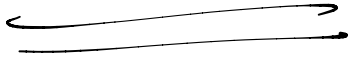


# Lecture-10

MBRL

$$P(s', r | s, a)$$

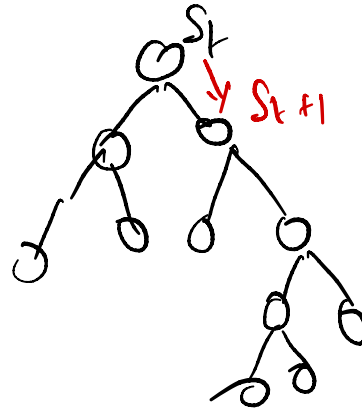


MCTS

use the model to learn the value fn / policy

- Dyna-Q.

online Planning to improve the policy



Stochastic optimization:

Model-Predictive Control (MPC):

$$S_t \rightarrow a_t$$

$$S_t \quad a_t, a_{t+1} \dots a_T$$

$$\underbrace{a_t, a_{t+1} \dots a_T}_A = \underset{a_t \dots a_T}{\operatorname{argmax}} J(a_t \dots a_T)$$

Return.

Random Shooting:

① uniform dists over actions.



4

6

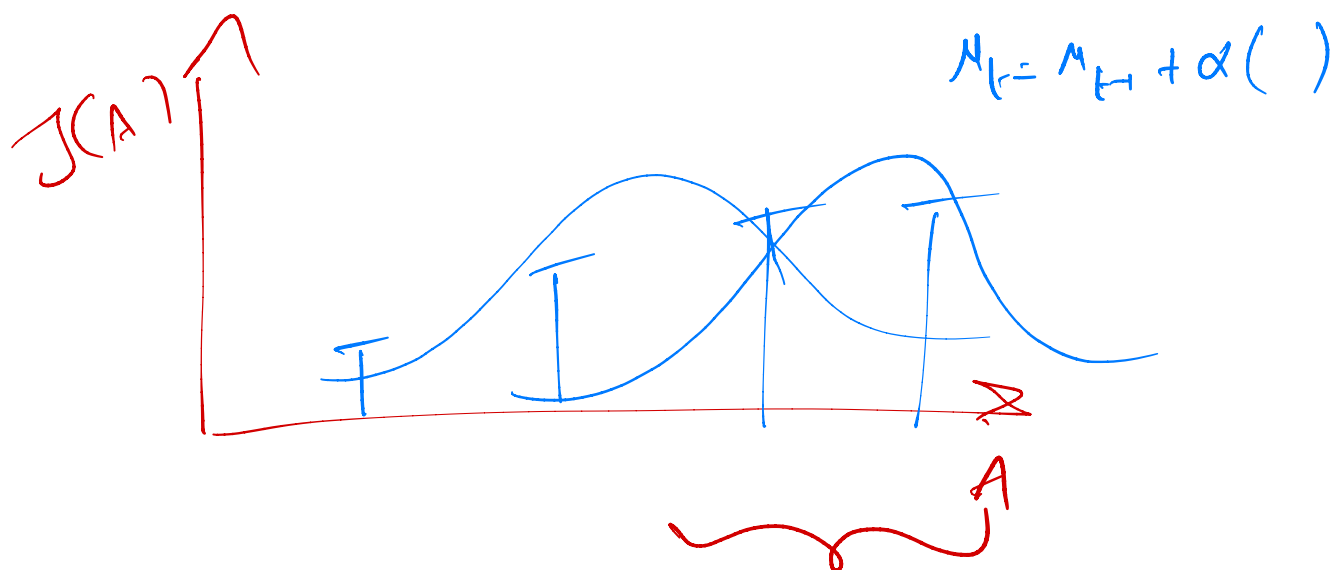
7

# Cross-entropy method:- (CEM)

Gaussian action.  $\mu, \sigma$

$$S_t \rightarrow a_t \quad a_{t+1} \quad a_{t+2} \quad a_{t+3}$$
$$\mu_1 \sigma_1 \quad \mu_2 \sigma_2 \quad \mu_3 \sigma_3 \quad \mu_4 \sigma_4$$

- ① Sample  $A_1 \dots A_N$  from  $p(A)$
- ② evaluate  $J(A_1) \dots J(A_N)$
- ③ Pick  $M$  elite candidates  $A_{i_1} \dots A_{i_M}$   
with highest value.
- ④ refit  $p(A)$  to  $A_1 \dots A_N$



J

