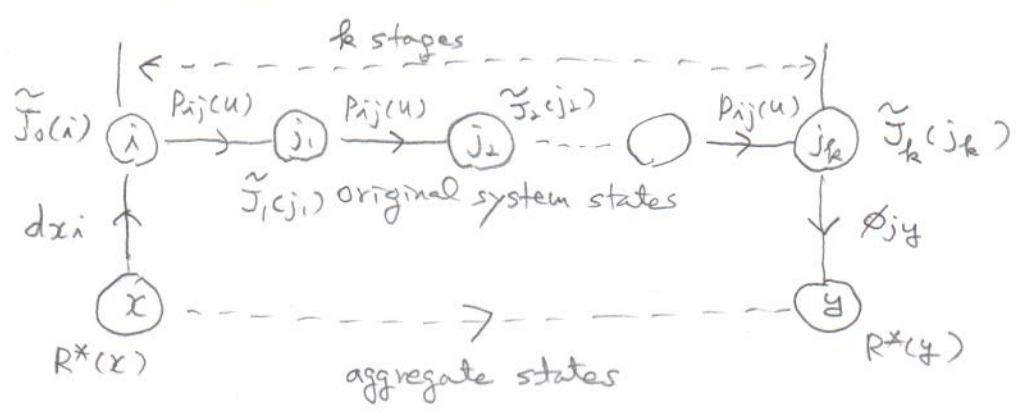


Multistep Aggregation



- Bellman equation for R^*

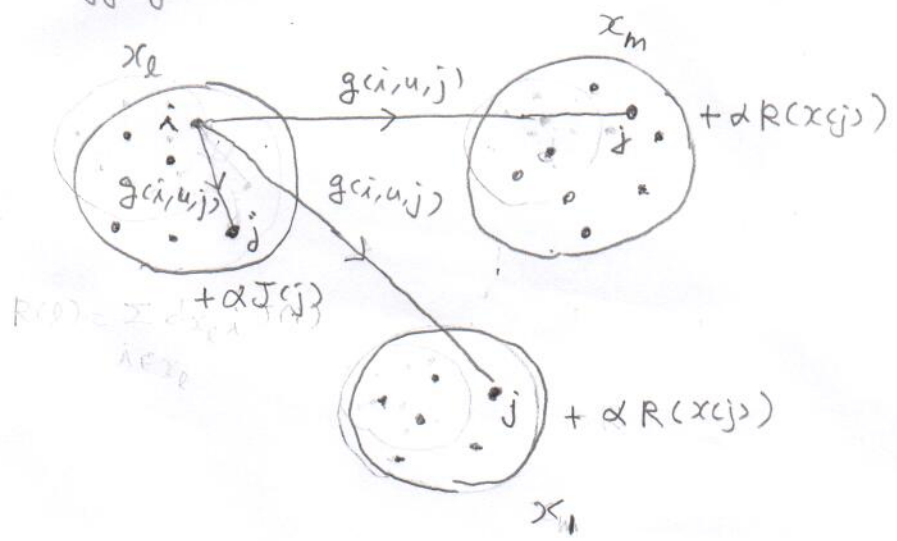
$$R^* = DT^k \Phi R^* \quad (56)$$

- Improved error bound

$$\| \tilde{J}_k - J^* \|_\infty \leq \frac{2\alpha^k \delta}{1 - \alpha^k}$$

Asynchronous Distributed Aggregation

- Partition original system states into a set of aggregate states $A = \{x_1, \dots, x_m\}$
 i.e., hard aggregation



- Bellman equation for each $i \in x_\ell$

$$\cong H_\ell(i, u, J, R)$$

$$J(i) = \min_{u \in U(i)} \left(\sum_{j=1}^n P_{ij}(u) g_{(i,u,j)} + \alpha \sum_{j \in x_\ell} P_{ij}(u) J(j) + \sum_{j \notin x_\ell} P_{ij}(u) R(x_j) \right)$$

where $R(x) = \sum_{i \in x_\ell} d_{xi} J(i)$

Synchronous distributed VI method

Each processor $l=1, \dots, m$ maintains/updates a local cost $J(i)$ for $i \in X_l$ and maintains/updates/communicates an aggregate cost

$$R(l) = \sum_{i \in X_l} d_{X_l i} J(i)$$

according to the following update rule:

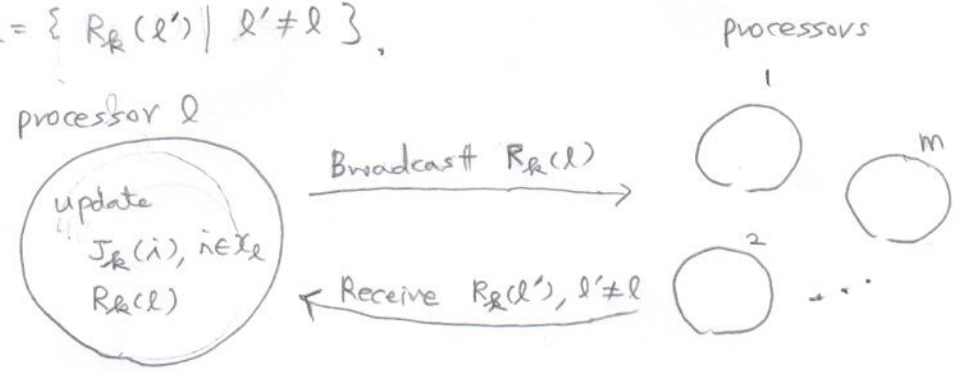
$$J_{k+1}(i) = \min_{u \in U(i)} H_l(i, u, J_k, R_k) \quad (58)$$

$$R_k(l) = \sum_{i \in X_l} d_{X_l i} J_k(i)$$

where $H_l(i, u, J_k, R_k) = \sum_{j=1}^n P_{ij}(u) g(i, u, j) + \alpha \sum_{j \in X_l} P_{ij}(u) J_k(j) + \alpha \sum_{j \notin X_l} P_{ij}(u) R_k(j)$

$$J_k = \{ J_k(i) \mid i \in X_l \}$$

$$R_k = \{ R_k(l') \mid l' \neq l \}$$



Asynchronous distributed VI method

$$J_{k+1}(i) = \begin{cases} \min_{u \in U(i)} H_i(i, u, J_k, R_{1,k}(1), \dots, R_{m,k}(m)), & \forall i \in I_l, k \\ J_k(i) & \text{otherwise} \end{cases}$$

where $0 \leq \tau_{l,k} \leq k$ for $l=1, \dots, m$,

$$R_{\tau}(l) = \sum_{i \in X_l} d_{X_l i} J_{\tau}(i), \quad \forall l=1, \dots, m$$

and

- $I_{l,k}$: the set of states $i \in X_l$ which are updated at time k
- $\tau_{l,k}$: the time when the corresponding aggregate cost, lastly reported to processor l by time k , was computed at the corresponding processor
- $k - \tau_{l,k}$: delays between the current time k and the time $\tau_{l,k}$ when the corresponding aggregate costs were computed at other processors.

- Convergence of the asynchronous distributed VC method

a) Total asynchronism assumption

Every $i \in X_l$ belongs to $I_{l,k}$ for infinitely many k (so each cost component is updated infinitely often) and $\lim_{k \rightarrow \infty} \tau_{l,k} = \infty$ for all $l = 1, \dots, m$ (so that processors eventually communicate more recently computed aggregate costs to other processors).

b) Asynchronous convergence theorem

The maximum-norm contraction property of the mapping underlying the synchronous iteration (58) establishes the total asynchronous convergence.

Proof) Refer to Chapter 2 of the text book "Dynamic programming and Optimal Control: Approximate Dynamic Programming", Vol. 2, 4th Edition, Athena Scientific 2012.

Refer to "Parallel and Distributed Computation: Numerical Methods", Athena Scientific 2015.

Q-Learning

(48)

- Used when there is no explicit model of the system, i.e., a model-free method
- Instead of approximating the cost function of a particular policy, it updates Q-factors with an optimal policy, thereby avoiding the multiple policy evaluation steps of the PI method.

Q-Learning: A stochastic VI Algorithm

- The optimal Q-factors are defined for all pairs (i, u) with $u \in U(i)$,

by

$$Q^*(i, u) = \sum_{j=1}^n P_{ij}(u) (g(i, u, j) + \alpha J^*(j)) \quad (60)$$

- Note that

optimal Q-factor = expected 1st stage cost when action u is applied
+ expected cost-to-go associated with an optimal policy when action u is applied.

- In contrast, the Bellman's equation for J_μ is defined for all i by

$$J_\mu(i) = \sum_{j=1}^n P_{ij}(\mu(i)) (g(i, \mu(i), j) + \alpha J_\mu(j)) \quad J_\mu = T_\mu J_\mu$$

i.e., expected 1st stage cost associated with a policy μ
+ expected cost-to-go associated with a policy μ

- The Bellman's equation for J is defined for all i by

$$J(i) = \min_{u \in U(i)} \sum_{j=1}^n P_{ij}(u) (g(i, u, j) + \alpha J(j)) \quad J = TJ \quad (61)$$

i.e., minimum $\left[\begin{array}{l} \text{expected 1st stage cost when action } u \text{ is applied} \\ \text{+ expected cost-to-go associated with an optimal policy} \\ \text{when action } u \text{ is applied} \end{array} \right]$

- By combining (60) and (61), we obtain the Bellman's equation for $Q^*(i, u)$

$$Q^*(i, u) = \sum_{j=1}^n P_{ij}(u) (g(i, u, j) + \alpha \min_{v \in U(j)} Q^*(j, v)) \quad (62)$$

and the relationship between $J(i)$ and $Q^*(i, u)$ as

$$J(i) = \min_{u \in U(i)} Q^*(i, u) \quad (63)$$

- The proof of existence and uniqueness of the fixed point of the Bellman equation (62) is the same as the proof of that of the Bellman equation for J .

- VI method for Q^*

$$Q_{k+1} = F Q_k$$

where

$$(FQ)(i, u) = \sum_{j=1}^n P_{ij}(u) (g(i, u, j) + \alpha \min_{v \in U(j)} Q(j, u)), \quad \forall (i, u)$$

- One can show that F is a maximum-norm contraction with modulus α so that the VI algorithm converges to Q^* from every starting point Q_0 .

- Simulation-based VI method for Q^*

Generate an infinitely long sequence of state-action pairs $\{(i_k, u_k)\}$

Given the pair (i_k, u_k) , generate next state j_k according to the probabilities $P_{ij}(u_k)$.

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

$$Q_{k+1}(i, u) = (1 - \delta_k) Q_k(i, u) + \delta_k (F_k Q_k)(i, u); \quad \forall (i, u)$$

where

$$(F_k Q_k)(i, u) = \begin{cases} g(i_k, u_k, j_k) + \alpha \min_{v \in U(j_k)} Q_k(j_k, v) & \text{if } (i, u) = (i_k, u_k) \\ Q_k(i, u) & \text{if } (i, u) \neq (i_k, u_k) \end{cases}$$

- Note that $(F_k Q_k)(i_k, u_k)$ is a single sample approximation of the expected value defining $(F Q_k)(i_k, u_k)$ in (64). — (65)

- The stepsize δ_k should satisfy

$$\delta_k > 0, \quad \forall k, \quad \sum_{k=0}^{\infty} \delta_k = \infty, \quad \sum_{k=0}^{\infty} \delta_k^2 < \infty$$

which are typical of stochastic approximation methods, as for example when $\delta_k = c_1 / (ck + c_2)$ where c_1 and c_2 are some positive constants.

- The rigorous convergence proof is beyond our scope; see Tsitsiklis, J. N., "Asynchronous Stochastic Approximation and Q-Learning", Vol. 16, pp. 185-202, which embeds Q-learning with a broad class of asynchronous stochastic approximation algorithms.

- The simulation-based VI method in (65), which belongs to a class of asynchronous stochastic approximation algorithms, is the original and most widely known Q-learning algorithm. This algorithm was originally reported in C. J. C. H. Watkins, "Learning from Delayed Rewards," Ph.D. Thesis, Cambridge University.

- In practice, Q-learning has some drawbacks, the most important of which is that the number of Q-factors/state-control pairs (i, u) may be excessive.

\Rightarrow Linear approximation architecture for Q and/or state aggregation