

150AML Spring 2011

Summary of part I
Problem Formulation and
Algorithms from [SB]

150AML

Roni Khardon, Tufts University

MDP Model

- Given by transitions $\Pr(s'|s,a)$
- reward $r(s,a)$
- Criterion: expected total discounted reward (discount factor gamma)

150AML

Roni Khardon, Tufts University

MDP Model

- Two problems:
 - calculate value of policy
 - calculate optimal value and policy
- We may or may not have a model and may or may not construct one during calculation

150AML

Roni Khardon, Tufts University

Backup Operators

Bellman Backup

$$[B(V)](s) = \max_a [r(s,a) + \sum_{s'} \Pr(s'|s,a)V(s')]$$

Bellman Backup restricted to policy

$$[B^\Pi(V)](s) = r(s,\Pi(s)) + \sum_{s'} \Pr(s'|s,\Pi(s))V(s')$$

Bellman Backup restricted to state

$$[B_{s^*}(V)](s^*) = \max_a [r(s^*,a) + \sum_{s'} \Pr(s'|s^*,a)V(s')]$$

$$[B_{s^*}(V)](s) = V(s) \text{ when } s \neq s^*$$

150AML

Roni Khardon, Tufts University

Policy Evaluation (calculate V^Π)

- Solve linear equations $V = B^\Pi(V)$
- Iterative Alg: Repeat $V \leftarrow B^\Pi(V)$
- Monte Carlo: $V(s_t) \leftarrow V(s_t) + \alpha[R_t - V(s_t)]$
- TD(0): $V(s_t) \leftarrow V(s_t) + \alpha[r_t + \gamma V(s_{t+1}) - V(s_t)]$
- TD(lambda): $\lambda=0 \Rightarrow$ TD(0) $\lambda=1 \Rightarrow$ MC

150AML

Roni Khardon, Tufts University

Optimization

- VI: Repeat $V \leftarrow B(V)$
- PI: Repeat $\Pi = \text{greedy}(V); V = V^\Pi$
- MPI: Repeat
 - $[\Pi = \text{greedy}(V); \text{Repeat } m \text{ times: } V = B^\Pi(V)]$
- Linear Programming
 - minimize $\sum_s V(s)$ subject to $V \geq B^a(V)$

150AML

Roni Khardon, Tufts University

On Line Optimization (RTDP)

Repeat:

pick start state s

Repeat: [in state s]

$$\forall a, Q(s,a) = r(s,a) + \gamma \sum_{s'} \Pr(s'|s,a) [\max_{a'} Q(s',a')]$$

$P =$ greedy or epsilon-greedy w.r.t. Q

Choose action a using policy P

Sample s' from $\Pr(s'|s,a)$

$s=s'$

Note: This performs $B_s(V)$ on the $Q()$ representation

150AML

Roni Khardon, Tufts University

On Line Optimization (SARSA)

Repeat:

[in state s] take action a ; observe r,s'

choose next action a' using policy P

$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma Q(s',a') - Q(s,a)]$$

$P =$ epsilon-greedy w.r.t. Q

$s=s'; a=a'$

150AML

Roni Khardon, Tufts University

On Line Optimization (SARSA lambda)

Repeat:

[in state s] take action a ; observe r,s'

choose next action a' using policy P

$$\delta = r + \gamma Q(s',a') - Q(s,a)$$

$$e(s,a) = e(s,a) + \delta$$

$$\text{for all } (s^*, a^*), Q(s^*, a^*) = Q(s^*, a^*) + \alpha \delta e(s^*, a^*)$$

$$e(s^*, a^*) = \gamma \lambda e(s^*, a^*)$$

$P =$ epsilon-greedy w.r.t. Q

$s=s'; a=a'$

150AML

Roni Khardon, Tufts University

On Line Optimization (Q learning)

Repeat:

[in state s] take exploration policy action a ;

observe r,s'

$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

150AML

Roni Khardon, Tufts University

State Aggregation

- Partition states into disjoint sets

$$p_1, p_2, \dots, p_n$$

- State s mapped to $p_i = p(s)$

$$[B_{p_i}(V)](p(s^*)) = \max_a [r(s^*, a) + \sum_{s'} \Pr(s'|s^*, a) V(p(s'))]$$

$$[B_{p_i}(V)](p(s)) = V(p(s)) \text{ when } p(s) \neq p(s^*)$$

- Parametric aggregation more general

150AML

Roni Khardon, Tufts University

On Line Optimization (SARSA)

Repeat:

[in state s] take action a ; observe r,s'

choose next action a' using policy P

$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma Q(s',a') - Q(s,a)]$$

$P =$ epsilon-greedy w.r.t. Q

$s=s'; a=a'$

150AML

Roni Khardon, Tufts University

On Line Optimization (SARSA)

(with linear function approximation)

Repeat:

[in state s] take action a ; observe r, s'

choose next action a' using policy P

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$$

$$\bar{w} \leftarrow \bar{w} + \alpha[r + \gamma Q_{\bar{w}}(s', a') - Q_{\bar{w}}(s, a)] \bar{\phi}(s, a)$$

$P = \text{epsilon-greedy w.r.t. } Q$

(P represented implicitly)

$s = s'$; $a = a'$

150AML

Roni Khardon, Tufts University

On Line Optimization (SARSA)

(with linear function approximation)

Repeat:

[in state s] take action a ; observe r, s'

choose next action a' using policy P

$$\bar{w} \leftarrow \bar{w} + \alpha[r + \gamma Q_{\bar{w}}(s', a') - Q_{\bar{w}}(s, a)] \bar{\phi}(s, a)$$

$P = \text{epsilon-greedy w.r.t. } Q$

(P represented implicitly)

$s = s'$; $a = a'$

150AML

Roni Khardon, Tufts University

Dyna-SARSA-RTDP

Repeat:

[in state s] take action a ; observe r, s'

choose next action a' using policy P

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$$

$$r(s, a) = \dots \quad p(s^* | s, a) = \dots$$

Update Q : RTDP with $r()$, $p()$ estimates

$P = \text{epsilon-greedy w.r.t. } Q$

$s = s'$; $a = a'$

150AML

Roni Khardon, Tufts University