





 Illustration

 Imagine agent wandering around in environment.

 • How does it learn utility values of each state?

 • (i.e., what are good / bad states? avoid bad ones...)

Reinforcement learning attempts to answer this question.









Naive updating

- Comes from adaptive control theory.
- (a) ``Sampling'' --- agent makes random runs through environment; collect statistics on final payoff for each state (e.g. when at (2,3), how often do you reach +1 vs. -1?)
- Learning algorithm keeps a running average for each state. Provably converges to true expected values (utilities).

Issues

• Main drawback: slow convergence.

See next figure.

- In relatively small world takes agent
- over 1000 sequences to get a reasonably small (< 0.1) root-mean-square error compared with true expected values.
- Question: Is sampling necessary?
- Can we do something completely different? Note: agent knows the environment (i.e. probability of state transitions) and final rewards.

CS-424 Gregory Dudek

Temporal difference learning

- Combine ``sampling" with ``calculation'' or stated differently: it's using a sampling approach to solve the set of equations.
- Consider the transitions, observed by a wandering agent.
- Use an observed transition to adjust the utilities of the observed states to bring them closer to the constraint equations.

Temporal difference learning

- U(i) : utility in state i
- R(i): Reward
- When observing a transition from i to j bring U(i) value closer to that of U(j)
- Use update rule:
- U(i) = U(i) + k(R(i) + U(j) U(i)) **
- k is the learning rate parameter.
- Rule is called the temporal-difference or TD equation. (note we take the difference between successive states.)

Mobile Robotics Research

- Mobile Robotics exemplifies these issues.
- The canonical problems:

Where am I

(position estimation)

How do I get there

(path planning)

Mapping and exploration

How do I explain what I saw, and describe what I need to remember

 $\begin{aligned} & \textbf{Kohonen learning} \\ & \textbf{We are given a set of nodes } \{n_i\} \text{ and a set of examples } \\ & \{s_i\}. \end{aligned}$ $\begin{aligned} & \textbf{For each training examples } s_i \\ & \textbf{Compute distance } d \text{ between } s_i \text{ and each } \\ & \textbf{node } n_j \text{ described by weights } w_{j,i}. \\ & \textbf{d} = \texttt{sqrt} ((s_i w_{j,i}^T)^2) \end{aligned}$ $\begin{aligned} & \textbf{Find node with min distance} \\ & \textbf{For each node } \underline{close} \text{ to n} \\ & where close means } D(n,n') < \mathbb{K} \\ & \textbf{update weights:} \\ & w_{j,i} = w_{j,i} + c(s_i - w_{j,i}) \\ & \textbf{Make } c \text{ and } K \text{ smaller.} \end{aligned}$

Kohonen: issues? Exploits both distribution of examples <u>and their frequency</u> to establish the mapping. Accomplishes clustering and recognition. May serve as a model of biological <u>associative memory</u>. Can project a high-dimensional space to a lower dimensional comparison space: <u>dimensionality reduction</u>. We had to choose the features in advance. A priori size and structure of the network.

• What guarantees do we have?

Planning & search

• Many planning problems can be viewed as search:

- 1.
 - States are world configurations.
 - Operations are actions that can be taken
 - Find a set of intermediate states from current to goal world configuration.
- 2.
 - States are plans
 - Operations are changes to the sequence of possible actions.
 - Find a path from the null plan (do nothing) to a plan that gets to the right goal.

CS-424 Gregory Dudek

Basic formalism

• Basic logical formalism derived from STRIPS.

- State variables determine what actions can or should be taken: in this context they are <u>conditions</u>
 - Shoe_untied()
 - Door_open(MC)
- An operator (remember those?) is now a <u>triple</u> Preconditions Additions Deletions