# PLANNING AS LEARNING

### Overview

- The last few lectures have considered techniques for planning.
- We have considered that an agent knows all bout its environment and plans by thinking hard about what it wants to do.
- Instead we can think of planning as a process of exploring the environment around the agent.
- We'll look at a number of approaches based on this idea.
- These are all types of *reinforcement learning*.
- They are also ideas that link closely back to the things we talked about when we talked about search.

cis32-spring2009-parsons-lect17

# • We'll start by considering that we have a state-space in which the agent is carrying out actions.

- We want to come up with a plan.
- Well, in fact what we end up with is a *policy*.
- That is matrix that tells us which action to carry out in which state.
- This is a *conditional plan*, which is much more robust than a linear plan as produced by STRIPS.
- One way to get this is by applying a heuristic the heuristic value for each state helps to tell us which action we should pick.

## Learning heuristics

- We will start by assuming that the agent knows the results and costs of each operation.
- We will also assume that it can build the whole search tree.
- This is just what we did for previous searches.
- We then set h(n) = 0 for all *n* and run an A<sup>\*</sup> search.
- When the agent has expanded node  $n_i$  to give a set of children  $\delta(n_i)$ , it updates its  $h(n_i)$  to be:

$$h(n_i) := \min_{n_j \in \delta(n_i)} \left[ h(n_j) + c(n_i, n_j) \right]$$

where  $c(n_i, n_j)$  is the cost of moving from  $n_i$  to  $n_j$ .

• We further assume that the agent can recognise the goal state and knows that *h*(*goal*) is 0.

cis32-spring2009-parsons-lect17

cis32-spring2009-parsons-lect17

- This won't do much for the agent the first time–it is just uniform cost search.
- However, subsequent searches will "zoom in" on the right solution faster and faster.
- This happens as the  $h_T(n)$  values propagate back from the goal.
- (There are few enough values that these can be stored in a table.)
- Each run propagates the true cost of getting to the goal further back through the search.
- Eventually, the minimal cost path can just be read off the tree.

#### cis32-spring2009-parsons-lect17



## Learning without a model of action

- As described this kind of search is a "thought experiment" that an agent carries out.
- In the case of the navigating robot, it is planning its route across the grid.
- Alternatively it would be possible for the agent to actually carry out the operations to see what happens.
- In the case of the robot it could move through the room randomly at first, working out over a number of runs what the outcomes of actions were, and which were most useful at which point.
- (To do this, the agent will have to build a model of the state space in its "head").

cis32-spring2009-parsons-lect17

- What we assume is that:
  - The agent can distinguish the states it visits (and name them).
  - The agent knows how much actions cost once it has taken them.
- The process starts at the start state *s*<sub>0</sub>.
- The agent then takes an action (maybe at random), and moves to another state. And repeats.
- As it visits each state, it names it and updates the heuristic value of this state as:

$$h(n_i) := [h(n_j) + c(n_i, n_j)]$$

- where  $n_i$  is the node in which an action is taken,  $n_j$  is the node the action takes the agent to, and  $c(n_i, n_j)$  is the cost of the action.
- $h(n_j)$  is zero if the node hasn't been reached before.

cis32-spring2009-parsons-lect17

• Whenever the agent has to choose an action *a*, it chooses it by:

$$a = \operatorname{argmin}_{a} \left[ h(\sigma(n_i, a)) + c(n_i, \sigma(n_i, a)) \right]$$

where  $\sigma(n_i, a)$  is the state reached from  $n_i$  after carrying out a.

- As before, the estimated minimum cost path to the goal is built up over repeated runs.
- However, allowing some randomness in the choice of actions increases the chance that the "estimated minimum cost path" really is the best path.

```
cis32-spring2009-parsons-lect17
```

- Potentially you could consider all the things it is possible to measure.
- Then:

$$h(n) = w_1 W(n) + w_2 P(n) + .$$

- We then learn which weights are best.
- One way to do this is as follows:
- After expanding  $n_i$  to  $\delta(n_i)$  we adjust the weights so that:

$$h(n_i) := h(n_i) + \beta \left( \min_{n_j \in \delta(n_i)} \left[ h(n_j) + c(n_i, n_j) \right] - h(n_i) \right)$$

• We modify *h*(*n<sub>i</sub>*) by adding some proportion of (controlled by β) of the difference between what we thought *h*(*n<sub>i</sub>*) was before the expansion, and what we think it is after.

# Learning without a search graph

- For many interesting problems, it is not possible to store all the states/nodes and build the entire search graph.
- Provided we have a model of the effects of actions, we can still search with an evaluation function.
- We start by assembling a heuristic as a linear combination of some set of plausible functions.
- For the 8-puzzle these might be:
  - W(n): number of tiles out of place.
  - P(n): sum of distance each tile is from home.
- Plus any additional functions you can think of.

```
cis32-spring2009-parsons-lect17
```

• We can rewrite this as:

$$h(n_i) := (1 - \beta)h(n_i) + \beta \min_{n_i \in \delta(n_i)} [h(n_j) + c(n_i, n_j)]$$

- $\beta$  controls how fast the agent learns—how much weight we give to the new estimate of the heuristic.
- If  $\beta = 0$  there is no adjustment.
- If  $\beta = 1$ ,  $h(n_i)$  is thrown away immediately.
- Low values of  $\beta$  lead to slow learning, and high values mean that performance is erratic.
- Note that this *temporal difference approach* can also work without a model of the effects of actions
  - Need modifications which we won't deal with here.

cis32-spring2009-parsons-lect17

11

cis32-spring2009-parsons-lect17

12

10

## Rewards not goals

- For many tasks agents don't have short term goals, but instead accrue *rewards* over a period of time.
- For such a case we definitely want a *policy*  $\pi$  which says what action should be carried out in a given state.
- We express the reward an agent gets as  $r(n_i, a)$ , and if doing *a* in  $n_i$  takes the agent to  $n_j$ , then:

 $r(n_i, a) = -c(n_i, n_j) + \rho(n_j)$ 

where  $\rho(n_j)$  is a reward for being in state  $n_j$ .

- We want an optimal policy  $\pi^*$  which maximises the reward at every node.
- (Rewards are typically discounted over time.)

cis32-spring2009-parsons-lect17

- Given a policy  $\pi$ , we can compute the value of each node—the reward the agent will get if it starts at that node and follows the policy.
- If the agent is at  $n_i$  and follows  $\pi$  to  $n_j$  then the agent will get reward:

$$V^{\pi}(n_i) = r(n_i, \pi(n_i)) + \gamma V^{\pi}(n_j)$$

where  $\gamma$  is the discount factor.

- Think of this as the opposite of bank interest.
- The optimum policy then gives us the action that maximises this reward:

$$V^{\pi^*}(n_i) = \max_a \left[ r(n_i, a) + \gamma V^{\pi^*}(n_j) \right]$$

- One way to find the optimum policy is by searching through all possible policies.
- Start with a random policy and calculate its reward.
- Then guess another policy and see if it has a better reward (kind of slow).
- Better would be to tweak the policy by swapping some *a* in *n<sub>i</sub>* for an *a'* with a higher *r*(*n<sub>i</sub>*, *a'*).
- A good heuristic for tweaking actions may help us find the best policy.

14

16

• But there are better approaches.

cis32-spring2009-parsons-lect17

13

15

• If we knew what the values of the nodes were under *π*<sup>\*</sup>, then we could easily compute the optimal policy:

$$\pi^*(n_i) = \operatorname{argmax}_a \left[ r(n_i, a) + \gamma V^{\pi^*}(n_j) \right]$$

- The problem is that we don't know these values.
- But we can find them out using *value iteration*.
- We start by guessing (randomly is fine) an estimated value V(n) for each node.

```
cis32-spring2009-parsons-lect17
```

cis32-spring2009-parsons-lect17



 $\operatorname{argmax}_{a} \left[ r(n_i, a) + \gamma V(n_j) \right]$ 

that is the best thing given what we currently know.

• We then update  $V(n_i)$  by:

$$V(n_i) := (1 - \beta)V(n_i) + \beta \left[r(n_i, a), \gamma V(n_j)\right]$$

- Progressive iterations of this calculation make V(n) a closer and closer approximation to  $V^{\pi^*}(n_i)$ .
- Intuitively this is because we replace the estimate with the actual reward we get for the next state (and the next state and the next state).

```
cis32-spring2009-parsons-lect17
```

## Summary

- This lecture has looked at a number of approaches to learning plans.
- We started by thinking about an agent exploring a search space and trying to learn a heuristic.
  - Since a good heuristic tells us which action to take, it is kind of like a plan generator.
- We started with a simple learning problem and progressively considered more complex scenarios.
- This created a battery of reinforcement learning methods that can be applied in a wide variety of situations.

18

cis32-spring2009-parsons-lect17

17