















<text><text><list-item><section-header><list-item><list-item>

Value Iterationfunction VALUE-ITERATION(mdp, ϵ) returns a utility functioninputs: mdp, an MDP with states S, transition model T, reward function R, discount γ ϵ , the maximum error allowed in the utility of any statelocal variables: U, U', vectors of utilities for states in S, initially zero δ , the maximum change in the utility of any state in an iterationrepeat $U \leftarrow U'; \delta \leftarrow 0$ for each state s in S do $U'[s] \leftarrow R[s] + \gamma \max_{a} \sum_{s'} T(s, a, s') U[s']$ if $|U'[s] - U[s]| > \delta$ then $\delta \leftarrow |U'[s] - U[s]|$ util $\delta < \epsilon(1 - \gamma)/\gamma$ return UFigure 17.4 The value iteration algorithm for calculating utilities of states. The termination condition is from Equation (17.8).







Policy Iteration

<u>Idea</u>: iteration works even with inaccurate utilities, when one action is clearly better than all others. Thus we can alternate two steps:

I. policy evaluation: calculate U(s) of each state under a policy Π

2. policy improvement: set new Π ' using one-step look-ahead from U(s)

<u>Algorithm</u>:

- Pick a policy Π at random
- Repeat:
 - I. Compute the utility of each state for Π $U_{t+1}(s) \leftarrow \mathbf{R}(s) + \gamma \sum_{s} T(s, \Pi(s), s') U_{t}(s')$

Simpler, linear equations because no ,,max" step. Can be solved in $O(n^3)$ for *n* states

- 2. Compute a new policy Π ' given these utilities Π '(s) = arg max_a $\sum_{s} T(s,a,s') U(s')$

- If $\Pi' = \Pi$ then return Π

Fixpoint of Bellman eq., optimal policy



Policy Iteration

Remarks:

- When no improvement, U is fixpoint and the policy it optimal
- There are finitely many policies and each step yields a better policy, so policy iteration must terminate
- Policy evaluation (step 1.) uses an iterative update to estimate the current U using a simplified Bellman update
- Possible to update U or policy only for a subset of states (instead of all), using either kind of updating (policy improvement, simplified value iteration): asynchronous policy iteration
 - could update only those states that are likely reached
 - more efficient, can be guaranteed to converge

Matlab example D. Barber (2012). Bayesian Reasoning and Machine Learning, Cambridge Univ. Press. \rightarrow BRML toolkit 17 Partially Observable Markov Decision Problems (POMDP) MDPs: fully observable environments and Markovian transition models "POMDP" account in addition for partially observable environments: Which state is the agent in? Utility of s? Optimal action? A POMDP is given by I. Initial state s₀ 2. Transition model $T(s,a,s') = \mathbf{P}(s'|a,s)$

- 3. Reward function R(s) or R(s,a,s'), additive
- 4. Observation (measurement) model: O(s,o) = P(o|s) sensing operation in state s, returns multiple observations o, with a probability distribution

POMDPs

<u>Idea in a nutshell:</u>

Belief state b(s) = prob. distribution over all possible states

Example: in 4x3-world = point in 11-dim continuous space

Define the agent's policy over its belief state: $\Pi^*(b)$ That is, actions depend on *beliefs*, not the state the agent is in!

POMDP decision cycle:

- I. Given current belief state b, execute $a=\Pi^*(b)$
- 2. Get new observations (measurements) o
- 3. Bayesian update of belief states:

 $b'(s') = \alpha O(s',o) \Sigma_s T(s,a,s')b(s)$

How to determine Π *(b) ?

19









Example

<u>Parameters</u>: The actions u1 and u2 are terminal actions. The action u3 is a sensing action that potentially leads to a state transition. The horizon is finite and γ =1.

$r(x_1, u_1)$	=	-100	$r(x_2, u_1)$	=	+100
$r(x_1, u_2)$	=	+100	$r(x_2, u_2)$	=	-50
$r(x_{1}, u_{3})$	=	-1	$r(x_2, u_3)$	=	-1
$p(x_1' x_1, u_3)$	=	0.2	$p(x_2' x_1,u_3)$	=	0.8
$p(x_1' x_2, u_3)$	=	0.8	$p(z_{2}' x_{2},u_{3})$	=	0.2
$p(z_1 x_1)$	=	0.7	$p(z_2 x_1)$	=	0.3
$p(z_1 x_2)$	=	0.3	$p(z_2 x_2)$	=	0.7

Payoffs in POMDPs

In MDPs, the payoff (or reward) depends on the state of the system. But in POMDPs, the true system state is not exactly known. Thus, we compute the expected payoff from belief state, by integrating over all

states:
$$r(b, u) = E_x[r(x, u)]$$

= $\int r(x, u)p(x) dx$
= $p_1 r(x_1, u) + p_2 r(x_2, u)$

Example:

 $r(b, u_1) = -100 p_1 + 100 p_2$ If certain we are in xI and execute $= -100 p_1 + 100 (1 - p_1)$ ul, we receive a reward of -100 If certain we are in x2 and execute $r(b, u_2) = 100 p_1 - 50 (1 - p_1)$ ul, the reward is +100. In between. linear combination weighted by $r(b, u_3) = -1$ probabilities 25



Resulting policy for T=1

Given we have a finite POMDP with T=1, we would use $V_1(b)$ to determine the optimal policy.

In our example, the optimal policy for T=1 is

$$\pi_1(b) = \begin{cases} u_1 & \text{if } p_1 \leq \frac{3}{7} \\ u_2 & \text{if } p_1 > \frac{3}{7} \end{cases}$$

This is the upper thick graph in the diagram.

27

Piecewise Linearity, Convexity

The resulting value function $\mathsf{V}_{\mathsf{I}}(\mathsf{b})$ is the maximum of the three functions at each point

$$V_{1}(b) = \max_{u} r(b, u)$$

=
$$\max \begin{cases} -100 p_{1} + 100 (1 - p_{1}) \\ 100 p_{1} - 50 (1 - p_{1}) \\ -1 \end{cases}$$

It is piecewise linear and convex.

Pruning

If we carefully consider $V_1(b)$, we see that only the first two components contribute.

The third component can therefore safely be pruned away from $V_1(b)$.

$$V_{1}(b) = \max \left\{ \begin{array}{cc} -100 \ p_{1} \ +100 \ (1-p_{1}) \\ 100 \ p_{1} \ -50 \ (1-p_{1}) \end{array} \right\}$$

29

Increasing the Time Horizon

If we go over to a time horizon of T=2, the agent can also consider the sensing action u3.

Suppose we perceive z1 for which p(z1 | x1)=0.7 and p(z1 | x2)=0.3. Given the observation z1 we update the belief using Bayes rule. Thus $V_1(b | z1)$ is given by

$$V_{1}(b \mid z_{1}) = \max \begin{cases} -100 \cdot \frac{0.7 p_{1}}{p(z_{1})} + 100 \cdot \frac{0.3 (1-p_{1})}{p(z_{1})} \\ 100 \cdot \frac{0.7 p_{1}}{p(z_{1})} & -50 \cdot \frac{0.3 (1-p_{1})}{p(z_{1})} \end{cases} \\ = \frac{1}{p(z_{1})} \max \begin{cases} -70 p_{1} + 30 (1-p_{1}) \\ 70 p_{1} & -15 (1-p_{1}) \end{cases} \end{cases}$$

Expected Value after Measuring

Since we do not know in advance what the measurement will be (z_1, z_2) , we have to compute the expected belief

$$\bar{V}_{1}(b) = E_{z}[V_{1}(b \mid z)] \\
= \sum_{i=1}^{2} p(z_{i}) V_{1}(b \mid z_{i}) \\
= \max \left\{ \begin{array}{c} -70 \ p_{1} \ +30 \ (1-p_{1}) \\ 70 \ p_{1} \ -15 \ (1-p_{1}) \end{array} \right\} \\
+ \max \left\{ \begin{array}{c} -30 \ p_{1} \ +70 \ (1-p_{1}) \\ 30 \ p_{1} \ -35 \ (1-p_{1}) \end{array} \right\}$$

31

Resulting Value Function

The four possible combinations yield the following function which again can be simplified and pruned.

$$\bar{V}_{1}(b) = \max \begin{cases} -70 \ p_{1} \ +30 \ (1-p_{1}) \ -30 \ p_{1} \ +70 \ (1-p_{1}) \\ -70 \ p_{1} \ +30 \ (1-p_{1}) \ +30 \ p_{1} \ -35 \ (1-p_{1}) \\ +70 \ p_{1} \ -15 \ (1-p_{1}) \ -30 \ p_{1} \ +70 \ (1-p_{1}) \\ +70 \ p_{1} \ -15 \ (1-p_{1}) \ +30 \ p_{1} \ -35 \ (1-p_{1}) \end{cases} \\ = \max \begin{cases} -100 \ p_{1} \ +100 \ (1-p_{1}) \\ +40 \ p_{1} \ +55 \ (1-p_{1}) \\ +100 \ p_{1} \ -50 \ (1-p_{1}) \end{cases} \end{cases}$$

State Transitions (Prediction)

When the agent selects u3 its state potentially changes.

When computing the value function, we have to take these potential state changes into account.

$$p'_{1} = E_{x}[p(x_{1} | x, u_{3})]$$

= $\sum_{i=1}^{2} p(x_{1} | x_{i}, u_{3})p_{i}$
= $0.2p_{1} + 0.8(1 - p_{1})$
= $0.8 - 0.6p_{1}$

33

Resulting Value Function

What's the resulting Value Function after executing u3 ?

• Taking also the state transitions into account, we finally obtain.

$$\bar{V}_1(b \mid u_3) = \max \left\{ \begin{array}{cc} 60 \ p_1 & -60 \ (1-p_1) \\ 52 \ p_1 & +43 \ (1-p_1) \\ -20 \ p_1 & +70 \ (1-p_1) \end{array} \right\}$$

Value Function for T=2

 Taking into account that the agent can either directly perform ul or u2, or first u3 and then ul or u2, we obtain (after pruning)

$$\bar{V}_2(b) = \max \left\{ \begin{array}{rrr} -100 \ p_1 & +100 \ (1-p_1) \\ 100 \ p_1 & -50 \ (1-p_1) \\ 51 \ p_1 & +42 \ (1-p_1) \end{array} \right\}$$





Why Pruning is Essential?

Each update introduces additional linear components to V. Each measurement squares the number of linear components. Thus, an unpruned value function for T=20 includes more than 10547,864 linear functions. At T=30 we have 10561,012,337 linear functions. The pruned value functions at T=20, in comparison, contains only 12 linear components. The combinatorial explosion of linear components in the value function are the major reason why POMDPs are impractical for most

applications and are still an active research area. So far only applied successfully to very small state spaces with small numbers of possible observations and actions. In other words: <u>Finding efficient state</u> <u>representations is a key to solve (PO)MDPs!</u>

37