

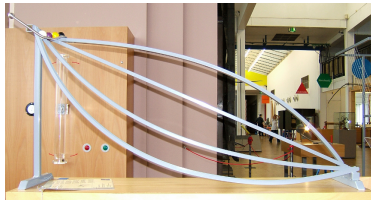
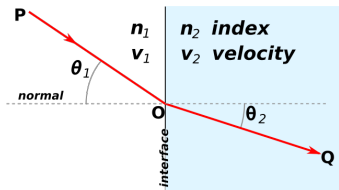
RL 6: The Bellman Equation

Michael Herrmann

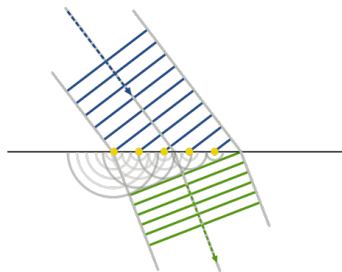
University of Edinburgh, School of Informatics

31/01/2014

Fermat's principle



Huygens–Fresnel principle



Brachistochrone curve
(Johann Bernoulli, 1696)
 \Rightarrow Calculus of variations

x : state

C : scalar cost

V : value

D : value of final state

Value at starting state

$$V(x(0), 0) = \min_u \left\{ \int_0^T C[x(t), u(t)] dt + D[x(T)] \right\}$$

Hamilton-Jacobi-Bellman equation

$$\dot{V}(x, t) + \min_u \{ \nabla V(x, t) \cdot F(x, u) + C(x, u) \} = 0$$

with $\dot{x}(t) = F[x(t), u(t)]$ determining the evolution of the state

N.B. This is copied from wikipedia and included here only for comparison. The important part begins on the next slide.

Bellman's Principle of Optimality: 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. (1957)

The reward hypothesis (Sutton): That all of what we mean by goals and purposes can be well thought of as the maximisation of the cumulative sum of a received scalar signal (reward).

Formulate learning problem such that the principle can be applied.

Value functions: Definition

Policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ or $a \sim \pi(\cdot|s)$ (means: $P(a_t = a) = \pi(a|s_t)$)

Value function $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$

$$V^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s \right]$$

assuming that the initial probability $\pi^0(s_0) > 0$

This is assuming an MDP with fixed $\pi: (\mathcal{S}, P^\pi, \pi^0)$ which is extended to $(\mathcal{S}, P^\pi, \pi^0, \mathcal{R})$. The latter is a *Markov reward process* which arises naturally by assigning a reward distribution $R(\cdot|s)$ to each state s or to each state-action pair according to

$$R^\pi(r|s) = \sum_{a \in \mathcal{A}} \pi(a|s) R(r|s, a)$$

Value functions for state-action pairs

Policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ or $a \sim \pi(\cdot|s)$

Value function $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

$$Q^\pi(s, a) = E \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid s_0 = s, a_0 = a \right]$$

assuming that the initial probability $\pi^0(s_0) > 0$ and that $\pi(s_0) = a_0$ (deterministic) or $\pi(s_0, a_0) > 0$ (stochastic).

First action a_0 is applied now, later actions are chosen by π .

[Note that the initial distribution π^0 and the policy π are different mathematical objects.]

Optimal value functions

For MDPs an optimal policy always exists ($s \in \mathcal{S}$, π fixed)

$$V^*(s) = \sup_{\pi} V^{\pi}(s)$$

For state-action pairs we have ($s \in \mathcal{S}$ and $a \in \mathcal{A}$)

$$\begin{aligned} V^*(s) &= \sup_{a \in \mathcal{A}} Q^*(s, a) \\ Q^*(s, a) &= r(s, a) + \gamma \sum_{u \in \mathcal{S}} P^{a=\pi(s)}(s, u) V^*(u) \end{aligned}$$

Suppose π satisfies

$$\sum_{a \in \mathcal{A}} \pi(a|s) Q^*(s, a) = V^*(s)$$

for all $s \in \mathcal{S}$. Then π is optimal.

Namely, $\pi(\cdot|s)$ selects the action(s) that maximise(s) $Q^*(s, \cdot)$.

So, optimality implies greediness and knowing $Q^*(s, a)$ allows us to act optimally.

Analogously, knowing V^* , r and P suffices to act optimally.

Bellman Equations for deterministic policies in an MDP

How to find the value of a policy?

$$V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{u \in \mathcal{S}} P(s, \pi(s), u) V^\pi(u)$$

This is the Bellman equation for V^π .

Define the Bellman operator for π as $T^\pi : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$ (maps value functions to value functions)

$$(T^\pi V)(s) = r(s, \pi(s)) + \gamma \sum_{u \in \mathcal{S}} P(s, \pi(s), u) V(u)$$

Then naturally,

$$T^\pi V^\pi = V^\pi$$

which is nothing but a compact formulation of the equation on top of this slide. This is a linear equation in V^π and T^π .

If $0 < \gamma < 1$ then T^* is a contraction w.r.t. the maximum norm.

Bellman optimality equations

How to characterise the optimal policy? Use the Bellman optimality principle.

$$V^*(s) = \sup_{a \in \mathcal{A}} \left(r(s, a) + \gamma \sum_{u \in \mathcal{S}} P(s, a, u) V^*(u) \right)$$

Bellman optimality operator $T^* : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$

$$(T^*V)(s) = \sup_{a \in \mathcal{A}} \left(r(s, a) + \gamma \sum_{u \in \mathcal{S}} P(s, a, u) V(u) \right)$$

Then naturally,

$$T^*V^* = V^*$$

which is nothing but a compact formulation of the equation on top of this slide.

If $0 < \gamma < 1$ then T^* is a contraction w.r.t. the maximum norm.

Bellman Operators for state-action value functions

$$T^\pi : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{S} \times \mathcal{A}}, \quad T^* : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$$

$$T^\pi Q(s, a) = r(s, a) + \gamma \sum_{u \in \mathcal{S}} P(s, a, u) Q(u, \pi(s))$$

$$T^* Q(s, a) = r(s, a) + \gamma \sum_{u \in \mathcal{S}} P(s, a, u) \sup_{b \in \mathcal{A}} Q(u, b)$$

T^π is a linear operator, but T^* is not. Both, T^π and T^* are contractions w.r.t. the maximum norm.

Defining $Q(s, \pi(s)) = Q^\pi$ we have $T^\pi Q^\pi = Q^\pi$ and Q^π is the unique solution of this equation. Similarly, we have $T^* Q^* = Q^*$ and Q^* is the unique solution of this equation.

Dynamic programming for solving MDPs

Value iteration

Starting from arbitrary V_0

$$V_{t+1} = T^* V_t$$

global convergence?

For state-action values functions

$$Q_{t+1} = T^* Q_t$$

Once V_t (or Q_t) is close to V^* (or Q^*) then the greedy policy is close to optimal.

Suppose Q is fixed and π is the greedy policy w.r.t. Q . Then

$$V^\pi(s) \geq V^*(s) - \frac{2}{1-\gamma} \|Q - Q^*\|_\infty$$

Singh and Yee, 1994

Dynamic programming for solving MDPs

Policy iteration

Fix an arbitrary initial policy π_0 .

Policy evaluation: At iteration $t > 0$ compute the action-value function underlying π_t

Policy improvement: Given Q^{π_t} define π_{t+1} as the policy that is greedy w.r.t. Q^{π_t} .

Works similar to value iteration, but policy evaluation is computationally more costly.

- *The* Bellman equation is the Bellman optimality equation. It characterises the optimal strategy based on the Bellman optimality principle
- It uses the transition probabilities
- Outlook: Use the actual process to estimate the transition probabilities or to directly sample the value function (or the state-action value function)
- Next: value prediction

- The Bellman (optimality) equation characterises an optimal value function
- In general this equation is not solvable
- Solution is possible by iterative schemes
- Need to take into account the embeddedness of the agent

Temporal Difference (TD) Learning for Value Prediction

Ideal value function

$$\begin{aligned}V_t &= \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r_{\tau} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \\&= r_t + \gamma (r_{t+1} + \gamma r_{t+2} + \dots) \\&= r_t + \gamma \sum_{\tau=t+1}^{\infty} \gamma^{\tau-(t+1)} r_{\tau} \\&= r_t + \gamma V_{t+1}\end{aligned}$$

Real value function is based on estimates of V_t and V_{t+1} , which may not obey this relation. Even if the estimates \hat{V}_t and \hat{V}_{t+1} are far from their true values we can at least require consistency, i.e. minimise the absolute value of the δ error (δ for $\delta\iota\alpha\varphi\omicron\rho\rho\acute{\alpha}$)

$$\delta_{t+1} = r_t + \gamma \hat{V}_{t+1} - \hat{V}_t$$

The simplest TD algorithm

Let \hat{V}_t be the t -th iterate of a learning rule for estimating the value function V .

Let s_t the state of the system at time step t .

$$\delta_{t+1} = r_t + \gamma \hat{V}_t(s_{t+1}) - \hat{V}_t(s_t)$$

$$\hat{V}_{t+1} = \begin{cases} \hat{V}_t(s) + \eta \delta_{t+1} & \text{if } s = s_t \\ \hat{V}_t(s) & \text{otherwise} \end{cases}$$

$$\begin{aligned} \hat{V}_{t+1}(s_t) &= \hat{V}_t(s_t) + \eta \delta_{t+1} = \hat{V}_t(s_t) + \eta (r_t + \gamma \hat{V}_t(s_{t+1}) - \hat{V}_t(s_t)) \\ &= (1 - \eta) \hat{V}_t(s_t) + \eta (r_t + \gamma \hat{V}_t(s_{t+1})) \end{aligned}$$

The update of the estimate \hat{V} is an exponential average over the cumulative expected reward.

Initialise η and γ and execute after each state transition

```
function TD0( $s, r, s1, V$ ) {  
     $\delta := r + \gamma * V[s1] - V[s];$   
     $V[s] := V[s] + \eta * \delta;$   
    return  $V;$   
}
```

Remarks on the TD(0) Algorithm

- If the algorithm converges it must converge to a value function where the expected temporal differences are zero for all states. This state satisfies the Bellman optimality equation.
- The continuous version of the algorithm can be shown to be globally asymptotically stable
- TD(0) is a stochastic approximation algorithm. If the system is ergodic and the learning rate is appropriately decreased, it behaves like the continuous version.

Robbins-Monro conditions

How to choose learning rates? If

$$\sum_{t=0}^{\infty} \eta_t = \infty, \quad \sum_{t=0}^{\infty} \eta_t^2 < \infty,$$

then $V_t(\cdot)$ will behave as the temporally continuous variant

$$\frac{dV(\cdot)}{dt} = r + (\gamma P - I) V(\cdot)$$

Choosing e.g. $\eta_t = c t^{-\alpha}$, the conditions hold for $\alpha \in (\frac{1}{2}, 1]$:

- $\alpha > 1$: goal is reachable even after been temporally trapped
- $\alpha = 1$: smallest step sizes, but still possible
- $\alpha \leq \frac{1}{2}$: large fluctuations can happen even after long time

Iterate-averaging (Polyak & Juditsky, 1992) gives best possible asymptotic rate of convergence

Practically: fixed step sizes or finite-time reduction (see earlier slide)

Many slides are adapted from web resources associated with Sutton and Barto's Reinforcement Learning book

... before being used by Dr. Subramanian Ramamoorthy in this course in the last three years.

... today's lecture followed mostly the book on Algorithms for Reinforcement learning by C. Szepesvari, Ch. 1.