\sigma}{k+1}$$

- TD(0) is much slower than LSPE or (31) but requires much less overhead
- TD(0) can be also viewed as a stochastic approximation method of the fixed point equation following

$$r = r - \sigma (Cr - d) \iff Cr = d.$$

# Multi-step PVI Methods

- Bellman's eq. for a stationary policy  $\mu$

$$J_\mu = T_\mu J_\mu = g_\mu + \alpha P_\mu J_\mu$$

- Replace  $T_\mu$  with a multi-step operator that has the same fixed point of Bellman's equation. For example,

$$T_\mu^\lambda, \lambda > 1$$

or 
$$T_\mu^\lambda = (1-\lambda) \sum_{l=0}^{\infty} \lambda^l T_\mu^{l+1}, \lambda \in (0, 1)$$

Why?

$$\begin{aligned}
J_\mu &= T_\mu J_\mu \\
\lambda J_\mu &= \lambda T_\mu^2 J_\mu \\
\lambda^2 J_\mu &= \lambda^2 T_\mu^3 J_\mu \\
&\vdots \\
\lambda^l J_\mu &= \lambda^l T_\mu^{l+1} J_\mu \\
&\vdots
\end{aligned}$$

$$\Rightarrow (1 + \lambda + \dots) J_\mu = \sum_{l=0}^{\infty} \lambda^l T_\mu^{l+1} J_\mu$$

$$\Rightarrow J_\mu = (1-\lambda) \sum_{l=0}^{\infty} \lambda^l T_\mu^{l+1} J_\mu$$

$$\begin{aligned}
&J_\mu = T_\mu^\lambda J_\mu \\
&\text{where } T_\mu^\lambda = (1-\lambda) \sum_{l=0}^{\infty} \lambda^l T_\mu^{l+1}
\end{aligned}$$

" $\lambda$ -weighted multi-step Bellman Equation"

————— (38)

- Note that

$$J_\mu = T_\mu^\lambda J_\mu = g_\mu^\lambda + \alpha P_\mu^\lambda J_\mu$$

where  $J_\mu(\lambda) = \sum_{l=0}^{\infty} \alpha^l \lambda^l P_\mu^l g_\mu = (I - \alpha \lambda P_\mu)^{-1} g_\mu$

$P_\mu(\lambda) = (1-\lambda) \sum_{l=0}^{\infty} \alpha^l \lambda^l P_\mu^{l+1}$  — (39)

- Projected Bellman equation for  $T_\mu(\lambda)$

$\Phi r = \Pi T_\mu(\lambda) \Phi r$

or

$c(\lambda) r = d(\lambda)$

where  $c(\lambda) = \Phi' M (I - \alpha P_\mu(\lambda)) \Phi$ ,  $d(\lambda) = \Phi' M g_\mu$ .

Proposition 15.

Let Assumptions 1 and 2 hold, and let  $\Pi$  be the projection w.r.t. the weighted Euclidean norm  $\|\cdot\|_2^\xi$  where  $\xi$  is the steady-state probability vector of the Markov chain corresponding to the given policy  $\mu$ . Then,

(a) The mappings  $T_\mu(\lambda)$  and  $\Pi T_\mu(\lambda)$  corresponding to  $\mu$  are contractions of modulus

$$\alpha_\lambda = \frac{\alpha(1-\lambda)}{1-\alpha\lambda}$$

w.r.t.  $\|\cdot\|_2^\xi$ .

(b) We have

$$\|J_\mu - \Phi r_\lambda^*\|_2^\xi \leq \frac{1}{\sqrt{1-\alpha_\lambda^2}} \|J_\mu - \Pi J_\mu\|_2^\xi$$

where  $\Phi r_\lambda^*$  is the  $\uparrow$  unique fixed point of  $\Pi T_\mu(\lambda)$ .

Proof) Similar to the proof of Proposition 13.  
Refer to the text book.

LSTD(x), LSPEC(x) and TDC(x)

LSTD(x) method

$$y_k = (C_k^{(x)})^{-1} d_k^{(x)}$$

LSPEC(x) method

$$y_{k+1} = y_k - \delta G_k (C_k^{(x)} y_k - d_k^{(x)}), \quad k=0, 1, 2, \dots$$

where  $\delta$  is the stepsize and  $G_k$  is a scaling matrix that converges to some  $G$  such that spectral radius of  $I - \delta G_k C^{(x)}$  is smaller than 1.

• One possibility is to choose  $\delta=1$  and

$$G_k = \left( \frac{1}{k+1} \sum_{t=0}^k \phi(i_t) \phi(i_t)' \right)^{-1}$$

where one can show that  $G_k \rightarrow G$  as  $C_k^{(x)} \rightarrow C^{(x)}$  and  $d_k^{(x)} \rightarrow d^{(x)}$  due to the law of large numbers argument.

• Introduce a vector called eligibility vector

$$z_t = \sum_{m=0}^t (\alpha \lambda)^{t-m} \phi(i_m) \quad \text{--- (40)}$$

which is the weighted sum of the present and past feature vectors  $\phi(i_m)$  from the simulations. Then, one can show that

$C_k^{(x)} \rightarrow C^{(x)}$  and  $d_k^{(x)} \rightarrow d^{(x)}$  as  $k \rightarrow \infty$  where

$$C_k^{(x)} = \frac{1}{k+1} \sum_{t=0}^k z_t ( \phi(i_t) - \alpha \phi(i_{t+1}) )'$$

and

$$d_k^{(x)} = \frac{1}{k+1} \sum_{t=0}^k z_t g(i_t, i_{t+1}) \quad \text{--- (41)}$$

In particular, (40) and (41) can be rewritten by means of recursive formulas in the form of stochastic approximation as

$$z_k = \alpha_k z_{k-1} + \phi(i_k)$$

$$c_k^{(\alpha)} = (1 - \delta_k) c_{k-1}^{(\alpha)} + \delta_k z_k (\phi(i_k) - \alpha \phi(i_{k+1}))'$$

$$d_k^{(\alpha)} = (1 - \delta_k) d_{k-1}^{(\alpha)} + \delta_k z_k g(i_k, i_{k+1})$$

with the initial conditions  $z_1 = 0, c_1 = 0, d_1 = 0$  and  $\delta_k = \frac{1}{k+1}, k=0, 1, \dots$

(42)

Alternatively, LSPE(x) can be written as

$$r_{k+1} = r_k - \frac{\gamma}{k+1} G_k \sum_{t=0}^k z_t g_{k,t} \tag{43}$$

where  $g_{k,t}$  is the tempore difference

$$g_{k,t} = \phi(i_t)' r_k - \alpha \phi(i_{t+1})' r_k - g(i_t, i_{t+1})$$

TD(x) method

- A single sample approximation of the LSPE(x) in (43) with  $G_k = I$ , i.e.,

$$r_{k+1} = r_k - \frac{\gamma}{k+1} G_k \sum_{t=0}^k z_t g_{k,t}$$

$$\approx r_k - \gamma_k z_k g_{k,k} \quad \text{with} \quad \gamma_k = \frac{\gamma}{k+1}$$