

1.

For DFS, with noise $w_t = 0$.

State: $x_0 \in X$, which is a finite space.

goal: $\min_{u_0: T-1} q(x_T) + \sum_{t=0}^{T-1} l_t(x_t, u_t)$

s.t. $x_{t+1} = f(x_t, u_t)$

$x_t \in X, u_t \in U$.

we can transfer it into a DSP problem:

state \Leftrightarrow node. $V := \left(\bigcup_{t=0}^T \{(t, x_t) \mid x_t \in X\} \right) \cup \{\tau\}$

initial start point $s = (0, x_0)$

τ is an artificial terminal state.

cost \Leftrightarrow edge.

$$C = \left\{ ((t, x_t), (t+1, x_{t+1}), c) \mid c = \min_{\substack{u \in U(x_t) \\ s.t. x_{t+1} = f(x_t, u)}} l_t(x_t, u) \right\} \cup \{(T, x_T), \tau, q(x_T)\}$$

a. Thus, for Dijkstra's algorithm, its backward version
is

Backward Dijkstra's Algorithm

$$OPEN \leftarrow \{\tau\}, g_\tau = 0, g_i = \infty \quad \forall i \in V \setminus \{\tau\}$$

while OPEN is not empty do

 Remove $i = \arg \min_{j \in OPEN} g_j$ from OPEN

 for $k \in \text{Children}(i)$

 if $g_k > g_i + c_{k,i}$

$g_k \leftarrow g_i + c_{k,i}$

$\text{Parent}(k) \leftarrow i$

 if $k \neq s$ then

$OPEN.append(k)$

It actually investigates the same states as Dynamic Programming(DP). For DP.

$$V_t(i) = \min_{j \in V \setminus \{t\}} c_{ij} + V_{t+1}(j)$$

$$\pi_t(i) = \arg \min_{j \in V \setminus \{t\}} c_{ij} + V_{t+1}(j)$$

will still investigate all children of i to find out the minimum.

(b) Yes, because of the equivalence between DFS and DSP.

Here, a heuristic means a lower bound estimation for cost that takes from state x_t to terminal state x_T .

Backward A* Algorithm

$$OPEN \leftarrow \{T\}, CLOSED = \{\}, \epsilon \geq 1$$

$$g_T = 0, g_i = \infty \text{ for } i \in V \setminus \{T\}$$

while $S \notin CLOSED$ do

Remove $i = \arg \min_{j \in OPEN} f_j = \arg \min_{j \in OPEN} g_j + h_j$ from OPEN

$CLOSED.append(i)$

for $k \in \text{Children}(i)$ and $k \notin CLOSED$ do

if $g_k > g_i + c_{ki}$ then

$g_k \leftarrow g_i + c_{ki}$

$\text{Parent}(k) \leftarrow i$

$OPEN.append(k)$

2.

(a) Because $h^{(1)}$ and $h^{(2)}$ are consistent, then

$$h(x_i) = \max \{ h^{(1)}(x_i), h^{(2)}(x_i) \} = 0$$

Also,

$$\begin{aligned} h(x_i) &= \max \{ h^{(1)}(x_i), h^{(2)}(x_i) \} \\ &\leq \max \{ h^{(1)}(x_j) + \text{cost}(x_i, x_j), h^{(2)}(x_j) + \text{cost}(x_i, x_j) \} \\ &= \text{cost}(x_i, x_j) + \max \{ h^{(1)}(x_j), h^{(2)}(x_j) \} \end{aligned}$$

Hence h is also consistent.

(b) If $h^{(1)}, h^{(2)}$ are heuristic, then

$$h(x_i) = h^{(1)}(x_i) + h^{(2)}(x_i) = 0.$$

Plus,

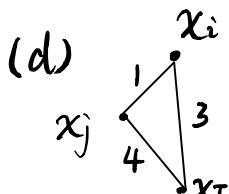
$$\begin{aligned} h(x_i) &= h^{(1)}(x_i) + h^{(2)}(x_i) \\ &\leq 2\text{cost}(x_i, x_j) + h^{(1)}(x_j) + h^{(2)}(x_j) \\ &= 2\text{cost}(x_i, x_j) + h(x_j) \end{aligned}$$

Let $\epsilon = 2$, then h is 2-consistent.

(c) h is consistent, then

$$\begin{aligned} h(x_i) &\leq \text{cost}(x_i, x_r) + h(x_r) \\ &\leq \text{dist}(x_i, x_r) \end{aligned}$$

Thus h is also admissible.



If h is consistent then

$$h(x_i) \leq 1 + h(x_r)$$

$$h(x_j) \leq 4 + 0 = 4$$

$$\begin{array}{ll}
 \text{If I set } h(x_i) = 3, \text{ then } h(x_i) > 1 + h(x_j) \\
 h(x_j) = 1 & h(x_j) \leq 4 + 0 = 4 \\
 h(x_t) = 0 & \text{however } h(x_i) \leq \text{dist}(x_i, x_t) = 3 \\
 & h(x_j) \leq \text{dist}(x_j, x_t) = 4
 \end{array}$$

Thus $h(x_i) = 3, h(x_j) = 1, h(x_t) = 0$ is admissible
but not consistent.

(e) Assume i is expanded and h is consistent, then
 $f_i = g_i + h_i \leq f_j \quad \forall j \in OPEN$.

Suppose g_i can be improved again, i.e. g_i is larger than
the least cost from s to i ($g_i > \text{dist}(s, i)$)

Then, there at least exists one state j on an optimal
path from s to i , s.t. $j \in OPEN, j \notin CLOSED, f_i \leq f_j$

However,

$$f_i = g_i + h_i > \text{dist}(s, i) + h_i = g_j + [\text{dist}(j, i) + h_i] \geq g_j + h_j = f_j$$

contradicts with $f_i \leq f_j$

Hence g_i cannot be approved.

(f) h is neither consistent, nor admissible.

h is not consistent.

For example, $x_i \in \mathbb{R}^2$, Suppose $x_i = (6, 8)^T, x_t = (0, 0)^T$

$$\text{then } h(x_i) = 8 + 0.4 \times 6 = 10.4.$$

$$\text{cost}(x_i, x_t) = 10, \quad h(x_t) = 0$$

$$\text{then } h(x_i) > h(x_t) + \text{cost}(x_i, x_t).$$

Therefore, h is not consistent.

h is not admissible.

Still suppose $x_i = (6, 8)^\top$, $h(x_i) = 10.4$

Then $\text{dist}(x_i, x_t) = \text{cost}(x_i, x_t) = 10 < h(x_i)$

Thus h is not admissible.

3.

(a) We can formulate this problem as follows:

State space: $X = [n] = \{1, 2, \dots, n\}$, state $x_t \in X$.

Control space: $U = \{1, 0\}$, where "1" means we choose to

recharge, "0" means we choose not to recharge,
i.e. browsing & chatting using phone.

Control input $u_t \in U$.

Motion model: $x_{t+1} = \begin{cases} x_t, & \text{if } u_t = 1, \text{ with prob. } 1-q \\ 1, & \text{if } u_t = 1, \text{ with prob. } q \\ j, & \text{if } u_t = 0, \text{ with prob. } P(x_t, j) \end{cases}$

Reward/Cost: $l(x_t = i, u_t = 0) = -r(i), i \in [n]$

$$l(x_t = i, u_t = 1) = \begin{cases} qC & \text{with prob. } q \\ 0 & \text{with prob. } 1-q \end{cases}$$

Then value function:

$$V_t(x_t) = \mathbb{E}_{x_{t+1}, T} \left[\sum_{\tau=t}^{T-1} \gamma^{\tau-t} l(x_\tau, u_\tau) \mid x_t \right]$$

As $T \rightarrow \infty$, it becomes a discounted Infinite-horizon

problem: $\min_{u \in U} V(x) = \min_{u \in U} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t l(x_t, u_t) \mid x_0 = x \right]$

$$x_{t+1} \sim p_f(\cdot \mid x_t, u_t) = \begin{cases} p_f(x_{t+1} \mid x_t, u_t = 0) = P(x_t, x_{t+1}) \\ p_f(x_t \mid x_t, u_t = 1) = 1 - q \\ p_f(1 \mid x_t, u_t = 1) = q \end{cases}$$

$$x_t \sim X = \{1, 2, 3, \dots, n\}$$

$$u_t \sim U = \{1, 0\}$$

Related Bellman equation:

$$V^*(x) = \min_{u \in U(x)} \{l(x, u) + \gamma \mathbb{E}_{x' \sim p_f(\cdot \mid x, u)} [V(x')]\} \quad \forall x \in X$$

where $l(x, u)$ is defined above.

It can also be written in a compacted version:

$$V^*(x) = \min_{u \in U(x)} \left\{ -r(x)(1-u) + qC + \gamma \sum_{x' \in X} p_f(x'|x, u) V^*(x') \right\}$$

$$= \min_{u \in U(x)} \begin{cases} -r(x) + \gamma \sum_{x'=l}^n p(x, x') V^*(x') & \text{if } u=0, l \in [n] \\ qC + r[qV^*(1) + (1-q)V^*(x)] & \text{if } u=1 \end{cases}$$

(b) Use Value Iteration (VI) to prove $V^*(i)$ increases in i :

① Assume $t=0$, $V_0(i) = 0$, then $t=1$,

$$V_1^*(i) = \min_{u \in \{1, 0\}} \{-r(i), qC\} = -r(i), \text{ since } r(i) > 0.$$

According to that $r(i)$ is decreasing as i increases,

$V_1^*(i) = -r(i)$ is increasing as i increases.

② When $t > 0$, suppose $V_t^*(i)$ increases as i increases, then we want to prove $V_{t+1}^*(i)$ also has the same property. (Mathematical Induction).

Let $Q_t(i, u) = V_t(i)$,

$$\text{then } V_{t+1}^*(i) = \min \{Q_{t+1}(i, u=1), Q_{t+1}(i, u=0)\}.$$

We want to prove

$$V_{t+1}^*(i+1) = \min \{Q_{t+1}(i+1, u=1), Q_{t+1}(i+1, u=0)\} < V_{t+1}^*(i) \quad (*)$$

If we can prove

$$Q_{t+1}(i+1, u=1) > Q_{t+1}(i, u=1)$$

$$\text{and } Q_{t+1}(i+1, u=0) > Q_{t+1}(i, u=0).$$

then $(*)$ is correct.

When $n=1$, based on VI:

$$\overline{V_{t+1}(i)} = qC + \gamma [qV_t(1) + (1-q)V_t(i)]$$

$$V_{t+1}(i+1) = qC + \gamma [qV_t(1) + (1-q)V_t(i+1)].$$

By induction assumption, $V_t(i) < V_t(i+1)$, $1-q > 0$

then $V_{t+1}(i+1) > V_{t+1}(i)$, i.e. $Q_{t+1}(i+1, u=1) > Q_{t+1}(i, u=1)$

When $u=0$, based on n VI:

$$\overline{V_{t+1}(i)} = -r(i) + \gamma \sum_{j=1}^n P(i,j) V_t(j).$$

we know $-r(i)$ is monotonously increasing over $i \in [n]$.

$$\sum_{j=1}^n P(i, j) V_t(j) = \underbrace{\left(\sum_{j=1}^n P(i, j) \right)}_{>0} \frac{V_t(1)}{} + \sum_{j=2}^n \underbrace{\left(\sum_{l=j}^n P(i, l) \right)}_{>0} \frac{(V_t(j) - V_t(j-1))}{}$$

\downarrow monotonously
 increasing monotonously
 increasing (Induction
 assumption)
 (Stochastic dominance)

Thus, $\sum_{j=1}^n P(i,j)V_t(j)$ is also monotonously increasing over $i \in [n]$.

Then $V_{t+1}(i)$ is monotonously increasing.

$$\text{so } V_{t+1}(i) < V_{t+1}(i+1)$$

Therefore, (*) holds and $V_{t+1}^*(i)$ increases in $i \in [n]$.

Based on ①②, the conclusion holds.

(c) $V^*(x)$

$$= \min_{u \in U(x)} \begin{cases} -r(x) + \gamma \sum_{x' \in \mathcal{S}} P(x, x') V^*(x') & \text{if } u=0, l \in [n] \\ qc + \gamma [qV^*(l) + (1-q)V^*(x)] & \text{if } u=1 \end{cases}$$

If $u=1$ is optimal, $V^*(i) = qc + \gamma [qV^*(l) + (1-q)V^*(i)]$,

then we can solve for $V^*(i)$, i.e. $V^*(i)$ is a constant. $V^*(i) = \frac{qc + \gamma q V^*(l)}{1 - \gamma(1-q)}$.

If $u=0$ is optimal policy,

$$V^*(i) = -r(i) + \gamma \sum_{j=1}^n P(i, j) V^*(j), l \in [n]$$

is increasing in $i \in [n]$.

Thus, if $-r(i) + \gamma \sum_{j=1}^n P(i, j) V^*(j) = qc + \gamma [qV^*(l) + (1-q)V^*(i)]$ exists zero-point, then the thresh nature holds.

Otherwise, $t=1$ or n is the threshold.

