| MDP's | POMDP's | Dec. POMDP's ↳ Decentralized | Idea from Self |
|---|---|---|---|
| Similarity: Structure ie Agent & environment Goals | | | A side note from an idea in Theory of Computation |
| Differences: 1 Agent: Environment is known *Goals are generalized reward f's. * Actions are Stochastic. ∘ Perfect perception | 1 Agent but environment is partially known ie partial observability | other agents + POMDP's | * We know that given a finite memory we could only perform if the problem could be expressed in finite States * Therefore if we could express the problem in a finite state space we could solve it ; even if it takes non-polynomial time . |

Atomic Model

↳ Exponential ⇒ finik horizon cases
↳ Undecidable ⇒ in infinite horizon cases

why Dec. MDP's exists ? ⇒ If multiple agents exist then they won't know about other agents.

Doubt Question ⇒ These are still in the markovian Domain. what about the Non-markovian? Ans ⇒ Non Markovian / Non-Markovian Dynamic still can be compiled into MDPs.

There are still many other types ⇒ SMDP's , MMDP's
⇓ ⇓
Time syn.    multiple Agents having apriori information

* Finding optimal policy in MDP is polynomial in state level only.

Non-trivial learning still could be done at atomic level.
Case in Point :- Atomic is there is no structure. Even Linear Regression has structure to it.

RL ⇒ Can be done at atomic level.
↳ Basically if I am in at a state then
to avoid falling into bad decisions.

A* search ← MDP differs

# Introduction to Gridworld
## Ideas from Russell Norvig

* The transitions are Markovian in a sense that probability of reaching a state s' from s depends on only s & not on the history.
* Rewards must be "bounded".

## MDP formally is

A sequential decision problem for a fully observable, stochastic environment with a Markovian transition model and additive rewards is called a MDP & consists of a set of states, set of Actions in each state a transition model $P(s'|s,a)$ & reward fn $R(s)$.

## Policy

↳ A solution must specify what the agent should do for any state that the agent might reach. { Think back to DFA in ToC }

Optimal Policy ⟹ Policy that yields the highest expected utility.

## Finite or Infinite Horizon

↳ After a fixed time N after which nothing matters

For ex stand at (3,1) & N=3 then head directly.



If N = 100 ⟹ take safer route from left. With finite horizon ⟹ the optimal action in a given state could change over time.

1) Additive Rewards :- Utility = $R(s_0) + R(s_1) ---$

2) Discount Rewards :→ $R(s_0) + \gamma R(s_1) +$
-- --

Let's see non-terminal

* Now our policy depends on the the $R(s)$.

Different $R(s)$ ⟹ Different behaviour

Moving forward we see that infinite horizon pose a natural question ⟹ 1) If agent does not reach a terminal state or if the agent never reaches one then ⟹ environment ∞ Utilities ∞

BUT

Infinite summation of Discount

$$= R_{max} / (1-\gamma)$$

## Points to note

* If agent agents up in a terminal state ⟹ no infinite sequences ↳ "Proper Policy".

* A way to deal with infinite seq would be the averaging argument.

## Bottom Line

↳ Discounted Rewards present the fewest difficulties in evaluating state sequences.

# Optimal Policy and Utilities of States

- **Recap**
- Utility ⇒ Sum of discounted rewards during the sequence
- **Compare Policies** ? ? ? ⇒ Expected Utilities obtained when executing

**Start**
↳ Assume an agent is in some state s & define $S_t$ {a variable}
↓
a state agent reaches at time t. [ $S_0 \Rightarrow s$ ; Simple enough ]

Probability distribution over state sequence $S_1, S_2, \_\_$ is determined by the initial state s, the policy $\pi$ & transition model for environment

Expected Utility $U^\pi(s) = E\left[\sum_{t=0}^{\infty} r^t R(S_t)\right]$

{ Logical Reasoning ;⇒ $E(x) = p \times R(p)$ }

$\pi_s^* \to max\ U^\pi(s)$
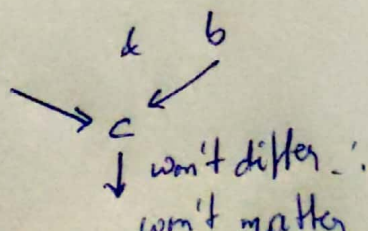
Basically although optimal, this optimality is based on the start state. Lie let's say i start from another state ⇒ other policy will be optimal ; quite obvious.}

Now in case of infinite horizon optimal policy is independent of the starting state.

**Short Proof**

Start with a ⟶ d ⟶ b
Intermediate ⟶ c
↓ won't differ ∴
won't matter

## POMDP's

Even the partially obse is different from value parameterized $f^n$ that w looked before. Even then the agent did not know about it's current location We have some partial informa about the states.