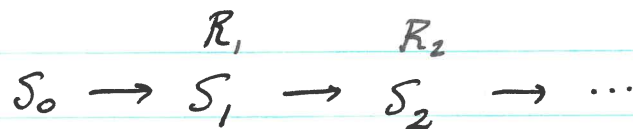# Chapter 2: Markov reward processes (MRP)

- State: $S_t$     $t = 0, 1, \ldots$

- Reward: $R_t \in \mathbb{R}$     $t = 0, 1, \ldots$
  - Cost / utility at time $t$
  - RV correlated with transitions $S_t \to S_{t+1}$
  - Sometimes: fct of $s, s'$     $r(s, s')$

- Transition probability (dynamics):

$$p(s', r \mid s) = P(S_{t+1} = s', R_{t+1} = r \mid S_t = s)$$

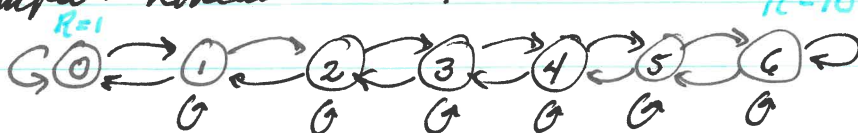$$S_0 \xrightarrow{R_1} S_1 \xrightarrow{R_2} S_2 \to \cdots$$

  - Reward depends on $s, s'$
  - Assumed homogeneous (time-independent, stationary)
  - No actions yet: dynamics + rewards

- Reward probability: $p(r \mid s) = \sum_{s'} p(s', r \mid s)$

- Transition probability: $p(s' \mid s) = \sum_{r} p(s', r \mid s)$    *pure dynamics*

- Expected state reward:

$$\rho(s) = E[R_{t+1} \mid S_t = s] = \sum_{r} r \, p(r \mid s)$$
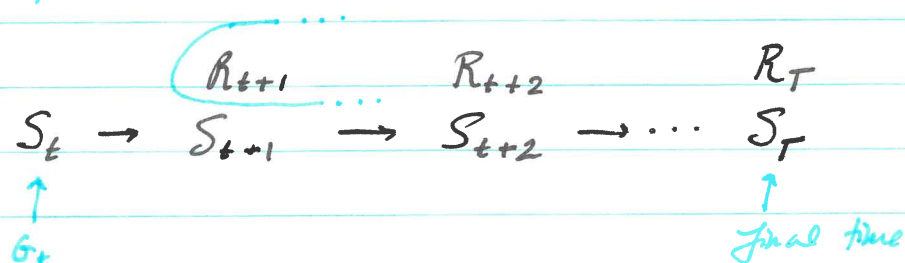
  - One-time reward from state $s$

- Example: linear model / Mars rover    $R = 10$    $R$ deterministic w/o $S_t$



$R = 1$

$\rho(0) = 1$    $\rho(6) = 10$

· Return : $G_t = R_{t+1} + R_{t+2} + \cdots + R_T$

present     future

$$S_t \rightarrow S_{t+1} \xrightarrow{R_{t+1} \quad \cdots \quad R_{t+2}} S_{t+2} \rightarrow \cdots \quad S_T \quad R_T$$

$\uparrow G_t$             final time

· Total reward / cost to go from time $t$
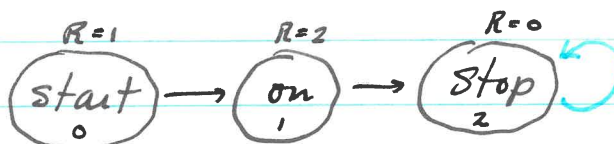
· Episodic process : $T$ finite

         Process stops / terminates

                                  Calculate backward

$\Rightarrow \quad G_T = 0 \quad G_{T-1} = R_T \quad G_{T-2} = R_{T-1} + R_T \quad \cdots$

· Continuing process : $T = \infty$

                  Process goes on and on

· Example :

         $R=1$      $R=2$      $R=0$

         (start) $\rightarrow$ (on) $\rightarrow$ (stop) $\circlearrowleft$

         0        1        2

         $0 \rightarrow 1 \rightarrow 2$     $1 \rightarrow 2$     $2$     Episodic

         $0 \rightarrow 1 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow \cdots$     Continuing

· Discounted return : $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 G_{t+3} + \cdots$

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

  · Discount rate $\gamma \in [0,1]$    Balances present / future rewards

  · $G_t < \infty$ for $\gamma \in [0,1)$
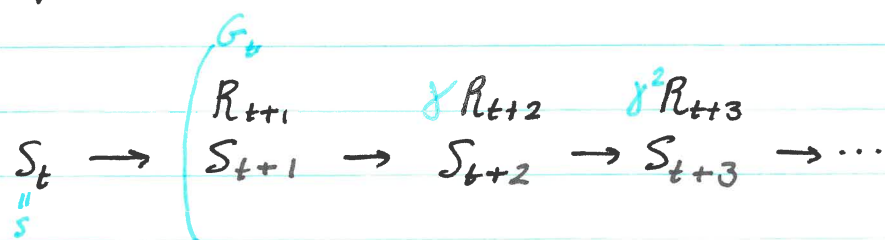
  · Myopic : $\gamma \approx 0$     $G_t = R_{t+1}$ for $\gamma = 0$ immediate reward

  · Far-sight : $\gamma \approx 1$

  · Iteration :

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

· Value function:  $v(s) = E[G_t \mid S_t = s]$

$$S_t \rightarrow \overbrace{\begin{array}{cccc} R_{t+1} & \gamma R_{t+2} & \gamma^2 R_{t+3} \\ S_{t+1} \rightarrow & S_{t+2} \rightarrow & S_{t+3} \rightarrow \cdots \end{array}}^{G_t}$$

- · State value fct : Expected return / cost to go from state $s$
- · Doesn't depend on time $t$
  - · Stationary MP
  - · Continuing process / infinite return
- · $\gamma = 0$ :   $v(s) = \rho(s)$    cost when leaving $S$
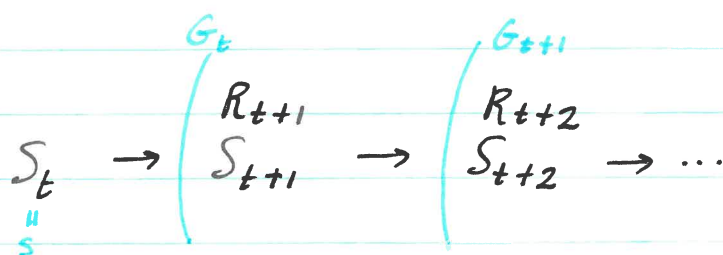
· Example : linear model
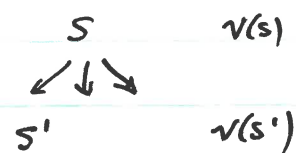


$R=1$        $R=10$

deterministic moves / transitions



0.4   0.4   0.4   0.2

$v(s)$ ?   See CW

· Bellman equation :

$$v(s) = E[R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$= E[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s]$$

$$S_t \rightarrow \overbrace{\begin{array}{c} R_{t+1} \\ S_{t+1} \rightarrow \end{array}}^{G_t} \overbrace{\begin{array}{c} R_{t+2} \\ S_{t+2} \rightarrow \cdots \end{array}}^{G_{t+1}}$$

$v(s)$    $v(S_{t+1})$

- Explicit expectation:

$$v(s) = \sum_r r \, p(r|s) + \gamma \sum_{s'} v(s') \, p(s'|s)$$

$$= \sum_{r,s'} r \, p(s', r|s) + \gamma \sum_{s',r} v(s') \, p(s', r|s)$$

$$= \rho(s) + \gamma \, (Pv)(s) \qquad P_{ij} = p(j|i)$$

- Matrix notation:  $v = \rho + \gamma Pv$  $\qquad v_i = v(i)$

  - $v_i = v(i)$  column vector  $|S'| \times 1$
  - $\rho_i = \rho(i)$  "  "
  - $\gamma$  scalar
  - $P_{ij} = p(j|i)$  $|S'| \times |S'|$ matrix

- Solution:  $v = (I - \gamma P)^{-1} \rho$

  - Always exists for $\gamma \in [0,1)$
  - Unique solution

- Bellman operator:  $T(v) = \rho + \gamma Pv$

$$v = T(v) \qquad \text{Value function} = \text{fixed pt of } T$$

- Example: linear model.
  - Write $\rho$, $P$
    - $7 \times 1$   $7 \times 7$
  - Solve numerically.

· Two ways to define MRPs:

    ①                                                ②

$p(s', r \mid s)$

     ↳ $p(r \mid s)$

            ↳ $\rho(s) = E[R_{t+1} \mid S_t = s]$                             $\rho(s)$

     ↳ $p(s' \mid s)$

         ↳ $P$                                                 $P$

$$v(s) = E[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s]$$

$$v = \rho + \gamma P v$$