

BACKTRACKING

Tree Search

```
func TREF-SEARCH(problem, frontier) returns a solution or failure
frontier ← INSERT (MAKE-NODE (INITIAL-STATES Problem)), frontier)
while not IS-EMPTY(frontier) do
    node ← POF(frontier)
    if problem.IS-GOAL(node.STATE) then return node
    for each child-node in EXPAND(problem, node) do
        add child-node to frontier
return frontier
```

Depth-First-Search (DFS)

Completeness: maybe not complete.
Optimality: find the leftmost, might not optimal

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

Breadth-First-Search (BFS)

Completeness: Complete

Optimality: not optimal

Time complexity: $O(b^s)$ Space complexity: $O(b^s)$

Uniform Cost Search

Completeness: complete.

Optimality: optimal

Time complexity: $O(b^{C^*})$ (C^* is the optimal path cost
minimal cost between 2 nodes: Σ)
space complexity: $O(b^{C^*})$

Admissibility: $b_n, 0 \leq h(n) \leq h^*(n)$

Theorem: for a given search problem, if the admissibility constraint is satisfied by a heuristic function h , using A^* search with n will yield an optimal solution.

```
function A*-GRAPH-SEARCH(problem, frontier)
reached ← an empty dict mapping nodes to the cost to each one
frontier ← INSERT (MAKE-NODE (INITIAL-STATE, problem))
while not IS-EMPTY(frontier) do
    node, node.COST-TO-NODE ← POF(frontier)
    if problem.IS-GOAL(node.STATE) then return node
    if node.STATE is not in reached or reached[node.STATE] >=
        node.COST-TO-NODE then
        reached[node.STATE] = node.COST-TO-NODE
        for each child-node in EXPAND(problem, node) do
            frontier ← INSERT (child-node, child-node.COST+node.COST-TO-NODE), frontier)
    return failure
```

Dominance: If heuristic a is dominant over heuristic b , then the estimated goal distance for a is greater than the estimated goal distance for b for every node in state space.
 $V_n | h(a) \geq h(b)$

```
function BACKTRACKING-SEARCH(csp) returns solution/failure
    return RECURSIVE-BACKTRACKING({ }, csp)
function RECURSIVE-BACKTRACKING(assignment, csp) returns soln/failure
    if assignment is complete then return assignment
    var ← SELECT-UNASSIGNED-VARIABLE(VARIABLES[csp], assignment, csp)
    for each value in ORDER-DOMAIN-VALUES(var, assignment, csp) do
        if value is consistent with assignment given CONSTRAINTS[csp] then
            add {var = value} to assignment
            result ← RECURSIVE-BACKTRACKING(assignment, csp)
            if result ≠ failure then return result
            remove {var = value} from assignment
    return failure
```

OPTIMIZATION 1. Filtering (Forward checking)

```
function AC-3(csp) returns the CSP, possibly with reduced domains
inputs: csp, a binary CSP with variables  $X_1, X_2, \dots, X_n$ 
local variables: queue, a queue of arcs, initially all the arcs in csp
while queue is not empty do
    tail ← REMOVE-FIRST(queue)
    if REMOVE-INCONSISTENT-VALUES( $X_i, X_j$ ) then
        for each  $X_k \in \text{NEIGHBORS}[X_i]$  do
            add ( $X_k, X_i$ ) to queue
```

```
function REMOVE-INCONSISTENT-VALUES( $X_i, X_j$ ) returns true iff succeeds
removed ← false
for each  $x \in \text{DOMAIN}[X_i]$  do
    if no value  $y$  in  $\text{DOMAIN}[X_j]$  allows  $(x, y)$  to satisfy the constraint  $X_i \leftrightarrow X_j$ 
    then delete  $x$  from  $\text{DOMAIN}[X_i]$ ; removed ← true
return removed
if tree structured with tops:  $O(nd^2)$ 
```

AC-3 has a worst-case time complexity of $O(ded^3)$

OPTIMIZATION 2. Ordering

MRV (Minimum Remaining Values) (choose variable)

LCV (Least Constraining Values) (choose value)

Topological Sort

Minimax

- \forall agent-controlled states. $V(s) = \max_{s' \in \text{successors}(s)} V(s')$
- \forall opponent-controlled states $V(s) = \min_{s' \in \text{successors}(s)} V(s')$
- terminal states $V(s) = \text{known}$

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

```
def max-value(state):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor))
    return v
```

EXPECTIMAX

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

```
def max-value(state):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor))
    return v
```

```
def exp-value(state):
    initialize v = 0
    for each successor of state:
        p = probability(successor)
        v += p * value(successor)
    return v
```

Alpha-Beta Pruning.

α : MAX's best option on path to root
 β : MIN's best option on path to root

```
def max-value(state,  $\alpha, \beta$ ):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor,  $\alpha, \beta$ ))
        if  $v \geq \beta$  return v
         $\alpha$  = min( $\alpha, v$ )
    return v
```

```
def min-value(state,  $\alpha, \beta$ ):
    initialize v = +∞
    for each successor of state:
        v = min(v, value(successor,  $\alpha, \beta$ ))
        if  $v \leq \alpha$  return v
         $\beta$  = max( $\beta, v$ )
    return v
```

Temporal Difference Learning

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

$$\text{Sample} = R(s, \pi(s, s')) + \gamma \text{Sample}$$

$$V_k^T(s) \leftarrow (1-\delta) V_k^T(s) + \delta \text{Sample}_k$$

The recurrence definition of $V_k^T(s)$:

$$V_k^T(s) \leftarrow \delta [(1-\delta)^{k-1} \text{sample}_1 + \dots + (1-\delta)^{k-2} \text{sample}_{k-1} + \text{sample}_k]$$

Q-Learning

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

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

$$\text{Sample} = R(s, a, s') + \gamma \max_{a'} Q_k(s', a')$$

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

Approximate Q-Learning

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s) = \vec{w} \cdot \vec{f}(s)$$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a) = \vec{w} \cdot \vec{f}(s, a)$$

where $\vec{f}(s) = [f_1(s), f_2(s), \dots, f_n(s)]^T$

and $\vec{f}(s, a) = [f_1(s, a), f_2(s, a), \dots, f_n(s, a)]^T$

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

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

$$Q(s, a) \leftarrow Q(s, a) + \delta \cdot \text{difference}$$

E-greedy

$0 < \epsilon \leq 1$. if ϵ is large, prefer randomly.

if ϵ is small, prefer learn by slowly.

ϵ greedy

Exploration Functions

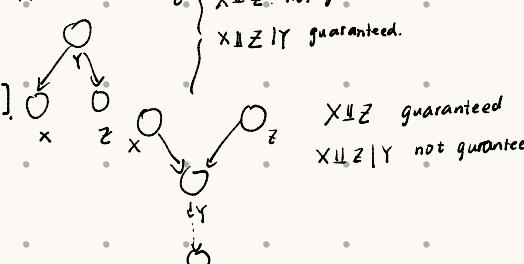
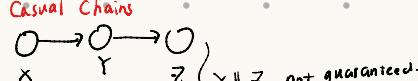
$$Q(s, a) \leftarrow (1-\delta) Q(s, a) + \delta [R(s, a, s') + \gamma \max_{a'} f(s', a')]$$

$$f(s, a) = \frac{k}{N(s, a)}$$

Bayesian Network

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{parents}(X_i))$$

Causal Chains



1 Shade all observed nodes $\{Z_1, \dots, Z_k\}$ in the graph.

2 Enumerate all undirected paths from X to Y .

3 For each path:

- a Decompose the path into triples (segments of 3 nodes).
- b If all triples are active, this path is active and d-connects X to Y .

4 If no path d-connects X and Y , then $X \perp\!\!\!\perp Y | \{Z_1, \dots, Z_k\}$.

Reinforcement Learning

{ Model-based Learning $\Rightarrow T(s, a, s')$
Model-Free Learning, $= \frac{\#(s, a, s')}{\#(s, a)}$

At convergence obtain an optimal policy-

$$\log(L(w_1 \dots w_k)) = \sum_{i=1}^n \sum_{k=1}^K t_{i,k} \log \left(\frac{e^{w_k^T f(x_i)}}{\sum_{l=1}^K e^{w_l^T f(x_i)}} \right) t_{i,k}$$

$$\log(L(w_1 \dots w_k)) = \sum_{i=1}^n \sum_{k=1}^K t_{i,k} \log \left(\frac{e^{w_k^T f(x_i)}}{\sum_{l=1}^K e^{w_l^T f(x_i)}} \right)$$

$$V_{\pi_j} \log(L(w)) = \sum_{i=1}^n V_{\pi_j} \sum_{k=1}^K t_{i,k} \log \left(\frac{1}{\sum_{l=1}^K e^{w_l^T f(x_i)}} \right)$$

$$= \sum_{i=1}^n (t_{i,j} - \frac{e^{w_j^T f(x_i)}}{\sum_{l=1}^K e^{w_l^T f(x_i)}}) f(x_i)$$

Neural Networks

$$L^{acc}(w) = \frac{1}{n} \sum_{i=1}^n (\text{sgn}(w \cdot f(x_i)) - y_i)^2$$

Sometimes, we want an output more expressive than a binary label. It then becomes useful to produce a probability for each of the (N) classes, reflecting our degree of certainty that the data point belongs to each class. As in multi-class logistic regression, we store a weight vector for each class (j), and estimate probabilities with the softmax function:

$$\sigma(x_i)_j = \frac{e^{f(x_i)^T w_j}}{\sum_{l=1}^N e^{f(x_i)^T w_l}} = P(y_i = j | f(x_i); w).$$

Given a vector output by f , softmax normalizes it to output a probability distribution. To derive a general loss function for our models, we can use this probability distribution to generate an expression for the likelihood of a set of weights:

$$\ell(w) = \prod_{i=1}^n P(y_i | f(x_i); w).$$

This expression denotes the likelihood of a particular set of weights explaining the observed labels and data points. We seek the set of weights that maximizes this quantity. This is equivalent to finding the maximum of the log-likelihood expression:

$$\log \ell(w) = \log \prod_{i=1}^n P(y_i | f(x_i); w) = \sum_{i=1}^n \log P(y_i | f(x_i); w).$$

Step function

$$f(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$b(x) = \frac{1}{1+e^{-x}} \text{ sigmoid.}$$

$$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases} \text{ ReLU}$$

Loss Functions and Multiplication

Optimization

$$\nabla_w L(w) = \left[\frac{\partial L(w)}{\partial w_1}, \dots, \frac{\partial L(w)}{\partial w_n} \right]$$

- Model-based learning of $T(s, a, s')$ and $R(s, a, s')$.
- Direct Evaluation to estimate $V(s)$.
- Temporal Difference learning to estimate $V(s)$.
- Learning to estimate $Q(s, a)$.

(b) In the limit of infinite timesteps, under which of the following exploration policies is Q -learning guaranteed to converge to the optimal Q-values for all state? (You may assume the learning rate α is chosen appropriately, and that the MDP is ergodic; i.e., every state is reachable from every other state with non-zero probability.)

- A fixed policy taking actions uniformly at random.
- A greedy policy.
- An ϵ -greedy policy
- A fixed optimal policy.

heuristic

graph search: must consistent

tree search: can be admissible

CNF & DNF

CNF: And of ORS: $(A \vee B) \wedge (C \vee D)$

DNF: OR of ANDS: $(A \wedge B) \vee (C \wedge D)$

Direct Evaluation

	Episode 1	Episode 2
B	east, C, -1 C, east, D, -1 D, exit, x, +10	B, east, C, -1 C, east, D, -1 D, exit, x, +10
C		
D		
	Episode 3	Episode 4
E	north, C, -1 C, east, D, -1 D, exit, x, +10	E, north, C, -1 C, east, A, -1 A, exit, x, -10

Walking through the first episode, we can see that from state D to termination we acquired a total reward of 10, from state C we acquired a total reward of $(-1) + 10 = 9$, and from state B we acquired a total reward of $(-1) + (-1) + 10 = 8$. Completing this process yields the total reward across episodes for each state and the resulting estimated values as follows:

s	Total Reward	Times Visited	$V^*(s)$
A	-10	1	-10
B	16	2	8
C	16	4	4
D	30	3	10
E	-4	2	-2

Multi-Class Logistic Regression

$$P(y_i | f(x); w) = \frac{e^{w_i^T f(x)}}{\sum_{k=1}^K e^{w_k^T f(x)}}$$

$$L(w_1 \dots w_k) = \sum_{i=1}^n p(y_i | f(x_i); w)$$

$$= \sum_{i=1}^n \prod_{k=1}^K \left(\frac{e^{w_k^T f(x_i)}}{\sum_{l=1}^K e^{w_l^T f(x_i)}} \right)^{t_{i,k}}$$

