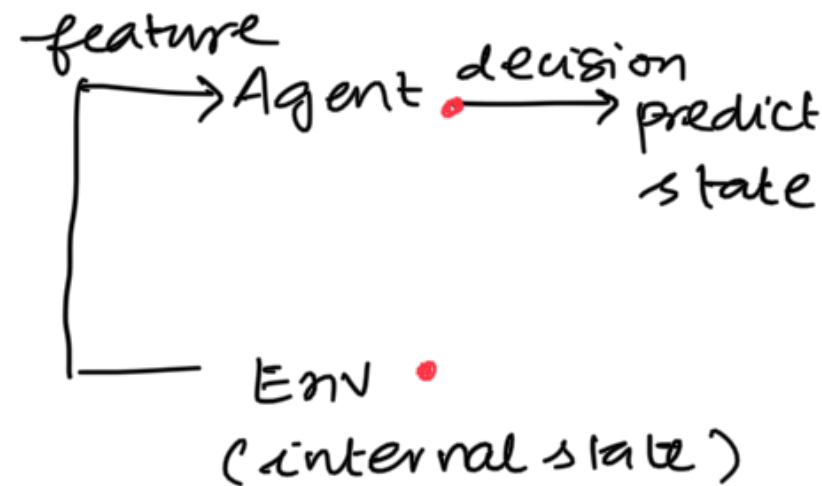
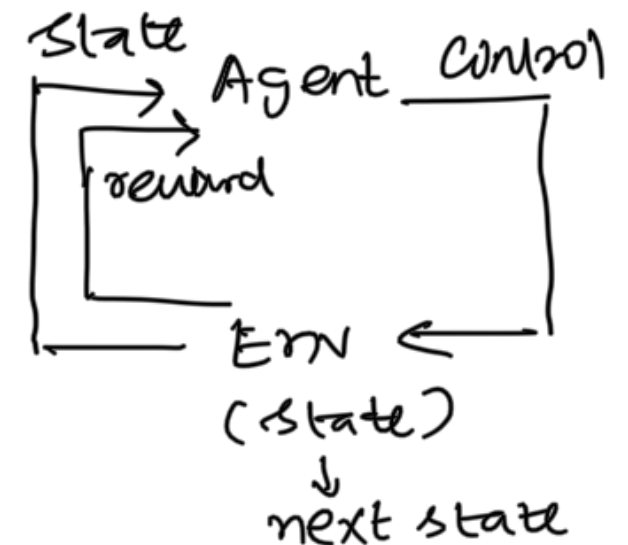


Descriptive
unsupervised
No decision

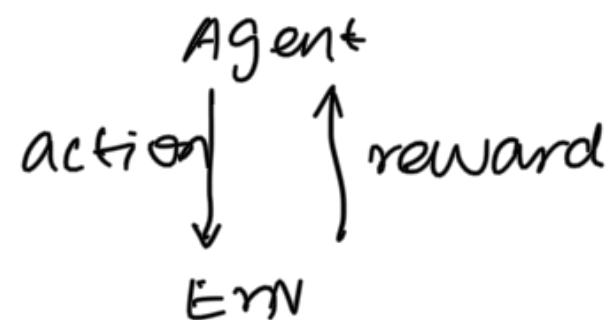


Predictive
Supervised
Predictive decision
No loop



Control
Reinforced
control decision
+ loop

Multi Armed Bandits



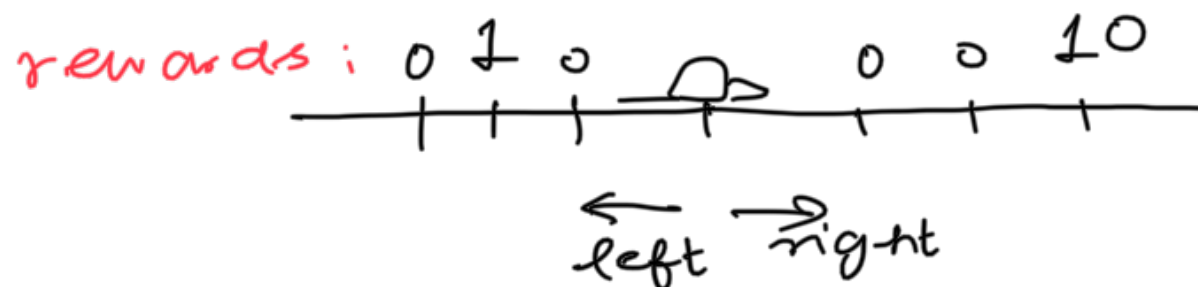
Goal: To look at the challenges in the control or reinforcement learning setting and slice out the multi-armed bandits

Application of : (learning to play)
control : playing chess/Gw, autonomous driving, robotics (Mujoco), play video games (Atari)

inventory control, traffic signal control, optimal investment, online ad placement

Challenge 1:

env: deterministic,



$$\sum_{t=1}^T r_t$$
$$\sum_{t=1}^T (0.9)^t r_t$$

$t=5$ on right

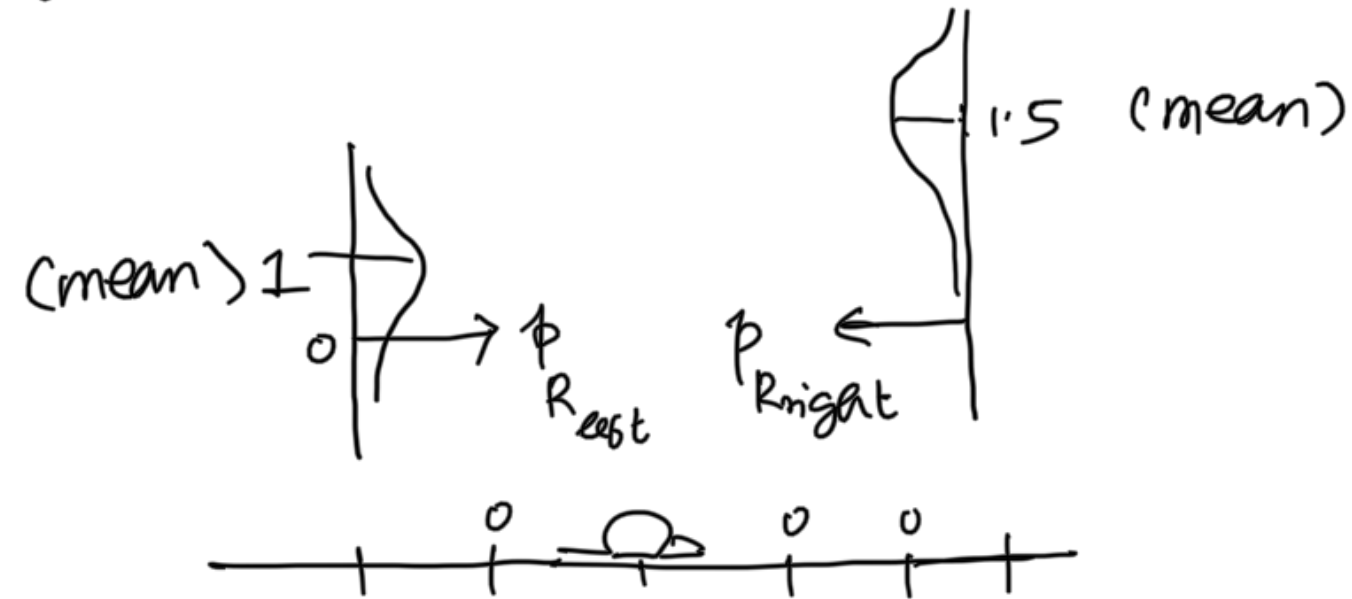
Say $T = 1$, go left ↑ switched
 $T = 3$, go right
 $T = 4$, go right
 $T = 8$, go right
 $T = 9$, go right for 3 steps then left

Moral: Immediate vs Future

challenge 2: Say I play a game of chess, and I win, is it possible to pinpoint the winning move?

Moral: Temporal credit assignment

challenge 3! Environment is not deterministic



$$\max \mathbb{E} \left[\sum_{t=1}^T r_t \right]$$

* you have to visit each of these places multiple times

* say I visit each location 10 times, collect samples, look at sample mean and decide

(few good samples at left location and few bad samples at right location can screw us)

Moral: Cannot explore and commit

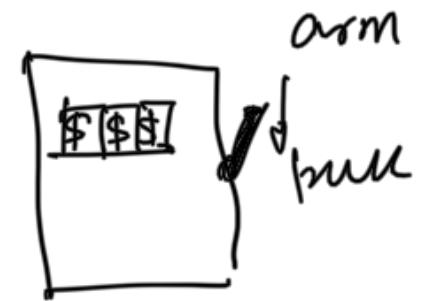
* in order to gain information we actually need ∞ samples (keep exploring)

* while we explore, we should not over explore bad choices, (exploit what we know)

Exploration VS Exploitation

Notation for Multi-Armed Bandits

- $A = \{1, \dots, k\}$ arms



- arm a is associated with distribution P_a

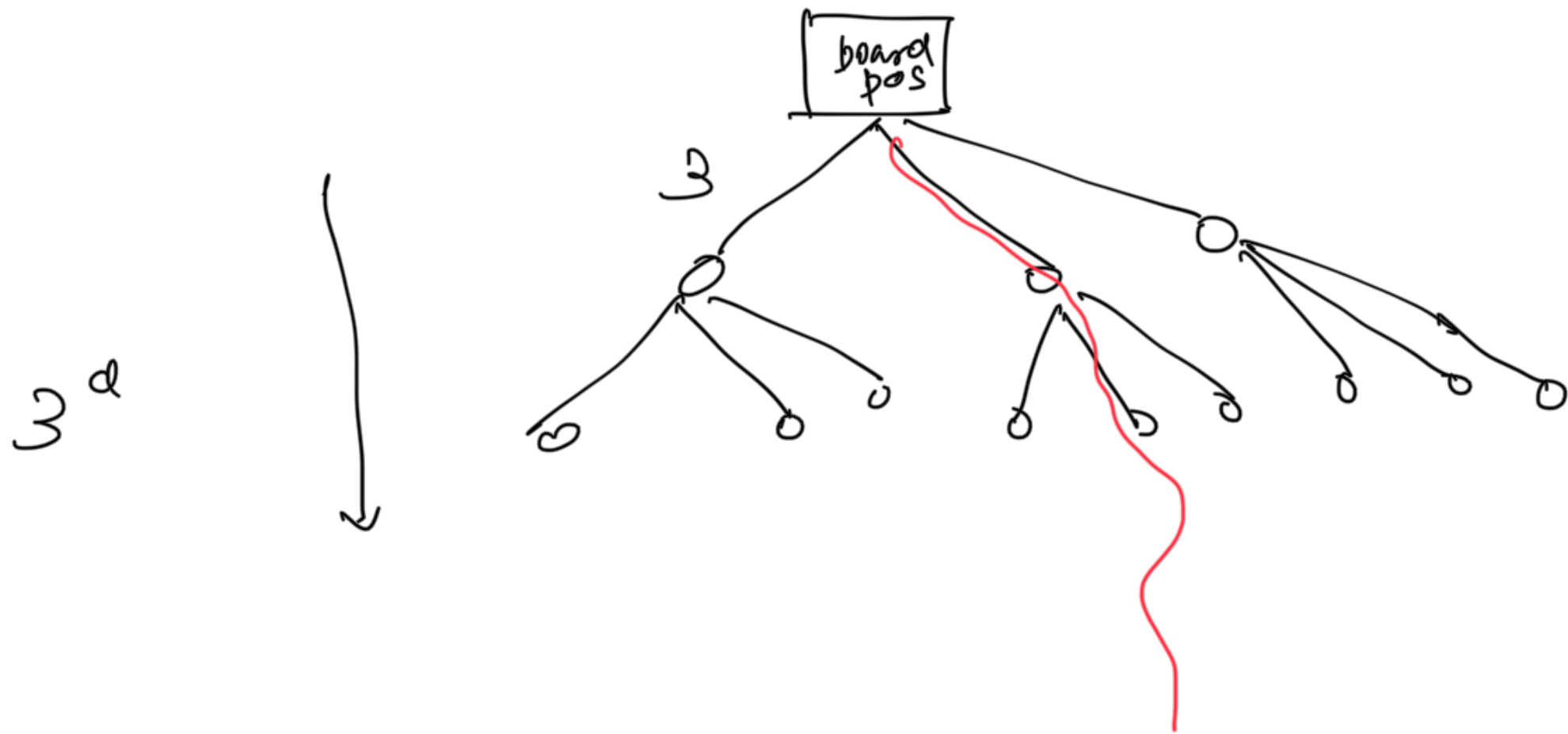
- at time t , we pick A_t (pulling arm)

- Reward $X_t \sim P_{A_t}$
- Mean rewards $\mu(a) = \mathbb{E}_P [X_t]$
- Best reward $\mu_* = \max_a \mu(a)$
- Best arm $a_* = \arg \max_a \mu(a)$
- Suboptimality Gaps $\Delta_a = \mu_* - \mu(a)$
- Goal ; Minimize

$$R_n = \text{Regret}_n = \mathbb{E} \left[\sum_{t=1}^n (\underbrace{\mu_* - X_t}_{\text{R.V.}}) \right] \quad n = 10$$

Why / where do we see the impact of
explore vs exploit

Example: Game



How to solve the MAB problem

• Sample mean \rightarrow True Mean
 Sequence /
 Convergence \rightarrow $\{x_i, y_{i,0}\}_{i=1}^n$ iid $\frac{\sum_{i=1}^n x_i}{n} \rightarrow \mathbb{E}[x_1]$
 of random variables

• Probability concentrates around true mean

Concentration \rightarrow
 σ
probability

