

uninformed search methods  
 DFS  
 BFS  
 VCF

Filtering: improve backtracking by forward checking

filtering: keep track of domains for unassigned variables and cross off bad options  
 Forward checking: cross off values that violate a constraint when added to the existing assignment

informed search  $A^*$

Admissible heuristics

$h$  is admissible if:  
 $0 \leq h(n) \leq h^*(n)$

$h^*(n)$  is the cost of node  $n$  optimal in tree search

Dominance:  $h_a \geq h_c$  if  
 $\forall n: h_a(n) \geq h_c(n)$

consistent heuristic

$\forall A, C: h(A) - h(C) \leq \text{cost}(A \text{ to } C)$

$A^*$  graph search is optimal

CSPs

Variables Domain Goal test is a set of Independent subproblems are connected components of Constraint graph  
 Constraints  
 Variables:  $WA, NT, Q, \dots$   
 Domain:  $D = \{\text{Red, green, blue}\}$   
 Implicit:  $WA \neq NT$   
 Explicit:  $(WA, NT) \in \{(Red, green), (Red, blue), \dots\}$   
 Solution:

Backtracking = DFS + variable ordering + fail-on-violation  
 0. for each value in order domain values  
 1. Only consider assignment to a single variable at each

Step 2. check assignment, if good, add to assignment, if no, remove value from assignment.

Consistency of a single arc, An arc  $X \rightarrow Y$  is consistent iff every  $X$ , there's some  $Y$  which could be assigned without violate a constraint  
 Delete from tail, neighbors of  $X$  need to be rechecked

Ordering = MRV: minimum Remaining values

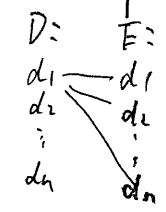
choose values with the fewest legal left values in its domain (variable) (the most constrained variable)

LCV least constraining value heuristic

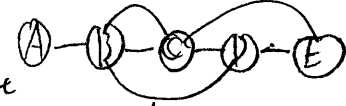
$B = (1,3), (2,3) A = (2,4), (2,5), (1,4) \dots E = (3,4), (4,5)$

We choose  $B$  by MRV, and for each  $(1,3), (2,3)$ , enforce arc consistency, if  $(1,3)$  leaves  $A = (2,4) E = \emptyset$  We chose  $(2,3)$   
 $(2,3)$  leaves  $A = (2,4), (2,5) E = (3,4)$  because  $3 > 2$   
 $D = (1,2)$

Tree-structured CSPs  $O(nd^2)$

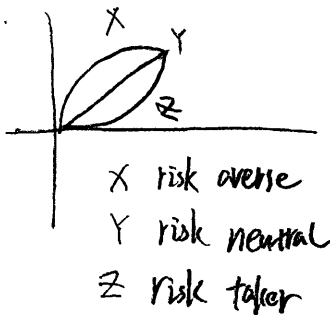


(Remove subtree see a tree)



Linear in  $n$ , and check both domains  $d^2$

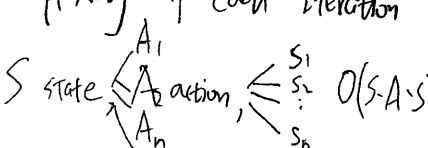
Cutset, local search (iterative improvements) min conflict heuristic

Lottery preference  $A \succ B$  policy Evaluation fix  $\pi$  policy extraction  
 $L = [p, A; (1-p), B]$   $V_0^\pi(s) = 0$   $\pi^*(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$   
 Utility of this lottery:  $V^{\pi_{k+1}}(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$   $\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$   
 $U(L) = pU(A) + (1-p)U(B)$   $OCS^2$  One step look-ahead  
  
 X risk averse  
 Y risk neutral  
 Z risk taker  
 Exploration vs Exploitation  
 Unknown state to explore unknown state vs Using known knowledge to act  
 Direct Evaluation input policy  

|   |   |
|---|---|
| A | D |
| B | C |

 observed Episodes  
 Episode 1: A, right, D, 1  
 Episode 2: ...  
 Waste information  
 Output value  

|   |    |
|---|----|
| 4 | 10 |
| 2 | 8  |

  
 Model-Based learning  
 1. Learn empirical MDP model  $\uparrow$  and  $\hat{R}$   
 2. solve learned MDP  
 Temporal difference learning (model free)  
 Fixed policy sample =  $R(s, \pi(s), s') + \gamma V^\pi(s')$   
 Update  $V(s)$   $V^\pi(s) = (1-\alpha)V^\pi(s) + \alpha \text{sample}$   
 Can turn values into a new policy  
 Q-learning  
 sample =  $R(s, a, s') + \gamma \max_{a'} Q(s, a', s')$   
 $Q(s, a) = (1-\alpha)Q(s, a) + \alpha [\text{sample}]$   
 Approximate Q-learning  
 $Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots$   
 adjust weights  
 Transition:  $(s, a, r, s')$   
 difference =  $[r + \gamma \max_{a'} Q(s', a')] - Q(s, a)$   
 $Q(s, a) = Q(s, a) + \alpha [\text{difference}]$   
 $w_i = w_i + \alpha [\text{difference}] f_i(s, a)$   
 Model-based learning: estimate Transition and Reward functions  
 Model-free learning: 1. off-policy learning: q-learning  
 2. on-policy learning: TD learning and direct evaluation  
 Minimax efficiency  
 Like DFS Time  $O(b^m)$   
 Space  $O(bm)$   
 $\alpha$  Max's best option on path to root  
 $\beta$  Min's best option on path to root  
 active Reinforcement learning  
 goal: learning the optimal policy/values  
 max-value if  $V \geq \beta$  return  $\beta$   
 min-value if  $V \leq \alpha$  return  $\alpha$   
 Random action (E-greedy)  
 with small prob  $\epsilon$ , act random  
 with large  $1-\epsilon$ , act on current policy  
 Better idea: explore area whose goodness is not established, eventually stop exploring  
 Visit count  $n$ ,  $f(a, n) = w + \frac{k}{n}$   
 $Q(s, a) = R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$   
 difference =  $[r + \gamma \max_{a'} Q(s', a')] - Q(s, a)$   
 $Q(s, a) = Q(s, a) + \alpha [\text{difference}]$   
 $w_i = w_i + \alpha [\text{difference}] f_i(s, a)$   
 Model-based learning: estimate Transition and Reward functions  
 Model-free learning: 1. off-policy learning: q-learning  
 2. on-policy learning: TD learning and direct evaluation  
 complexity of each iteration  $O(S^2A)$   
  
 $S$  state  $\swarrow$  action  $\searrow$   $O(S \cdot A \cdot S)$