

$$\begin{cases} \dot{x} = f(x(t), u(t)) & t > 0 \\ x(0) = x_0 \end{cases}$$

$$u(t): [0, +\infty) \rightarrow \mathbb{R}^m$$

$$J[u(\cdot)] = \int_0^T r(x(t), u(t)) dt + g(x(T))$$

PMP:

$$(ODE) \quad \dot{x}^*(t) = \nabla_p H(x^*(t), p^*(t), u^*(t))$$

$$(ADI) \quad \dot{p}^*(t) = -\nabla_x H(x^*(t), p^*(t), u^*(t))$$

$$(PMP) \quad H(x^*(t), p^*(t), u^*(t)) = \max_{u \in A} H(x^*(t), p^*(t), u)$$

(hp 5 动态规划 (Dynamic programming))

核心思想: 把一个控制问题放入到一族“从任意时刻、任意状态出发”的问题中。再用这一族的最优函数刻画最优控制。

Hamilton-Jacobi-Bellman 方程 (HJB 方程, Bellman PDE)

$$\text{考虑} \quad \int_0^{\infty} \frac{\sin x}{x} dx$$

$$\text{定义参数化积分} \quad I(\alpha) = \int_0^{\infty} e^{-\alpha x} \frac{\sin x}{x} dx \quad \alpha > 0$$

$$\begin{aligned} \text{对 } \alpha \text{ 求导} \quad I'(\alpha) &= \frac{d}{d\alpha} \int_0^{\infty} e^{-\alpha x} \frac{\sin x}{x} dx \\ &= \int_0^{\infty} \frac{d}{d\alpha} e^{-\alpha x} \frac{\sin x}{x} dx = \int_0^{\infty} (-x) e^{-\alpha x} \frac{\sin x}{x} dx \end{aligned}$$

$$-I(\alpha) = -\int_0^{\infty} e^{-\alpha x} \sin x dx$$

$$\text{令 } u = e^{-\alpha x}, \quad du = -\alpha e^{-\alpha x} dx \quad \text{则} \quad du = -\alpha e^{-\alpha x} dx, \quad v = -\cos x$$

$$\int u dv = uv - \int v du$$

$$J(\alpha) = \int_0^{\infty} e^{-\alpha x} \sin x dx \quad J(\alpha) = 1 - \alpha I(\alpha)$$

$$= \left[ e^{-\alpha x} (-\cos x) \right] \Big|_0^{\infty} - \int_0^{\infty} (-\cos x) (-\alpha e^{-\alpha x}) dx$$

$$= \left[ -e^{-\alpha x} \cos x \right] \Big|_0^{\infty} - \alpha \int_0^{\infty} e^{-\alpha x} \cos x dx$$

$$= 0 - (-1) - \alpha \int_0^{\infty} e^{-\alpha x} \cos x dx$$

$$I(\alpha) = \int_0^{\infty} e^{-\alpha x} \cos x dx \quad I(\alpha) = \alpha J(\alpha)$$

$$= \left[ e^{-\alpha x} \sin x \right] \Big|_0^{\infty} + \alpha \int_0^{\infty} e^{-\alpha x} \sin x dx$$

$$= \alpha \int_0^{\infty} e^{-\alpha x} \sin x dx$$

联立可得  $J(\alpha) = 1 - \alpha(\alpha J(\alpha)) = 1 - \alpha^2 J(\alpha)$

$$\therefore J(\alpha) = \frac{1}{1+\alpha^2}$$

因此  $I'(\alpha) = -\frac{1}{1+\alpha^2}$

$$I(\alpha) = -\arctan \alpha + C$$

$$I(\alpha) = \int_0^{\infty} e^{-\alpha x} \frac{\sin x}{x} dx \quad \alpha > 0$$

当  $\alpha \rightarrow \infty$  时, 有  $I(\alpha) \rightarrow 0$ , 且  $\arctan \alpha \rightarrow \frac{\pi}{2}$

$$I(\infty) = -\frac{\pi}{2} + C = 0 \Rightarrow C = \frac{\pi}{2}$$

因此  $\int_0^{\infty} \frac{\sin x}{x} dx = I(0) = \frac{\pi}{2}$

**Remark:** 通过引入参数  $\alpha$ , 对  $\{I(\alpha)\}$  做整体分析, 更容易得到原来的值.

**动态规划:** 不只看一个固定的初始状态和初始时间, 而是考虑所有起始点  $(x, t)$ .

$$\text{考虑} \quad \begin{cases} \dot{x}(t) = f(x(t), u(t)) & t \in [0, T] \\ x(0) = x_0 \end{cases}$$

控制  $u(\cdot) \in A \subset \mathbb{R}^m$   $A$  是紧集

$$\text{收益函数} \quad J[u(\cdot)] = \int_0^T r(x(t), u(t)) dt + g(x(T))$$

扩展到所有起始时刻  $t$  的子问题

$$\begin{cases} \dot{x}(s) = f(x(s), u(s)) & s \in [t, T] \\ x(t) = x \in \mathbb{R}^n \end{cases}$$

$$\text{收益函数} \quad J_{x,t}[u(\cdot)] = \int_t^T r(x(s), u(s)) ds + g(x(T))$$

Defn (价值函数, value function) 对每一个起点  $(x, t) \in \mathbb{R}^n \times [0, T]$ , 价值函数为

$$v(x, t) = \sup_{u(\cdot) \in \mathcal{U}} J_{x,t}[u(\cdot)] = \sup_{u(\cdot) \in \mathcal{U}} \left\{ \int_t^T r(x(s), u(s)) ds + g(x(T)) \right\}$$

当  $t = T$  时,

$$v(x, T) = g(x).$$

Hamilton-Jacobi-Bellman 方程

假设: ① 控制集  $A$  为紧集

② 对每个  $(x, t)$ ,  $\sup_{u(\cdot) \in \mathcal{U}} J_{x,t}[u(\cdot)]$  可得 (存在最优控制)

③  $v \in C^1$  ( $v$ -阶连续可微)

Thm (HJB 方程) 基于上述假设, 价值函数  $v$  满足

$$v_t(x, t) + \max_{\alpha \in A} \{ f(x, \alpha) \cdot \nabla_x v(x, t) + r(x, \alpha) \} = 0$$

且  $v(x, T) = g(x)$ .

Remark: 对应 Hamiltonian

$$H(x, p) := \max_{\alpha \in A} \{ f(x, \alpha) \cdot p + r(x, \alpha) \}$$

$$\text{HJB} \Rightarrow v_t(x, t) + H(x, \nabla_x v(x, t)) = 0$$

Proof: 给定任意  $(x, t)$ , 取步长  $h > 0$ , 且  $t+h < T$ . 在区间  $[t, t+h]$  取一个常控制, 即

$$\alpha(s) \equiv a \quad s \in [t, t+h]$$

对应状态轨迹

$$\begin{cases} \dot{x}(s) = f(x(s), a) \\ x(t) = x \end{cases}$$

在  $s \in [t+h, T]$  采用最优控制, 最优收益为

$$v(x(t+h), t+h)$$

整体策略总收益为

$$\int_t^{t+h} r(x(s), a) ds + v(x(t+h), t+h)$$

有

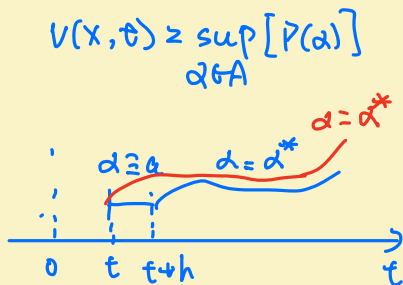
$$v(x, t) \geq \int_t^{t+h} r(x(s), a) ds + v(x(t+h), t+h) \quad (1)$$

令  $h \rightarrow 0$ , (1) 式可写为

$$\frac{1}{h} \int_t^{t+h} r(x(s), a) ds + \frac{v(x(t+h), t+h) - v(x, t)}{h} \leq 0$$

由于  $\lim_{h \rightarrow 0} \frac{v(x(t+h), t+h) - v(x, t)}{h} = \nabla_x v(x, t) \cdot \dot{x}(t) + v_t(x, t)$

$$\lim_{h \rightarrow 0} \frac{1}{h} \int_t^{t+h} r(x(s), a) ds = r(x, a)$$



代回原式, 可得

$$\dot{x}(t) = f(x, a)$$

$$r(x, a) + \nabla_x V(x, t) \cdot f(x, a) + V_t(x, t) \leq 0$$

$$\Rightarrow r(x, a) + \nabla_x V(x, t) \cdot f(x, a) + V_t(x, t) \leq 0 \quad \forall a \in A$$

$$\Rightarrow \max_{a \in A} \{ r(x, a) + \nabla_x V(x, t) \cdot f(x, a) + V_t(x, t) \} \leq 0 \quad (2)$$

令  $\alpha^*$  是从  $(x, t)$  出发的最优控制, 考虑  $SG[t, t+h]$ .

$$V(x, t) = \int_t^{t+h} r(x^*(s), \alpha^*(s)) ds + V(x^*(t+h), t+h)$$

类似的, 除以  $h$  并令  $h \rightarrow 0$ , 可得

$$V_t(x, t) + \nabla_x V(x, t) \cdot f(x, \alpha^*(t)) + r(x^*, \alpha^*(t)) = 0$$

即存在  $\alpha^* = \alpha^*(t) \in A$  使得.

$$V_t(x, t) + \nabla_x V(x, t) \cdot f(x, \alpha^*) + r(x^*, \alpha^*) = 0$$

那么, 可以找到一最优值, 使得.

$$\max_{a \in A} \{ V_t(x, t) + \nabla_x V(x, t) \cdot f(x, a) + r(x, a) \} \geq 0 \quad (3)$$

综合 (2)、(3), 可得

$$\max_{a \in A} \{ V_t(x, t) + \nabla_x V(x, t) \cdot f(x, a) + r(x, a) \} = 0$$

### 动态规划求最优控制的两步法

简单来说: 先解 HJB 方程, 再构造控制

Step 1: 求解 HJB 方程, 得到  $v$  的表达式:

$$\begin{cases} V_t(x, t) + \max_{a \in A} \{ \nabla_x V(x, t) \cdot f(x, a) + r(x, a) \} = 0 & 0 \leq t < T \\ V(x, T) = g(x) \end{cases}$$

Step 2: 由  $v$  构造最优反馈控制.

$$\alpha^* = \alpha(x, t) \in \arg \max_{a \in A} \{ f(x, a) \cdot \nabla_x V(x, t) + r(x, a) \}$$

然后考虑系统

$$\begin{cases} \dot{x}^*(s) = f(x^*(s), \alpha(x^*(s), s)) & s \in [t, T] \\ x^*(t) = x \end{cases}$$

$$\alpha^*(s) := \alpha(x^*(s), s)$$

Thm (最优性验证定理) 设  $v$  为 HJB 方程的解, 则上述控制是从  $(x, t)$  出发的最优控制.

Proof:  $v(x, t)$  为起点的子问题收益函数

$$P_{x,t}[\alpha^*(\cdot)] = \int_t^T r(x^*(s), \alpha^*(s)) ds + g(x^*(T))$$

根据 HJB 方程

$$v_t(x^*(s), s) + f(x^*(s), \alpha^*(s)) \cdot \nabla_x v(x^*(s), s) + r(x^*(s), \alpha^*(s)) = 0$$

$$\Rightarrow v_t(x^*(s), s) + \nabla_x v(x^*(s), s) \cdot \dot{x}^*(s) = -r(x^*(s), \alpha^*(s)).$$

$$\text{由于 } \frac{d}{ds} v(x^*(s), s) = \nabla_x v(x^*(s), s) \cdot \dot{x}^*(s) + v_t(x^*(s), s)$$

$$\Rightarrow \frac{d}{ds} v(x^*(s), s) = -r(x^*(s), \alpha^*(s))$$

两边在  $[t, T]$  上积分

$$\int_t^T \frac{d}{ds} v(x^*(s), s) ds = v(x^*(s), s) \Big|_t^T$$

$$= v(x^*(T), T) - v(x^*(t), t)$$

$$= - \int_t^T r(x^*(s), \alpha^*(s)) ds$$

根据  $v(x, T) = g(x)$ , 可得

$$g(x^*(t)) - v(x^*(t), t) = - \int_t^T v(x^*(s), \alpha^*(s)) ds$$

化简可得

$$P_{x,t}[\alpha^*(\cdot)] = v(x, t)$$

E.g.: 三个速度的 - 维系统.

$$\begin{cases} \dot{x}(s) = \alpha(s) & 0 \leq t \leq s \leq 1 \\ x(t) = x \end{cases}$$

控制集  $A = \{-1, 0, 1\}$

$$\text{收益函数 } P_{x,t}(\alpha(\cdot)) = - \int_t^1 |x(s)| ds$$

$$\text{价值函数 } v(x, t) = \sup_{\alpha \in A} P_{x,t}(\alpha(\cdot)) = - \inf_{\alpha \in A} \int_t^1 |x(s)| ds$$

剩余时间为  $1-t$  ① 以速度  $-1$  运动, 位置变化  $-(1-t)$

② 以速度  $+1$  运动, 位置变化  $(1-t)$

给出边界直线  $x = t-1$  及  $x = 1-t$ , 可将  $(x, t) \in \mathbb{R} \times [0, 1]$  分为三块

$$\text{Region I} = \{(x, t) \mid x < t-1, 0 \leq t \leq 1\}$$

$$\text{Region II} = \{(x, t) \mid |x| \leq 1-t, 0 \leq t \leq 1\}$$

$$\text{Region III} = \{(x, t) \mid x > 1-t, 0 \leq t \leq 1\}$$

Region III: 最优控制  $\alpha(s) \equiv -1, s \in [t, 1]$

$$x(s) = x - (s-t) > 0$$

$$v(x, t) = - \int_t^1 x(s) ds = - \int_t^1 (x - (s-t)) ds$$

$$= - \frac{1-t}{2} (2x + t - 1)$$

Region I: 最优控制  $\alpha(s) \equiv 1, s \in [t, 1]$

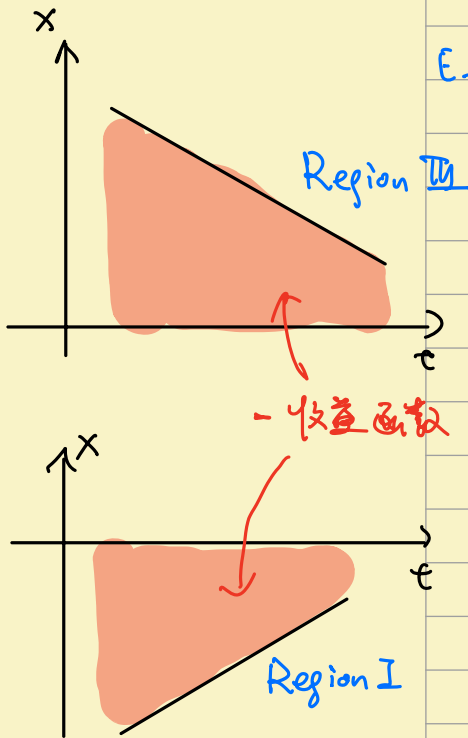
$$x(s) = x + (s-t) < 0$$

$$v(x, t) = \int_t^1 x(s) ds = - \frac{1-t}{2} (-2x + t - 1)$$

Region II: 跑到原点的时间  $|x|$

$$\int_t^1 |x(s)| ds = \int_0^{|x|} (|x| - u) du = \frac{x^2}{2}$$

$$v(x, t) = - \frac{x^2}{2}$$



验证 HJB 方程.

$$\text{HJB: } V_t(x, t) + \max_{a \in \{-1, 0, 1\}} \{a V_x(x, t) - |x|\} = 0$$

$$\max_{a \in \{-1, 0, 1\}} a V_x \begin{cases} V_x & V_x > 0 \\ 0 & V_x \geq 0 \\ -V_x & V_x < 0 \end{cases} = |V_x|$$

因此

$$V_t(x, t) + \max_{a \in \{-1, 0, 1\}} \{a V_x(x, t) - |x|\} = 0$$

$$\Rightarrow V_t(x, t) + \max_{a \in \{-1, 0, 1\}} \{a V_x(x, t)\} - |x| = 0$$

$$\Rightarrow V_t(x, t) + |V_x(x, t)| - |x| = 0$$

在 Region I 时,  $V(x, t) = -\frac{1-t}{2} (-2x + t - 1)$

$$V_t = -(-x + t - 1)$$

$$V_x = 1 - t$$

$$|V_x| = 1 - t$$

} 代入 HJB 方程成立.

同理, Region II, 和 III 成立.

练习:

$$\begin{cases} \dot{x}(s) = a(s) \\ x(t) = x \in \mathbb{R} \end{cases}$$

$$0 \leq t \leq s \leq 2$$

控制集  $A = \{-1, 0, 1\}$

收益函数  $r(x, a) = -|x|$

① 仿照 5.2.1 求出 Region I, II, III.

② 求出三种情形下  $V(x, t)$  表达式

③ 三种情形下, 验证 HJB 方程满足