Reinforcement Learning

Reinforcement Learning Lectures 4 and 5

Gillian Hayes

18th January 2007

informatics

Gillian Hayes

```
RL Lecture 4/5
```

18th January 2007

- informatics

- Framework
- Rewards, Returns
- Environment Dynamics
- Components of a Problem
- $\bullet\,$ Values and Action Values, V and Q
- Optimal Policies
- Bellman Optimality Equations

Goal: maximise total reward received

Gillian Hayes

RL Lecture 4/5

18th January 2007

a informatics

State/ Situation s t Reward r AGENT

Action a_t Where is boundary between agent and environment?

Task: one instance of an RL problem – one problem set-up Learning: how should agent change policy? Overall goal: maximise amount of reward received over time

r t+1

s t+1



Goals and Rewards

Immediate reward r at each step. We must maximise expected cumulative reward:

Return = Total reward $R_t = r_{t+1} + r_{t+2} + r_{t+3} + \cdots + r_{\tau}$

Discounted Reward

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots$$
$$= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

 $0 \leq \gamma < 1$ discount factor \rightarrow discounted reward finite if reward sequence $\{r_k\}$ bounded

 $\gamma=0:$ myopic $~~\gamma\rightarrow1:$ agent far-sighted. Future rewards count for more

Dynamics of Environment

Choose action a in situation s: what is the probability of ending up in state s'? Transition probability

$$P_{ss'}^{a} = \Pr\{s_{t+1} = s' \mid s_t = s, a_t = a\}$$

BACKUP DIAGRAM



 $R^a_{ss'}$ and $\rho(s,a)$

Reward functions

 $R^a_{ss'}$ expected next reward given current state s and action a and next state s' $\rho(s,a)$ expected next reward given current state s and action a

$$\rho(s,a) = \sum_{s'} P^a_{ss'} R^a_{ss'}$$

Sometimes you will see $\rho(s,a)$ in the literature, especially that prior to 1998 when S+B was published.

Sometimes you'll also see $\rho(s)$. This is the reward for being in state s and is equivalent to a "bag of treasure" sitting on a grid-world square (e.g. computer games – weapons, health).

7 Informatics

Sutton and Barto's Recycling Robot 1

RL Lecture 4/5

If action a chosen in state s and subsequent state reached is s' what's the

 $R^a_{aa'} = E\{r_{t+1} \mid s_t = s, a_t = a, s_{t+1} = s'\}$

If we know P and R then have complete information about environment – may

- At each step, robot has choice of three actions:
 - go out and search for a can

Dynamics of Environment

expected reward?

need to learn them

Gillian Hayes

- wait till a human brings it a can
- go to charging station to recharge
- Searching is better (higher reward), but runs down battery. Running out of battery power is very bad and robot needs to be rescued
- Decision based on current state is energy high or low
- Reward is no. cans (expected to be) collected, negative reward for needing rescue

This slide and the next based on an earlier version of Sutton and Barto's own slides from a previous Sutton web resource.

informatics

18th January 2007

Sutton and Barto's Recycling Robot 2

 $\begin{array}{ll} \mathsf{S}{=}\{\mathsf{high}, \mathsf{low}\} & \mathsf{A}(\mathsf{high}) = \{\mathsf{search}, \mathsf{wait}\} & \mathsf{A}(\mathsf{low}) = \{\mathsf{search}, \mathsf{wait}, \mathsf{recharge}\} \\ \mathsf{R}^{\mathrm{search}} & \mathsf{expected} \text{ no. cans when searching} & \mathsf{R}^{\mathrm{wait}} & \mathsf{expected} \text{ no. cans when waiting} \\ \mathsf{R}^{\mathrm{search}} > R^{\mathrm{wait}} \end{array}$



10 informatics

Action Values Q

Q-Action Value of taking action a in state s under policy π is $Q^{\pi}(s, a) =$ expected return starting in s, taking a and then following policy π

$$Q^{\pi}(s,a) = E_{\pi}\{R_t \mid s_t = s, a_t = a\}$$

= $E_{\pi}\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a\}$

What is the backup diagram?



Values V

Policy π maps situations $s \in S$ to (probability distribution over) actions $a \in A(s)$ V-Value of s under policy π is $V^{\pi}(s)$ = expected return starting in s and following policy π

$$V^{\pi}(s) = E_{\pi}\{R_t \mid s_t = s\} = E_{\pi}\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s\}$$

BACKUP DIAGRAM FOR V(s)



11 Informatics

Recursive Relationship for V

$$V^{\pi}(s) = E_{\pi} \{R_{t} \mid s_{t} = s\}$$

= $E_{\pi} \{\sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} \mid s_{t} = s\}$
= $E_{\pi} \{r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} \mid s_{t} = s\}$
= $\sum_{a} \pi(s, a,) \sum_{s'} P^{a}_{ss'} [R^{a}_{ss'} + \gamma E_{\pi} \{\sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} \mid s_{t+1} = s'\}]$
= $\sum_{a} \pi(s, a,) \sum_{s'} P^{a}_{ss'} [R^{a}_{ss'} + \gamma V^{\pi}(s')]$

This is the BELLMAN EQUATION. How does it relate to backup diagram?

Gillian Hayes

13 informatics

Recursive Relationship for Q

$$Q^{\pi}(s,a) = \sum_{s'} P^{a}_{ss'}[R^{a}_{ss'} + \gamma \sum_{a'} \pi(s',a')Q(s',a')$$

Relate to backup diagram

Grid World Example

Check the V's comply with Bellman Equation From Sutton and Barto P. 71, Fig. 3.5



3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

Gillian Hayes RL Lectu	e 4/5 18th January 2007	Gillian Hayes	RL Lecture 4/5	18th January 2007
	14 Informatics			15 informatics
Relating (Q and V		Relating V and Q	
$Q^{\pi}(s,a) = E_{\pi} \{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} \}$ $= E_{\pi} \{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} \}$ $= \sum_{s'} P^{a}_{ss'} [R^{a}_{ss'} + \gamma^{b}_{ss'}]$ $= \sum_{s'} P^{a}_{ss'} [R^{a}_{ss'} + \gamma^{b}_{ss'}]$	$ s_{t} = s, a_{t} = a \}$ $\gamma^{k} r_{t+k+2} s_{t} = s, a_{t} = a \}$ $E_{\pi} \{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} s_{t+1} = s' \}]$ $V_{\pi}(s')]$		$V^{\pi}(s) = E_{\pi} \{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} \mid s_{t} = \sum_{a} \pi(s, a) Q^{\pi}(s, a) $	<i>s</i> }



Optimal Policies π^*

An optimal policy has the highest/optimal value function $V^*(s)$ It chooses the action in each state which will result in the highest return

Optimal Q-value $Q^*(s,a)$ is reward received from executing action a in state s and following optimal policy π^* thereafter

$$V^{*}(s) = \max_{\pi} V^{\pi}(s)$$
$$Q^{*}(s,a) = \max_{\pi} Q^{\pi}(s,a)$$
$$Q^{*}(s,a) = E\{r_{t+1} + \gamma V^{*}(s_{t+1}) \mid s_{t} = s, a_{t} = a\}$$

RL Lecture 4/5

Bellman Optimality Equations 1

Bellman equations for the optimal values and Q-values

$$V^{*}(s) = \max_{a} Q^{\pi^{*}}(s, a)$$

= $\max_{a} E_{\pi^{*}} \{ R_{t} \mid s_{t} = s, a_{t} = a \}$
= $\max_{a} E_{\pi^{*}} \{ r_{t+1} + \gamma \sum_{k} \gamma^{k} r_{t+k+2} \mid s_{t} = s, a_{t} = a \}$
= $\max_{a} E\{ r_{t+1} + \gamma V^{*}(s_{t+1}) \mid s_{t} = s, a_{t} = a \}$
= $\max_{a} \sum_{s'} P^{a}_{ss'}[R^{a}_{ss'} + \gamma V^{*}(s')]$

Gillian Hayes

RL Lecture 4/5

18th January 2007

Bellman Optimality Equations 1 18 informatics

$$Q^{*}(s,a) = E\{r_{t+1} + \gamma \max_{a'} Q^{*}(s_{t+1},a') \mid s_{t} = s, a_{t} = a\}$$
$$= \sum_{s'} P^{a}_{ss'}[R^{a}_{ss'} + \gamma \max_{a'} Q^{*}(s',a')]$$

Value under optimal policy = expected return for best action from that state.

19 informatics

Bellman Optimality Equations 2

If dynamics of environment $R^a_{ss^\prime}, P^a_{ss^\prime}$ known, then can solve equations for V^* (or $Q^*).$

Given V^* , what then is optimal policy? I.e. which action a do you pick in state s?

The one which maximises expected $r_{t+1} + \gamma V^*(s_{t+1})$, i.e. the one which gives the biggest

 $\sum_{s'}$ (instant reward + discounted future maximum reward)* $P^a_{ss'}$

So need to do one-step search

Bellman Optimality Equations 2 20 Informatics	21 informatics			
There may be more than one action doing this \rightarrow all OK All GREEDY actions	Assumptions for Solving Bellman Optimality Equations			
Given Q^* , what's the optimal policy? The one which gives the biggest $Q^*(s, a)$, i.e. in state s , you have various Q values, one per action. Pick (an) action with largest Q .	 Know dynamics of environment P^a_{ss'}, R^a_{ss'} Sufficient computational resources (time, memory) BUT Example: Backgammon OK 10²⁰ states ⇒ 10²⁰ equations in 10²⁰ unknowns, nonlinear equations (max) Often use a neural network to approximate value functions, policies and models ⇒ compact representation Optimal policy? Only needs to be optimal in situations we encounter - some very rarely/never encountered. So a policy that is only optimal in those states we encounter may do 			
Gillian Hayes RL Lecture 4/5 18th January 2007	Gillian Hayes RL Lecture 4/5 18th January 2007			
22 informatics	Components of an RL Problem 23			
Components of an RL Problem	• value – V or Q			
Agent, task, environment	• policy			
States, actions, rewards	• poncy			
Policy $\pi(s, a) \rightarrow$ probability of doing a in s	• model			
Value $V(s) \rightarrow$ number – Value of a state	sometimes subject to conditions, e.g. learn best policy you can within given time			
Action value $Q(s, a)$ – Value of a state-action pair Model \mathcal{D}^a_{a} – probability of going from $a \to a'$ if do a	Learn to maximise total future discounted reward			
Nodel $\Gamma_{ss'} \rightarrow$ probability of going from $s \rightarrow s$ if do <i>a</i>				
Return $R \rightarrow sum$ of future rewards				
Total future discounted rewards $r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \cdots = \sum_{k=0}^{\infty} r_{t+k+1} \gamma^k$				
Learning strategy to learn (continued)				

Agent, task, environment	RL Buzzwords
--------------------------	---------------------

Actions, situations/states, rewards Policy Environment dynamics and model

Return, total reward, discounted rewards Value function V, action-value function Q Optimal value functions and optimal policy Complete and incomplete environment information Transition probabilities and reward function Model-based and model-free learning methods

Temporal and spatial credit assignment

Exploration/exploitation tradeoff

Gillian Hayes

RL Lecture 4/5

18th January 2007