

Comp 135 Introduction to Machine Learning and Data Mining

Fall 2016

Professor: Roni Khardon

Computer Science
Tufts University

MDP Model

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

comp135

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

comp135

Roni Khardon, Tufts University

Backup Operators

Bellman Backup

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

Extracting a policy

$$[Greedy(V)](s) = \operatorname{argmax}_a [r(s, a) + \sum_{s'} \Pr(s' | s, a) V(s')]$$

Bellman Backup restricted to policy

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

comp135

Roni Khardon, Tufts University

Policy Evaluation (calculate V^π)

- Solve linear equations $V = B^\pi(V)$
- Iterative Alg: Repeat $V \leftarrow B^\pi(V)$
- The solution is V^π

comp135

Roni Khardon, Tufts University

Planning / Optimization

- VI: Repeat $V \leftarrow B(V)$
- PI: Repeat $\pi \leftarrow \text{greedy}(V); V \leftarrow V^\pi$
- Another view of PI:
 - Repeat $Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} p(s' | s, a) V^\pi(s')$
 - $\pi(s) = \operatorname{argmax}_a Q^\pi(s, a)$

comp135

Roni Khardon, Tufts University

Learning

- Transition and reward model not given
- Learn model and plan, or use model free method
- Bandits: are "1 state MDPs"
- MC: Monte Carlo: evaluate $Q(s,a)$ using independent random rollouts
- TD: Temporal Difference: $Q(s,a)$ estimate uses previous value of next state in the rollout

comp135

Roni Khardon, Tufts University

Exploration policy

- is crucial so that active RL does not get trapped with good estimate of bad policy
- Epsilon-exploration: pick optimal action with prob= $[1-\epsilon]$ and random action with prob= ϵ
- Softmax exploration

$$p(a_i) = \frac{e^{Q(a_i)/T}}{\sum_k e^{Q(a_k)/T}}$$

comp135

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

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

comp135

Roni Khardon, Tufts University

On Line Optimization (Q learning)

Repeat:

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

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

comp135

Roni Khardon, Tufts University

RL in practice

- Cannot afford to enumerate states
- In some problems cannot afford to enumerate actions
- Must use generalization.
- The $V()$, $Q()$, $\pi()$ are explicitly represented as functions of state/action
- (e.g. decision tree; neural network)
- Adapt algorithms to learn these representation

comp135

Roni Khardon, Tufts University