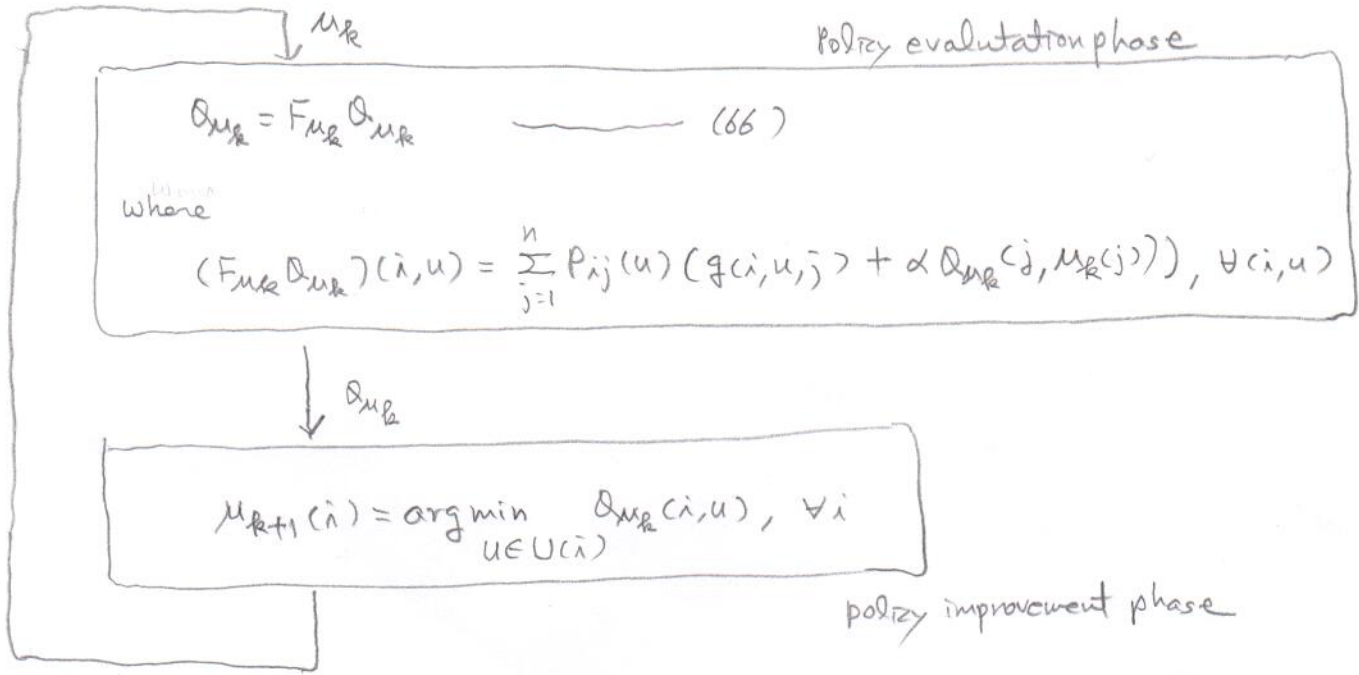


Q-Learning and PI



- Optimistic PI

$Q_{\mu_k} = F_{\mu_k}^{M_k} Q_{\mu_{k-1}}$  instead of solving (66) exactly

Simulation-based Extreme Optimistic PI

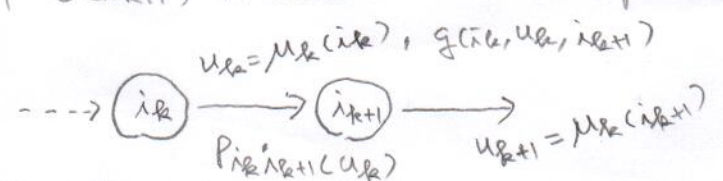
- Known as SARSA (state - Action - Reward - state - Action)
- Use a single sample between policy updates
- On-policy learning

At the start of iteration  $k$ , we have  $Q_{\mu_{k-1}}$ ,  $\mu_k$ , we are at the current state  $i_k$ , and we have chosen a control  $u_k$ . Then,

- (1) Simulate the next transition  $(i_k, i_{k+1})$  using the transition probabilities  $P_{ik+1}(u_k)$
- (2) Generate the control  $u_{k+1}$  from the minimization

$$u_{k+1} = \mu_k(i_{k+1}) = \operatorname{argmin}_{u \in U(i_{k+1})} Q_{\mu_{k-1}}(i_{k+1}, u)$$

Note that in some schemes,  $u_{k+1}$  is chosen with a small probability to be a random element of  $U(i_{k+1})$  in order to enhance exploration, e.g.)  $\epsilon$ -greedy



(3) Update Q-factor via single-sample approximation of extreme  
optimistic PI, i.e.,  $m_R = 1$

$$Q_{u_R}(i_R, u_R) = F_{u_R} Q_{u_{R-1}}(i_R, u_R)$$

$$\Rightarrow Q_{u_R}(i_R, u_R) = (1 - \gamma_R) Q_{u_{R-1}}(i_R, u_R) + \gamma_R (g(i_R, u_R, i_{R+1}) + \alpha Q_{u_{R-1}}(i_{R+1}, u_{R+1})) \quad (67)$$

$$\Rightarrow Q_{u_R}(i_R, u_R) = Q_{u_{R-1}}(i_R, u_R) + \underbrace{\gamma_R (g(i_R, u_R, i_{R+1}) + \alpha Q_{u_{R-1}}(i_{R+1}, u_{R+1}) - Q_{u_{R-1}}(i_R, u_R))}_{\text{TD}} \quad (68)$$

on-policy  $\leftarrow$

Comparison of SARSA in (67) with Q-learning in (65).

- Both SARSA and Q-learning algorithms are model-free method.
- SARSA is on-policy TD learning and Q-learning is off-policy TD learning.
- SARSA is an optimistic PI method and Q-learning is a VI method.
- Q-learning outperforms SARSA, in particular, due to the off-policy learning capability.
- The convergence of Q-learning can be proven using the theory of asynchronous distributed stochastic approximation whereas the convergence properties of SARSA are unclear.
- Both SARSA and Q-learning can be easily generalized to multistep algorithms, i.e., SARSA( $\lambda$ ) and Q-learning( $\lambda$ ), to reduce bias.
- Both SARSA and Q-learning can use linear approx. architecture or state aggregation to overcome the curse of the excessive number of Q-factors/state-control pairs  $(i, u)$ .

# Q-Factor Approximation and Projected Equations

- Introduce a linear parametriz architecture for Q

$$\tilde{Q}_\mu(\bar{i}, u, r) = \phi(\bar{i}, u)' r$$

or 
$$\tilde{Q}_\mu(\bar{i}, u, r) = \phi(\bar{i})' \delta_u, \forall u \in U(\bar{i})$$

- The projected Bellman equation for a policy  $\mu$

$$\mathbb{E}r = \Pi F_\mu(\mathbb{E}r)$$

where

$$F_\mu(\mathbb{E}r)(\bar{i}, u) = \sum_{j=1}^n P_{ij}(u) (g(\bar{i}, u, j) + \alpha(\mathbb{E}r)(j, \mu(j))) \quad (69)$$

- Simulation-based Implementation

$$\sum_{t=0}^k \phi(\bar{i}_t, u_t) g_{k,t} = 0 \quad \text{"LSTD(0)"}$$

$$r_{k+1} = r_k - \left( \sum_{t=0}^k \phi(\bar{i}_t, u_t) \phi(\bar{i}_t, u_t)' \right)^{-1} \sum_{t=0}^k \phi(\bar{i}_t, u_t) g_{k,t} \quad \text{"LSPE(0)"}$$

$$r_{k+1} = r_k - \delta_k \phi(\bar{i}_k, u_k) g_{k,k} \quad \text{"TD(0)"}$$

where  $g_{k,t}$  are the corresponding TD errors i.e.

$$g_{k,t} = \phi(\bar{i}_t, u_t)' r_k - \alpha \phi(\bar{i}_{t+1}, u_{t+1})' r_k - g(\bar{i}_t, u_t, \bar{i}_{t+1}).$$

## SARSA for approximate Q

At the start of iteration  $k$ , we have the current  $r_k$  and  $u_k$ , we are at some state  $\bar{i}_k$ , and we have chosen a control  $u_k$ . Then,

- (1) Simulate the next transition  $(\bar{i}_k, \bar{i}_{k+1})$  using the transition probabilities  $P_{ik}(u_k)$
- (2) Generate the control  $u_{k+1}$  from the minimization

$$u_{k+1} = u_k(\bar{i}_{k+1}) = \arg \min_{u \in U(\bar{i}_{k+1})} \tilde{Q}_{u_k}(\bar{i}_{k+1}, u)$$

Note that in some schemes,  $u_{k+1}$  is chosen with a small probability to be a random element of  $U(\bar{i}_{k+1})$  in order to enhance exploration, e.g.)  $\epsilon$ -greedy

(3) Update parameter vector via TD( $\phi$ )

$$r_{k+1} = r_k - \delta_k \phi(i_k, u_k) g_{k,k}$$

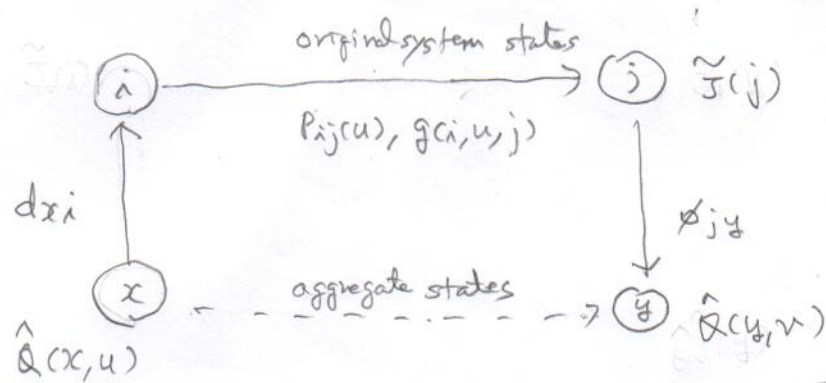
$$\Rightarrow r_{k+1} = r_k + \delta_k \phi(i_k, u_k) \left( g(i_k, u_k, i_{k+1}) + \alpha \phi(i_{k+1}, u_{k+1})' r_k - \phi(i_k, u_k)' r_k \right) \quad (10)$$

- Note that approximate SARSA in (10) can be viewed as an single-sample approximation of extreme PI ( $\mu_k=1$ ) for the projected Bellman equation in (64).
- Compare SARSA in (68) and approximate SARSA in (90).

Q-Learning and Aggregation

- Assume that the control constraint set  $U(i)$  is independent of the state and is denoted by  $U$ .
- The aggregate Bellman equation for  $\hat{Q}(x, u), x \in A, u \in U$

$$\hat{Q}(x, u) = \hat{g}(x, u) + \alpha \sum_{y \in A} \hat{P}_{xy}(u) \min_{v \in U} \hat{Q}(y, v) \quad (11)$$



$$= \sum_{i=1}^n dx_i \sum_{j=1}^n P_{ij}(u) g(i, u, j) + \alpha \sum_{y \in A} \left( \sum_{i=1}^n dx_i \sum_{j=1}^n P_{ij}(u) \phi_{jy} \right) \min_{v \in U} \hat{Q}(y, v)$$

$$= \sum_{i=1}^n dx_i \sum_{j=1}^n P_{ij}(u) \left( g(i, u, j) + \alpha \sum_{y \in A} \phi_{jy} \min_{v \in U} \hat{Q}(y, v) \right)$$

- VI method for Q

$$\hat{Q}_{k+1} = F \hat{Q}_k \tag{12}$$

where

$$(F \hat{Q}_k)(x, u) = \sum_{i=1}^n dx_i \sum_{j=1}^n P_{ij}(u) (g(i, u, j) + \alpha \sum_{y \in A} \phi_{jy} \min_{v \in U} \hat{Q}_k(y, v)), \forall (x, u)$$

- One can show that F is a maximum-norm contraction with modulus  $\alpha$  so that the VI algorithm converges to  $\hat{Q}^*$  from every starting point  $\hat{Q}_0$ .

- Simulation-based VI method for Q

Generate an infinite sequence of state-control pairs  $\{(x_k, u_k)\} \subset A \times U$ . Given the pair  $(x_k, u_k)$ , generate an original system state  $i_k$  according to the disaggregation probabilities  $dx_{k,i}$ , and then next state  $j_k$  according to the probabilities  $P_{ij}(u_k)$ . Finally, generate an aggregate state  $y_k$  according to the aggregation probabilities  $\phi_{j_k y}$ .

For each  $(x_k, u_k, i_k, j_k, y_k)$  generated, update the Q-factor of  $(x_k, u_k)$  as follows:

$$\hat{Q}_{k+1}(x, u) = (1 - \gamma_k) \hat{Q}_k(x, u) + \gamma_k (F_k \hat{Q}_k)(x, u), \forall (x, u) \tag{13}$$

where

$$(F_k \hat{Q}_k)(x, u) = \begin{cases} g(i_k, u_k, j_k) + \alpha \min_{v \in U} \hat{Q}_k(y_k, v) & \text{if } (x, u) = (x_k, u_k) \\ \hat{Q}_k(x, u) & \text{if } (x, u) \neq (x_k, u_k) \end{cases}$$

- Note that  $(F_k \hat{Q}_k)(i_k, u_k)$  in (13) is a single sample approximation of the expected values defining  $(F \hat{Q}_k)(i_k, u_k)$  in (12).

- After solving the Q-factors  $\hat{Q}$  for the aggregate system, the Q-factors  $\tilde{Q}$  for the original system are approximated by

$$\tilde{Q}(j, v) = \sum_{y \in A} \phi_{jy} \hat{Q}(y, v) \tag{14}$$

Then, the optimal cost-to-go function  $\tilde{J}(j)$  of the original system is approximated by

$$\tilde{J}(j) = \min_{v \in U} \tilde{Q}(j, v).$$

and the corresponding suboptimal policy  $\tilde{\mu}(i)$  can be obtained by

$$\tilde{\mu}(i) = \arg \min_{v \in U}$$