

Depth-First Search Optimal vs. complete planning
Breadth-First Search Planning vs. replanning
Uniform-Cost Search

Search Problems

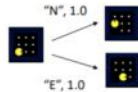
용어!

○ A **search problem** consists of:

○ A state space



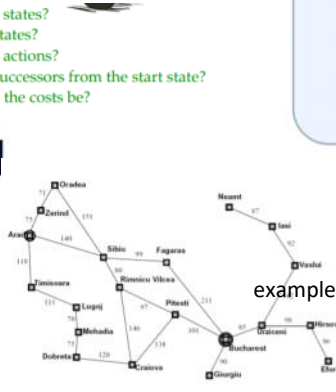
○ A successor function (with actions, costs)



○ A start state and a goal test

○ A **solution** is a sequence of actions (a plan) which transforms the start state to a goal state

What are the states?
How many states?
What are the actions?
How many successors from the start state?
What should the costs be?



Problem이 무엇인지에 따라 달라지는 것을 볼 수 있다.

A **search state** keeps only the details needed for planning (abstraction)

○ Problem: Pathing

- States: (x,y) location
- Actions: NSEW
- Successor: update location only
- Goal test: is (x,y)=END

○ Problem: Eat-All-Dots

- States: ((x,y), dot booleans)
- Actions: NSEW
- Successor: update location and possibly a dot boolean
- Goal test: dots all false

○ State space:

○ Cities

○ Successor function:

○ Roads: Go to adjacent city with cost = distance

○ Start state:

○ Arad

○ Goal test:

○ Is state == Bucharest?

○ Solution?

Search Trees

Depth-First Search

Breadth-First Search

Uniform Cost Search



No, **optimal**

$O(b^m)$



Only if costs are all 1

$O(b^s)$



optimal

$O(b^{C^*/\epsilon})$

Tree A*

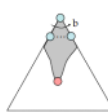
Combining UCS and Greedy

$$f(n) = g(n) + h(n)$$

Uniform-cost Greedy

Admissible Heuristics

$$0 \leq h(n) \leq h^*(n)$$



never **expand** a state twice

A* Graph Search

노드 확장 순서,

Consistency of Heuristics

결과, 남아있는 노드

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

takes time
space does the fringe take?

LIFO stack

FIFO queue

priority queue
cheapest node first.

```
function TREE-SEARCH(problem, fringe) return a solution, or failure
  fringe ← INSERT(MAKE-NODE(INITIAL-STATE[problem]), fringe)
  loop do
    if fringe is empty then return failure
    node ← REMOVE-FRONT(fringe)
    if GOAL-TEST(problem, STATE[node]) then return node
    for child-node in EXPAND(STATE[node], problem) do
      fringe ← INSERT(child-node, fringe)
    end
  end
```

No information about goal location

-> 단점, 방향성x

greedy



UCS

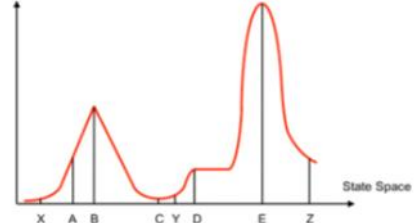


and

A*

```
function GRAPH-SEARCH(problem, fringe) return a solution, or failure
  closed ← an empty set
  fringe ← INSERT(MAKE-NODE(INITIAL-STATE[problem]), fringe)
  loop do
    if fringe is empty then return failure
    node ← REMOVE-FRONT(fringe)
    if GOAL-TEST(problem, STATE[node]) then return node
    if STATE[node] is not in closed then
      add STATE[node] to closed
      for child-node in EXPAND(STATE[node], problem) do
        fringe ← INSERT(child-node, fringe)
      end
    end
  end
```

Objective Function



QUIZ >

Starting from X, where do you end up ?

Starting from Y, where do you end up ?

Starting from Z, where do you end up ?

풀기

Constraint Satisfaction Problems



Variables: WA, NT, Q, NSW, V, SA, T

Domains: $D = \{\text{red, green, blue}\}$

Constraints: adjacent regions must have different colors



n-queens

Variables: Q_k

Domains: $\{1, 2, 3, \dots, N\}$

Constraints:

Backtracking Search

백트래킹: 더 이상 전진불가

Idea 1: One variable at a time

Idea 2: Check constraints as you go

Forward checking



Arc Consistency in a CSP



-> efficient

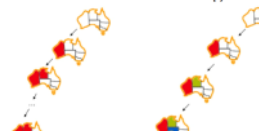
Ordering: Minimum Remaining Values, Least Constraining Value

Filtering

Structure – turns out trees are easy.

BFS 아주 비효율

DFS vs. Backtracking



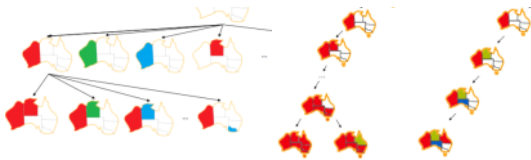
백트래킹에서 진화

-> Filtering : 포워드 체킹

상대가 반드시

최선의 선택을 할 것임

Expectimax 출현



백트레이킹에서 진화
-> Filtering : 포워드 체킹

상대가 반드시
최선의 선택을 할 것임

Expectimax 출현

Worst-Case vs. Average Case

Minimax Values

def max-value(state):

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

def min-value(state):

```
initialize v = +∞
for each successor of state:
    v = min(v, value(successor))
return v
```

Just like (exhaustive) DFS

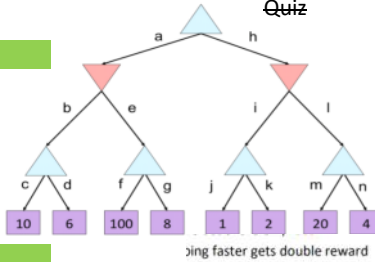
Time: $O(b^m)$

Space: $O(bm)$

Depth Matters

Depth 지점에 따라

결과 다름



Alpha-Beta Pruning

def max-value(state, α, β):

```
initialize v = -∞
for each successor of state:
    v = max(v, value(successor, α, β))
    if v ≥ β return v
    α = max(α, v)
return v
```

def min-value(state, α, β):

```
initialize v = +∞
for each successor of state:
    v = min(v, value(successor, α, β))
    if v ≤ α return v
    β = min(β, v)
return v
```

Expectimax

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
```

An MDP is defined by:

- o A set of states $s \in S$
- o A set of actions $a \in A$
- o A transition function $T(s, a, s')$
 - o Probability that a from s leads to s' , i.e., $P(s' | s, a)$
 - o Also called the model or the dynamics
- o A reward function $R(s, a, s')$
 - o Sometimes just $R(s)$ or $R(s')$
- o A start state
- o Maybe a terminal state policy $\pi^*: S \rightarrow A$

Markov decision processes:

- o Set of states S
- o Start state s_0
- o Set of actions A
- o Transitions $P(s' | s, a)$ (or $T(s, a, s')$)
- o Rewards $R(s, a, s')$ (and discount γ)

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

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

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

현실(don't know T or R)

A set of states $s \in S$

A set of actions (per state) A

A model $T(s, a, s')$

A reward function $R(s, a, s')$

Il looking for a policy $\pi(s)$

Model-Based Learning Model-Free Learning

Model-Free (temporal difference) Learning

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

Update to $V(s)$: $V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)sample$

Same update: $V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$

$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + \alpha [R(s, \pi(s), s') + \gamma V^\pi(s')] \quad Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

Policy Iteration

0단계: policy를 임의로 정한다.

1단계: 현재의 policy에 대해 (not optimal) utility를 계산

2단계: 계산한 utility를 이용하여 policy를 개선

개선한 policy를 현재의 policy로 한다.

3단계: 1~2단계를 policy가 수렴할 때까지 반복한다.

상태집합 $S \ni s$

시작상태 (start state) 종료상태 (terminal state)

행동집합 $A \ni a$

전이함수(transition function) $T(s, a, s') = P(s' | s, a)$

현재 상태 s에서 행동 a를 했을 때, 다음상태가 s'이 될 확률

보상함수(reward function) $R(s, a, s')$

Utility 보상의 합

reward sequence가 주어졌을 때 구함

$$Utility = r_0 + r_1 + r_2 + r_3 + \dots + r_{end}$$

$$utility = \sum_{t=0}^{t=end} r_t$$

$$Utility = \gamma^0 r_0 + \gamma^1 r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \dots + \gamma^{end} r_{end}$$