

Two Issues in Policy Iteration

- (exploration in policy evaluation step
- oscillation in policy improvement step

Exploration

- To evaluate μ , need to generate samples using μ which biases the simulation by underrepresenting states (states with small β_i) that are unlikely to occur under μ .
- As a result, the cost-to-go estimates of these underrepresented states will be highly inaccurate, and seriously impact the next improved policy $\bar{\mu}$.
- The irreducibility of P_μ of μ may be hard or impossible to guarantee.
- Need to change the sampling mechanism and the simulation formulas.
- Consider a modified transition probability matrix which is irreducible

$$\bar{P}_\mu = (I - B)P_\mu + BQ$$

where B is a diagonal matrix with diagonal components $\beta_i \in [0, 1]$ and Q is another transition probability matrix.

(Why? \bar{P}_μ could address both the underrepresented state issue and the difficulty with multiple recurrent classes and transient states)

- At state i , the next state is generated with probability $1 - \beta_i$ according to $P_{ij}(\mu(i))$, and with probability β_i according to Q_{ij} .

↳ Artificial transition \Rightarrow "exploration"
 (c.f.) exploration vs exploitation

- Consider projected Bellman equation for \bar{P}_M

$$\bar{\Phi}r = \bar{\Pi} \bar{T}_M(\alpha) \bar{\Phi}r$$

where $\bar{T}_M(\alpha)j = \bar{g}_M(\alpha) + \alpha \bar{P}_M(\alpha)j$ and $\bar{\Pi}$ denotes the projection on S w.r.t. $\|\cdot\|_{\bar{\xi}}$ where $\bar{\xi}$ is the steady-state distribution corresponding to \bar{P}_M .

- However, this projected Bellman equation is problematic because it solves r^* such that $\bar{\Phi}r^* \sim \bar{T}_M(\alpha) \bar{\Phi}r^*$, not $\bar{\Phi}r^* \sim T_M(\alpha) \bar{\Phi}r^*$.
- Instead, consider the following modified projected Bellman eq. for \bar{P}_M

$$\bar{\Phi}r = \bar{\Pi} T_M(\alpha) \bar{\Phi}r \quad \text{--- (4*)} \quad \text{"exploration-enhanced projected Bellman equation"}$$

where $T_M(\alpha)j = g_M(\alpha) + \alpha P_M(\alpha)j$ and $\bar{\Pi}$ denotes the projection on S w.r.t. $\|\cdot\|_{\bar{\xi}}$ where $\bar{\xi}$ is the steady-state distribution corresponding to \bar{P}_M .

Exploration using Extra Transitions

- Apply only to $\lambda=0$ (single-step Bellman equation) cases.
- Generate a state sequence $\{\hat{i}_0, \hat{i}_1, \dots\}$ according to any steady-state distribution $\bar{\xi}$ (more balanced one than ξ such as the uniform distribution).
- Generate a sequence of "independent" transitions $\{(i_0, j_0), (i_1, j_1), \dots\}$ according to the original transition matrix P_M .
- The exploration-enhanced projected Bellman equation (4*) yields

$$Cr = d$$

$$\text{where } C = \bar{\Phi}' M_{\bar{\xi}} (I - \alpha \bar{P}_M) \bar{\Phi}, \quad d = \bar{\Phi}' M_{\bar{\xi}} g_M.$$

- One can derive the estimators for C and d by invoking the law of large numbers argument

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

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

- Thus, we have

LSTD(0):

$$\sum_{t=0}^k \phi(i_t) \bar{\delta}_{k,t} = 0$$

where $\bar{\delta}_{k,t} = \phi(i_t)' r_k - \alpha \phi(j_t)' r_k - g(i_t, j_t)$

LSPE(0):

$$r_{k+1} = r_k - \frac{\gamma}{k+1} \left(r_k - \sum_{t=0}^k \phi(i_t) \bar{\delta}_{k,t} \right)$$

TD(0):

$$r_{k+1} = r_k - \frac{\gamma}{k+1} \phi(i_k) \bar{\delta}_{k,k}$$

Exploration using Modified Temporal Differences

- An alternative exploration approach that works for all $\lambda \geq 0$.
- Use explicit knowledge of the transition probabilities $P_{ij}(u_i)$ and \bar{P}_{ij} , so it is not a model-free approach.
- Generate a single state sequence $\{i_0, i_1, \dots\}$ according to the exploration-enhanced transition matrix \bar{P} (i.e., it is a off-policy approach).

- Define a modified exploration-enhanced temporal differences

$$\bar{g}_{k,t} = \phi(\hat{i}_t) r_k - \frac{P_{\hat{i}_t \hat{i}_{t+1}}}{\bar{P}_{\hat{i}_t \hat{i}_{t+1}}} (\alpha \phi(\hat{i}_{t+1}) r_k + g(\hat{i}_t, \hat{i}_{t+1})) \quad (45)$$

where P_{ij} and \bar{P}_{ij} denote the ij th components of P_M and \bar{P} , respectively.

- Note that the approximation of an expected value w.r.t. a given distribution (induced by the transition matrix P_M) by sampling w.r.t. a different distribution (induced by the exploration-enhanced transition matrix \bar{P}) is reminiscent of "importance sampling". The probability ratio $\frac{P_{\hat{i}_t \hat{i}_{t+1}}}{\bar{P}_{\hat{i}_t \hat{i}_{t+1}}}$ in (45) provides the necessary correction.

- One can derive the estimators for C and d using the modified temporal differences in a recursive form

$$C_k^{(t)} = (1 - \delta_k) C_k^{(t-1)} + \delta_k z_k \left(\phi(\hat{i}_k) - \alpha \frac{P_{\hat{i}_k \hat{i}_{k+1}}}{\bar{P}_{\hat{i}_k \hat{i}_{k+1}}} \phi(\hat{i}_{k+1}) \right)$$

$$d_k^{(t)} = (1 - \delta_k) d_k^{(t-1)} + \delta_k z_k \frac{P_{\hat{i}_k \hat{i}_{k+1}}}{\bar{P}_{\hat{i}_k \hat{i}_{k+1}}} g(\hat{i}_k, \hat{i}_{k+1})$$

where z_k are modified eligibility vectors given by

$$z_k = \alpha \left(\frac{P_{\hat{i}_{k-1} \hat{i}_k}}{\bar{P}_{\hat{i}_{k-1} \hat{i}_k}} z_{k-1} + \phi(\hat{i}_k) \right)$$

with $z_{-1} = 0$, $C_{-1} = 0$, $d_{-1} = 0$ and $\delta_k = \frac{1}{k+1}$, $k=0, 1, \dots$

Contraction Properties of Exploration-Enhanced Methods

$$\bar{P}r = \bar{\pi} T(\lambda) \bar{P}r$$

Proposition 16.

Assume that \bar{P} is irreducible and $\bar{\xi}$ is its invariant distribution. Then $T(\lambda)$ and $\bar{\pi} T(\lambda)$ are contractions w.r.t. $\|\cdot\|_{\bar{\xi}}$ for all $\lambda \in [0, 1)$ provided $\bar{\alpha} < 1$ where

$$\bar{\alpha} = \frac{\alpha}{\sqrt{1 - \max_{\lambda=1, \dots, n} \beta_{\lambda}}}$$

The associated modulus of contraction is at most equal to

$$\frac{\bar{\alpha}(1-\lambda)}{1-\bar{\alpha}\lambda}$$

Proof) Refer to the text book.