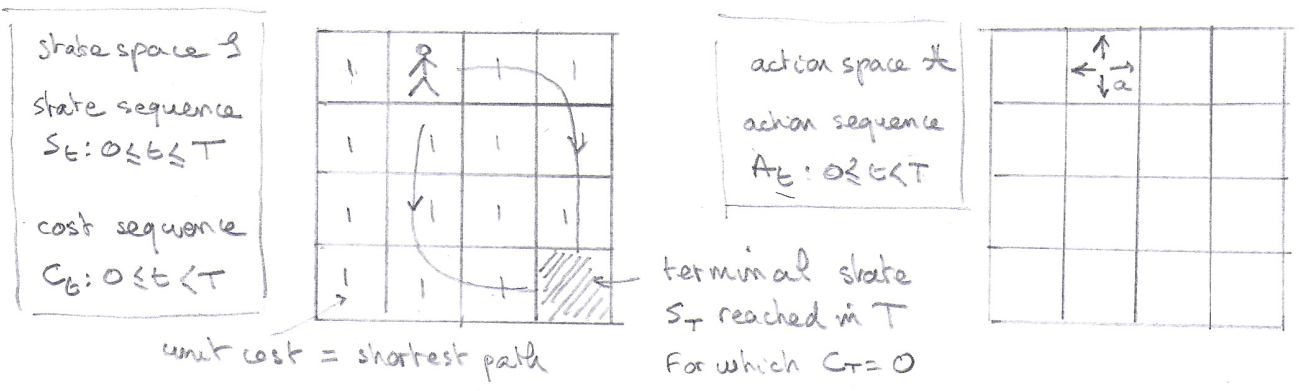


Reinforcement learning (so far)

Example: shortest path problem to terminal state in grid world



- * MDP dynamics is specified by a transition kernel: $P_{s's'}^a = \mathbb{P}[S_{t+1}=s' | S_t=s, A_t=a]$
- * A policy is a decision strategy $\pi(a|s) = \mathbb{P}[A_t=a | S_t=s]$

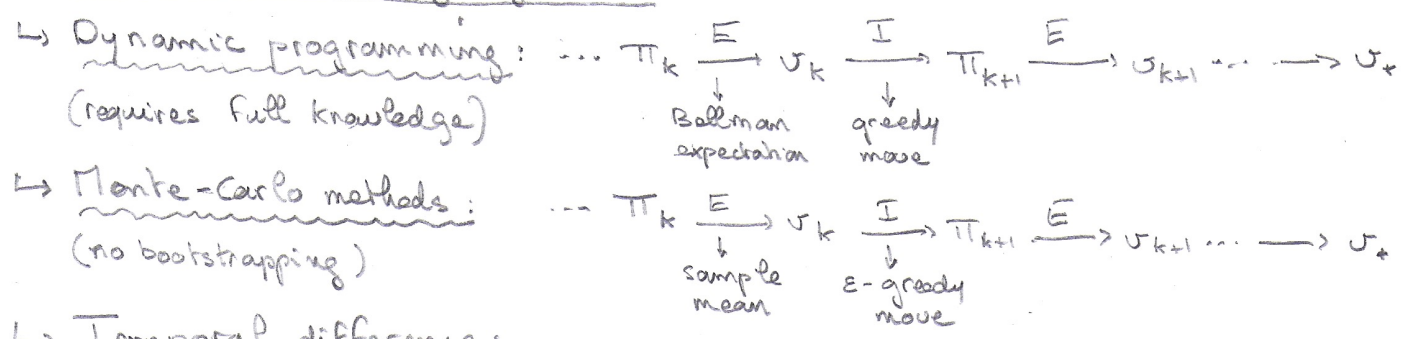
Cost Formulation of the MDP (as in control theory)

The optimal value function (cost-to-go function) minimizes the expected cost (length):
$$J_*(s) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^T C_{t+k+1} \mid S_t=s \right]$$

Characterization via Bellman optimality

$$\begin{aligned}
 J_*(s) &= \min_a \left\{ \mathbb{E}[C_{t+1} | S_t=s, A_t=a] + \mathbb{E} \left[\mathbb{E} \left[\sum_{k=0}^{T-1} C_{t+k+2} \mid S_{t+1} \right] \mid S_t=s, A_t=a \right] \right\} \\
 &= \min_a \left\{ c(s,a) + \mathbb{E} [J_*(S_{t+1}) \mid S_t=s, A_t=a] \right\} = \sum_{s'} P_{s's'}^a J_*(s')
 \end{aligned}$$

Reinforcement learning algorithms



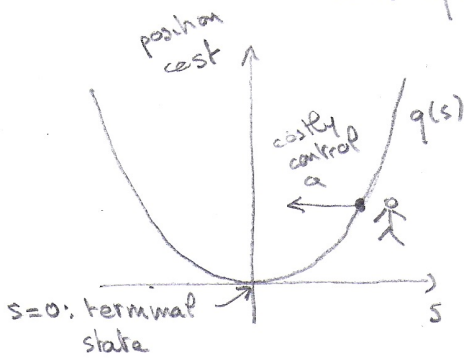
cheapest move from successor

$$q_{t+1}(s_t, A_t) \leftarrow (1-\alpha_t) q_t(s_t, A_t) + \alpha_t [C_{t+1} + \min_a q_t(s_{t+1}, a)]$$

⇒ Computational cost (and learning time) remains fundamentally limited by nonlinearities (presence of the "min")

Analytical solution from control theory

Example: An agent is moving in a continuous state space incurring a cost $q(s)$ at each position. The agent want to minimize its long term cost by exerting control on its dynamics (reaching zero cost state as soon as possible). However exerting control also come at a cost (the faster, the more expensive).



Optimal control theory resolves the above tradeoff.

Controlled dynamics: $\dot{x} = a \leftarrow$ control is a time-dependent function

Infinitesimal cost: $c(s, u) = q(s) + \frac{a^2}{2} = \frac{s^2 + a^2}{2} \leftarrow$ specifies tradeoff

Cost-to-go Function: $V_*(s) = \min_a \left[\int_0^T c(s_t, a_t) dt \right]$, T : time to reach zero.

Bellman optimality: $V_*(s) = \min_a \left[c(s, a) dt + V_*(s + a dt) \right]$

$$V_*(s) = \min_a \left[\frac{s^2 + a^2}{2} dt + V_*(s) + V'_*(s) a dt \right]$$

$$\boxed{-\frac{s^2}{2} = \min_a \left[\frac{a^2}{2} + V'_*(s) a \right]}$$

⚠ Minimization can be carried out explicitly (quadratic function of a)

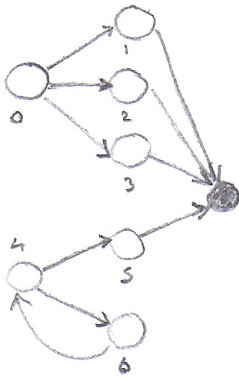
$$a_* = -V'_*(s) \Rightarrow -\frac{s^2}{2} = -\frac{V'_*(s)^2}{2} \Rightarrow \boxed{V_*(s) = \frac{s^2}{2}, s = s_0 e^{-t}}$$

Key ingredients allowing resolution: * special functional form for cost
* minimization in a continuous setting

Expecting full analytical resolution is unreasonable in general. However one can hope to simplify the problem if the two above ingredients can be adapted to MDP settings.

Continuous Framework for MDPs

Example: shortest path on a graph with terminal state.

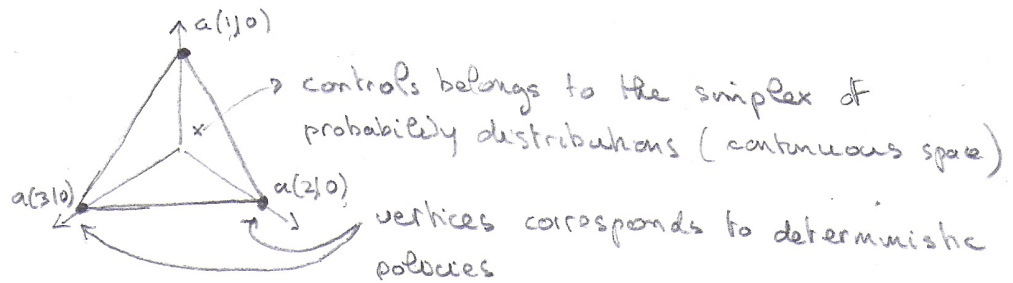


1) First idea: space of control = space of transition probabilities

$$a(s'|s) = P[S_{t+1} = s' | S_t = s, A_t = a] = P_{ss'}^a$$

↓
control can generate any transition probability as long as the transitions are allowed by the graph

For $s=0$, 3 possible transitions to 1, 2, 3.



2) Second idea: passive dynamics = reference dynamics to quantify the cost of controls

Passive dynamics is naturally defined as a random walk on the graph via a transition kernel $p(s'|s) = P_{ss'}^{\emptyset}$ ← no control

3) Third idea: control decision cost = distance from the passive dynamics

$$c(s, a) = \underbrace{q(s)}_{\text{site specific cost}} + \underbrace{D_{KL}(a(\cdot|s) || p(\cdot|s))}_{\substack{\text{control/action cost} \\ = \text{Kullback Leibler divergence}}} = q(s) + \sum_{s'} a(s'|s) \log \frac{a(s'|s)}{p(s'|s)}$$

⇒ Bellman equation:

$$V_*(s) = \min_a \left\{ c(s, a) + \sum_{s'} P_{ss'}^a V_*(s') \right\}$$

$$V_*(s) = q(s) + \min_a \left\{ \sum_{s'} a(s'|s) \left[\log \left(\frac{a(s'|s)}{p(s'|s)} \right) + V_*(s') \right] \right\}$$

Reduction to linear Bellman equation

Key insight: performing a nonlinear change of variable by introducing

the "desirability" function: $z_*(s) = e^{-v_*(s)}$.

Goal: using the well-known KL divergence property to perform the minimization explicitly in Bellman equation.

Bellman equation:
$$-\log z_*(s) = \min_a \left\{ q(s) + \sum_{s'} a(s'|s) \log \left(\frac{a(s'|s)}{p(s'|s) z_*(s')} \right) \right\}$$

almost a KL divergence: normalization!

To normalize we introduce: $Z(s) = \sum_{s'} p(s'|s) z_*(s')$. Then:

$$-\log z_*(s) = \min_a \left\{ q(s) + \sum_{s'} a(s'|s) \log \left(\frac{a(s'|s) Z(s)}{p(s'|s) z_*(s')} \right) - \sum_{s'} a(s'|s) \log Z(s) \right\}$$

\swarrow $D_{KL}(a(\cdot|s) \parallel \frac{p(\cdot|s) z_*(\cdot)}{Z(s)})$ \searrow

Minimization: $D_{KL} = 0 \iff a(\cdot|s) = \frac{p(\cdot|s) z_*(\cdot)}{Z(s)}$

↑
KL divergence

optimal control

Linear Bellman equation:

$$-\log z_* = q - \log Z \iff z_* = e^{-q} Z \iff \begin{cases} z_*(s) = e^{-q(s)} \sum_{s'} p(s'|s) z_*(s') \\ z_*(s) = e^{-q(s)} \mathbb{E} [z_*(S_{t+1}) | S_t = s] \end{cases}$$

linear equation about z_* !

Concretely under matrix form

$$\begin{bmatrix} \underline{z}_T \\ \underline{z}_N \end{bmatrix} = \begin{bmatrix} Id & 0 \\ 0 & e^{-q_N} \end{bmatrix} \begin{bmatrix} Id & 0 \\ \Pi_{NT} & \Pi_{NN} \end{bmatrix} \begin{bmatrix} \underline{z}_T \\ \underline{z}_N \end{bmatrix}$$

\leftarrow terminal nodes $q_T = 0$
 (zero cost) $z_T = 1$

\leftarrow non terminal nodes

$$\left(\begin{bmatrix} e^{q_N} \end{bmatrix} - \Pi_{NN} \right) \underline{z}_N = \Pi_{NT}$$

↑ diagonal matrix

\leftarrow directly solvable equation as opposed to eigenvalue equation (solvable by iteration)

Embedding of classical MDP

Example: shortest path on a graph with terminal state.

In MDP, this corresponds to deliver a unit cost at each site except for the terminal node. Actions are free and deterministic!

By contrast, actions have inherent costs in linear Bellman setting.

Idea: solving MDP via approximation. Specifically let assign $q(s) = p$ for all non terminal nodes. Taking $p \rightarrow +\infty$ makes the action cost small (negligible) compared to the residence cost p , while the linear problem remains solvable:

$$\text{shortest path} = \lim_{p \rightarrow +\infty} \frac{v_+^p(s)}{p} \leftarrow \text{value function obtained for cost } p.$$

How general is this approximation method? Very!

Classical MDP specified by:

$$\Downarrow P_{ss'}^a = P[S_{t+1}=s' | S_t=s, A_t=a]$$

$$\Downarrow c(a,s) = E[C_{t+1} | S_t=s, A_t=a]$$

Linear Bellman Framework specified by:

$$\Downarrow q(s) \leftarrow \text{site specific cost}$$

$$\Downarrow p(s'|s) = P_{ss'}^{\phi} \leftarrow \text{passive dynamics}$$

A MDP can be represented (embedded) in the linear framework if there $q(s), p(s'|s)$ yielding the MDP costs $c(a,s)$ for the control $a(s'|s) = P_{ss'}^a$, i.e.:

$$q(s) + \sum_{s'} P_{ss'}^a \log \left(\frac{P_{ss'}^a}{p(s'|s)} \right) = c(a,s) \rightarrow |X(s)| \text{ equations at fixed } s \text{ about } q(s) \text{ and } p_{s'} = \log p(s'|s)$$

$$q(s) - \sum_{s'} P_{ss'}^a \log p(s'|s) = c(a,s) - \sum_{s'} P_{ss'}^a \log P_{ss'}^a = b(a)$$

In matrix form: $q \underline{1} - P \underline{1} = \underline{b}$ where P is $|X(s)| \times |S|$ matrix with $\sum_{s'} e^{ls'} = \sum_{s'} p(s'|s) = 1 \leftarrow \text{normalization}$

If P is generally row-rank deficient ($|S| \gg |X(s)|$), thus there is \underline{c} such that $P \underline{c} = \underline{b}$. Choosing $q = -\log \sum_{s'} \exp(-c_{s'})$ guarantees that $\underline{1} = q \underline{1} - \underline{c}$ solve the system of equations.

Applications

1) Dynamic programming:

- A) Given a MDP, build a linear embedding
- B) Then optimize the policy by solving linear Bellman equation
- C) Find an approximate optimal strategy for the original MDP by picking the most likely action (greedy move)

2) Monte Carlo method:

The desirability function z_* can be learnt from sampled experiences. The estimator is build recursively from the linear Bellman equation.

$$z_*(s) = e^{-q(s)} \mathbb{E} [z_*(s_1) | s_0=s] = e^{-q(s)} \mathbb{E} [e^{-q(s_1)} \mathbb{E} [z_*(s_2) | s_1=s_1] | s_0=s] \dots$$

$$\hat{z}(s) = \mathbb{E} \left[e^{-\sum_{t=0}^{\infty} q(s_t)} \mid s_0=s \right]$$

 where s_t is sampled according to the passive dynamics

3) Temporal difference method:

A naive α -constant MC scheme (α : learning rate) is

$$\hat{z}(s_t) \leftarrow \hat{z}(s_t) + \alpha [z^*(s_t) - \hat{z}(s_t)] \leftarrow \text{no bootstrapping}$$

↑ target ↑ averaging estimate

The linear Bellman equation suggests a target with bootstrapping

$$z_*(s) = \mathbb{E} [e^{-q(s_t)} z_*(s_{t+1}) \mid s_t=s]$$

↑ expectation = MC averaging ↑ evaluation at successor state = bootstrapping

$$\hat{z}(s_t) \leftarrow \hat{z}(s_t) + \alpha_t (e^{-q(s_t)} \hat{z}(s_{t+1}) - \hat{z}(s_t))$$

 where the learning rates satisfy usual conditions:
 $\sum_{t \geq 0} \alpha_t = +\infty$
 $\sum_{t \geq 0} \alpha_t^2 < +\infty$