

## Search Problem

State space [location (dots, booleans)] / Successor function / Start & Goal state

World state: Agent pos  $\times 2^{\text{Food counts}} \times \text{Ghost pos} \times \text{# of levels} \times \text{actions}$

## Tree Search

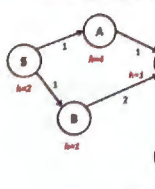
	Strategy	Fringe	is Complete	is Optimal	Time complexity	Space
DFS	a deepest node	LIFO stack	Yes/no cycle	No/leftmost	$O(b^m)$	$O(bm)$
BFS	a shallowest node	FIFO queue	Yes	Yes/same costs	$O(b^b)$	$O(b^m)$
UCS	a cheapest node	priority queue	Yes	Yes	$O(b^{b^2})$	
Greedy	a node that is closest to a goal		badly guide DFS			
A*	Orders by $f(n) = g(n) + h(n)$		Admissible (optimistic) solution $\leq h^*(n)$	Yes		

## \* Heuristic

A function that estimates how close a state is to a goal

Graph search never expand a node twice A\* with consistent heuristic

State space graph



Search tree



Main idea: estimated heuristic costs  $\leq$  actual costs

- Admissibility: heuristic cost  $\leq$  actual cost to goal  
 $h(A) \leq \text{actual cost from A to G}$
- Consistency: heuristic "arc" cost  $\leq$  actual cost for each arc  
 $h(A) - h(C) \leq \text{cost}(A \text{ to } C)$
- Consistency implies admissibility
- Consequences of consistency:
  - The f value along a path never decreases, because:  
 $h(A) \leq \text{cost}(A \text{ to } C) + h(C)$

1 node  
b nodes  
 $b^2$  nodes  
 $b^3$  nodes  
 $b^m$  nodes

NP hard to find smallest cutset min

d: Domain Size

Constraint Satisfaction Problem

Variables/Domains/constraints

Suppose a graph of n variables can be broken into subproblems of only c variables  $\rightarrow O(n^c)(d^c)$

Basic solution: Backtracking ① One variable at a time ② Check constraints as you go

= DFS + variable-ordering + Fail-on-violation

Filtering: Forward Checking: Cross out values that violates a constraint when added to assignment; whenever any variables has no value left, backtrack.

Arc Consistency: An arc  $X \rightarrow Y$  is consistent iff for every x in the tail there is some y in the head could be assigned without violating a constraint.

\* Delete from the tail \* If Y loses a value, arc  $X \rightarrow Y$  needs to be re-checked.

Ordering: Minimum Remaining Value: the variable with the fewest legal left values in domain

Least Constraining Value: the one rules out the fewest values in the remaining vars.

Structure: Tree structure: Choose a root variable, order variable that parents precede child

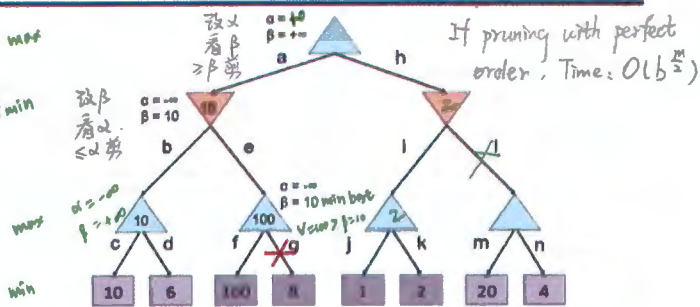
Remove backward:  $n \rightarrow 2$ , remove inconsistent + Assign forward  $1 \rightarrow n$  assign x

cutset structure: a set of variables s.t. the remaining constraint is a tree.

c: cutset size  
 $O((d^c)(n-c)d^2)$

## Alpha-Beta Example 2

Time  $O(b^m)$   
Space  $O(bm)$



## Adversarial search

Resource Limit: ① Depth-limited ② Limiting branching factor

Evaluation Function:  $\text{Eval}(u) = w_1 f_1(u) + w_2 f_2(u) + \dots$

## Minimax

def value(state):

if the state is a terminal state: return the state's utility  
if the next agent is MAX: return max-value(state)  
if the next agent is MIN: return min-value(state)

def max-value(state):

Initialize  $v = -\infty$   
for each successor of state:  
 $v = \max(v, \text{value(successor)})$   
return v

def min-value(state):

Initialize  $v = +\infty$   
for each successor of state:  
 $v = \min(v, \text{value(successor)})$   
return v

$\alpha$ : MAX's best option on path to root

$\beta$ : MIN's best option on path to root

## Alpha-beta

def max-value(state,  $\alpha$ ,  $\beta$ ):

Initialize  $v = -\infty$   
for each successor of state:  
 $v = \max(v, \text{value(successor, } \alpha, \beta))$   
if  $v \geq \beta$  return v  
 $\alpha = \max(\alpha, v)$   
return v

def min-value(state,  $\alpha$ ,  $\beta$ ):

Initialize  $v = +\infty$   
for each successor of state:  
 $v = \min(v, \text{value(successor, } \alpha, \beta))$   
if  $v \leq \alpha$  return v  
 $\beta = \min(\beta, v)$   
return v

## Expectimax

no pruning!

def value(state):

if the state is a terminal state: return the state's utility  
if the next agent is MAX: return max-value(state)  
if the next agent is EXP: return exp-value(state)

def max-value(state):

Initialize  $v = -\infty$   
for each successor of state:  
 $v = \max(v, \text{value(successor)})$   
return v

def exp-value(state):

Initialize  $v = 0$   
for each successor of state:  
 $p = \text{probability(successor)}$   
 $v += p * \text{value(successor)}$   
return v

Variable Elimination (VE) on Polytree BN = Sum-product message passing algorithm or belief propagation algorithm. (from leaf to root)

General: group nodes to have a polytree (junction/join) tree, Intractable (不易处理) with large cliques (large tree width)

Approximate, tractable but not guaranteed to work well.

## Propositional Logic

语法

语义

conjunction:  $\wedge$   
disjunctive:  $\vee$   
Valid/satisfiable/un-  
 $(\alpha \wedge \beta) \equiv (\beta \wedge \alpha)$  commutativity of  $\wedge$   
 $(\alpha \vee \beta) \equiv (\beta \vee \alpha)$  commutativity of  $\vee$   
 $((\alpha \wedge \beta) \wedge \gamma) \equiv (\alpha \wedge (\beta \wedge \gamma))$  associativity of  $\wedge$   
 $((\alpha \vee \beta) \vee \gamma) \equiv (\alpha \vee (\beta \vee \gamma))$  associativity of  $\vee$   
 $\neg(\neg\alpha) \equiv \alpha$  double-negation elimination  
 $(\alpha \Rightarrow \beta) \equiv (\neg\beta \Rightarrow \neg\alpha)$  contraposition  
 $(\alpha \Rightarrow \beta) \equiv (\neg\alpha \vee \beta)$  implication elimination  
 $(\alpha \Leftrightarrow \beta) \equiv ((\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha))$  biconditional elimination  
 $\neg(\alpha \wedge \beta) \equiv (\neg\alpha \vee \neg\beta)$  de Morgan  
 $\neg(\alpha \vee \beta) \equiv (\neg\alpha \wedge \neg\beta)$  de Morgan  
 $(\alpha \wedge (\beta \vee \gamma)) \equiv ((\alpha \wedge \beta) \vee (\alpha \wedge \gamma))$  distributivity of  $\wedge$  over  $\vee$   
 $(\alpha \vee (\beta \wedge \gamma)) \equiv ((\alpha \vee \beta) \wedge (\alpha \vee \gamma))$  distributivity of  $\vee$  over  $\wedge$   
 $B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1})$  Convert to CNF

1. Eliminate  $\Leftrightarrow$ , replacing  $\alpha \Leftrightarrow \beta$  with  $(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)$ .

$(B_{1,1} \Rightarrow (P_{1,2} \vee P_{2,1})) \wedge ((P_{1,2} \vee P_{2,1}) \Rightarrow B_{1,1})$

2. Eliminate  $\Rightarrow$ , replacing  $\alpha \Rightarrow \beta$  with  $\neg\alpha \vee \beta$ .

$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge (\neg(P_{1,2} \vee P_{2,1}) \vee B_{1,1})$

3. Move  $\neg$  inwards using de Morgan's rules and double-negation:

$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge ((\neg P_{1,2} \wedge \neg P_{2,1}) \vee B_{1,1})$

4. Apply distributivity law ( $\wedge$  over  $\vee$ ) and flatten:

$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge (\neg P_{1,2} \vee B_{1,1}) \wedge (\neg P_{2,1} \vee B_{1,1})$

Prove KB  $\models \alpha \rightarrow$  show KB  $\wedge \neg\alpha$  is unsatisfiable  $\rightarrow$  to.cnf (KB  $\wedge \neg\alpha$ )  $\rightarrow$  resolve to empty clause

Horn Logic ①  $P_1 \wedge P_2 \wedge \dots \wedge P_n \Rightarrow Q$ , ②  $P_1 \vee P_2 \vee \dots \vee P_n \vee Q$

$\Rightarrow$  Forward chain: Add new clauses into KB until Q is added  $\Rightarrow$  data-driving

Backward chain: goal-driving \* avoid loop (check if new subgoal is on the goal stack)

# avoid repeated work (has already be proved or fail)

Resolution:

AVX, BVX  
AVB

only sound and complete for Horn Logic

run in linear time  
Complete search algorithm can be used to produce a complete inference algorithm







# I. Probabilistic Temporal Model.

1. **Stationarity** assumption: same transition probabilities at all time steps.

Infinitely many variables.

**Markov** assumption: past and future independent given present. **first-order**. k-th order: dependent k earlier steps

2. Forward Algorithm:  $P(X_t) = \sum_{x_{t-1}} P(X_{t-1} = x_{t-1}) P(X_t | X_{t-1} = x_{t-1})$   
For most chains,  $P_{00} \perp P_0$ .

Stationary distribution:  $P_{00}(X) = \sum_{\pi} P(X|\pi) P_{00}(\pi)$ .

Application: Web Link Analysis, PageRank, Gibbs Sampling.

3. HMM. Initial distribution  $P(X_0)$ ; Transition Model  $P(X_t | X_{t-1})$ ; Emission model  $P(E_t | X_t)$ .

evidence is independent of everything else given current state.

4. **Fitting**:  $P(X_t | e_{1:t})$ . (belief state) recursive

$$P(X_{t+1} | e_{1:t+1}) = \alpha P(e_{t+1} | X_{t+1}) \sum_{x_t} P(x_t | e_{1:t}) P(X_{t+1} | x_t)$$

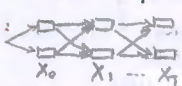
normalize update                      Predict

$$\alpha = \frac{1}{P(e_{t+1} | e_{1:t})}$$

$J_{1:t+1}$  = Forward ( $f_{1:t}, e_{t+1}$ ), start with  $f_{1:0} = P(X_0)$

Cost  $O(|X|^2)$ ,  $|X|$ : # of states.

**State trellis**:



each arc:  $P(x_t | x_{t-1}) P(e_t | x_t)$   
init arc:  $P(X_0)$   
product of weight on a path:  
state sequence probability

Forward algorithm:

sum over all possible paths:  $P(X_{t+1} | e_{1:t+1}) = \sum_{x_t} P(X_{t+1} | e_{1:t+1})$

DP:  $J_{1:t+1} = \alpha P(e_{t+1} | X_{t+1}) \sum_{x_t} P(X_{t+1} | x_t) J_{1:t}(x_t)$

keep track of the total probability of all path to it.

5. **Most likely explanation**:  $\arg\max_{x_{1:t}} P(x_{1:t} | e_{1:t})$

**Viterbi Algorithm**:

keep track of the maximum probability of any path to it.

$$m_{1:t+1} = \text{Viterbi}(m_{1:t}, e_{t+1})$$

$$= P(e_{t+1} | X_{t+1}) \max_{x_t} P(X_{t+1} | x_t) m_{1:t}(x_t)$$

Time:  $O(|X|^2 T)$ , Space:  $O(|X| T)$ .

6. Dynamic Bayes Net (DBN).

Every HMM is a DBN, Every discrete DBN can be represented HMM.  
exponentially fewer parameters.

Exact Infer: unroll as we go, eliminate all variables from prev time step to find  $P(X_T | e_{1:T})$ .

7. Particle Filtering.

N particles. Generally,  $N \ll |X|$ . Many  $x$ :  $P(x) = 0$

propagate forward:  $x_{t+1} \sim P(x_{t+1} | x_t)$

observe (weight):  $w = P(e_t | x_t)$

resample.

## II. Markov Decision Process. (Offline)

1. non-deterministic search problems.

discounting helps to converge.

2. The bellman equation:

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

3. Time-Limited values.

$V_k(s)$ : optimal value of s if game ends in k more steps.

- (DP) Value Iteration:  $V_0(s) = 0$

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

until convergence. each iteration  $O(S^2 A)$ . Total:  $O(S^2 A H)$

will converge to unique optimal values. Compared to Expectimax.

Policy extraction:

$$\pi^*(s) = \arg\max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$= \arg\max_a Q^*(s, a)$$

4. Q-value Iteration:  $Q_0(s, a) = 0$

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

3. Policy Iteration (DP)

Optimal. converge much faster under some conditions.

Policy evaluation

a. Iterative updates.  $V_0^\pi(s) = 0$

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

$O(S^2)$  per iteration.

b. Solve Linear system.  $O(S^3)$  per step

$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

Policy improvement.

$$\pi_{k+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_k}(s')]$$

## III. Reinforcement Learning. (Online)

1. Model-Based Learning.

Learn empirical MDP Model.

Count  $s'$  for each  $s, a$  normalise  $\hat{T}(s, a, s')$

Solve the learnt MDP.

2. Model-Free Learning.

- 2-1. Passive Reinforcement Learning. (Direct Evaluation).

act according to  $\pi$ , average the samples.

easy to understand

don't need T, R

correct a.v.g. value using sample transitions

waste info about state

connection, state learnt

separately, takes a long time

- 2-2. Temporal Difference Learning. (Evaluates a policy  $\pi$ )

sample =  $R(s, \pi(s), s') + \gamma V^\pi(s')$

update:  $V^\pi(s) \leftarrow (1-\alpha)V^\pi(s) + \alpha \cdot \text{sample}$

$$V^\pi(s) \leftarrow V^\pi(s) + \alpha (\text{sample} - V^\pi(s))$$

Decreasing  $\alpha$  can give converging averages.

- 2-3. Q-Learning (Active Reinforcement Learning) (Full RL)

sample =  $R(s, a, s') + \gamma \max_{a'} Q(s', a')$

$$Q(s, a) \leftarrow (1-\alpha) Q(s, a) + \alpha \cdot \text{sample}$$

converges to optimal policy even if acts suboptimally.

\* off-policy learning.

have to explore enough, eventually  $\alpha \downarrow$  small enough and not too quickly.

3. Explore.

$\epsilon$ -greedy. - randomly

- exploration function.

n: visit count.

$$f(u, n) = u + k/n$$

$$Q(s, a) \leftarrow \alpha R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$$

4. Regret.

a measure of total mistake cost.

minimizing regret > learning to be optimal.

5. Approximate Q-learning

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

$$\text{difference} = [r + \gamma \max_{a'} Q(s', a')] - Q(s, a)$$

$$w_i \leftarrow w_i + \alpha \cdot \text{difference} \cdot f_i(s, a)$$

$$(\frac{dw_i}{dt}(s, a))$$

6. Policy Search.

Real priority, get the ordering of Q-values right (action prediction)

start with an OK solution.

change weight up and down to see if it's better.

problem: may need to run many episodes.

better method exploit lookahead structure, sample wisely, change multiple parameters.



## IV. Supervised learning.

### 1. Classification

learning  $f$  with discrete output value.

#### 1-1 Model-Based Classification.

Naïve Bayes: Assume all features are independent effects of label

# of parameters is linear in  $n$

tree structure: linear inference time.

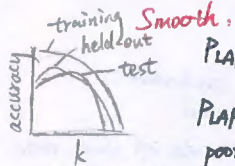
$$P(Y | w_1, w_2, \dots, w_n) \propto P(Y) \prod_i P(w_i | Y)$$

#### 1-2 Training.

$$\text{MLE: } P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}}$$

#### 1-3 Generalization and Overfitting.

why: too few training samples / noisy training data / too many attributes / too expressive. (store all spaces)



$$P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|} \quad k: \text{strength of prior}$$

$$P_{LAP,k}(x|y) = \frac{c(x,y) + k}{c(y) + k|X|}$$

poor performance when  $|X|$  or  $|Y|$  is large.

Linear interpolation ( $P(x|y)$  isn't too diff from  $P(x)$ )

$$P_{LIN}(x|y) = \alpha \hat{P}(x|y) + (1-\alpha) \hat{P}(x)$$

#### 1-4 Perceptron

$$\text{activation } w(x) = w \cdot f(x) \quad \begin{matrix} > 0 & \rightarrow & 1 \\ < 0 & \rightarrow & -1 \end{matrix}$$

$$\text{learn: } w = w + y^* \cdot f \quad y^* = \pm 1$$

If training set separable, perceptron converges.

#### 2. Regression.

learning  $f$  with real-valued output value.

$$L(w) = \sum_i (y_i - w^T x_i)^2, \quad w^* = (X^T X)^{-1} X^T y$$

Regularization: alleviate overfitting.

$$\text{LASSO: } L(w) = \sum_i (y_i - w^T x_i)^2 + \lambda \sum_k |w_k|$$

$$\text{Ridge Regression: } L(w) = \sum_i (y_i - w^T x_i)^2 + \lambda \sum_k w_k^2$$

"Occam's razor": prefer the simplest hypothesis consistent with the data.

## V. Unsupervised Machine Learning.

#### 1. Clustering what does similar mean?

one option:  $\text{dist}(x_i, y) = \sum_j (x_i - y_j)^2$  is small.

$$2. \text{K-means } \{ \{x_i\}, \{a_j\}, \{q_k\} \} = \sum_i \text{dist}(x_i, a_j)$$

pick  $k$  means

assign data to closest mean  $a_j = \arg\min_k \text{dist}(x_i, q_k)$

assign mean to average assigned points  $q_k = \frac{1}{|\{i: a_i = k\}|} \sum_{i: a_i = k} x_i$

stop when no points' assignments change.

Problem: Local optima; Equally Sized Clusters; Circular Clusters.

#### 3. Expectation-Maximization (EM).

pick  $k$  random cluster models (Gaussian)

assign data proportionately to different models

revise each cluster model on its assigned points

stop when no changes.

EM degrades to K-means if

1. All Gaussians are spherical and have identical weights and covariances  $\mu$  is only parameter

2. Label distribution in E step are point-estimations  $\sigma^2 \rightarrow 0$

$$L(\theta; D) \geq F(\theta; Q) = \sum_{j=1}^m \sum_i Q(i|j) \log \frac{P(z_i, x_i | \theta)}{Q(z_i | x_i)}$$

EM for HMM.

E Step: compute the distribution of hidden states given each training instance.

Infeasible to enumerate. But can compute expected counts of transitions and emissions using the forward and backward algorithms.

M Step: Update the parameters to maximize expected log likelihood based on distributions over hidden states. Closed-form solution: simply normalize the expected counts of transitions and emissions. Known as Baum-Welch algorithm.

## VI. Natural Language Parsing.

#### 1. CNF.

$$A \rightarrow BC$$

$$A \rightarrow w$$

conversion to CNF:

$$A \rightarrow B, B \rightarrow C \Rightarrow A \rightarrow C$$

$$S \rightarrow ABC \Rightarrow S \rightarrow XC, X \rightarrow AB$$

#### 2. Regular Grammar

$A \rightarrow aB$  or  $A \rightarrow a$ ; Probabilistic RG = HMM.

#### 3. Dependency Grammar. siblings, grandparent

The tree score is the sum of arc scores.

#### 4. Parsing (DP)

Probabilistic Parsing

$$P_{A,i,j} = \max_{B,C,k} P(A \rightarrow BC) \times P_{B,i,k} \times P_{C,k,j}$$

Run CYK on a HMM

CYK on HMMs  $\approx$  Viterbi

#### 5. Learning Grammars

Supervise Methods. (treebank)

MLE  $\times$  poor performance

增加额外标记.

Unsupervised Methods

EM.

#### 6. Context-free Grammar

① a set  $\Sigma$  of terminals (words)

② a set  $N$  of nonterminals (phrases)

③ a start symbol  $S \in N$

④ a set  $R$  of production rules

specifies how a nonterminal can produce

a string of terminals and/or nonterminals.

#### 7. CFG focus on constituents.