



Aalto University
School of Electrical
Engineering

ELEC-E8119 Reinforcement Learning Planning in Discrete Space Introduction to Optimal Control

Ville Kyrki

24.9.2019

Today

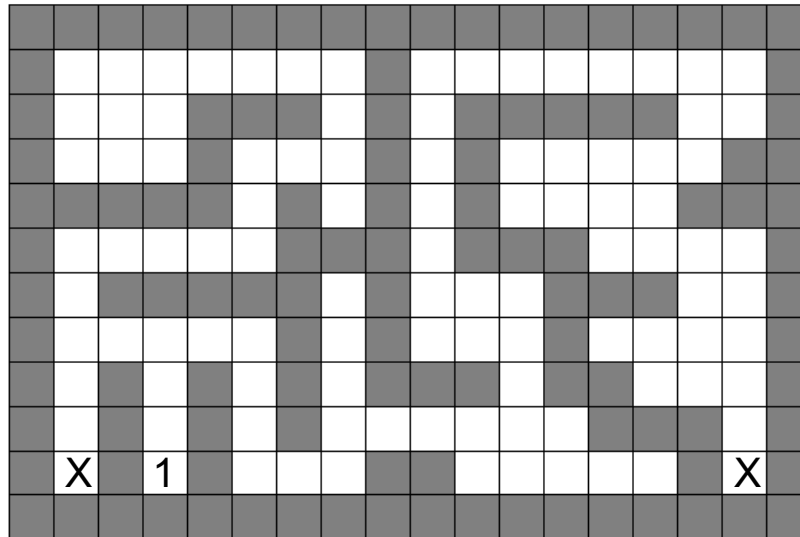
- Traditional (discrete) planning
- How to plan with complete knowledge?
- Connection to optimal control

Learning goals

- Formulate discrete planning problems.
- Use value iteration as a way to solve them.
- Explain relationship between discrete planning and optimal control.

Maze

- What are: state space, initial state, goal states set, action space, state transition (dynamics) function?



Handwritten notes and formula:

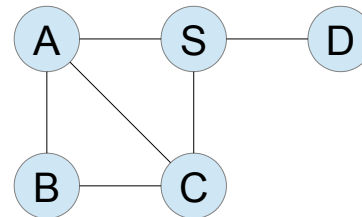
$f((x, y) + \vec{u}) = (x', y')$
if possible

Arrows indicate the action space \vec{u} (up, down, left, right) and the resulting state (x', y') .

Planning as graph search – General search algorithm

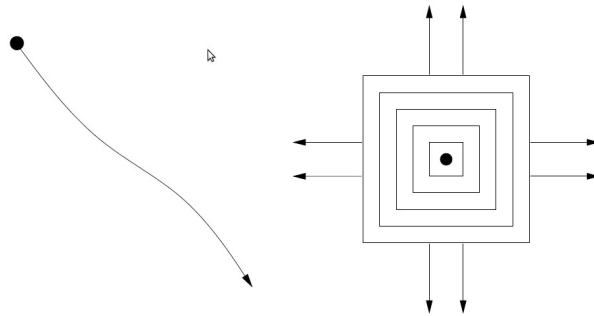
- Keep track of visited states (dead and alive).
 - Pick an alive state and visit one of the unvisited states (reachable and not dead).
 - Only difference between algorithms how to choose the state to visit.

Propose algorithms!



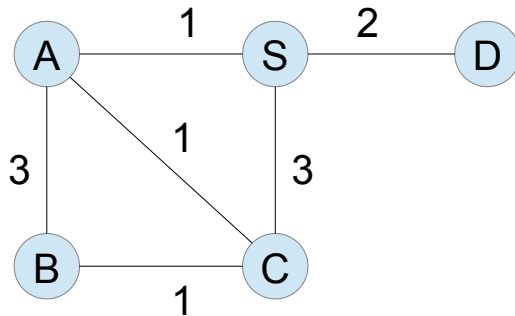
Search strategies

- Are all search strategies good?
 - What properties we want a strategy to have?
- *Systematic* graph search
 - All states will be visited (for finite graphs)
 - If solution exists, it will be found in finite time (for infinite graphs)



Shortest paths

- Associate each edge of graph with nonnegative cost
- Cost of plan is the sum of costs



Shortest path from S to B?

G	1	2
S	0	
A	∞	1
B	∞	∞ 4 3
C	∞	3 2
D	∞	2

Dijkstra

Dijkstra's algorithm

- Idea:
 - Keep track of “cost-to-come”, i.e. accumulated minimum cost to reach a particular state
 - Expand state with minimum cost-to-come
 - Update “cost-to-come” estimates until optimum is reached.
- Computational complexity $O(|V| \log |V| + |E|)$
- Exercise: Solve shortest paths from S in the previous graph using Dijkstra.
- A* is an extension of Dijkstra with heuristic lower bound.

Optimal planning (fixed-length plans)

Like feasible planning, plus

- cost functional *any kind of costs*

$$L(\tau_K) = \sum_{k=1}^K l(x_k, u_k) + l_F(x_{K+1})$$

$$\tau_K = (u_1, \dots, u_K)$$

$$l_F(x) = \begin{cases} 0, & x \in X_G \\ \infty, & x \notin X_G \end{cases}$$

- Goal: $\min_{\tau} L(\tau)$

Solving optimal planning

- Principle of Optimality (Bellman, 1957): An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision.
- Value-function:
 - cost-to-go $G_k^*(x_k)$

$$G_k^*(x_k) = \min_{u_k, \dots, u_K} \left\{ \sum_{i=k}^K l(x_i, u_i) + l_F(x_F) \right\}$$

Backward value iteration

- Assume we know $G_{k+1}^*(x)$
how to compute $G_k^*(x_k)$?

$$G_k^*(x_k) = \min_{u_k, \dots, u_K} \left\{ \sum_{i=k}^K l(x_i, u_i) + l_F(x_F) \right\}$$

$$\begin{aligned} G_k^*(x_k) &= \min_{u_k} \left\{ l(x_k, u_k) + G_{k+1}^*(x_{k+1}) \right\} & G_{K+1}^*(x) &= l_F(x) \\ &= \min_{u_k} \left\{ l(x_k, u_k) + G_{k+1}^*(f(x_k, u_k)) \right\} \end{aligned}$$



Value iteration, unknown length plans

- Iterating recursion until value function stationary: optimal cost plans have been found from all states that can reach a goal state

$$G^*(x) = \min_u \{l(x, u) + G^*(f(x, u))\}$$

Bellman equation

- Using G^* , optimal actions can be found from

$$u^* = \operatorname{argmin}_{u \in U(x)} \{l(x, u) + G^*(f(x, u))\}$$

- Complexity $O(K|X||U|) = O(|X|^2|U|)$

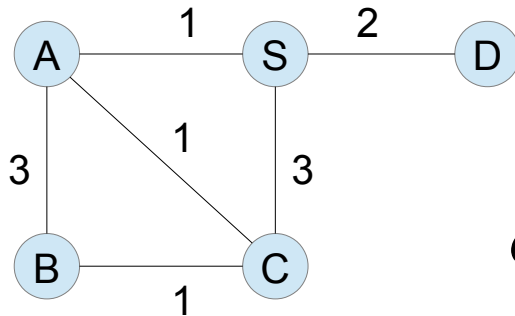
Why? Compare to Dijkstra!

Exercise

- Use backward value iteration for

$$x_I = S$$

$$X_G = \{B\}$$



	$k+1$	k	$k-1$	$k-2$	$k-3$
S	∞	∞	4	3	3
A	∞	3	2	2	2
B	0	0	0	0	0
C	∞	1	1	1	1
D	∞	∞	∞	6	5

Handwritten red annotations:
 - A red arrow points from the value 3 in the S row, $k-2$ column to the value 3 in the A row, $k-2$ column.
 - Below the value 6 in the D row, $k-2$ column, is the text "2+4".
 - Below the value 5 in the D row, $k-3$ column, is the text "2+3".

Reminder:

$$G^*(x) = \min_u \{l(x, u) + G^*(f(x, u))\}$$

What about continuous state variables?

- Simple systems can often be modeled as (locally) linear

$$x_{k+1} = A x_k + B u_k$$

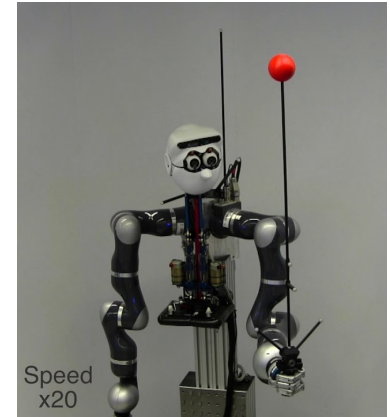
- Goal is often to reach a particular state with minimum effort.
- Goal can be modeled with a cost function similar to

x
if vector of length 2

$$J(u_1, \dots, u_K) = \sum_{k=1}^{K-1} (x_k^T Q x_k + u_k^T R u_k) + x_K^T Q_f x_K$$

for final state

$Q_{xx}x^2 + 2Q_{xy}xy + Q_{yy}y^2$



*quadratic cost
symmetric*

Linear quadratic regulator (LQR)

- Minimize $J(u_1, \dots, u_K) = \sum_{k=1}^{K-1} (x_k^T Q x_k + u_k^T R u_k) + x_K^T Q_f x_K$
for system $x_{k+1} = A x_k + B u_k$
can be solved by value iteration.

- Value function

$$V_t^*(x_t) = \min_{u_t, \dots, u_K} \sum_{k=t}^{K-1} (x_k^T Q x_k + u_k^T R u_k) + x_K^T Q_f x_K$$

- Bellman / Hamilton-Jacobi equation

$$V_t^*(x_t) = x_t^T Q x_t + \min_{u_t} (u_t^T R u_t + V_{t+1}^*(A x_t + B u_t))$$

Solution in the readings for this lecture (there's a minor error, though).

Find the error(s)

$$V_t^*(x_t) = x_t^T Q x_t + \min_{u_t} (u_t^T R u_t + V_{t+1}^*(A x_t + B u_t))$$

$$V_{t+1}^*(x) = x^T P_{t+1} x$$

$$\begin{aligned}
 0 &= \frac{\partial}{\partial u} \left(u_t^T R u_t + \underbrace{(A x_t + B u_t)^T}_{x^T A^T + u^T B^T} P_{t+1} (A x_t + B u_t) \right) \\
 &= 2 u^T R + 2 x^T A^T P_{t+1} B + \underbrace{2 u^T B^T P_{t+1} B}_{2 u^T B^T P B} \\
 &\quad \underbrace{x^T A^T P A x}_{2 x^T A^T P B} + \underbrace{x^T A^T P B u}_{2 x^T A^T P B} + \underbrace{u^T B^T P A x}_{2 x^T A^T P B} + \underbrace{u^T B^T P B u}_{2 u^T B^T P B}
 \end{aligned}$$

Find the error(s)

$$V_t^*(x_t) = x_t^T Q x_t + \min_{u_t} (u_t^T R u_t + V_{t+1}^*(A x_t + B u_t))$$

$$V_{t+1}^*(x) = x^T P_{t+1} x$$

$$0 = \frac{\partial}{\partial u} (u_t^T R u_t + (A x_t + B u_t)^T P_{t+1} (A x_t + B u_t))$$

$$= 2 u^T R + 2 x^T A^T P_{t+1} B + 2 u^T B^T P_{t+1} B$$

$$u = -(R + B^T P B)^{-1} (B^T P A) x$$

Now solve for u!

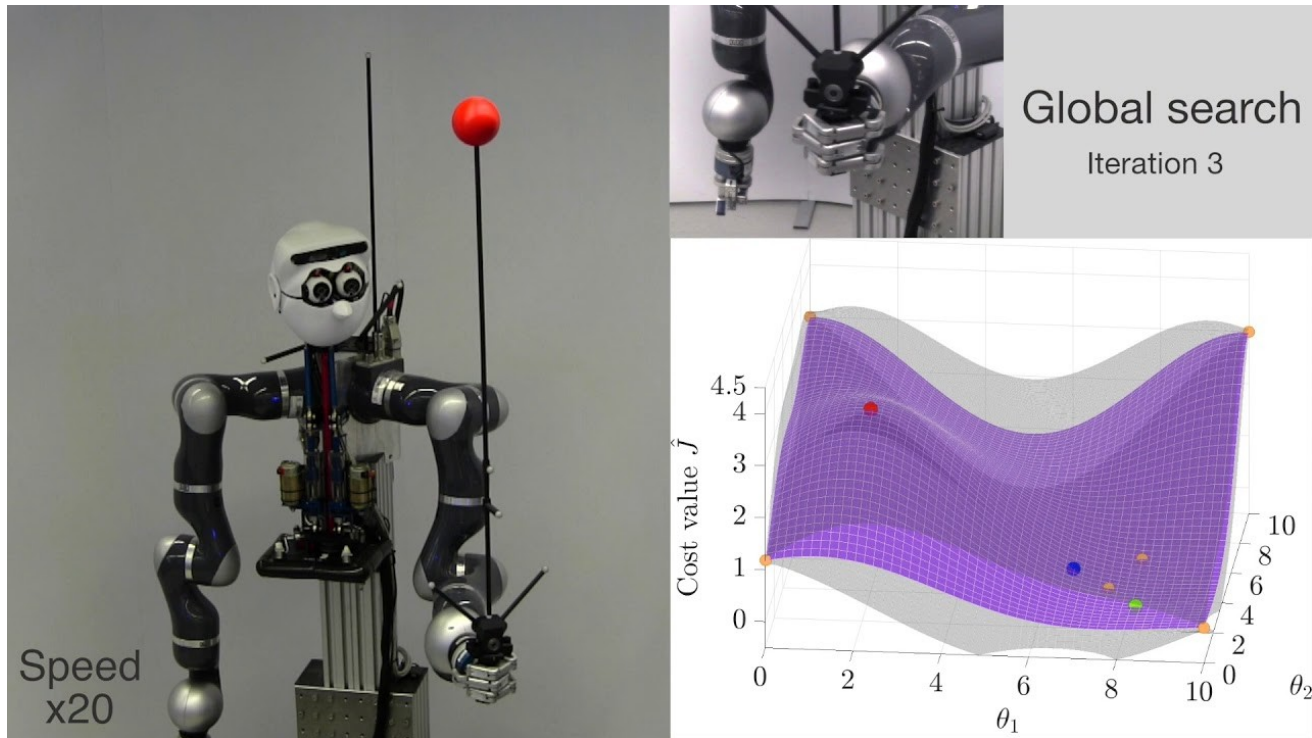
$$\cancel{2} u^T (R + B^T P B) = - \cancel{2} x^T A^T P B$$

$$u^T = - \frac{x^T A^T P B}{R + B^T P B} = - (R + B^T P B)^{-1} x^T A^T P B$$

Note:

$$\frac{\partial}{\partial x} x^T A x = x^T (A + A^T)$$

LQR example



Summary

- Feasible planning can be understood as search.
- Optimal planning can be solved by dynamic programming using a recursive formulation of value function.
- Optimal control of linear systems with quadratic costs can be done using Bellman equation similar to discrete planning.

Next time: Stochastic discrete state environment

- Readings:
 - Sutton & Barto, chapters 2-2.3, 2.5-2.6, 3-3.8 due week from now