

## Model in RL

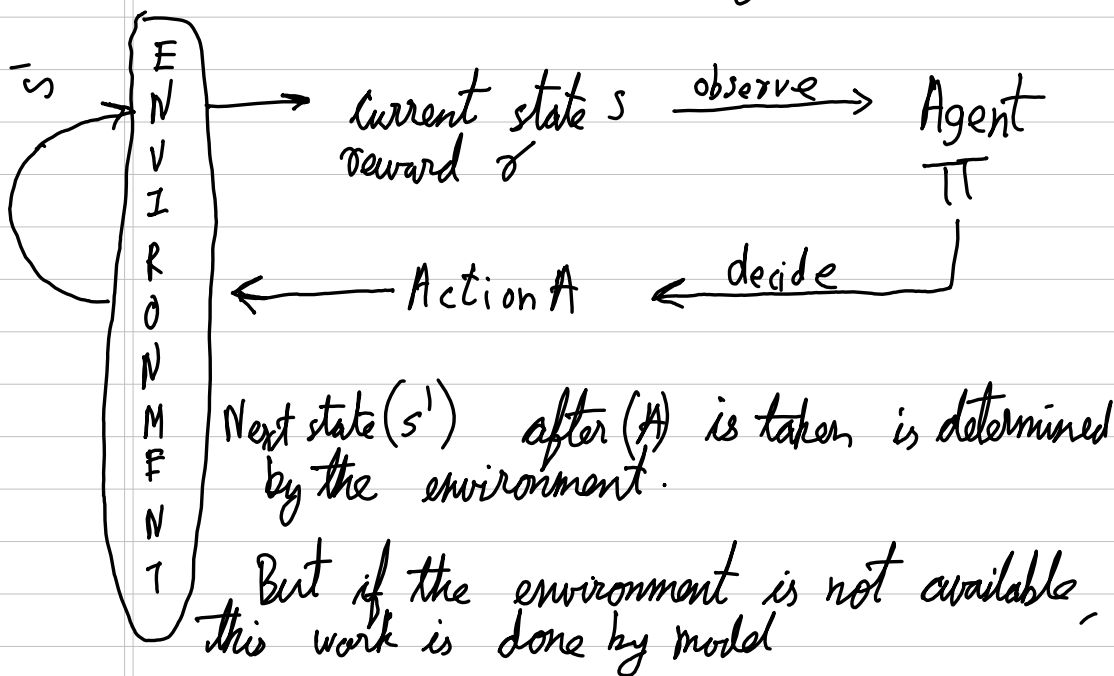
A model is anything that agent can use to predict the environment's behaviour

Given an state & action, a model predicts the next state & next reward

stochastic  $\rightarrow$  several possible states with probability (distribution model)

sample model  $\rightarrow$  give only 1 state with highest probability

determining what will be next state



Model can just replace environment.

## Planning in RL

Planning is any computational process that takes a model as input and produces or improves policy for interacting with modelled environment



Learning is when we actually interact with the environment



Planning uses simulated experiences generated by model

Learning uses real experiences

Model based (D.P., heuristic search) rely on Planning  
Model free (MC, TD, etc) rely on Learning

Learning	Planning
Real experience	simulated experience
Model free	Model based
May be slower depending on environment	faster

In many cases, learning algorithm can be used for planning.

eg Planning & learning where everything is the same except environment is replaced by model

# Model free vs Model based

Model based methods work well when the model is known.

In such cases (toy problems) the model and environment are the same

But in real life cases, if the knowledge about the exact environment is not known then  $\text{model} \neq \text{env.}$

This may be due to incorrect assumptions

(D.P.) so when  $P(s'|s, a)$  &  $R(s, a)$  are estimates and not true values,  $\text{model} \neq \text{env.}$

Also in some applications like recommendation systems and financial markets, the environment may change over time.

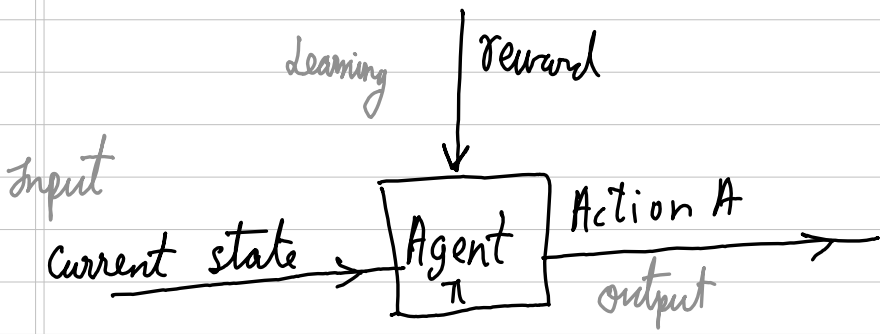
In such cases, the model used is just an estimate or approximation of the environment.

Model free methods on other hand learn directly from the environment

Model may be incorrect because  $\rightarrow$

- ① Environment is stochastic and only a limited amount of experience is available
- ② Model has learned using a function approximator that failed to generalize
- ③ Environment has changed and new behaviour has not been observed yet.

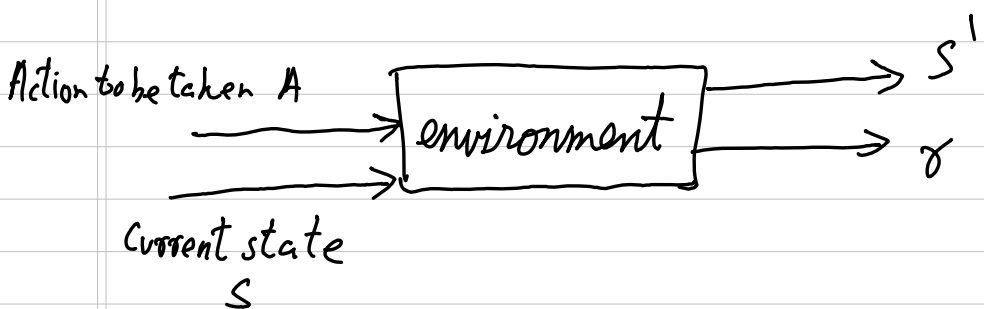
# Planning vs learning



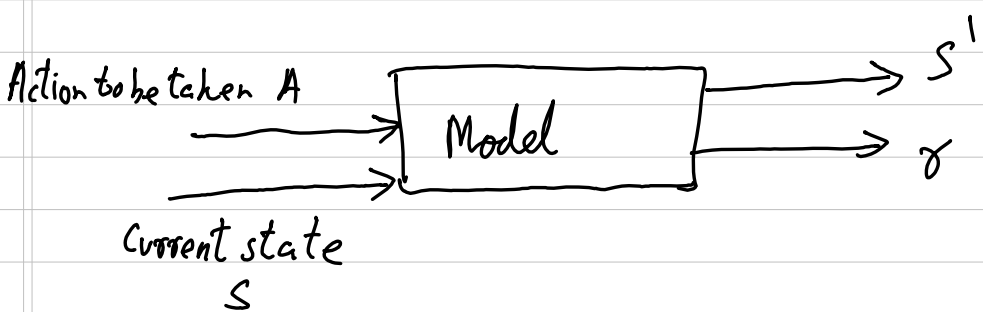
Agent is a black box with input  $S$  output  $A$ .  $\pi$  is used for learning

$S$  &  $\delta$  need to be supplied by environment or someone else

In learning RL  $\rightarrow$



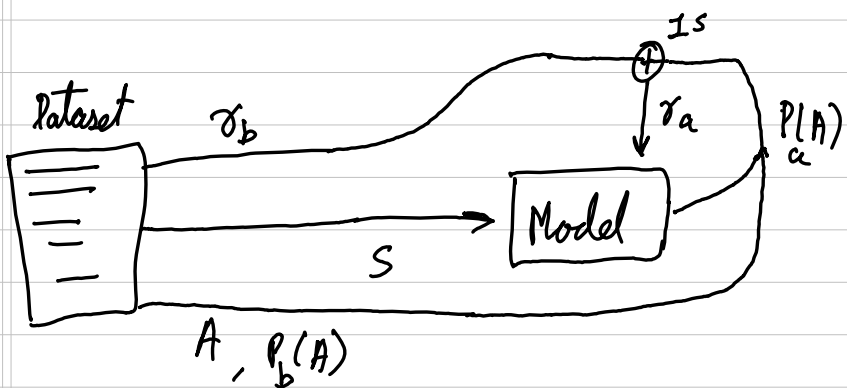
In planning RL  $\rightarrow$



A model is used to simulate the environment

Note  $\rightarrow$  for off policy historical learning

There is no environment, just dataset



$P_a(A)$  is used only to normalize the reward via importance sampling

Action has already been predecided by  $b$

# Direct & Indirect RL

\* What we can do with real experience

- ① Model learning  $\rightarrow$  Improve model to match it with the environment. (to cope environment changes)
- ② Direct RL  $\rightarrow$  Improve value functions and policy from environment (learning)
- ③ Indirect RL  $\rightarrow$  Improve value function and policy from model (Planning)

Direct methods are not affected by bias in model design

But real world interactions may be costly

With a model, agent can "Imagine" i.e. simulate transitions between states without needing to interact directly with the environment

If all these 3 methods are combined, then there will be increase in efficiency

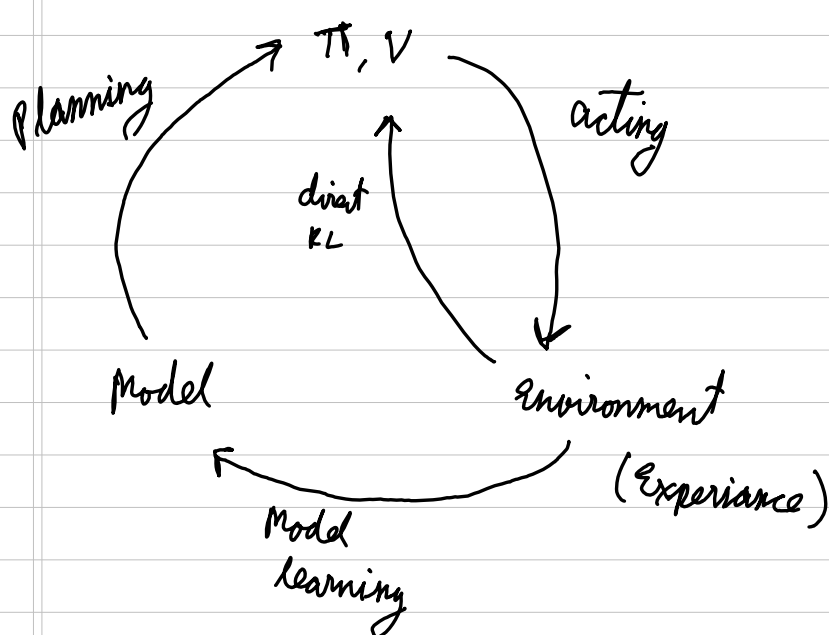
Agent will learn from real environment as well as simulated models

So the agent can update policy multiple times on single real update  $\rightarrow$  faster learning

Dyna - Q

Simple Architecture integrating major functions of an online planning agent

It combines model free & model based techniques.

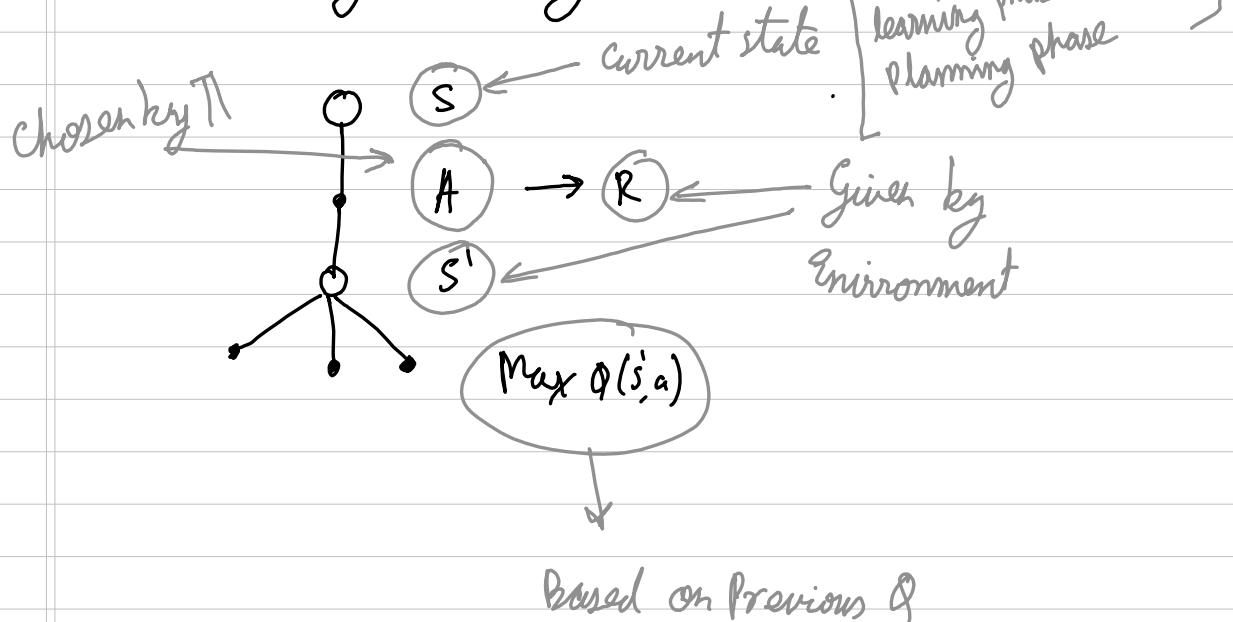


step 1) Take Action & learn from real experience

Agent selects action based on  $\pi$  ( $\epsilon$ -greedy)

Take action observe reward & new state

Update  $\pi$  using Q learning



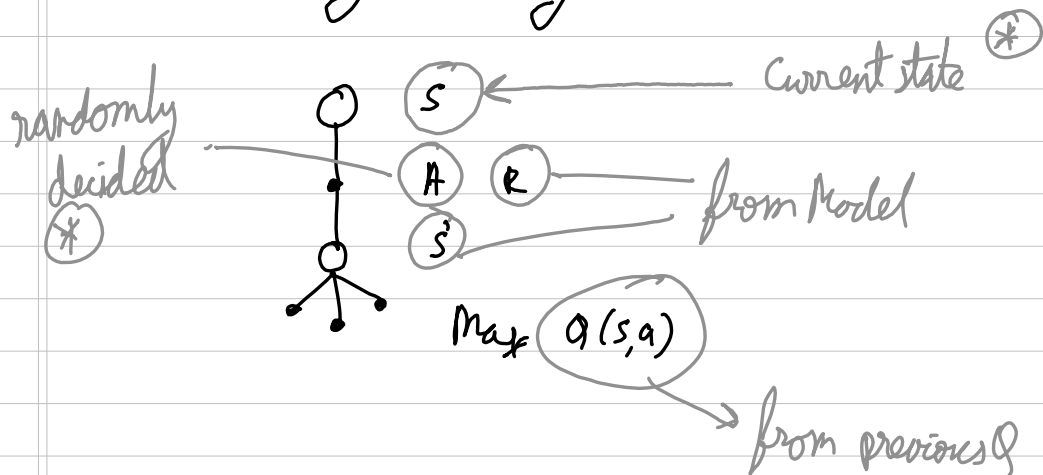
step 2) update the Model

record  $s'$  &  $R$  to update  $P(s, s', a)$   
 &  $R(s, a)$   
 ( $a$  is action taken)

step 3) simulate experience (Planning)

Agent generates number of simulated experiences using Model  
for each simulated experience, update Q-value as if it were a real experience

This step allows the agent to make rapid progress by learning from hypothetical experiences without actually interacting with the environment.



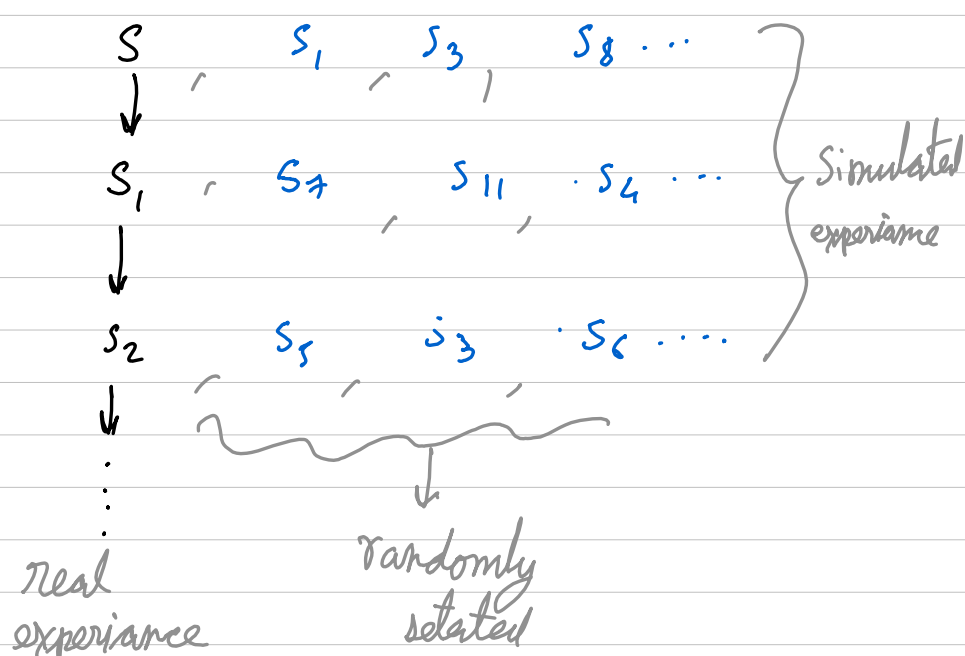
④ → Randomly selected from past state action pairs  $(S, a)$

Repeat Planning  $N$  times

Then go back to step 1

Note  $\rightarrow$  The current state for learning remains same as it was after the most recent real (environmental) experience, not the state from any of the planning steps.

The planning steps only modify the Q-values and not agents current state



# Dyna Q +

Designed to improve exploration

for Non stationary environments where  $\gamma, p$  changes over time

Introduces exploration bonus to agent to revisit the states and actions that it has not visited recently in simulated experience

Model keeps track of each state-action pair of how many time steps have elapsed since it was last tried in a real interaction.

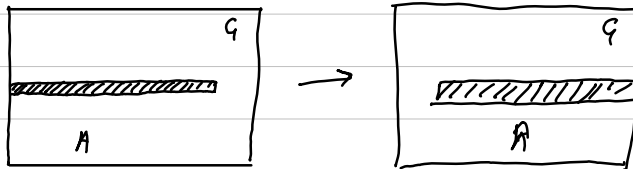
special bonus is assigned  $\gamma \rightarrow \gamma + k \sqrt{T}$

Parameter  
↑

Has cost for revisit, but worth it for changing environment

Transition not tried for  $T$  rewards in real world

eg Gridworld barriers moves



In such cases, revisiting states is helpful

## Prioritized Sweeping

Instead of updating random state action pairs during planning, prioritized sweeping selectively updates only those pairs with highest potential for learning.

These state action pairs are prioritized based on how much their  $Q$  values could change

By focusing only on these high impact updates, prioritized sweeping achieves faster learning with fewer updates.

$$\text{TD error} \rightarrow \left| \underbrace{\gamma + \gamma \max_a Q(s', a')}_{\text{estimated goodness}} - \underbrace{Q(s, a)}_{\text{current goodness}} \right|$$

High T.D. error means updating the  $Q$  value could significantly improve the agents knowledge

Once a  $(s, a)$  is visited,

calculate TD error for  $(s, a)$

Calculate TD error for predecessors

Add all to priority queue with  
priority = TD error (Thresholding can be done)

In planning phase, instead of randomly selecting actions, take actions from priority queue

For each predecessor, when they come in planning phase for update, calculate the  $Q$  of their predecessors too.

A predecessor state are the states that could transition into the current state under some action.

Priority sweeping has been found to dramatically increase the speed of convergence.