Using Q-Learning: Asterix and Obelix

Reinforcement Learning Lecture 13

Gillian Hayes

19th February 2007

informatics

Gillian Hayes	RL Lecture 13	19th February 2007	Gillian Hayes	RL Lecture 13	19th February 2007	
		a informatics			2 informatics	
A Q-Learning Example			Obelix Robot			
	C .		 Actions: move forw 	ard, turn left, turn hard left, turn ri	ght, turn hard right	
Obelix – Mahadevan and Connell (IBM)		 Sensors: sonar (35ft, 20degrees, time to echo), infrareds (4 inches), motorcurrent (increases if robot stuck) Can measure distance to object to nearest 1/4 inch, but use range bins: near/far Blind spot: object permanence a problem – if the robot turns through a small angle, small object disappear until the next sensor picks them up 				
• Asterix – John Hoar (DAI MSc 1996)						
• A behaviour-based robot learning to push a box						
• Learns behaviours by trial and error: Q-learning			• State space $2^{18} = 262144$ states: is this too many?			
 Rewards/punishments for doing right/wrong thing 			• Find box, push it, don't get stuck.			
• Simulator and real robot			• How to organise the RL problem? In one monolithic block? In submodules that			
• Reimplemented (he	re) on Asterix		learn parts of the task	</td <td></td>		
How to frame a prob	lem so that we can use reinforceme	nt learning with it	See diagrams of Obel	ix, its environment, its sonars, its se	ensor quantisation	
<u></u>						

- Asterix and Obelix robots
- Q-learning
- \bullet Statistical Clustering

Architecture for Obelix: Subsumption Architecture

• Behaviour-based modules: unwedger, pusher, finder



- Each module gets own reward
- \bullet 3 copies of Q-learning algorithm

	School of a
-	
Б	

- Priorities: Unwedger > Pusher > Finder (= default behaviour)
- Modules have applicability predicates
- So combine BB modules with RL, rather than using a monolithic controller

```
Gillian Hayes
                                  RL Lecture 13
                                                                 19th February 2007
                                                                                              Gillian Hayes
                                                                                                                                                              19th February 2007
                                                                                                                                RL Lecture 13
                                                               f informatics
                                                                                                                                                             -
                                                                                              • Default behaviour = box_finder (so no applicability predicate needed)
                     Applicability Predicates
                                                                                              • Perceptual aliasing - same state needs different actions depending on context,
                                                                                              e.g. BUMP
Under what sensor conditions is a module applicable?
e.g.
                                                                                              See diagram of control flow
   PUSHER_APPLICABLE
```

return TRUE else if BUMP recently activated (perceptual return TRUE else return FALSE

Current state + recent activation allows perceptual aliasing to be dealt with (a bump = false state often occurs after pushing the box; box bounces away from bump sensor)

 $\bullet~\mbox{If}>1$ applicable, priority ordering decides which gets a chance to learn



¹² informatics

Task

So, given

- states
- actions
- behaviour-based task modules and priority ordering
- applicability and reward functions

learn

 \bullet Transfer function from state \Rightarrow action which maximises cumulative expected reinforcement

- 1. Initialise Q(s,a) for all x,a to 0, $\alpha = 0.5$, $\gamma = 0.9$
- 2. Do forever:
 - Observe current state s_t
 - 90% of time choose an action that has maximum $Q(s_t,a_t)$, 10% choose some other action. Action is a_t .
 - Carry out action a_t in world. New state is s_{t+1}
 - Immediate reward for carrying out a in \boldsymbol{s}_t is \boldsymbol{r}_{t+1}
 - Update $Q: \Delta Q = \alpha [r_{t+1} + \gamma max_a Q(s_{t+1}, a) Q(s_t, a_t)]$ for state s_t and all states within weighted Hamming distance of 2 from s_t

Gillian Hayes	RL Lecture 13	19th February 2007

14 informatics

Alternative Generalisation Method – Statistical Clustering

Weighted Hamming distance OK for noisefree simulations, but use statistical clustering for noisy real world environment

For each action

A set of clusters of states with similar $\mathsf{Q}\mathsf{s}$

Cluster: $< p_{bit1}, p_{bit2}, \dots p_{bitn}, Q_c >$

 p_{bit1} is $p(bit1 \mid s \in c)$, i.e. probability that bit 1=1 given that the state s is in the cluster c – you can get this from counting over all the states that are already in the cluster and from sensor statistics. Q_c is the Q value of the cluster

To be in the cluster, a state's bits must match those of the cluster closely and the state must have a Q value very close to that of the cluster and differing from it by no more than δ

Gillian Hayes

RL Lecture 13

19th February 2007

15 informatics

Updating Clusters and Action Selection

New cluster: if difference between current state and existing clusters is too big, start a new cluster

Merge two clusters: if close enough merge and update new cluster statistics

Action selection: match state against all clusters, all actions. Calculate its Q-value for each action in turn:

$$Q(s,a) = \frac{\sum_{\text{clusters for that action } Q_{\text{cluster}} \times \text{weight}}{\sum_{\text{weight}}}$$

where weight = match probability for the state being in that cluster. Pick action with the biggest $Q(\boldsymbol{s},\boldsymbol{a})$

Scales better than Hamming

16 informatics

Experiments

Hamming compared with:

- random agent: chooses actions randomly
- \bullet hand-coded: give it the reward for each action that the learning agent would have got
- statistical clustering

on

- simulator
- robots

Results

- learning curves
- learning agents often as good as hand-coded

Gillian Hayes

RL Lecture 13

19th February 2007

18 informatics

RL for Real World Problems

Slow

Noisy sensors – can states be reidentified?

Generalisation – how to do it? 40% of processor time spent matching states to clusters in Asterix

Actions - reliability and repeatability

Monolithic vs. modular architecture

Specifying reward function difficult

But may be easier than writing controller

Most of effort goes into setting up the problem so that the RL algorithm can be applied.

BUT: if we get the right level of abstraction it can be done



John Hoar: Reinforcement Learning Applied to a Real Robot Task. DAI MSc dissertation, University of Edinburgh (1996).

S.Mahadevan and J.Connell: Automatic Programming of Behavior-Based Robots Using Reinforcement Learning. Artificial Intelligence 55, 311–365 (1992).

Jeremy Wyatt, John Hoar, Gillian Hayes: Design, Analysis and Comparison of Robot Learners. Robotics and Autonomous Systems 24(1-2), 17-32 (1998).

Gillian Hayes

Issues

Asterix: a good internal learning curve doesn't always correspond with good behaviour in world

push walls, not boxes

Gillian Hayes

- find boxes, then move away
- zoom around world as quickly as possible

Difficulty: designing a good reward function. Learner will ALWAYS take advantage....

So always evaluate on real task, even if learning curves increase properly

RL Lecture 13

19th February 2007

10 Informatics

nformatics