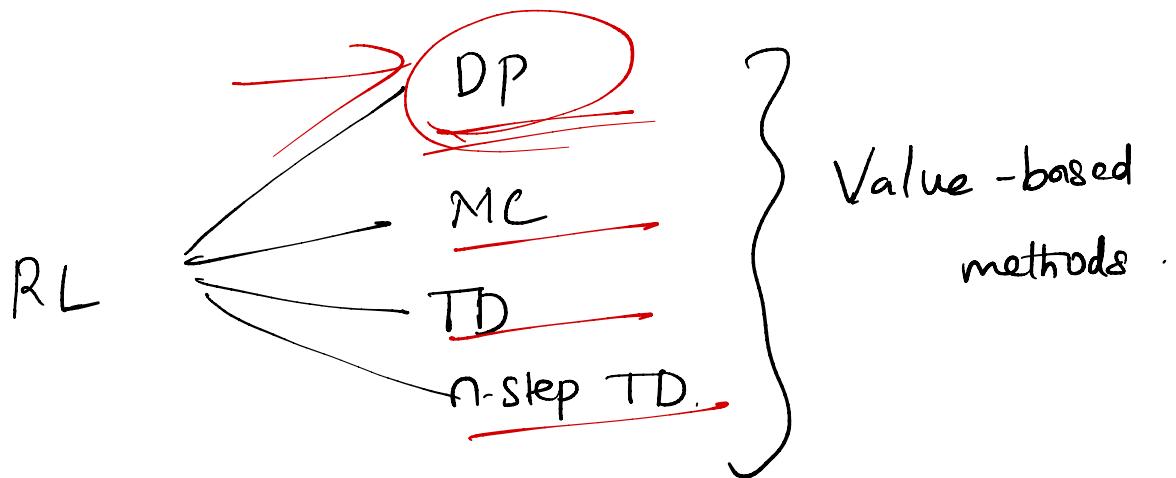


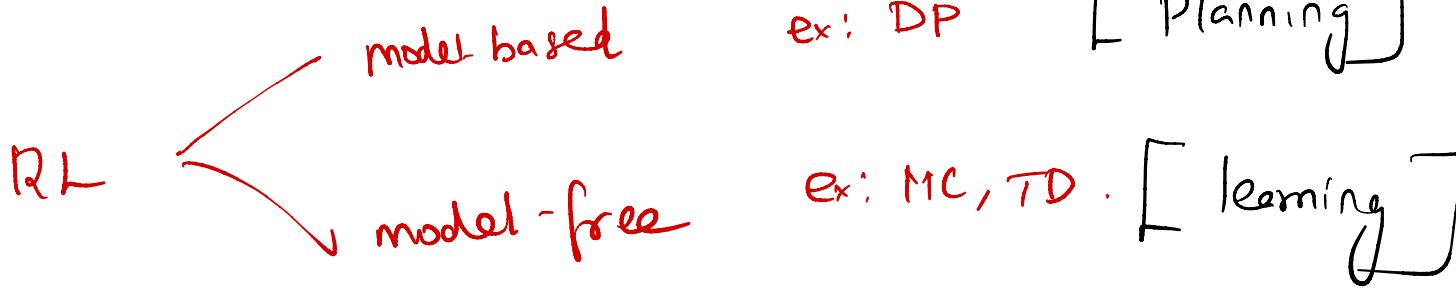
Lecture - 09



Model-based RL

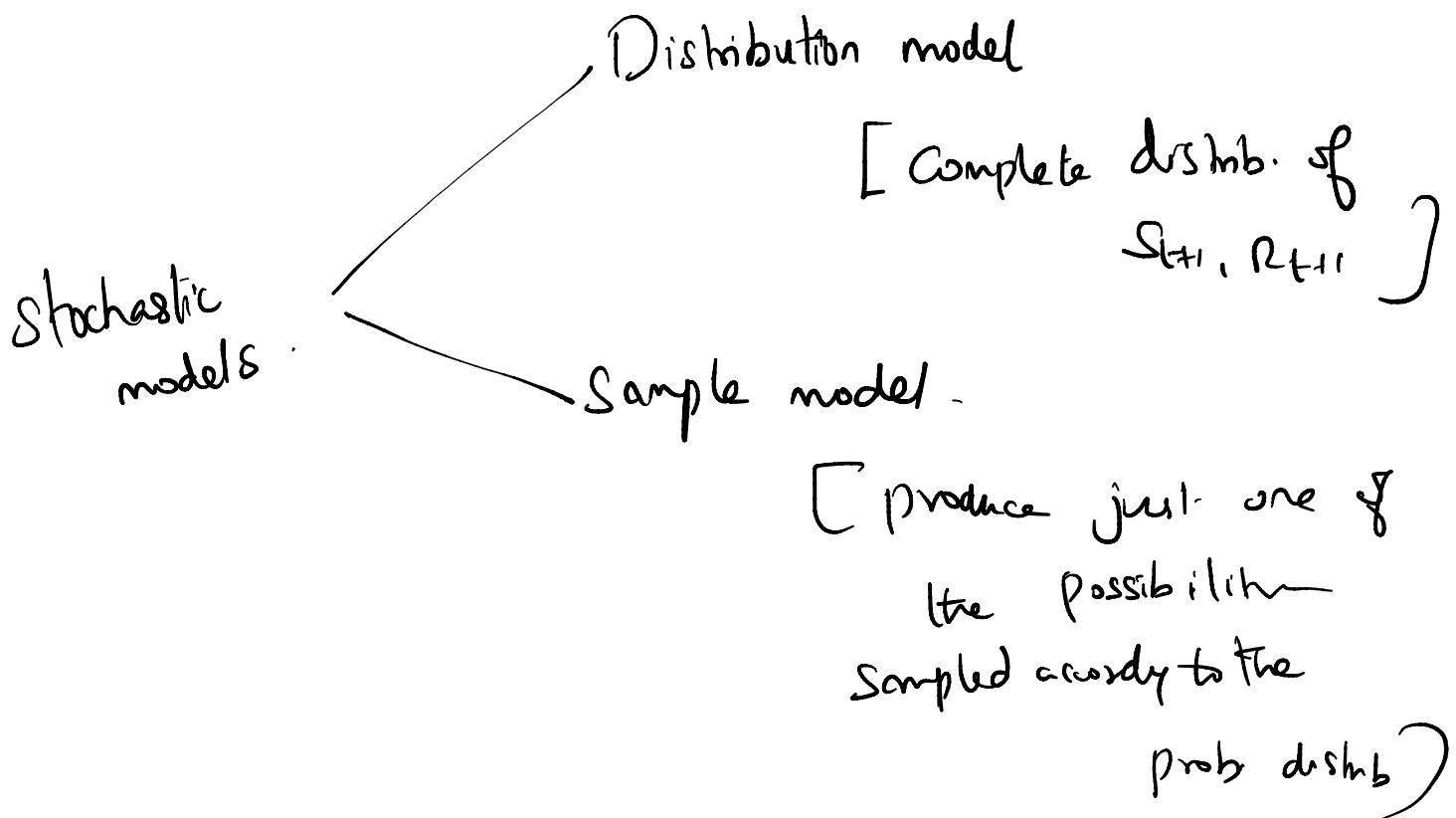
Model of the world \rightarrow dynamics fn.
given $s_t, a_t \rightarrow s_{t+1}, r_{t+1}$

$$P(s_{t+1}, r_{t+1} | s_t, a_t)$$

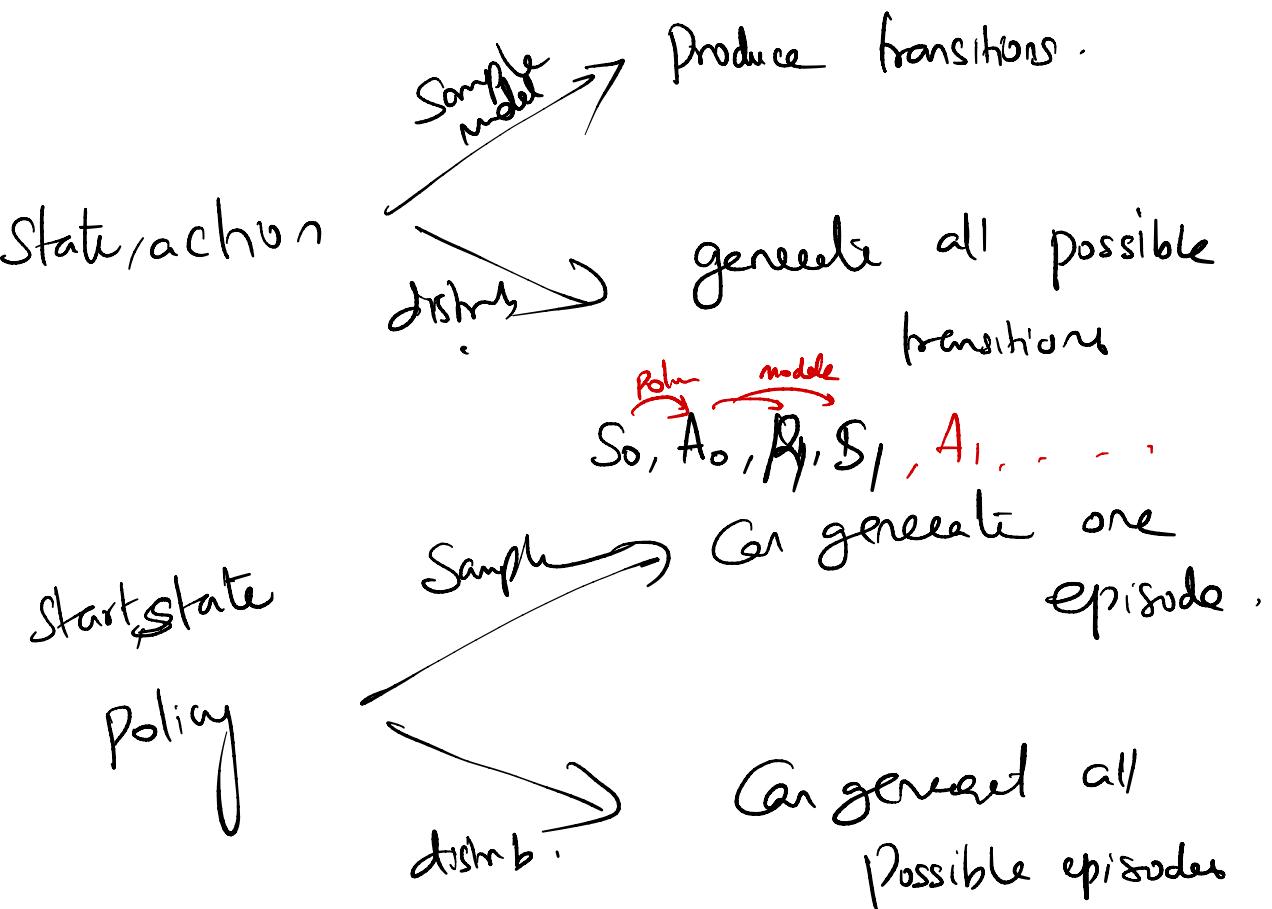


Model :

$$P(S_{t+1}, R_{t+1} | S_t, R_t)$$

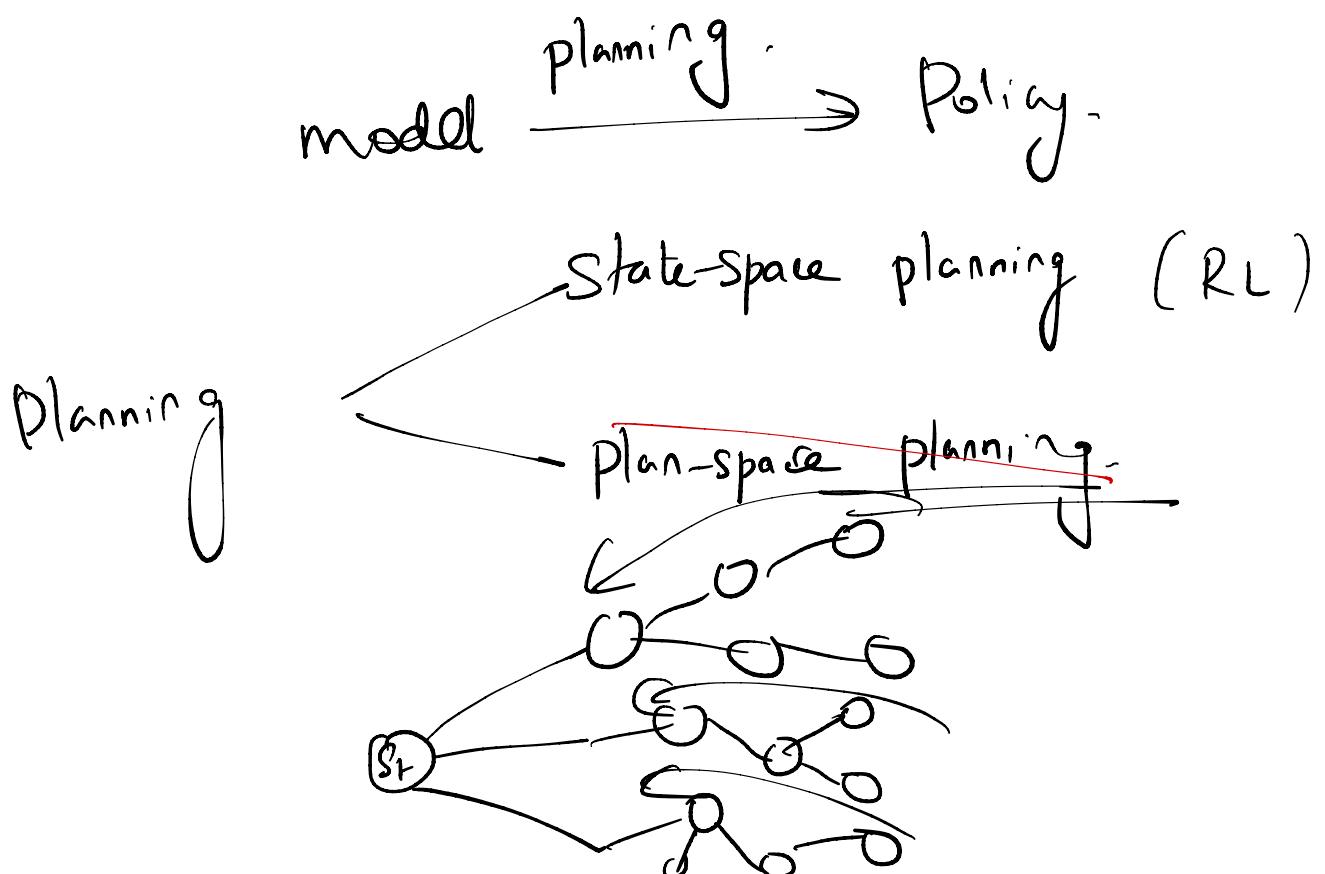


→ ex: modelling the sum of a dozen dice.



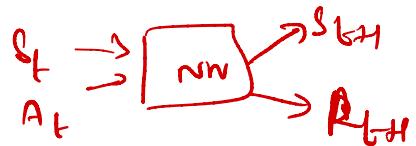
How to use the model?

— Simulate the env.



Learning the model:-

- tabular model.
- Linear model.
- Neural networks.



Planning and Learning :

model → learn better
value fn, Policy

model → to improve the policy
during decision time.
(Decision time planning)

MBRL



MERL



Random-sample one-step tabular Q-planning

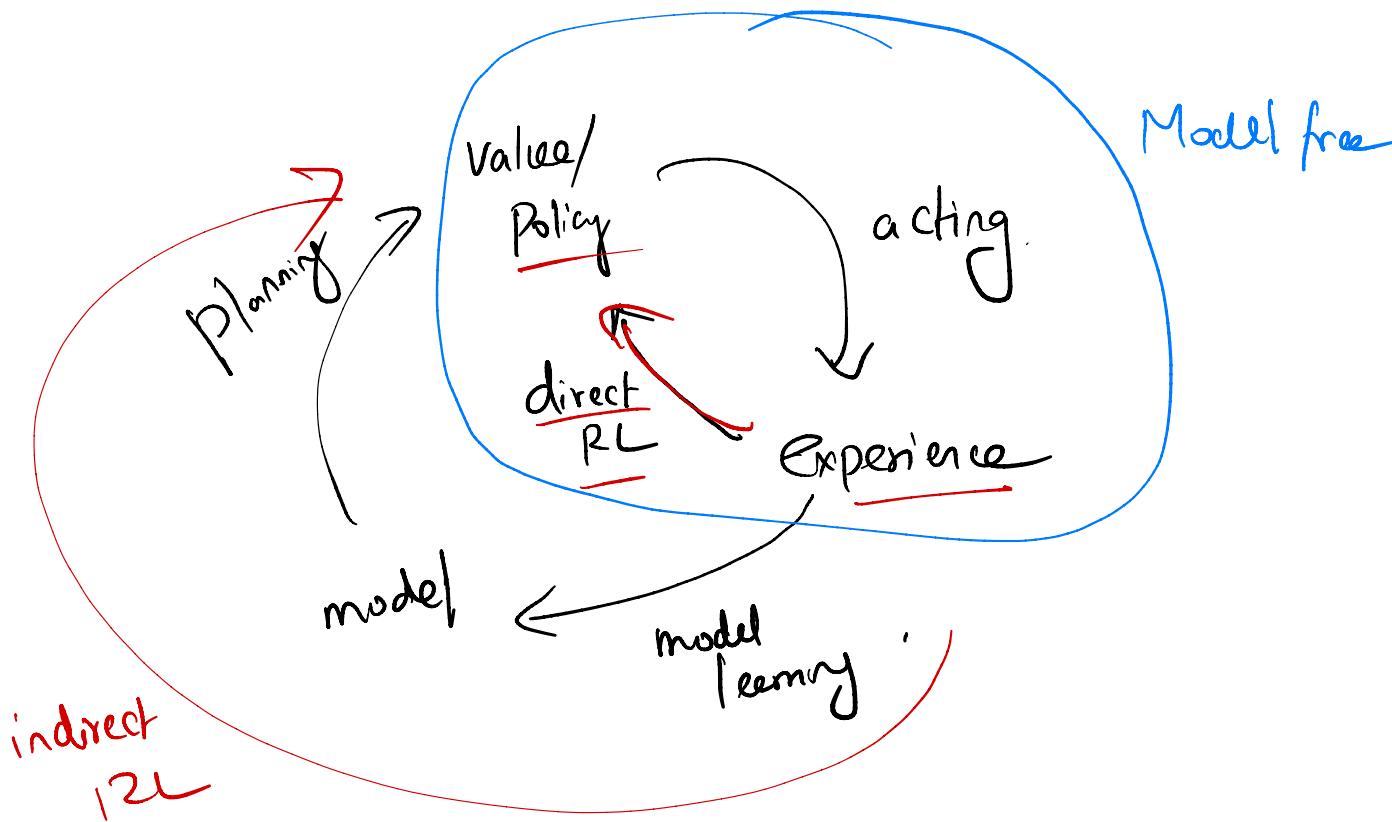
Loop forever:

1. Select a state, $S \in \mathcal{S}$, and an action, $A \in \mathcal{A}(S)$, at random
2. Send S, A to a sample model, and obtain
a sample next reward, R , and a sample next state, S'
3. Apply one-step tabular Q-learning to S, A, R, S' :

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

Dyna : Integrated planning, acting, and learning

Dyna-Q : simple framework for online planning agent

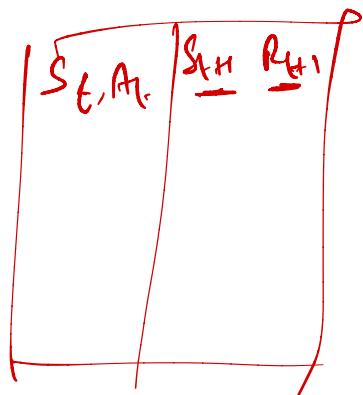


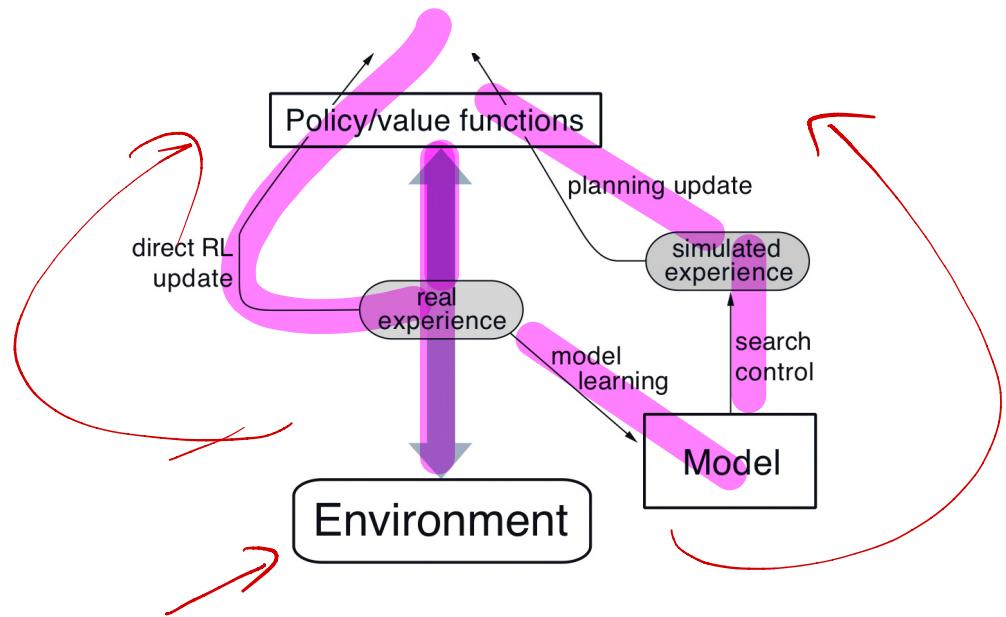
Dyna-Q:

Planning - random Sample one-step tabular
Q-Planning.

direct RL - one-step tabular Q-Learning.

model-learning - tabular model.





Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

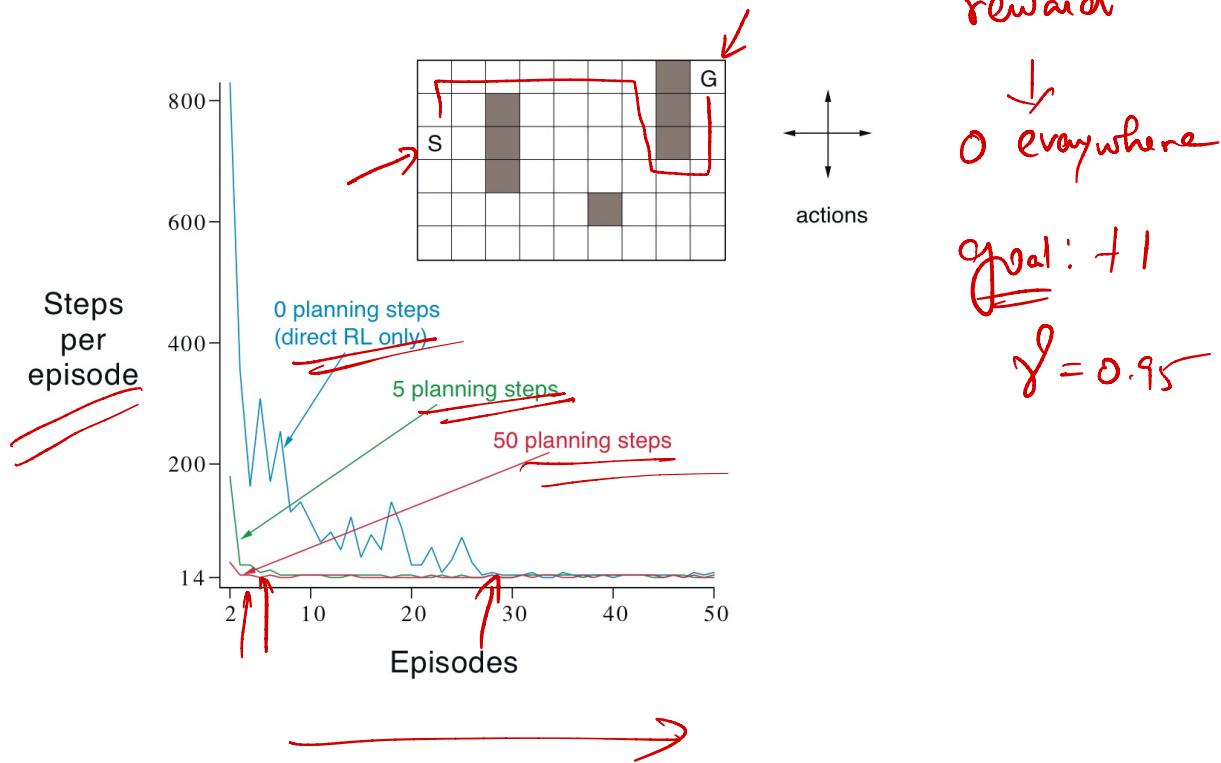
- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \epsilon\text{-greedy}(S, Q)$
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- (f) Loop repeat n times:
 $S \leftarrow$ random previously observed state
 $A \leftarrow$ random action previously taken in S
 $R, S' \leftarrow Model(S, A)$
 $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

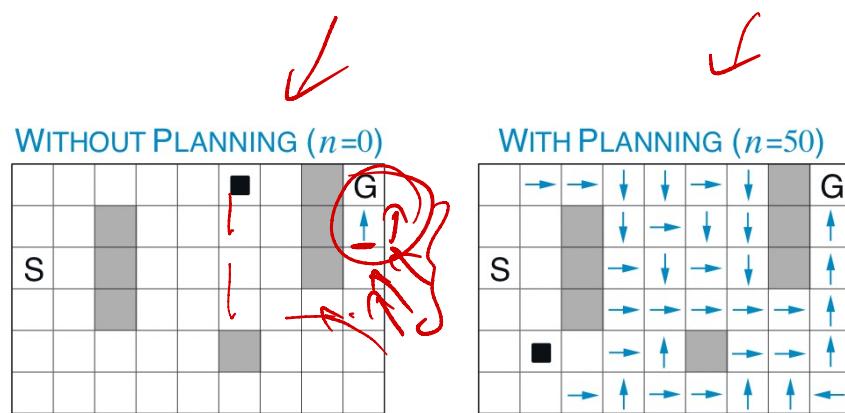
update for
NN model.

direct RL update

model learning

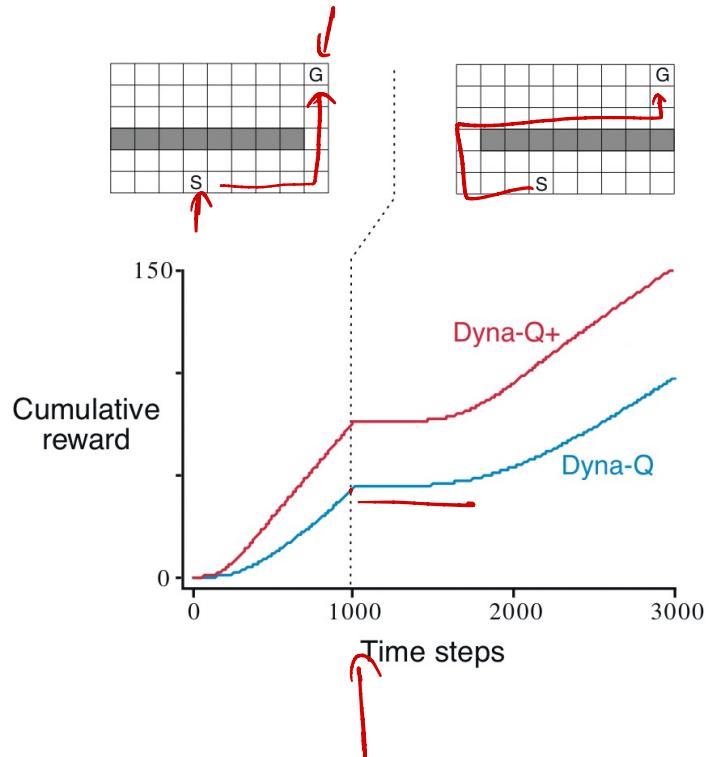
indirect RL update

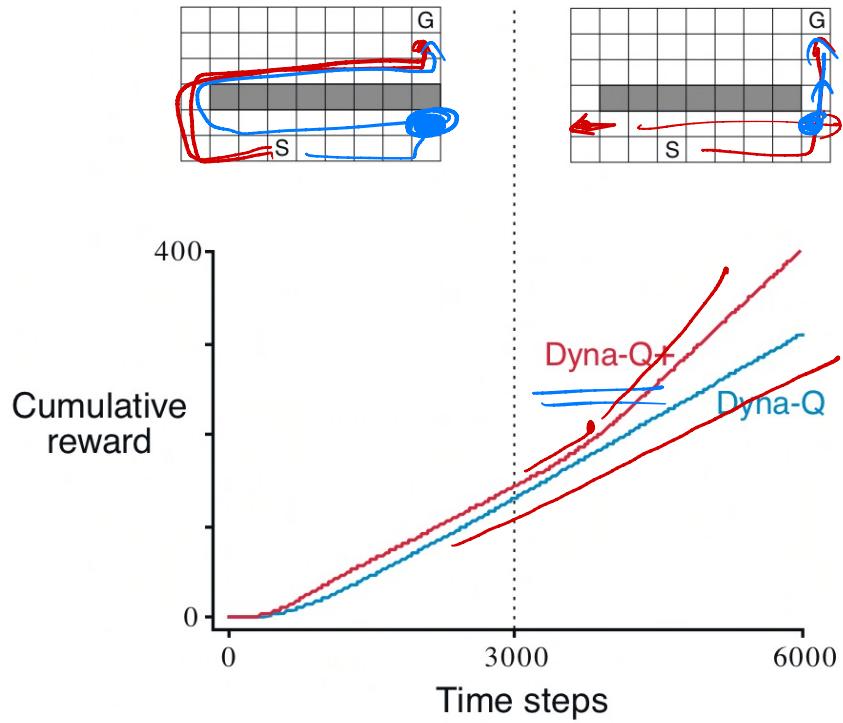




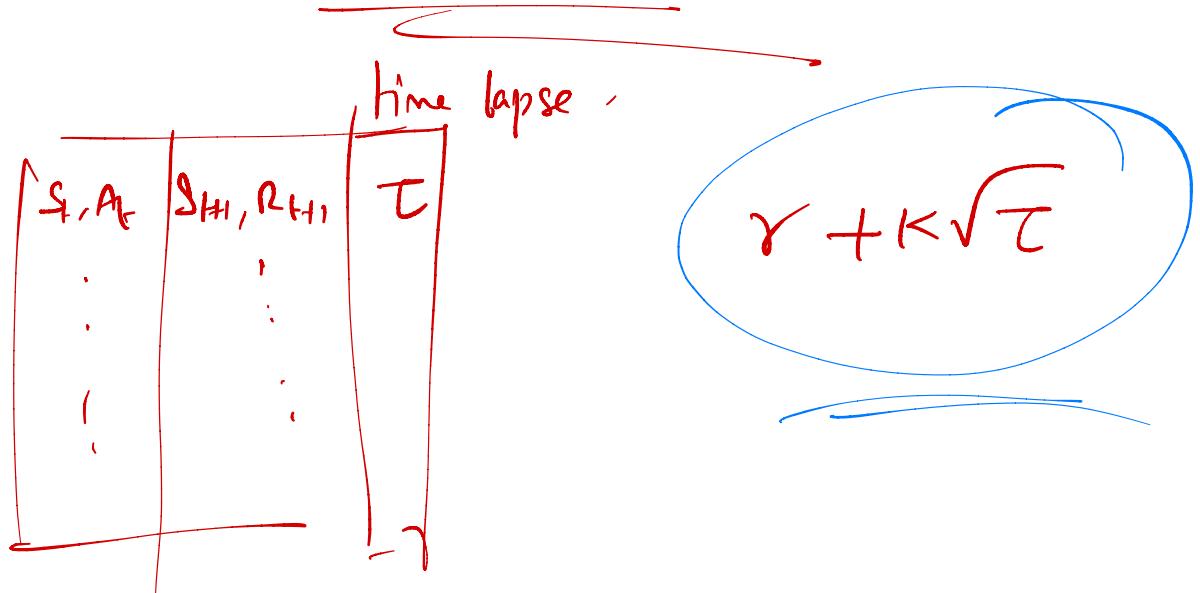
What if the model is wrong?

- ① world is stochastic.
- ② Function approximation

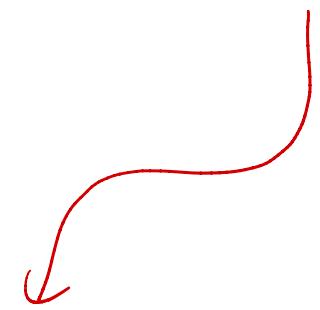




Exploration vs. Exploitation

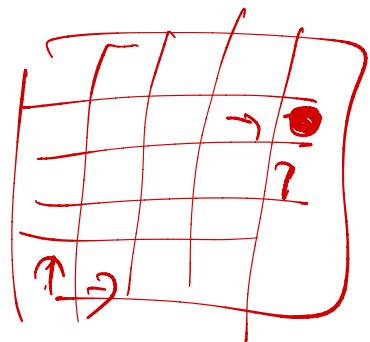


Prioritized Sweeping :-



backward focusing
Computation.

Priority queue

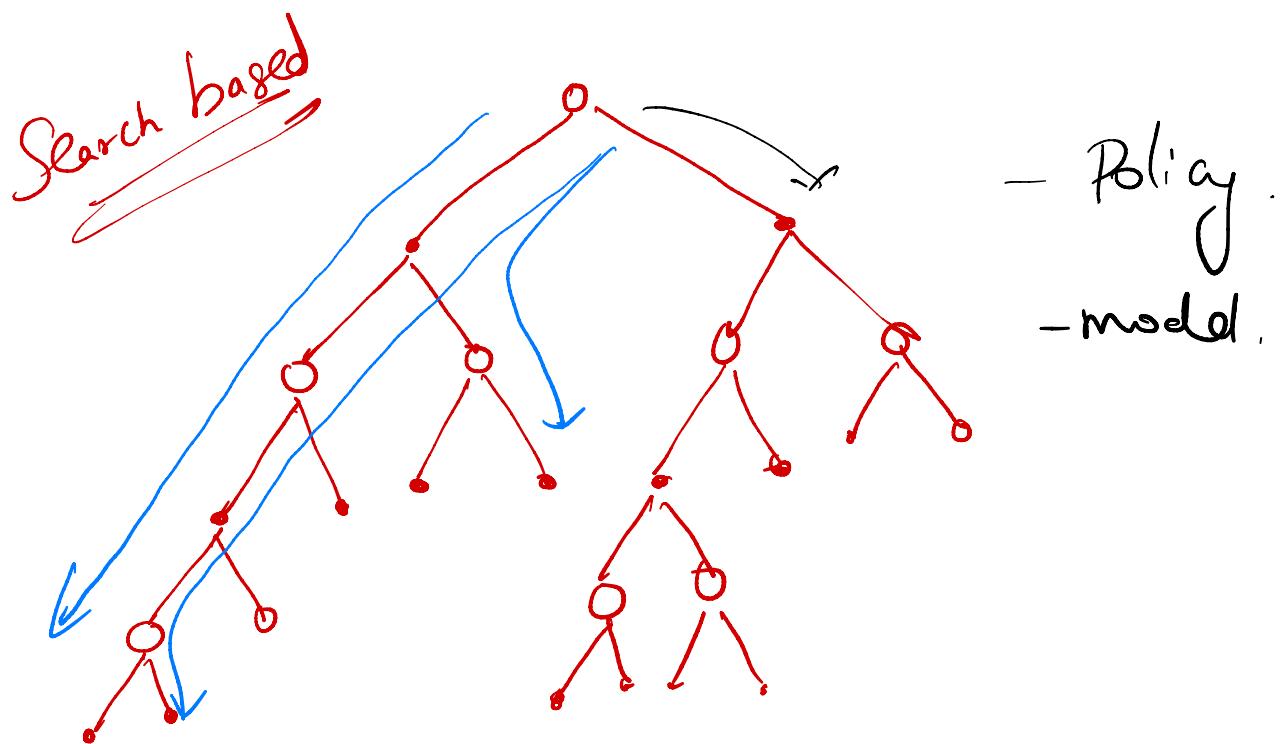
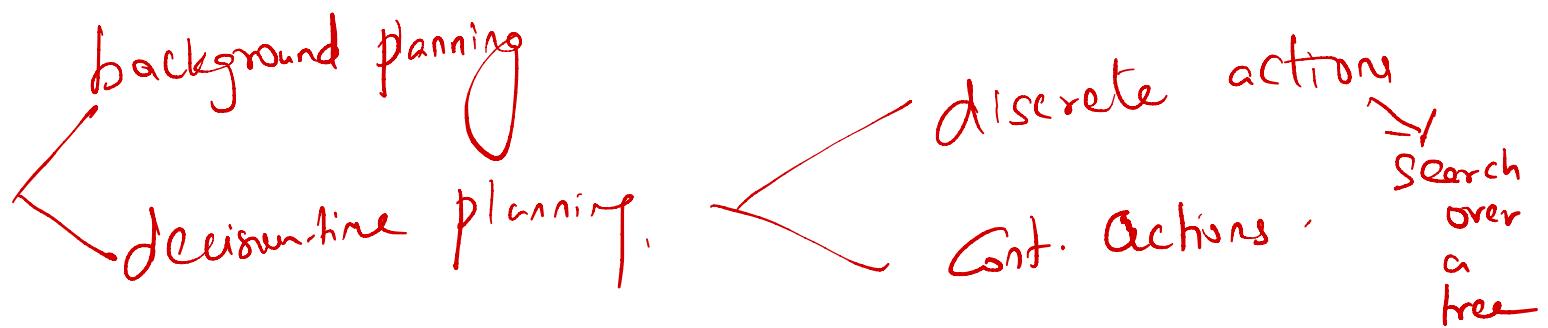


Prioritized sweeping for a deterministic environment

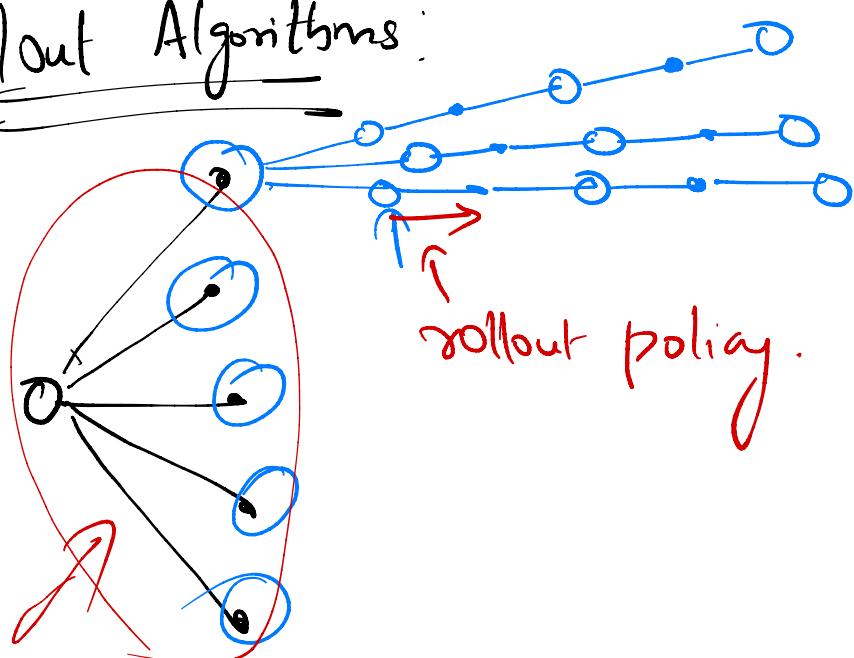
Initialize $Q(s, a)$, $Model(s, a)$, for all s, a , and $PQueue$ to empty
Loop forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \underline{policy(S, Q)}$
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Model(S, A) \leftarrow R, S'$
- (e) $P \leftarrow |R + \gamma \max_a Q(S', a) - Q(S, A)|.$
- (f) if $P > \theta$, then insert S, A into \underline{PQueue} with priority P
- (g) Loop repeat n times, while $PQueue$ is not empty:
 $S, A \leftarrow \underline{first(PQueue)}$
 $R, S' \leftarrow Model(S, A)$
 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
Loop for all \bar{S}, \bar{A} predicted to lead to S :
 $\bar{R} \leftarrow$ predicted reward for \bar{S}, \bar{A}, S
 $P \leftarrow |\bar{R} + \gamma \max_a Q(S, a) - Q(\bar{S}, \bar{A})|.$
if $P > \theta$ then insert \bar{S}, \bar{A} into $PQueue$ with priority P

Planning at decision time :-



Rollout Algorithms



Monte-Carlo Tree Search (MCTS)

Tree Policy

Rollout policy

Selection, expansion, simulation, backup.

