

Introduction, Markov, Bellman Equation

Wen Yu

Departamento de Control Automático
CINVESTAV-IPN

Fundamental challenge in AI and machine learning is:

Learn to make good decisions under uncertainty

- 1 Supervised learning: regression, classification, modeling (Teaching)
- 2 Unsupervised learning: clustering, dimension reduction
- 3 Reinforcement Learning: learn what to do; how to map situations to actions by maximizing a numerical reward signal. (Punishment)

Application of RL

- Video game
- Robotics
- Education
- Health care
- Solve optimization problems

Fundamental problem: Differences with the other machining learnings

- Optimization
- Delayed consequence
- Exploration: experiences
- Generalization: *Policy is mapping from past experience to action, why not just pre-program a policy? We never see before*

History of reinforcement learning

There are two main roads (1980s)

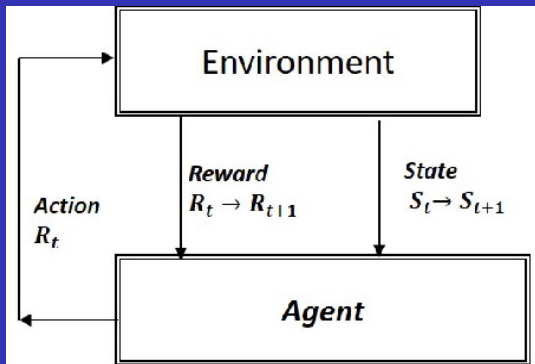
- 1 Concerns learning by trial and error, originated in the psychology of animal learning
- 2 Concerns the problem of optimal control, did not involve learning, dynamic programming

Concerning temporal-difference methods

1989, The temporal-difference and optimal control were fully brought together with: Q-learning

Reinforcement Learning vs other learning methods

- Reinforcement learning: uses training information to evaluate the actions
- Other learning methods: use training information to give correct actions
- Reinforcement learning: active exploration for searching good action
- Other learning methods: the action is the best or the worst

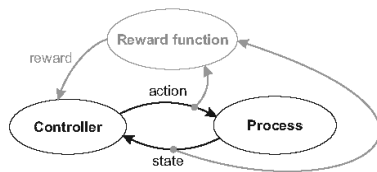


Learning problems facing a decision-maker interacting with its environment

Goal: Select actions to maximize total expected future reward

immediate/long term reward

Reinforcement Learning/dynamic programming vs control



In dynamic programming (DP) and reinforcement learning (RL), a controller (agent, decision maker) interacts with a process (environment), by means of three signals:

- 1 state signal, which describes the state of the process
- 2 action signal, which allows the controller to influence the process
- 3 reward signal, which provides the controller with feedback on its immediate performance.

- R.S. Sutton, A.G. Barto, *Reinforcement Learning: An Introduction*, 2nd edition, The MIT Press, 2018
- L.Busoni, R.Babuska, Bart.Schutter, D.Ernst, *Reinforcement Learning and Dynamic Programming Using Function Approximators Technology*, CRC Press, 2010

Notation in RL/DP vs control

RL vs DP

- RL uses Max/Value, DP uses Min/Cost
- State value \rightarrow State cost.
- Value (or state-value) function \rightarrow Cost function.

Controller

- Agent \rightarrow controller
- Action \rightarrow control
- Environment \rightarrow Dynamic system

The probability of the variable A takes a is

$$0 \leq P(A = a) \leq 1$$

Alternatives -> "add"

$$P(A = a_1 \text{ or } A = a_2) = P(A = a_1) + P(A = a_2)$$

Normalisation

$$\sum_{\text{all possible } a} P(A = a) = 1$$

Joint probability

$P(A = a, B = b)$ the probability that both $A = a$ and $B = b$ occur

Marginalisation

$$P(A = a) = \sum_{\text{all possible } b} P(A = a, B = b)$$

Conditional Probability

$P(A = a \mid B = b)$ the probability that $A = a$ occurs given the knowledge B

Product Rule

$$\begin{aligned} P(A = a, B = b) \\ &= P(A = a) P(B = b \mid A = a) \\ &= P(B = b) P(A = a \mid B = b) \end{aligned}$$

Independence, iff A and B are independent:

$$\begin{aligned} P(A = a \mid B = b) &= P(A = a) \\ P(B = b \mid A = a) &= P(B = b) \\ P(A = a, B = b) &= P(A = a) P(B = b) \end{aligned}$$

Bayes Rule

$$P(A = a | B = b) = \frac{P(B=b|A=a)P(A=a)}{P(B=b)}$$
$$P(A | B) = \frac{P(B|A)P(A)}{P(B)}$$

Continuous variables, "Sum" \rightarrow "Integral", the probability that X lies between x and $(x + dx)$ is $p(x)dx$, $p(x)$ is a probability density function

$$P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} p(x) dx$$
$$\int_{-\infty}^{+\infty} p(x) dx = 1$$
$$\int_{-\infty}^{+\infty} p(x, y) dy = p(x)$$

Expectations: Averages over a time series X

$$E[A] = \sum_a P(A = a) A$$

$$E[X] = \int_{-\infty}^{+\infty} p(x) x dx$$

If X and Y are independent

$$E[X \times Y] = E[X] \times E[Y]$$

Series Theory

Geometric series

$$1 + \alpha + \alpha^2 + \dots = \sum_{i=1}^n \alpha^i = s_n$$
$$\lim_{n \rightarrow \infty} = \frac{1}{1-\alpha}$$

Moving average

$$\bar{X} = \frac{1}{m} \sum_{i=n-m+1}^n x_i$$

Weighted moving average

$$\bar{X} = \frac{1}{\sum_{i=n-m+1}^n \alpha_i} \sum_{i=n-m+1}^n \alpha_i x_i = \frac{1}{\sum_{i=1}^m \alpha_i} \sum_{i=1}^m \alpha_i x_i$$

Exponentially weighted average

$$\bar{X} = \frac{1}{\sum_{i=1}^{\infty} \alpha_i} \sum_{i=1}^{\infty} \alpha_i x_i = (1 - \alpha) \sum_{i=1}^{\infty} \alpha^i x_i$$

Markov Processes (MP)

Sequence of past observations, actions and rewards

$$h_t = \{a_1, o_1, r_1, \dots, a_t, o_t, r_t\}$$

Markov assumption:

- Each event depends only on in the previous event, this process requires the nature of "memoryless"
- Or in order to predict the future, you only need to know the current state of the environment

The s_t is Markov iff

$$P \{s_{t+1} \mid s_t, a_t\} = P \{s_{t+1} \mid h_t, a_t\}$$

Prior examples h_t

No Markov Processes

The current blood pressure is the state, the action is whether to take medication or not, it is Not Markov, because there are many other states to influence the action and also all history blood pressures.

To be Markov

$$s_t = h_t$$

In practice: recent observation

$$s_t = o_t$$

Notation

MP is the sets of states and actions

$$MDP = \{\mathcal{S}, \mathcal{A}, \mathcal{P}\}$$

- \mathcal{P} is random variable defined discrete probability distributions of \mathcal{S} and \mathcal{A}
- s and \bar{s} are values of the random variables S_{t-1} and S_t , $s, \bar{s} \in \mathcal{S}$
- a is the value of the action A_t , $a \in \mathcal{A}(s)$

The process of $S_{t-1} - A_{t-1} \rightarrow S_t - A_t - R_t$ is

$$p\{S_t = \bar{s}, R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$$

or $p(\bar{s}, r \mid s, a)$

where $p(\cdot)$ is the dynamics of the MP, p specifies a probability distribution for each choice of s and a , so

$$\sum_{\bar{s} \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(\bar{s}, r \mid s, a) = 1$$

Components of RL

- Model: how world changes in response to agent's action
- Policy: mapping agent's states to action
- Value function: rewards from state and action

1) Transition or dynamic model, to predict next agent state

$$p(s_{t+1} = s' \mid s_t = s, a_t = a)$$

2) Reward model

$$r(s_t = s, a_t = a) = \mathbb{E}[(r_t \mid s_t = s, a_t = a)]$$

Reward, Return and Value

Horizon: number of time step in each episode, it can be infinite or finite

Return G : Discounted sum of rewards from t to horizon is

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

Value V : sum of future reward r under a particular policy π

$$\begin{aligned} V^\pi(s_t = s) &= \mathbb{E}\{G_t \mid s_t = s\} \\ &= \mathbb{E}_\pi[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \mid s_t = s] \end{aligned}$$

where $0 \leq \gamma \leq 1$ is discount factor

Infinite horizon and continuing task

$$T = \infty$$

$$G_{\infty} = \sum_{i=0}^{\infty} \gamma^k R_{t+1+i}$$

Although the sum of an infinite number of terms, it is still finite if the reward is nonzero and constant, and $\gamma < 1$. For example $R_t = 1$

$$G_{\infty} = \sum_{i=0}^{\infty} \gamma^k = \frac{1}{1 - \gamma}$$

Or $\gamma < 1$, the infinite sum has a finite value as long as the reward sequence R_t is bounded

$\gamma = 0$, immediate rewards

$\gamma = 1$, normal reward.

continuing task

$$G_t = \int_0^t \gamma^{\tau} R_{\tau} d\tau$$

Why discount

- Convenience of mathematical expression, this is also the most important
- Avoid infinite returns in MP
- Uncertainty in long-term rewards
- Immediate returns can get more benefits than delayed rewards
- Human/animal behaviour show immediate rewards are more important

Type of agents

- Model based:
 - 1 Explicit: model
 - 2 may not have policy or value function
- Model-free
 - 1 Explicit: policy or value functions
 - 2 No model

Exploration and Exploitation

Agent only experiences what happens for the action it has tried

- Exploration: try new things that enable the agent to make better decision in the future
- Exploitation: choose actions that have good reward from past experience

Exploration-Exploitation trade-off: Sacrifice reward to explore and learn potentially better policy

Exploration and Exploitation

Watch movies

- Exploration: watch a new movie
- Exploitation: watch a favorite movie which you have seen before

Drive a car

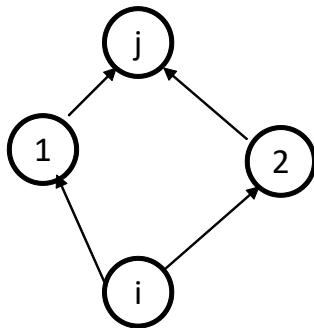
- Exploration: try a different route
- Exploitation: try fastest route given prior experience

- Evaluation: Estimate/predict the rewards from a given policy
- Control (Optimization): find the best policy

Markov Chains

When $n = 2$, 2-step transition probabilities (2 episodes)

$$P_{ij} = P_{i1}P_{1j}P_{i2}P_{2j}$$



Markov Chains: model of the state dynamics in the Markov property.
If it is time invariant, it is stationary, the transition probability from the state i to the state j is

$$P(x_{k+1} = j \mid x_k = i) \\ = P(x_{k+1} = j \mid x_k = i, x_{k-1} = i - 1, \dots, x_0 = 0)$$

n -step transition probability

$$p_{ij}^{(n)} = p(x_{k+n} = j \mid x_k = i) * p(x_n = j \mid x_0 = i)$$

Define transition matrix

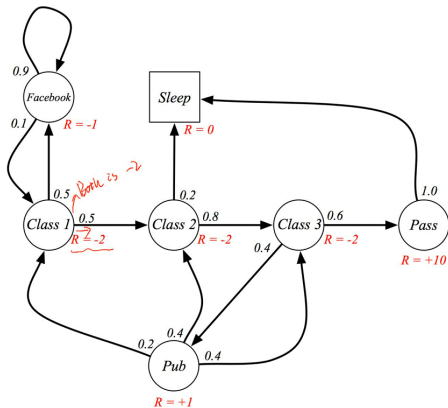
$$P = \begin{bmatrix} p(s_1 | s_1) & \cdots & p(s_N | s_1) \\ p(s_1 | s_2) & p(s_2 | s_2) & \vdots \\ p(s_1 | s_N) & & p(s_N | s_N) \end{bmatrix}$$

$p(s_1 | s_1)$ is the probability of $s_1 \rightarrow s_1$

The next state is s'

$$s' = Ps$$

Example: Student MRP



Example

Return G_t , $\gamma = \frac{1}{2}$. C_1 has 4 returns.

We can see that Class 1 (C_1) actually has 4 paths, Each path has a corresponding probability, $G_t (C_1)$ are

$C_1 C_2 C_3$ Pass Sleep	$G_1 = -2 + (-2)*\frac{1}{2} + (-2)*\frac{1}{4} + 10*\frac{1}{8} + 0*\frac{1}{16} = -2.25$
C_1 FB FB $C_1 C_2$ Sleep	$G_1 = -2 + (-1)*\frac{1}{2} + (-1)*\frac{1}{4} + (-2)*\frac{1}{8} + (-2)*\frac{1}{16} + 0*\frac{1}{32} = -3.125$
$C_1 C_2 C_3$ Pub $C_2 C_3$ Pass Sleep	$G_1 = -2 + (-2)*\frac{1}{2} + (-2)*\frac{1}{4} + (1)*\frac{1}{8} + (-2)*\frac{1}{16} + \dots = -3.41$
C_1 FB FB $C_1 C_2 C_3$ Pub C_1 FB FB FB $C_1 C_2 C_3$ Pub C_2 Sleep	$G_1 = -2 + (-1)*\frac{1}{2} + (-1)*\frac{1}{4} + (-2)*\frac{1}{8} + (-2)*\frac{1}{16} + (-2)*\frac{1}{32} + \dots = -3.20$

The value function needs "expected return" when evaluating the value of the state C_1 .

$$V(C_1) = \frac{1}{4} [(-2.25 + (-3.125) + (-3.41) + (-3.20))] = 2.996$$

Computing the value function

Average return

$$\begin{aligned} V^\pi(s_t = s) &= \mathbb{E}\{G_t \mid s_t = s\} \\ &= \mathbb{E}_\pi[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \mid s_t = s] \end{aligned}$$

Large number of episodes

Markov property yields additional structure:

Define the total discounted reward from time t as the return G_t ,

$$G_t = \sum_{i=0}^N \gamma^i r_{t+i}$$

where $0 \leq \gamma \leq 1$.

Recursive form

$$G_t = r_t + \gamma G_{t+1}$$

Bellman equation for MRP(Markov Reward Process)

The state-value function $V(s)$ is defined as

$$V(s_t) = \mathbb{E} \{G_t \mid S_t = s\}$$

Because the recursive form $G_t = r_t + \gamma G_{t+1}$,

$$\begin{aligned} V(s_t) &= \mathbb{E} \{G_t \mid S_t = s\} \\ &= \mathbb{E} \{r_t + \gamma G_{t+1} \mid s\} \\ &= \mathbb{E} \{r_t + \gamma V(S_{t+1}) \mid s\} \\ &= r_t + \gamma \mathbb{E} \{V(s_{t+1})\} \end{aligned}$$

where r_t is immediate reward, $\gamma \mathbb{E} \{V(s_{t+1})\}$ is discounted sum of future rewards

Bellman equation for MRP

Bellman equation is

$$V(s) = r(s) + \gamma \mathbb{E} \{ V(s') \}$$

where s' is the state in next time.

Because the expectations: Averages over a time series X is

$$E[A] = \sum_a p(A = a) A$$

The value expectation of next state $\mathbb{E} \{ V(s') \}$ can be obtained according to the probability distribution of s_t

$$\mathbb{E} \{ V(s') \} = \gamma \sum p(s' | s) V(s')$$

Bellman equation is

$$V(s) = r(s) + \gamma \sum p(s' | s) V(s')$$

Bellman equation for MRP

Bellman equation in matrix form

$$\begin{bmatrix} V(s_1) \\ \vdots \\ V(s_N) \end{bmatrix} = \begin{bmatrix} R(s_1) \\ \vdots \\ R(s_N) \end{bmatrix} + \gamma \begin{bmatrix} P(s'_1 | s_1) & \dots & P(s'_N | s_1) \\ P(s'_1 | s_2) & P(s'_2 | s_2) & \vdots \\ P(s'_1 | s_N) & & P(s'_N | s_N) \end{bmatrix} \begin{bmatrix} V(s_1) \\ \vdots \\ V(s_N) \end{bmatrix}$$

For infinite horizon, the value function is stationary,

$$\begin{bmatrix} V(s_1) \\ \vdots \\ V(s_N) \end{bmatrix} = \begin{bmatrix} R(s_1) \\ \vdots \\ R(s_N) \end{bmatrix} + \gamma \begin{bmatrix} P(s_1 | s_1) & \dots & P(s_N | s_1) \\ P(s_1 | s_2) & P(s_2 | s_2) & \vdots \\ P(s_1 | s_N) & & P(s_N | s_N) \end{bmatrix} \begin{bmatrix} V(s_1) \\ \vdots \\ V(s_N) \end{bmatrix}$$

So

$$V = R + \gamma PV$$
$$V = (I - \gamma P)^{-1} R$$

The matrix inverse require, computational complexity $O(N^3)$, N states

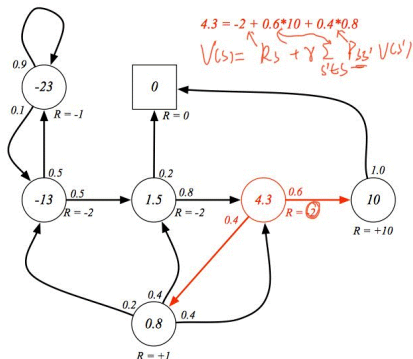
Value function calculation with Bellman equation

$$V(s) = r(s) + \gamma \sum p(s' | s) V(s')$$

when $\gamma = 1$

$$V(C3) = -2 + 0.6 * 10 + 0.4 * 0.8 = 4.3$$

Example: Bellman Equation for Student MRP



Markov Decision Processes (MDP)

MDP is the sets of states, actions, and rewards

$$MDP = \{S, A, \mathcal{R}, \mathcal{P}\}$$

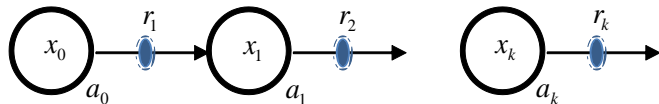
r is the value of the random variable \mathcal{R} , $r \in \mathcal{R}$

Because Markov property

$$p(s', r | s, a) = p\{s', r | S_{t-1} = s, A_{t-1} = a\}$$

The probability of each possible value for s_t , r_t depends only on the immediately preceding state s_{t-1} , a_{t-1}

Markov decision process



Policy at step t : a mapping from states to action probabilities

$$\pi_k : x_k \rightarrow a_k$$

or the probability that $a_k = a$ when $x_k = s$

Markov decision process

	Without Action	Consider Action (A)
Completely observable state	Markov Chain (MC)	Markov decision process (MDP)
Partially observable states	Hidden Markov Model (HMM)	Partially observable MDP (POMDP)

- Factored MDPs: add planning graph style heuristics use goal regression to generalize better
- Hierarchical MDP: hierarchy of sub-tasks and/or actions to scale better
- Partially Observable Markov Decision Processes, noisy sensors, partially observable environment, popular in robotics

Deterministic policy

$$\pi(s) = a$$

Stochastic policy

$$\pi(a | s) = (A_t = a | S_t = s)$$

Reinforcement learning: how the agent changes its policy π_k (changing the transition probabilities), such that a better MDP (with lower cost/higher expected reward) is reached.

MRP(Markov Reward Process)

$$MRP(S, R^\pi, P^\pi, \gamma) : MDP + \pi(a | s)$$

Because

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \{G | s\} = \mathbb{E}_\pi \{r_t + \gamma V^\pi(s') | s\} \\ &= (\sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a)) [r_t + \gamma V^\pi(s')] \end{aligned}$$

where $\sum_{s'} \sum_r$ are sums over all the possible values, they can be merged to $\sum_{s', r}$.

$$\begin{aligned} V^\pi(s) &= \sum_{a, s', r} \pi(a | s) p(s', r | s, a) [r_t + \gamma V^\pi(s')] \\ &= \mathbb{E}_\pi \{ \pi(a | s) [r + \gamma V^\pi(\bar{s})] \} \end{aligned}$$

where

$$\begin{aligned} \mathbb{E}_\pi(r^\pi(s)) &= \sum_{a \in \mathcal{A}} p(a | s) r(s, a) \\ \mathbb{E}_\pi(p^\pi(s' | s)) &= \sum_{a \in \mathcal{A}} \pi(a | s) p(s' | s, a) \end{aligned}$$

Policy evaluation

Value of a policy

Initialize $V_0(s) = 0$

For $k = 1$ until convergence ($|V_k - V_{k-1}| < \epsilon$)

For all s

$$V_k^\pi(s) = r_t(s, \pi(s)) + \gamma \sum p(s' | s, \pi(s)) V_{k-1}^\pi(s')$$

end

This is Bellman backup for a particular policy

Optimal policy

$$\pi^*(s) = \max_{\pi} [V^{\pi}(s)] \text{ for all } s \in \mathcal{S}$$

$v^*(s)$ is There exists solution

Optimal policy for a MDP in an infinite horizon problem is deterministic

Policy Iteration (PI)

Initialize $\pi_0(s)$ randomly for all states s

While $i == 0$ or the policy changes for any state ($|\pi_k - \pi_{k-1}| > 0$)

Policy Evaluation

$$MRP \rightarrow V_i^\pi$$

Policy Improvement

$$Policy \rightarrow \pi_{i+1}$$

$$i = i + 1$$

State-Action function: Q-function

This form can be also obtained from the discounted return

$$Q^{\pi}(s, a) = r(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^{\pi}(s')$$

Compute $Q^\pi(s, a)$ with a policy π_k

$$Q^{\pi_k}(s, a) = r(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^{\pi_k}(s') = V^\pi(s)$$

Compute new policy π_{k+1} for all $s \in S$

$$\pi_{k+1}(s) = \max_a Q^{\pi_k}(s, a), \quad \forall s \in S$$

so

$$\pi_{k+1}(s) \geq Q^{\pi_k}(s, \pi_k(s))$$

Optimal policies

For finite MDPs, an optimal policy is: a policy π is better than or equal to a policy $\bar{\pi}$, if its expected return is greater than or equal to that of $\bar{\pi}$ for all states

$$\pi \geq \pi', \text{ if and only if } V_{\pi}(s) \geq V_{\pi'}(s) \text{ for all } s$$

Monotonic improvement in policy

Theorem

For any MDP, the following points are true:

- There is an optimal strategy that is better or at least equal to any other strategy
- All optimal strategies have the same optimal value function
- All optimal strategies have the same value function.

If a policy does not change, it will never change again

There a maximum number of interaction of policy iteration

Bellman optimality equation

$V^*(s)$ also satisfies Bellman equation,

$$\begin{aligned} V^*(s) &= \max_{\pi} [V^{\pi}(s)] \text{ for all } s \in \mathcal{S} \\ &= \max_a \mathbb{E}_{\pi} \{r + \gamma V^*(s')\} \\ &= \max_a \sum_{\bar{s}, r} p(s', r | s, a) [r + \gamma V^*(s')] \end{aligned}$$

Bellman optimality equations

$$\begin{aligned} Q^*(s, a) &= \max \mathbb{E}_{\pi} \{r(s, a) + \gamma Q(s', a')\} \\ &= \mathbb{E}_{\pi} \{\max r(s, a) + \gamma \max_{a'} Q(s', a')\} \\ &= \sum_{\bar{s}, r} p(\bar{s}, r | s, a) [\max r(s, a) + \gamma \max_{a'} Q(s', a')] \\ &= \max r(s, a) + \gamma \sum p(s', r | s, a) V_{\pi}^*(s') \end{aligned}$$

Value Iteration (VI)

Initialize $V_0(s) = 0$ randomly for all states s

Loop until finite horizon or convergence

for each state

$$V_{k+1}(s) = \max_a r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V_k(s')$$

Bellman backup on value function

$$\pi_{k+1}(s) = \arg \max_a \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V_k(s') \right\}$$