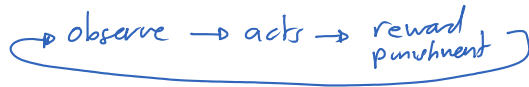


- Incorporate Time to a decision network



- long but unknown amount of time.
- indefinitely (infinite horizon process)

Problem #1: Utility

Consider :

- A: \$1,000, \$1,000, \$1,000, ...
- B: \$1, \$2, \$1, \$1, ...
- C: \$1,000,000, \$0, \$0, \$0, \$0, \$0, ...
- D: \$0, \$0, \$0, ... \$1,000,000, \$0, \$0, \$0, ...

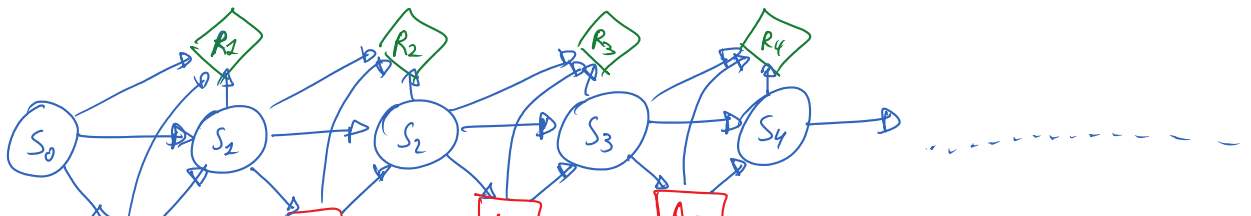
"Value"

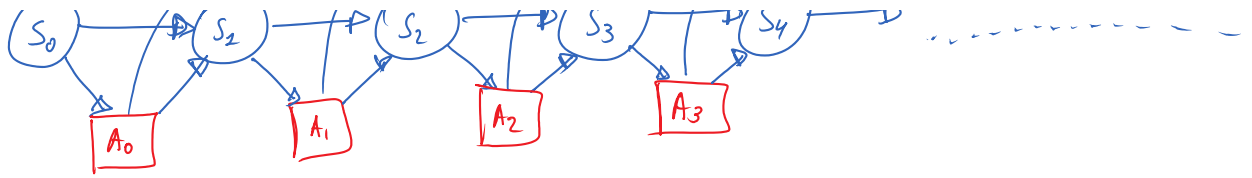
- $V = \sum_{i=0}^{\infty} r_i$
- $V = \lim_{n \rightarrow \infty} (r_1 + \dots + r_n) / n$
- $V = r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 r_4 + \gamma^4 r_5 \dots$
 γ : discount rate $0 \leq \gamma \leq 1$
 $= r_1 + \gamma (r_2 + \gamma (r_3 + \gamma (r_4 \dots)))$
- $V_t = r_t + \gamma V_{t+1}$

If we have the following properties

- the next state depends only on the current action.
- the effects of actions are stochastic, but do not change over time

Belief Network: Markov Decision Process (MDP)





an MDP consists of:

- a set S of states
- a set A of actions
- the dynamics $P(S_{i+1} | S_i, A_i)$
- Reward table $R(S, a, S')$

Sometimes we use expected value of R $R(S, a) = \sum_{S'} R(S, a, S') \cdot P(S' | S, a)$

- γ the discount factor

Example:

Agent Bob $S = \{\text{healthy}, \text{sick}\}$
 $A = \{\text{relax}, \text{party}\}$

What should bob do each weekend?

$P(S' | S, a)$

S	A	$P(S' = \text{healthy})$
healthy	relax	0.95
healthy	party	0.70
sick	relax	0.50
sick	party	0.10

$R(S, A)$		
S	A	
healthy	relax	70
healthy	party	100
sick	relax	0
sick	party	20

• A Policy:

A policy $\pi: S \rightarrow A$ what to do at each state.

Among policies, π^* optimal policy is the one with maximum expected reward.

An MDP with stationary dynamics always has an optimal policy.

expected reward.

An MDP with stationary dynamics always has an optimal policy.

- Partially Observable Markov Decision Processes (POMDP)
Markov Decision Process on a Hidden Markov Model

