

Search method: 1. Uniformed Search 2. Informed Search

* Uniformed search: 不知終點在哪 \Rightarrow 盲目尋找

Depth-First Search — Complete: ? , Optimal: No

以深度為主, fringe 乃 a LIFO stack

Breadth-First Search — Complete: Yes, Optimal: 不一定

以廣度為主, fringe 乃 a FIFO queue

但DFS會有'cycle走不出來'的問題 \Rightarrow 限制每次iteration的深度

Depth-First Search - Iterative Deepening

get DFS's space advantage with BFS's time advantage

但DFS & BFS都無考慮到cost

Uniform Cost Search — Complete: yes, Optimal: yes.

找累積到該node後最小的cost的node展開

Search method: 'Uniform Search' ² informed search

* Informed Search: 知道目標在什麼位置

Heuristics (啟發式) - 一種工具

有某個 'concept' 目標像 'concept' 去做 (concept 是自己定義的)

Greedy Search - Complete: ?, Optimal: No.

每次做的決策都是對當下最有利的 (by heuristics)

A* Search (UCS + greedy) - Complete: yes, Optimal: with admissible heuristic (tree)

因為 UCS 像 BFS 會一直尋找, 無目標; greedy 知道目標但 cost 可能會太大)

⇒ UCS 考慮 '走的成本', greedy 考慮 '到終點的成本' ($h(x)$)

但 $h(x)$ 是自己估計的, 需要一些條件, 否則不是 optimal

⇒ Admissible Heuristic (可接受的 h) $0 \leq h(n) \leq h^*(n)$, $h^*(n)$ 是 true cost to goal

但 tree search 會有重複 node 的問題 ⇒ extra work

Graph Search

將走過的 node 存起來才不會重複走

但因為有 heuristic 所以又會有一些問題

⇒ Consistency of heuristics $h(A) - h(C) \leq \text{cost}(A \rightarrow C)$

Summary

Tree A*: is optimal with admissible heuristics

Graph A*: is optimal with consistent heuristics

Constraint Satisfaction Problem 條件滿足的問題
達成某些條件才是重點 (目標在滿足條件的狀況下完成)

* Constraint Satisfaction Problem

state defines by variables x_i
with domain D
some constraints

- Unary: involves one variable
- Binary: involves two variable
- high-order: involves three or more

* Constraint Satisfaction Method

也可用 DFS, BFS... 但效果非常差

Backtracking Search = DFS + variable-ordering + fail-on-violation

is the basic uninformed algorithm

一個一個加, 再不違反條件的情況下繼續走

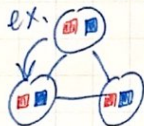
但等到發現錯誤時通常來不及了 \Rightarrow detect failure early

Filtering

assign 一個 variable 時將和它有相關 constraint 的也刪掉, 但它也無法全部檢查到

\Rightarrow Arc consistency: 双向看是否有違規 (if x loss a value, neighbors needs to be checked)

但也不是都可 fixed



Ordering

1. 將剩餘值少的先處理 minimum remaining value (MRV)

2. 選會造成較少限制的 least constraining value (LCV)

但也是還有 limit arc consistency 的問題, 且若是 k -consistency 會要大量計算

Structure

重新排列, 把不相關的 & 中間點切開

Tree-Structured

切開某部分使它成為一個 tree

Iterative Algorithm

將所有地塗色, 再慢慢修改到無不合規定的

* No fringe, Not complete & Not Optimal

Local Search

improve a single option until you can't make it better. 會再找到更好的解

* No fringe, incomplete but suboptimal

Adversarial method: 1. Deterministic 2. non-deterministic

希望演算法計算一個 strategy (policy) 可推薦採取的動作

Types of game: deterministic or stochastic? / one or more play? / zero sum?

Single Agent Tree

* value of a state: 在該位置有可能拿到的最高分

Adversarial Game Tree - Minimax search

計算每個 node 的 minimax value, Just like DFS

但若是一開始直接做, 則有可能為了 optimal 就放棄一些机会 (對手可能不是 optimal)

且 like DFS 則計算量大

Game tree Pruning - Alpha-Beta Pruning

在 minimax search 中, 已知拿到最高分了, 則某些子樹也不須要經過

α : max's best option on path to root

β : min's best option on path to root

但在部分過程中算的 node 的 value 可能是錯的, 不過對結果不會影響

但在真實狀況中, 無法計算到最後

⇒ depth-limited (限制計算 minimax, α, β 的深度)

用 evaluation function 取代 terminal utilities (依照現在的情況給分)

要這個是因為若 α 限制太淺會有 thrash 情況

但上面的演算法是取最安全的情況下, 但若是事情發生有概率...

Expectimax Search

計算平均而非 minimax, 且無法 pruning (因為計算平均)

Markov Decision Process

- Utilities: 將喜好轉為數字來表現

- Rational Preferences: $(A > B) \wedge (B > C) \Rightarrow (A > C)$

- Maximum expected utility (MEU): choose the action that maximizes expected utility

- Utility Scales: 測量 utility 的方式

Adversarial Method: 'Deterministic' non-deterministic

當有些 action 發生時不是固定的, 是有機率的

Markov Decision Process - Expectimax

接下來做的決定, 只和現在的狀態有關

在 deterministic 時, 我們要的是一個計畫

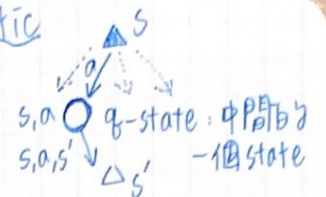
在 non-deterministic 中, 我們要的是一個 "policy"

one way to solve them is with "expectimax" but it didn't compute entire policies

但若是 reward 先拿到會較好? (因有些沒拿到 reward 前就結束了)

⇒ discounting: 得到的 rewards 隨時間下降 (⇒ 愈早拿到愈好) → 助於收斂

限制深度: 似 depth-limited search



做某個 action 有部分機率會跳到另一個 state

Markov decision Problem

- set of states $s \in S$
- set of actions $a \in A$
- transition function $T(s, a, s')$
- reward function $R(s, a, s')$

Optimal Quantities

- $V^*(s)$: 最佳動作得到的值
- $Q^*(s, a)$: 期望做 a 後拿到最好的值
- $\pi^*(s)$: 最佳的 policy

☆ V^* 不一定和 Q^* 相同, 都是 optimal 才會相同

但用 expectimax 會花太多時間, 且有太多重複的 state, tree 會無限深

Time-limited Value

計算 k 步的每個 state 的值就好, 不用一直無限算 (it's what a depth-k expectimax)

⇒ 愈深則 state 的值會慢慢收斂 (due to discounting)

Value Iteration

從底層開始做, 由累積的方式慢慢加上, repeat until convergence

⇒ will converge to unique optimal value

到目前為止

- policy: 每個 state 採取的 action

- utility: reward (with discounting) 的總和

- value: 在該 state 期望得到的最大 reward (max node)

- q-value: 在該 state 期望得到的最大 reward (chance node)

The Bellman Equations: 定義 "optimal utility" via expectimax

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

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

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

但要怎麼知道什麼時候會收斂?

Policy Iteration: policy evaluation + policy extraction

- evaluation: 依照你的 policy 給一個評分 → fixed policy π : 給固定的 policy

- extraction: 把好的 policy 解出來

⇒ 先給 (部分) 固定的 policy, 一直 evaluation 再慢慢更新修改到 optimal, 加快收斂

Reinforcement Learning - online planning

前面都是已經知道所有的規則，但若是不知道，只能一直try，然後修好又在不知T、R、P的情況下，慢去學習

* Reinforcement Learning

還是MDP的問題，但只是不知道T、R、P

Model-based Learning: 經過多次try後可計算T、R、P

Model-free Learning: 不去管T、R、P，直接算state的avg

Passive reinforcement learning (固定policy去採樣學習)

direct evaluation: 直接計算每個state的avg under π

“但每個state都是独立的學習 → 要花較久時間

→ 那為什麼不用policy evaluation呢? 因為我們沒有T、R

sample based policy evaluation: 在某個state不斷repeat做採樣

“但我們不一定能回到上一片

Temporal Difference Learning

不斷的回到原點重複做後學習再取動態平均

加上exponential moving average: 較晚學習到的資料比例占越大

(因為晚學到的較好)

“但我們的policy還是固定的... → learn Q-value instead of value”

Active reinforcement Learning (學習 optimal 的 policy / values)

Q-value iteration → Q-learning

學習 Q-value 再和前面的合併

“但在學習過程中或許還有別的路，要不要走走看? → explore

- off-policy learning

- converges to optimal policy even if you're acting suboptimally

“要不要去試試其它路?

Exploration: 加上 ϵ 讓state可能亂跳

“但有的 ϵ 不管走多久都會亂走 → 令 ϵ 隨時間 \downarrow exploration function

exploration function: 去過的地方傾向用原來的action, or explore new place

“但沒辦法通用? (改變門的位置就出不去了)

→ feature-based representation - linear value function

describes a state using a vector of features, 好處是經驗會被新整起來

壞處是雖然share同一個feature但值會很不一樣

policy search: start with an ok solution (e.g. Q-learning) then fine-tune by

hill climbing on feature weights
→ learn the policy that maximize the rewards, not value