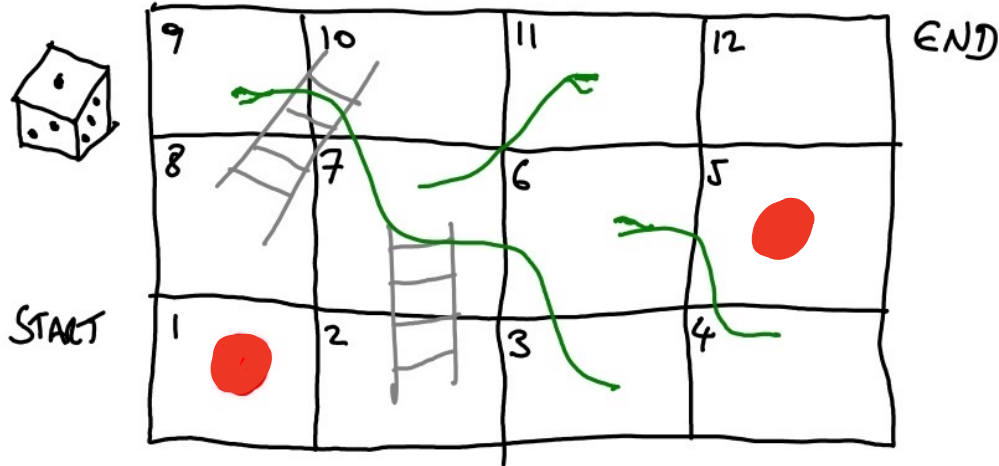


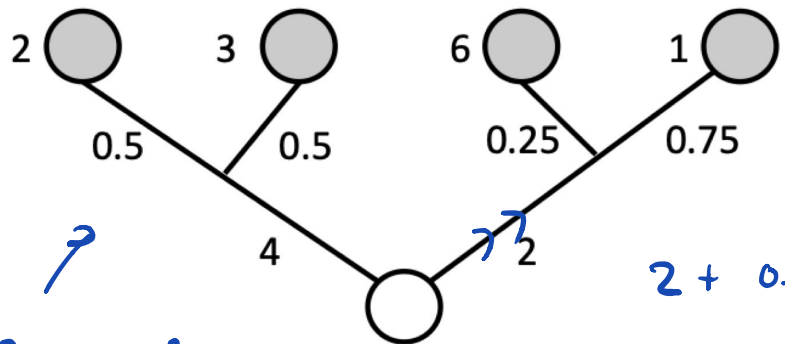
3. MARKOV DECISION PROCESSES

A MARKOV CHAIN WITH DECISIONS :



 [CHOOSE WHICH TO SOLVE]

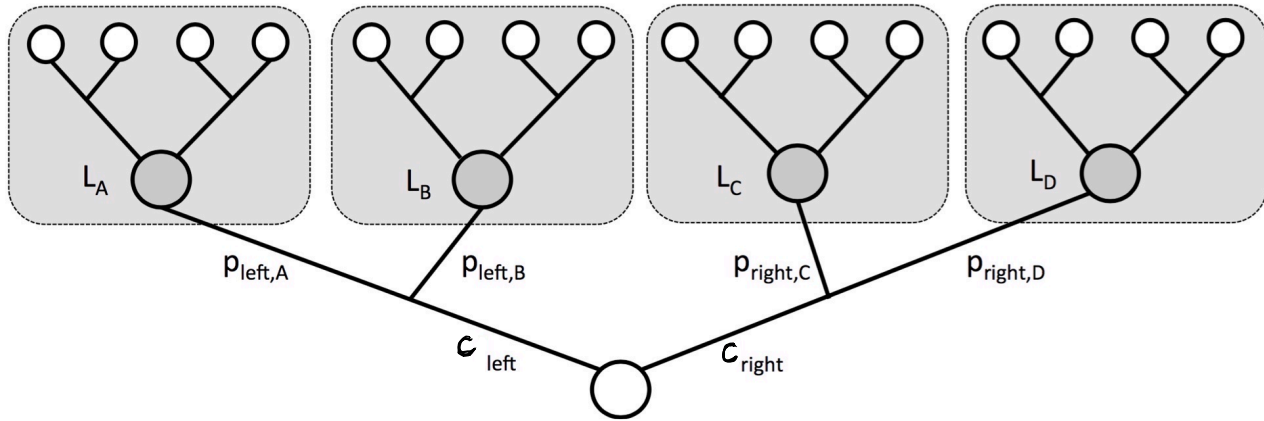
SOLVING AN MDP:



$$\begin{aligned} & 4 + 2 \times 0.5 + 3 \times 0.5 \\ &= 4 + 2.5 \\ &= 6 \end{aligned}$$

$$\begin{aligned} & 2 + 0.25 \times 6 + 0.75 \times 1 \\ &= 2 + 2.25 \\ &= 4.25 \end{aligned}$$

SOLVING AN MDP:



$$L(x) = \min_{a \in \{\text{LEFT}, \text{RIGHT}\}} \left\{ c(a) + \mathbb{E}_{x,a} [L(\hat{x})] \right\}$$

ABSTRACT DEFINITION:

STATE $x \in \mathcal{X}$
ACTION $a \in \mathcal{A}$
REWARD $r(x, a)$
NEXT STATE

$$\hat{x} = f(x, a, u)$$

POLICY $\pi(x)$ OR $\pi_t(x)$

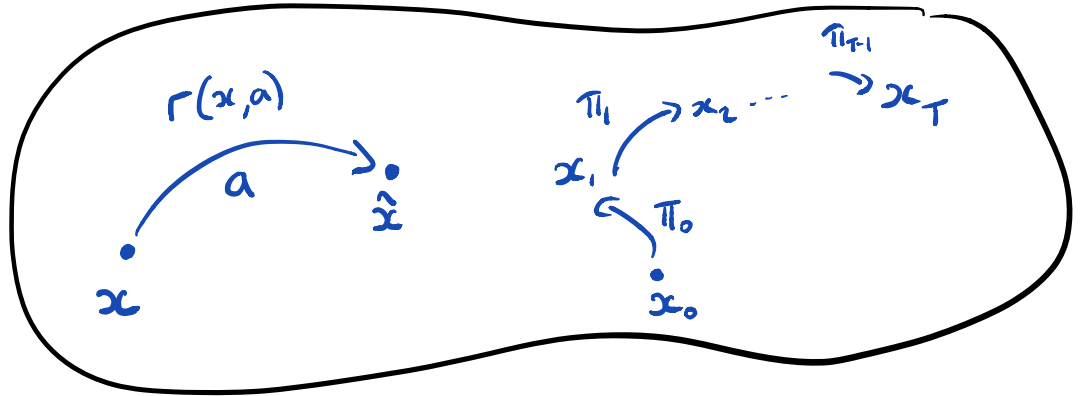
OBJECTIVE: [MAXIMIZE SUM OF REWARDS]

$$V_T(x_0) := \max_{\pi} \underbrace{\mathbb{E} \left[\sum_{t=0}^{T-1} r(x_t, \pi_t) + r(x_T) \right]}_{R_T(x, \pi)} \text{ OVER } \pi \in \mathcal{P}_t$$

↑
VALUE FUNCTION

[POLICY OVER t STEPS]

STATE SPACE \mathcal{X}



THE BELLMAN EQUATION

$$\mathbb{E}_{x,a}[V(\hat{x})] = \mathbb{E}[V(x_1) | X_0 = x, A_0 = a]$$

$$\downarrow = \mathbb{E}[V(f(x,a;u))]$$

THEOREM: $V_t(x) = \max_{a \in A} \left\{ r(x,a) + \mathbb{E}_{x,a}^{t-1}[V(\hat{x})] \right\}, \quad V_0(x) = r(x)$

HERE $\mathbb{E}_{x,a}[V(\hat{x})] = \mathbb{E}[V(x_1) | X_0 = x, A_0 = a]$

$$= \mathbb{E}[V(f(x,a,u))]$$

PROOF:

$$V_t(x) = \max_{\pi_{T-t}, \dots, \pi_{T-1}} \mathbb{E}_{x, \pi_{T-t}} \left[\sum_{s=T-t}^{T-1} r(X_s, \pi_s) + r(X_T) \right]$$

$\max_{\pi_{T-1}} \max_{\pi_{T-2}, \dots, \pi_{T-1}} \mathbb{E}_{x_{T-t}, \pi_{T-t}} \mathbb{E}_{x_{T-t+1}, \pi_{T-t+1}} \left[r(x_{T-t}, \pi_{T-t}) + \sum_{s=T-t+1}^{T-1} r(X_s, \pi_s) + r(X_T) \right]$

MOVE AS FAR TO RIGHT AS POSSIBLE \rightarrow

$$= \text{MAX}_{\frac{\pi_{T-1}}{a}} \left\{ r(x_{T-t}, \pi_{T-1}) + \mathbb{E}_{\frac{x_{T-t}, \pi_{T-t}}{x}} \left[\text{MAX}_{\frac{\pi_{T-2}, \dots, \pi_{T-1}}{\pi_{T-2}, \dots, \pi_{T-1}}} \mathbb{E}_{\frac{X_{T-t+1}, \pi_{T-t+1}}{X_{T-t+1}, \pi_{T-t+1}}} \left[\sum_{s=T-t+1}^{T-1} r(X_s, \pi_s) + r(X_T) \right] \right] \right\}$$

$V_{t-1}(\hat{X})$

$$= \text{MAX}_a \left\{ r(x, a) + \mathbb{E}_{x, a} [V_{t-1}(\hat{X})] \right\}$$

□

Ex 20. You need to sell a car. At every time $t = 0, \dots, T - 1$, you set a price p_t and a customer then views the car. The probability that the customer buys a car at price p is $D(p)$. If the car isn't sold by time T then it is sold for fixed price V , $V < 1$. Maximize the reward from selling the car and find the recursion for the optimal reward when $D(p) = (1 - p)_+$.

ANSWER:

SOLD
OR NOT
} r(p)
r(x, a)

$$V_t(x) = \begin{cases} V_t & \text{IF } x=0 \text{ NOT SOLD} \\ 0 & \text{IF } x=1 \text{ SOLD} \end{cases}$$

$$V_t = \text{MAX}_p \left\{ p D(p) + (1 - D(p)) \overbrace{V_{t-1}} \right\} \quad V_0 = V$$

$$= \text{MAX}_p \left\{ p(1-p) + (1 - (1-p))V_{t-1} \right\}$$

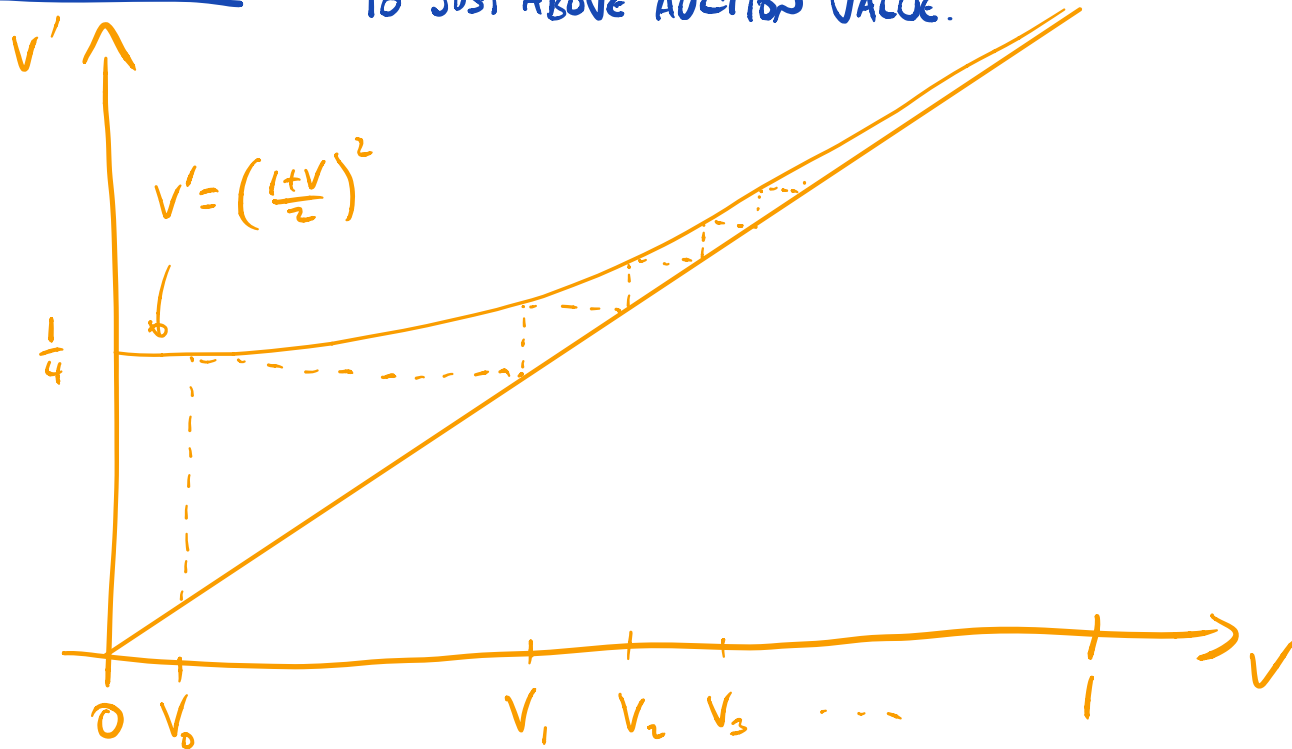
DIFFERENTIATE & SOLVE FOR p GIVE

$$V_t = \left(\frac{1 + V_{t-1}}{2} \right)^2$$

NOTE $V_t > V_{t-1}$ & ONLY FIXED POINT IS $V^* = 1$

So $V_t \rightarrow 1$ AS $t \rightarrow \infty$

A PICTURE: START HIGH & SLOWLY DECREASE PRICE
TO JUST ABOVE AUCTION VALUE.



Ex 21 (Call Option). You own a call option with strike price p . Here you can buy a share at price p making profit $X_t - p$ where x_t is the price of the share at time t . The share must be exercised by time T . The price of stock X_t satisfies

$$X_{t+1} = X_t + \epsilon_t$$

for ϵ_t IIDRV with finite expectation. Show that there exists a decreasing sequence $\{a_t\}_{0 \leq t \leq T}$ such that it is optimal to exercise whenever $X_s \geq a_s$ occurs.

ANSWER: BELLMAN EQN IS

$$V_t(x) = \text{MAX} \{ x - p, \mathbb{E}[V_{t+1}(x + \epsilon)] \}$$

1). POLICIES WITH t STEPS LEFT CONTAIN ALL POLICIES WITH $t-1$ STEPS

So $V_t(x) \geq V_{t-1}(x)$.

2) $V_t(x) - x = \text{MAX} \{ -p, \mathbb{E}[V_{t+1}(x + \epsilon)] - x \}$ IS DECREASING CONTINUOUS IN x .
DECREASING

BECAUSE $V_0(x) - x = \text{MAX} \{ -p, -x \}$ HAS THESE PROPERTIES. & MAX & EXPECTATION PRESERVE THESE PROPERTIES.

\therefore DECISION SWITCHES AT x_{T-t}^* .
 $-p = \mathbb{E}[V_{t+1}(x_{T-t}^* + \epsilon) - x_{T-t}^*]$ x_t DECREASES BECAUSE
 V_t INCREASES

A PICTURE:

