

# 强化学习与博弈论

Reinforcement Learning and Game Theory

陈旭

计算机学院



中山大學  
SUN YAT-SEN UNIVERSITY

---

# **Chapter 1: Simple Decisions Models**

# Some definitions

## Definition

*argmax* is defined by the following equivalence:

$$x^* \in \operatorname{argmax}_{x \in X} f(x) \iff f(x^*) = \max_{x \in X} f(x)$$

## Definition

A choice of behavior in a single-decision problem is called an *action*. The set of alternative actions available will be denoted as **A**. This will either be discrete set, e.g.,  $\{a_1, a_2, \dots, \}$ , or a continuous set, .e.g., the unit interval  $[0, 1]$ .

# More...

## Definition

A *payoff* is a function  $\pi : \mathbf{A} \rightarrow R$  that associates a numerical value with every action  $a \in \mathbf{A}$ .

## Definition

An action  $a^*$  is an *optimal action* if

$$\pi(a^*) \geq \pi(a) \quad \forall a \in \mathbf{A}.$$

or equivalently,  $a^* \in \operatorname{argmax}_{a \in \mathbf{A}} \pi(a)$ .

## Definition

An *affine transformation* changes payoff  $\pi(a)$  into  $\pi'(a)$  as

$$\pi'(a) = \alpha\pi(a) + \beta$$

where  $\alpha, \beta$  are constants independent of  $a$  and  $\alpha > 0$ .

## Theorem

*The optimal action is unchanged if payoffs are altered by an affine transformation.*

## Proof

because  $\alpha > 0$ , we have

$$\begin{aligned} \operatorname{argmax}_{a \in \mathbf{A}} \pi'(a) &= \operatorname{argmax}_{a \in \mathbf{A}} [\alpha\pi(a) + \beta] \\ &= \operatorname{argmax}_{a \in \mathbf{A}} \pi(a). \end{aligned}$$

## Example

The Convent Fields Soup Company needs to determine the price  $p$ . The demand function is:

$$Q(p) = \begin{cases} Q_0 \left(1 - \frac{p}{p_0}\right) & \text{if } p < p_0, \\ 0 & \text{if } p \geq p_0. \end{cases}$$

The payoff is  $\pi(p) = (p - c)Q(p)$  where  $c$  is the unit production cost.

- Solving, we have  $p^* = \frac{1}{2}(p_0 + c)$ .
- Now, let say we need to consider a fixed cost to build the factory, the payoff function is  $\pi(p) = (p - c)Q(p) - B$ , where  $B$  is a constant. What is  $p^*$ ?

# Uncertainty

## Modeling uncertainty

- If uncertainty exists, we compare the expected outcome for each action.
- Let  $X$  be the set of states with  $P(X = x)$ .
- Payoff for adopting action  $a$  is:

$$\pi(a) = \sum_{x \in X} \pi(a|x)P(X = x)$$

- An optimal action is

$$a^* \in \operatorname{argmax}_{a \in A} \sum_{x \in X} \pi(a|x)P(X = x).$$



## Example

- An investor has \$1000 to invest in one year. The available actions (1) put the money in the bank with 7% interest per year; (2) invest in stock which returns \$1500 if the stock market is good or returns \$600 if the stock market is bad.
- $P(\text{Good}) = P(\text{Bad}) = 0.5$ .
- Expected payoff:
  - ①  $\pi(a_1) = \$1070$ ;
  - ②  $\pi(a_2) = 1500/2 + 600/2 = \$1050$ .
- So  $a_1^*$  and we should put the money in the bank.

## Definition

Let  $\Omega = \{\omega_1, \omega_2, \dots\}$  be the set of possible outcomes.

- We say  $\omega_i \succ \omega_j$  if an individual *strictly prefers* outcome  $\omega_i$  over  $\omega_j$ .
- If the individual is indifferent:  $\omega_i \sim \omega_j$ .
- Either prefer or indifferent:  $\omega_i \succeq \omega_j$ .

## Definition

An individual will be called **rational under certainty** if his preference for outcomes satisfy the following conditions:

- (Completeness) Either  $\omega_i \succeq \omega_j$  or  $\omega_j \succeq \omega_i$ .
- (Transitivity) If  $\omega_i \succeq \omega_j$  and  $\omega_j \succeq \omega_k$ , then  $\omega_i \succeq \omega_k$ .

## Definition

A **utility function** is a function  $u : \Omega \rightarrow \mathbf{R}$  such that

$$u(\omega_i) > u(\omega_j) \iff \omega_i \succ \omega_j$$

$$u(\omega_i) = u(\omega_j) \iff \omega_i \sim \omega_j$$

The immediate consequence of this definition is an individual who is rational under certainty should seek to maximize his utility.

What happens when an action does not produce a definite outcome and instead, we allow each outcome to occur with a known probability?

### Definition

A **simple lottery**,  $\lambda$ , is a set of probabilities for the occurrence of every  $\omega \in \Omega$ . The probability that outcome  $\omega$  occurs is  $p(\omega|\lambda)$ . The set of all possible lotteries is denoted as **A**.

## Theorem

**Expected Utility Theorem:** *If an individual is rational, then we can define a utility function  $u : \Omega \rightarrow \mathbf{R}$  and the individual will maximize the payoff function  $\pi(a)$  (or the expected utility) given by*

$$\pi(a) = \sum_{\omega \in \Omega} p(\omega | \lambda(a)) u(\omega)$$

## Definition

An individual whose utility function satisfies

- $E(u(w)) < u(E(w))$ , it is said to be **risk averse**,
- $E(u(w)) > u(E(w))$ , it is said to be **risk prone**,
- $E(u(w)) = u(E(w))$ , it is said to be **risk neutral**.

## Example

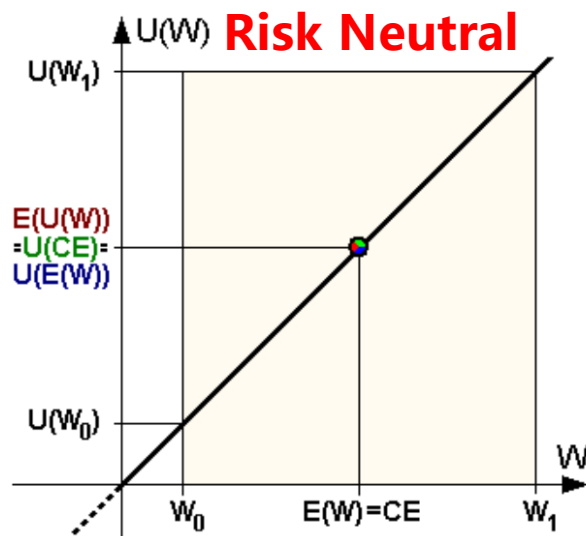
Someone flip a coin. If it is head (tail), you get \$1 (\$1M). Your utility function can be:

- $u(x) = x$ ,
- $u(x) = x^2$ ,

Classify the above as risk averse, risk prone and risk neutral utility function.

[https://en.wikipedia.org/wiki/Risk\\_aversion](https://en.wikipedia.org/wiki/Risk_aversion)

## Concave



Up to now, we assume finding an optimal action  $a^*$  from a given set  $\mathbf{A}$ . But the selection can be *randomized*. Does this allow one to achieve a higher payoff?

## Definition

We specify a **general behavior**  $\beta$  by giving a list of probabilities with which each available action is chosen. We denote the probability that action  $a$  is chosen by  $p(a)$  and  $\sum_{a \in \mathbf{A}} p(a) = 1$ . The set of all randomizing behavior is denoted by  $\mathbf{B}$ . The payoff of using behavior  $\beta$  is

$$\pi(\beta) = \sum_{a \in \mathbf{A}} p(a)\pi(a).$$

An **optimal behavior**  $\beta^*$  is one for which

$$\pi(\beta^*) \geq \pi(\beta) \quad \forall \beta \in \mathbf{B}.$$

or  $\beta^* \in \operatorname{argmax}_{\beta \in \mathbf{B}} \pi(\beta)$ .



## Definition

The support of a behavior  $\beta$  is the set  $\mathbf{A}(\beta) \subseteq \mathbf{A}$  of all the actions for which  $\beta$  specifies  $p(a) > 0$ .

## Theorem

*Let  $\beta^*$  be an optimal behavior with support  $\mathbf{A}^*$ . Then*

$$\pi(a) = \pi(\beta^*) \quad \forall a \in \mathbf{A}^*.$$

The consequence of the above theorem is that if a randomized behavior is optimal, then two or more actions are optimal as well. So randomization is not necessary but it may be used to break a tie.

## Example

A firm may make one of the marketing actions  $\{a_1, a_2, a_3\}$ . The profit for each action depends on the state of the economy  $\mathbf{X} = \{x_1, x_2, x_3\}$ :

	$x_1$	$x_2$	$x_3$
$a_1$	6	5	3
$a_2$	3	5	4
$a_3$	5	9	1

If  $P(X = x_1) = 1/2$ ,  $P(X = x_2) = P(X = x_3) = 1/4$ . What *are* the optimal behaviors?

## Example

A firm may make one of the marketing actions  $\{a_1, a_2, a_3\}$ . The profit for each action depends on the state of the economy  $\mathbf{X} = \{x_1, x_2, x_3\}$ :

	$x_1$	$x_2$	$x_3$
$a_1$	6	5	3
$a_2$	3	5	4
$a_3$	5	9	1

If  $P(X = x_1) = 1/2$ ,  $P(X = x_2) = P(X = x_3) = 1/4$ . What *are* the optimal behaviors?

## Answer

Because  $\pi(a_1) = \pi(a_3) = 5$  and  $\pi(a_2) = 3.75$ , optimal randomizing behaviors have support  $\mathbf{A}^* = \{a_1, a_3\}$  with  $p(a_1) = p$  and  $p(a_3) = 1 - p$  ( $0 < p < 1$ ). Using either  $a_1$  or  $a_3$  with probability 1 is also an optimal behavior.