

## 5.6 Q-Learning

Saturday, February 13, 2021 1:52 PM

Q-Learning directly updates estimates of Q-factors associated with an optimal policy.

Let us define the optimal Q-factor,

$$Q^*(i, u) = \sum_{j=0}^n p_{ij}(u) (g(i, u, j) + \gamma^* j) \quad i=0, \dots, n.$$

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

Combining the above two, we get,

$$Q^*(i, u) = \sum_{j=0}^n p_{ij}(u) (g(i, u, j) + \min_{v \in U(j)} Q^*(j, v)).$$

Using the uniqueness of Bellman's eq.,

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

In terms of Q-factors,

$$Q(i, u) = \sum_{j=0}^n p_{ij}(u) (g(i, u, j) + \min_{v \in U(j)} Q(j, v)).$$

More generally,

$$Q(i, u) = (1-\gamma) Q(i, u) + \gamma \sum_{j=0}^n p_{ij}(u) (g(i, u, j) + \min_{v \in U(j)} Q(j, v)).$$

Q-Learning is an approximate version of the above update where the expectation is replaced by a single sample,

$$Q(i, u) = (1-\gamma) Q(i, u) + \gamma (g(i, u, j) + \min_{v \in U(j)} Q(j, v)).$$

The Convergence of Q-Learning -

We first write a general statement.

$$Q_{\gamma t}(i, u) = (1-\gamma) Q(i, u) + \gamma (g(i, u, j) + \min_{v \in U(j)} Q(j, v))$$

$$Q_{\gamma t}(i, u) = 0 \text{ for } t \notin T^u. \text{ and } Q_{\gamma t}(0, u) = 0 \quad \forall t$$

Proposition-5: Suppose that -

$$\sum_{t=0}^{\infty} \gamma^t Q_{\gamma t}(i, u) = \infty, \quad \sum_{t=0}^{\infty} \gamma^t Q_{\gamma t}(i, u) < \infty, \quad \forall i, u \in U(i).$$

Then  $Q(i, u)$  converges to  $Q^*(i, u)$  with probability 1, in each of the following cases -

(i) If all policies are proper.

(ii) If Assumptions 1 and 2 hold and if  $Q(i, u)$  is bounded with probability 1.

Proof: Consider a mapping  $H$ ,

$$HQ(i, u) = \sum_{j=0}^n p_{ij}(u) (g(i, u, j) + \min_{v \in U(j)} Q(j, v)),$$

Then Q-Learning is of the form,

$$Q_{\gamma t+1}(i, u) = (1-\gamma) Q_{\gamma t}(i, u) + \gamma (g(i, u, j) + \min_{v \in U(j)} Q(j, v)).$$

$$\text{where } w_{\gamma t}(i, u) = g(i, u, j) + \min_{v \in U(j)} Q_{\gamma t}(j, v) = \sum_{j=0}^n p_{ij}(u) (g(i, u, j) + \min_{v \in U(j)} Q(j, v)).$$

Note that  $E(w_{\gamma t}(i, u) | F_t) = 0$  and  $E(w_{\gamma t}^2(i, u) | F_t) \leq \lambda (1 + \max_{j, v} Q^*(j, v))$ .

(i) We first consider the case when all policies are proper,

According to Chapter-2 there exist  $\zeta(i) > 0, \forall i$ , and a scalar  $\beta \in (0, 1)$ ,

$$\sum_{j=1}^n p_{ij}(u) \zeta(j) \leq \beta \zeta(i), \quad \forall i, u \in U(i).$$

For any two vectors  $Q$  and  $\bar{Q}$ ,

$$|HQ(i, u) - H\bar{Q}(i, u)| \leq \sum_{j=1}^n p_{ij}(u) \left| \min_{v \in U(j)} Q(j, v) - \min_{v \in U(j)} \bar{Q}(j, v) \right|,$$

$$\leq \sum_{j=1}^n p_{ij}(u) \max_{v \in U(j)} |Q(j, v) - \bar{Q}(j, v)|,$$

$$\leq \sum_{j=1}^n p_{ij}(u) \|Q - \bar{Q}\|_1 \zeta(j),$$

$$\leq \beta \|Q - \bar{Q}\|_1 \zeta(i).$$

Dividing both sides by  $\zeta(i)$  and taking min over  $\forall i$  and  $u \in U(i)$ ,

$$\|HQ - H\bar{Q}\|_1 \leq \beta \|Q - \bar{Q}\|_1$$

$\therefore H$  is a weighted max norm contraction.

Convergence of Q-Learning follows from

Proposition-4 of Chapter-4.

(ii) In the second case, we impose assumptions 1 and 2 and remove the case of all proper policies,

Clearly  $H$  is a monotone mapping,

$$Q \leq \bar{Q} \Rightarrow HQ \leq H\bar{Q},$$

Also, for  $\gamma$  a positive scalar and  $\zeta$  being a vector of all components equal to 1,

$$H(\zeta - \gamma \zeta) \leq H(\zeta + \gamma \zeta) \leq H\zeta + \gamma \zeta.$$

Finally,  $H\zeta - \gamma \zeta$  has a unique solution and converges to  $H\zeta$  with probability 1 based on Proposition-6 from Chapter-4.

Q-Learning for Discounted Problems

We may not convert the discounted problem to a stochastic shortest path problem. This can be handled directly,

$$Q(i, u) = (1-\gamma) Q(i, u) + \gamma (g(i, u, j) + \gamma \min_{v \in U(j)} Q(j, v))$$

Here,  $\gamma$  is the discount factor.

Similar convergence results are observed for the stochastic shortest path case.