

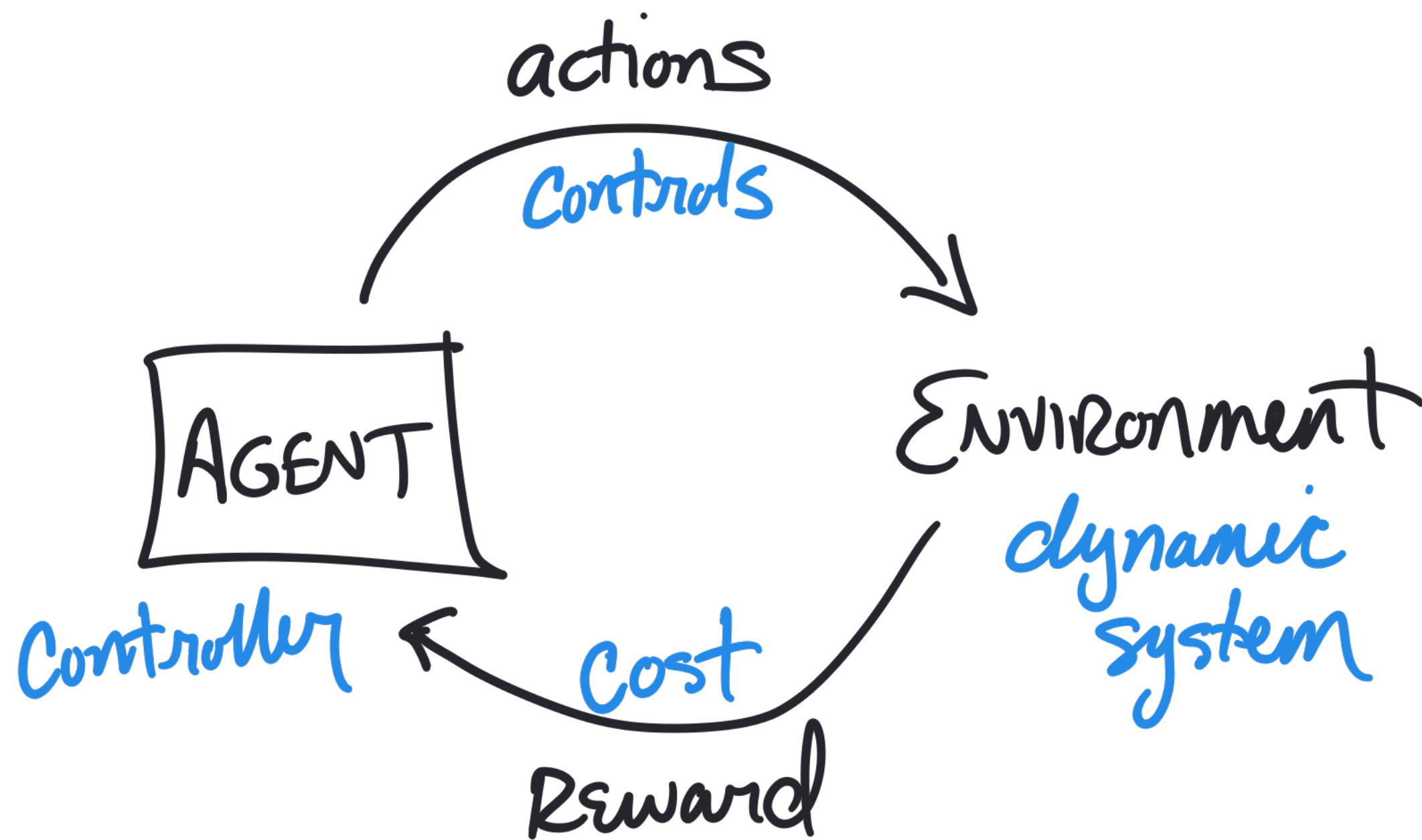
Introduction to Reinforcement Learning

BIOE 498/598 P1

RL

Computer sci
AI

optimal
control



OBJECTIVE: $\max \sum \text{Rewards}$
 $\min \sum \text{costs}$

objectives

[1. Exploration vs. Exploitation]

2. Policies

3. Deep RL

Multi-arm (k-arm) Bandit

Exploration

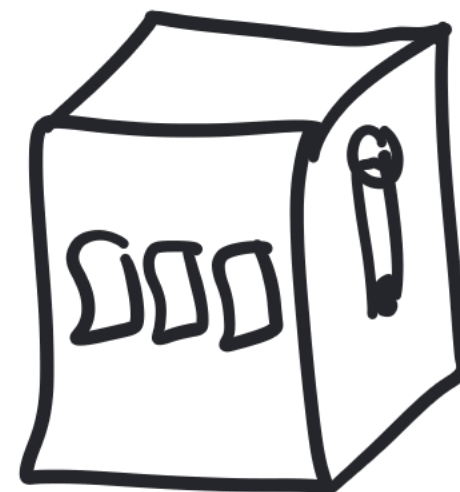
1. Determine the payout of all machines

2. Play the best machine.

Exploitation

In: 25¢

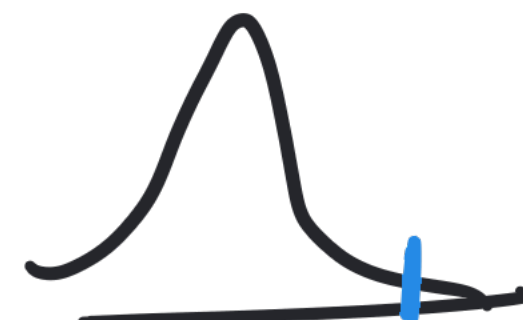
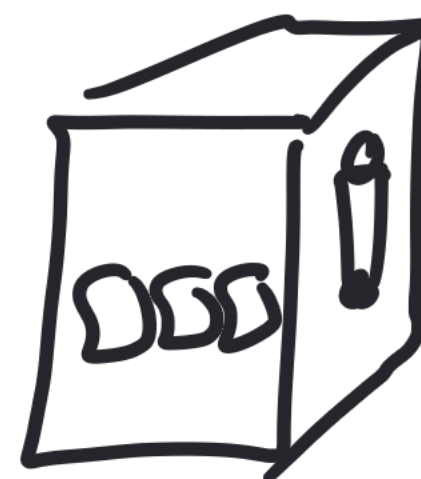
0



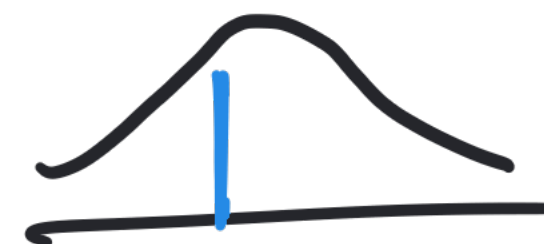
Payout



50¢



25¢



What is ~~the~~^a solution to Exploration vs. Exploitation?

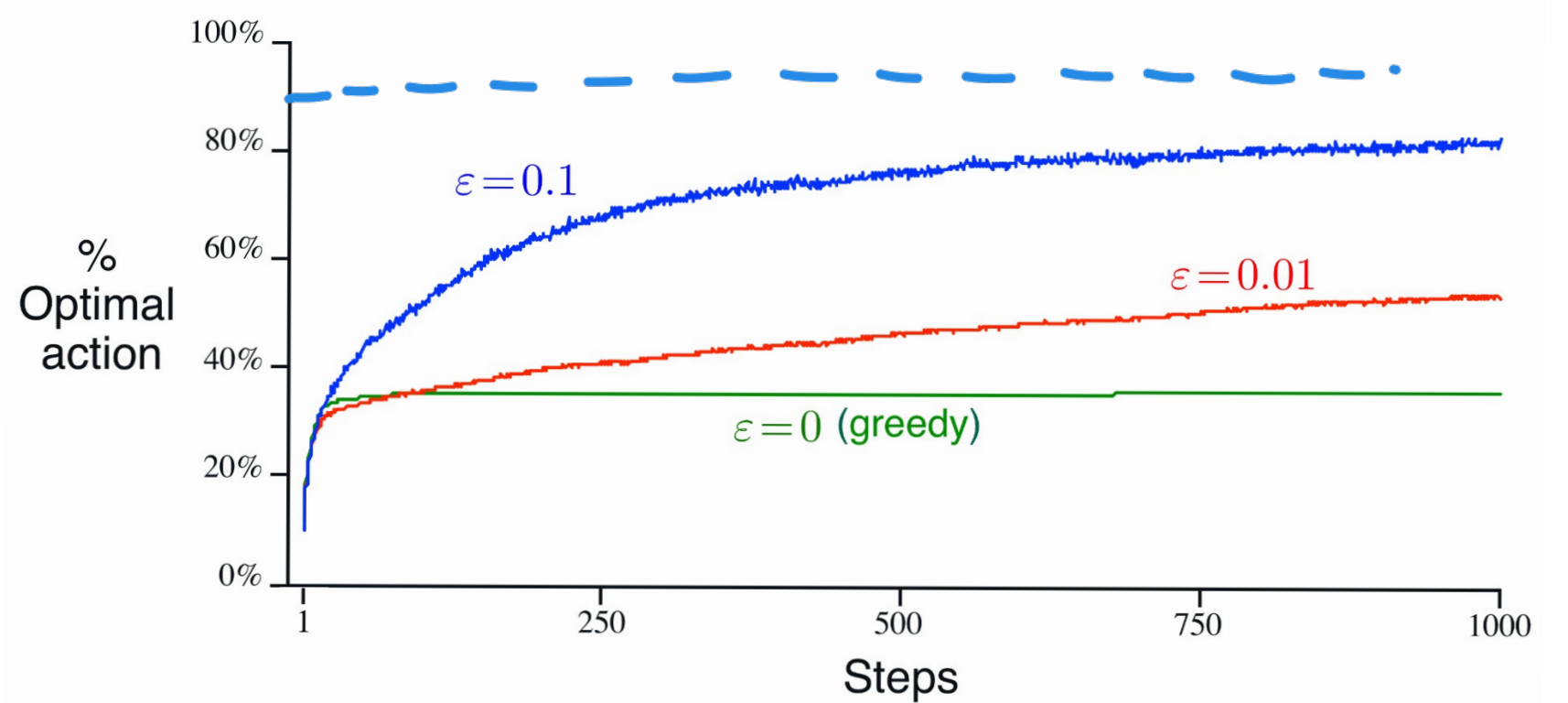
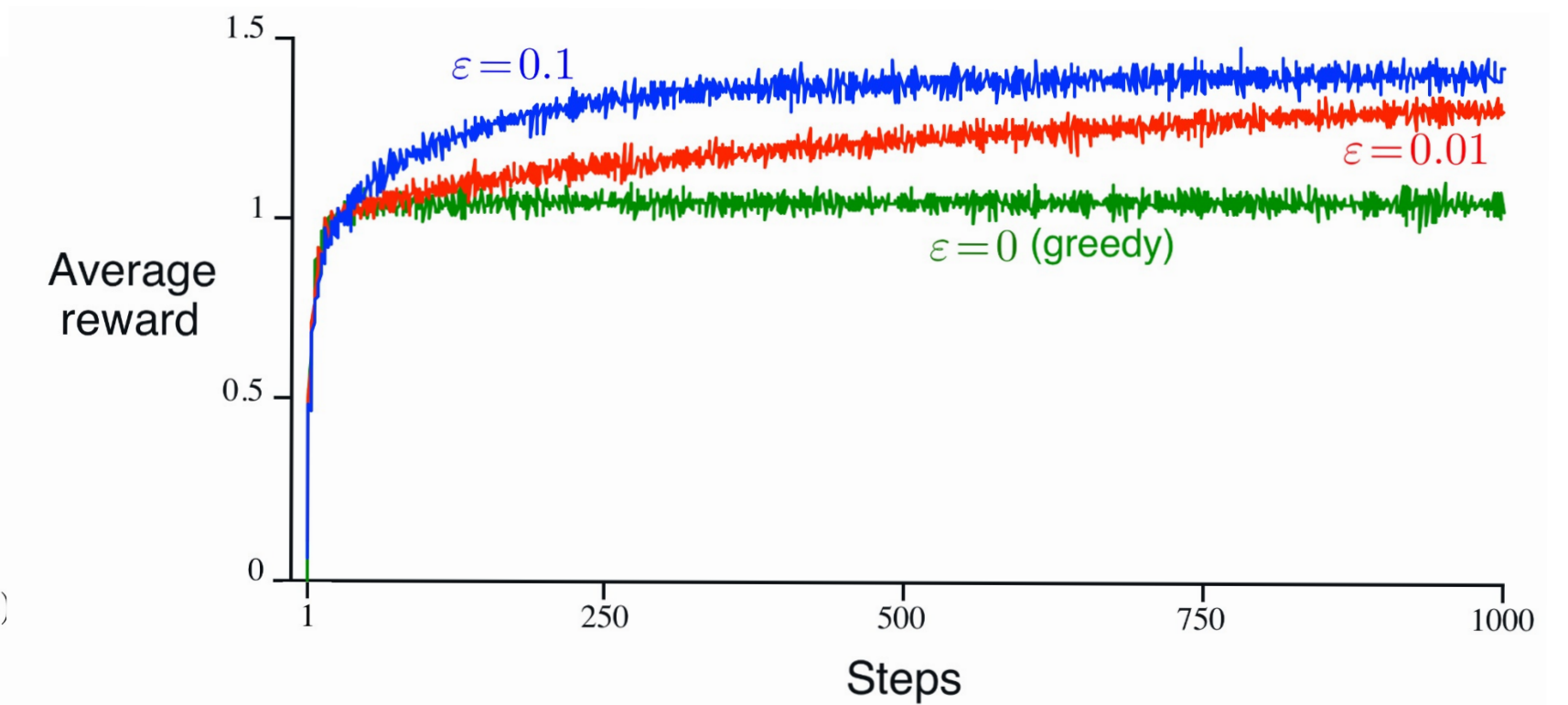
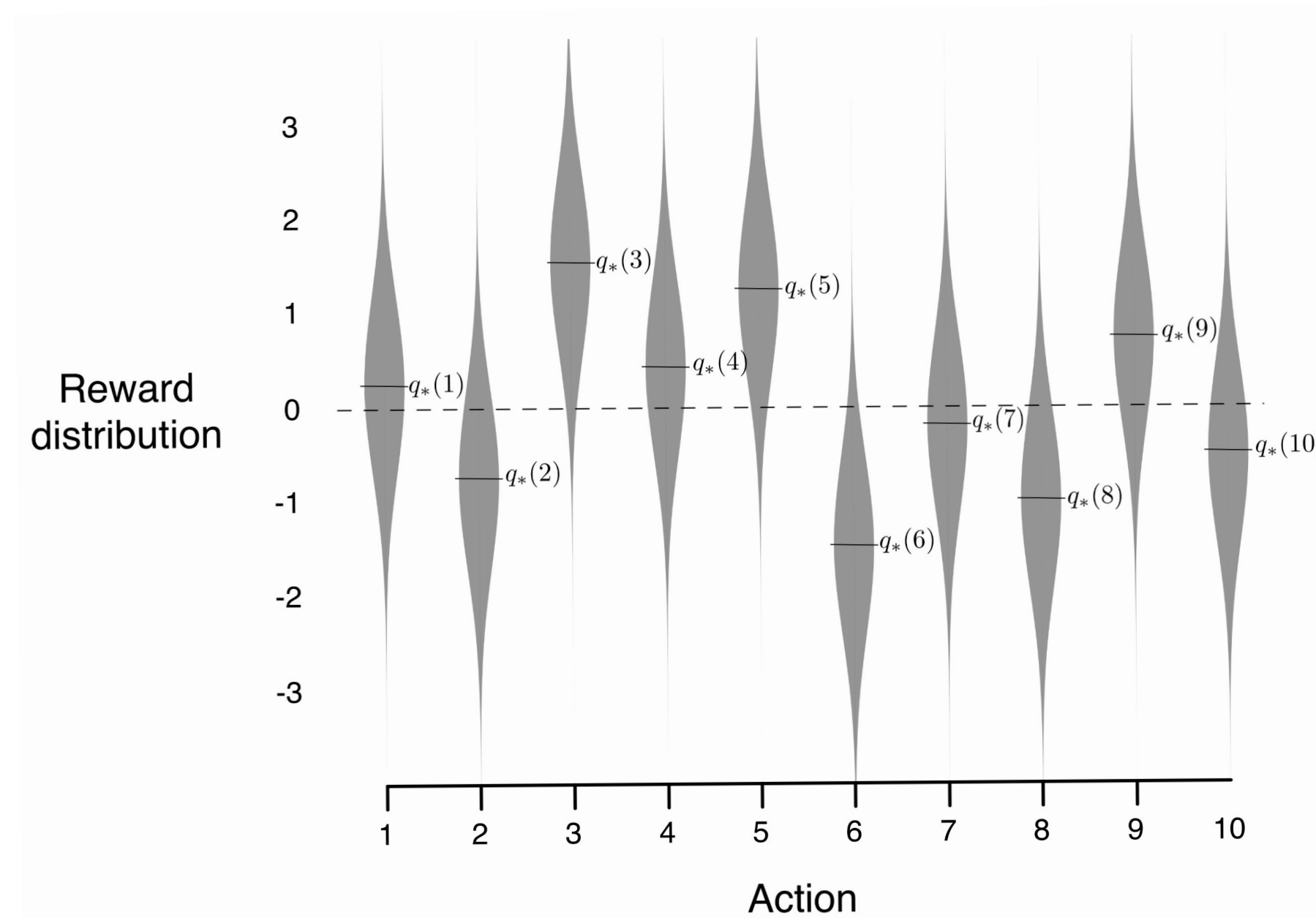
1. Do the optimal action most of the time. Exploit
2. Some of the time, do random actions. Explore

ϵ -greedy policy

fraction of times
where we explore

90% exploit, 10% explore $\rightarrow \epsilon = 0.1$

ϵ -greedy policies for 10-arm bandits



Why does this look familiar?

Explore

Resolution III : Screening, lots of factors
confounded effects



Exploit

RSM / CCD : Very intensive, few factors
clean, predictive