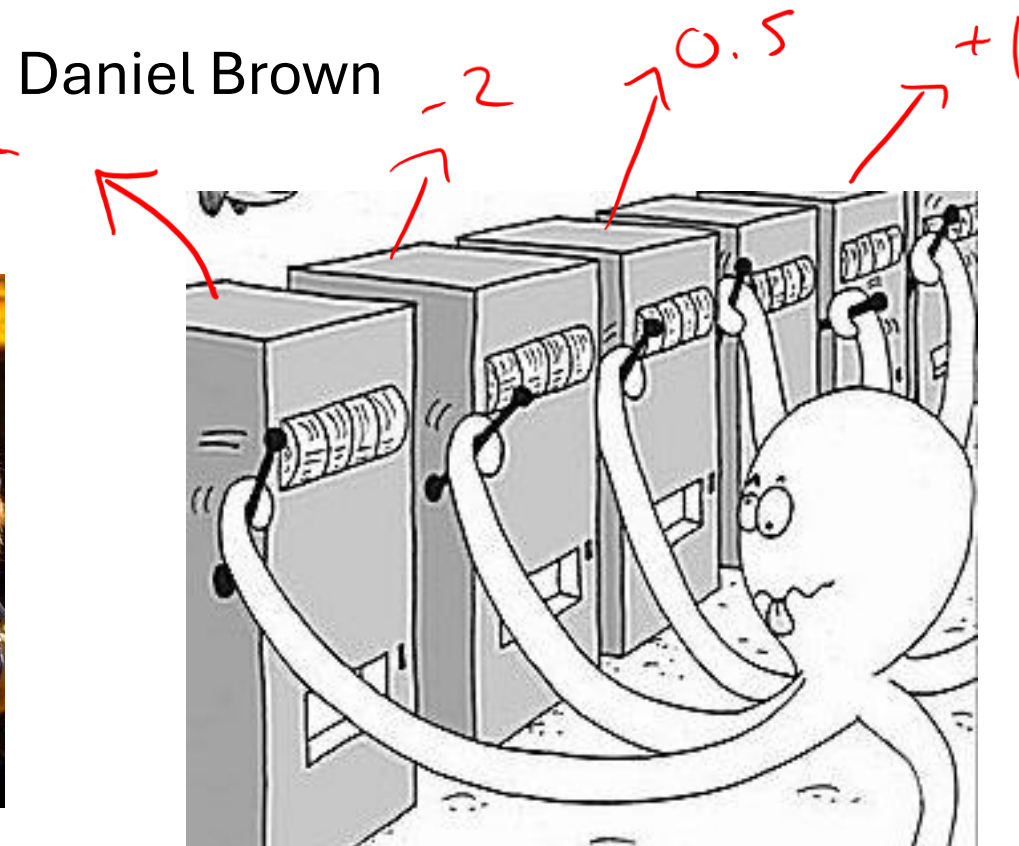


# Multi-Armed Bandits



# Evaluative feedback

Not supervised learning  
- We've not given labels



## REPORT CARD

Reading	B
Writing	C-
Mathematics	D
Science	C-
History	B+
Art	B-
P.E.	B

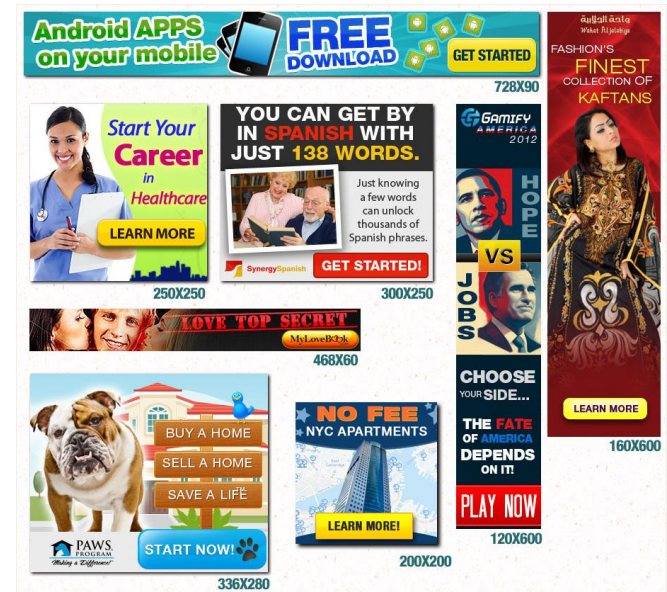
## Vital Signs Basics



MAB - reward  
- action space

# Applications

- Online Advertising and Recommendation
- Clinical Trials
- Robotics
- Dynamic Pricing
- Search Engine Optimization
- Education and Learning Platforms



$$\pi \rightarrow \operatorname{argmax}_i \mu_i \quad \mu_i = \mathbb{E}[r(a_i)]$$

# Problem formalism

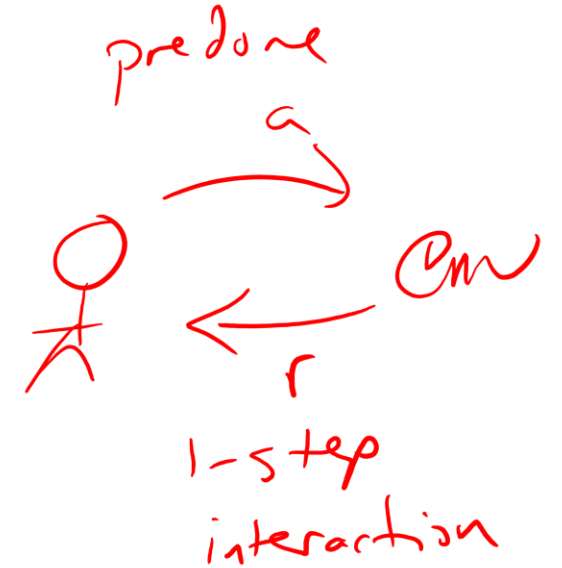
*Actions*

- Arms  $\mathcal{A} = \{a_1, \dots, a_k\}$ 
  - Each arm is associated with an unknown reward distribution
- Rewards  $r_t(a_i)$
- Possible Goals
  - Maximize cumulative reward (Minimize regret)
  - Best arm identification
- Standard Assumptions
  - Independence: Rewards from each arm are independent
  - Stationarity: Reward distributions don't change over time

$$r_t(a_1) \sim \mathcal{N}(0, 1) \quad \forall t$$

$$r(a_2) \sim \mathcal{U}[-2, 2]$$

$$r(a_3) = \mathcal{N}(10, 100)$$



stochasticity is allowed

# How should we solve this problem?

1) initial exploration of actions  
- figure out which has highest mean reward and take that action

median, std, full dist.

"Smart" exploration

Exploitation

Balance exploration & exploitation

Random

= maybe good at first

pure exploration

bad exploitation

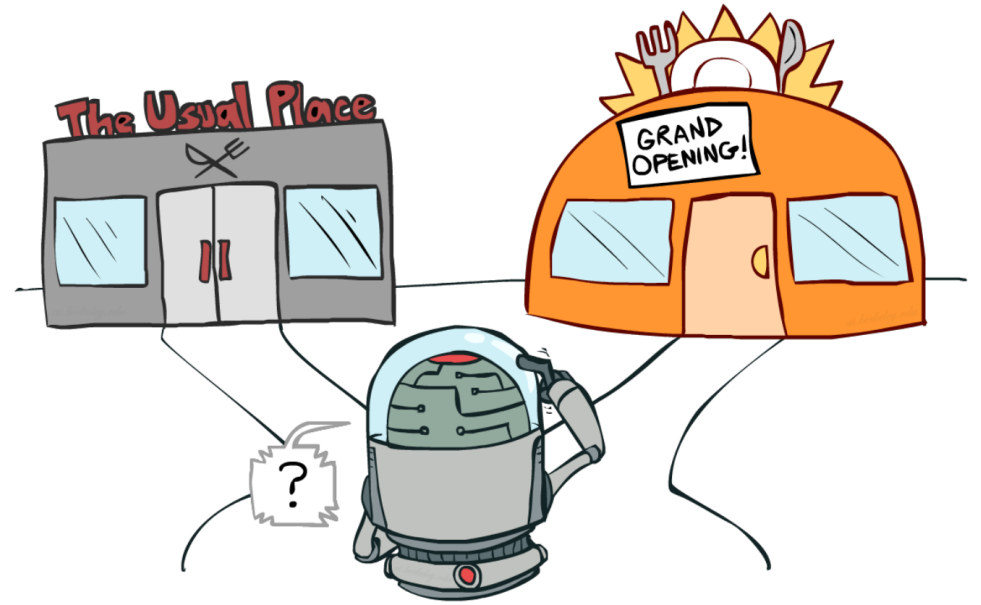
# Greedy

$\bar{\mu}_i = \text{sample ave}$

	$a_1$	$a_2$	$a_3$
$t=1$	0	0	1
$t=2$	0	0	1
$t=3$			1
$t=4$			-1
	$\bar{\mu}_1 = 0$	$\bar{\mu}_2 = 0$	$\bar{\mu}_3 = 0$

Not Good

# Exploration





# $\epsilon$ -Greedy

$\epsilon \in [0, 1]$  determines the prob  
of taking a rand  
action

pick  $\epsilon$

randomly generate  $x \in [0, 1]$

if  $x \leq \epsilon$

take rand action uniformly

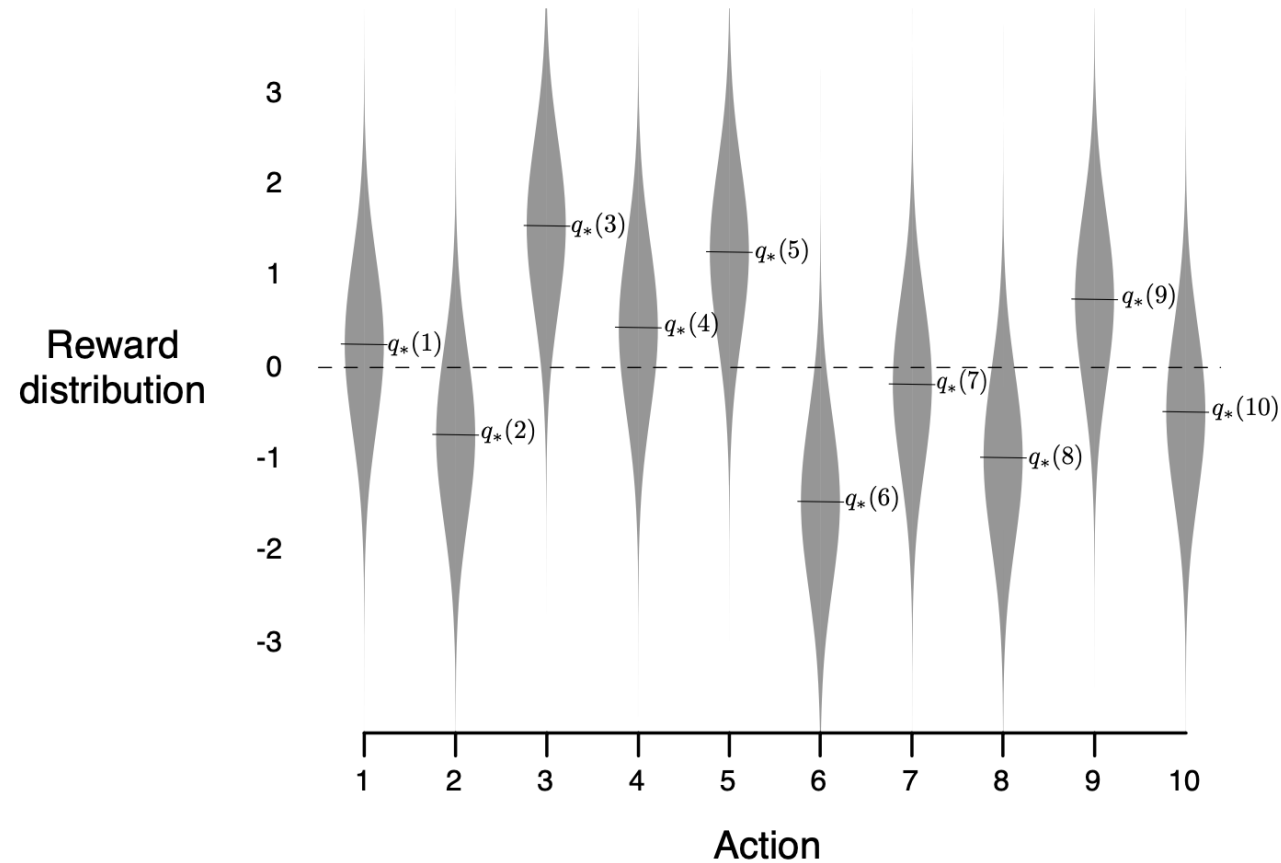
else

take greedy action

\* Anneal  $\epsilon$  : start  $\epsilon = 1$  then  $\epsilon \rightarrow 0$  as  $t \rightarrow \infty$

## 2.3 The 10-armed Testbed

To roughly assess the relative effectiveness of the greedy and  $\varepsilon$ -greedy action-value methods, we compared them numerically on a suite of test problems. This was a set of 2000 randomly generated  $k$ -armed bandit problems with  $k = 10$ . For each bandit problem, such as the one shown in Figure 2.1, the action values,  $q_*(a)$ ,  $a = 1, \dots, 10$ ,

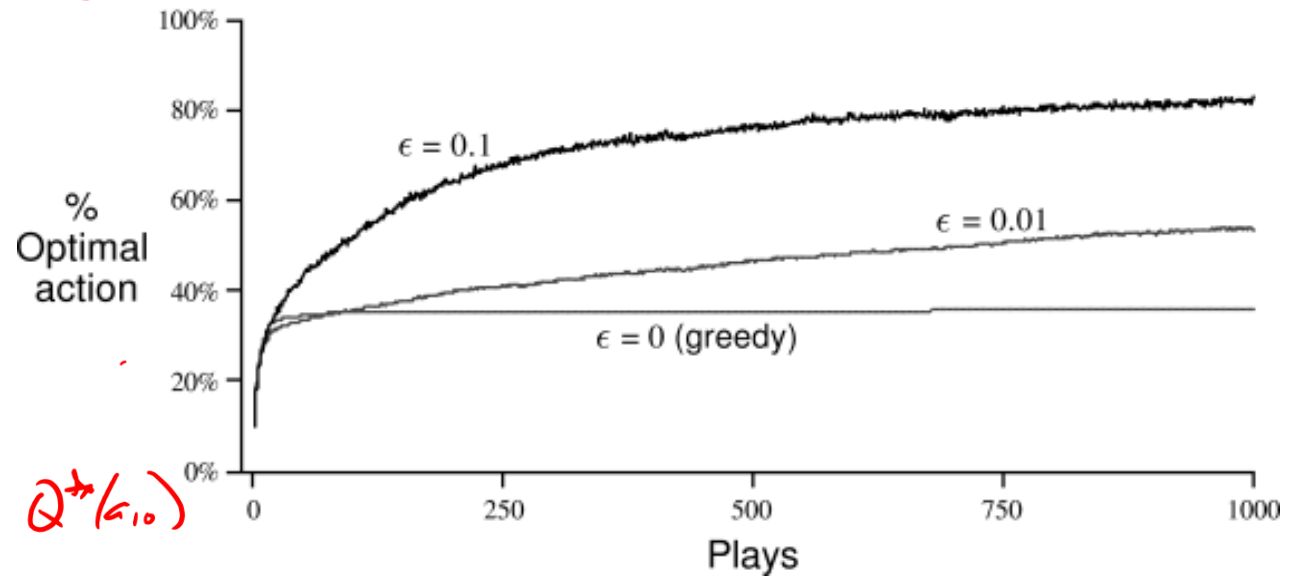
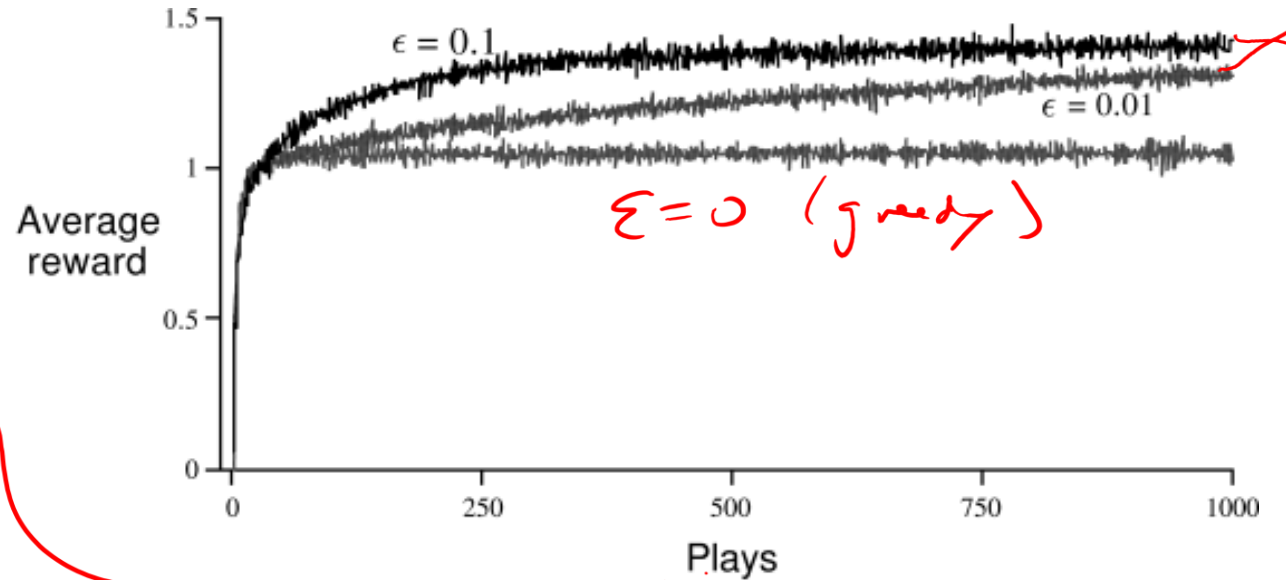
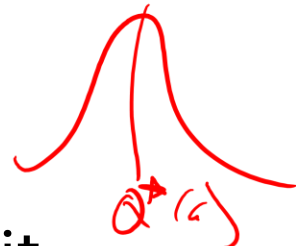


**Figure 2.1:** An example bandit problem from the 10-armed testbed. The true value  $q_*(a)$  of each of the ten actions was selected according to a normal distribution with mean zero and unit variance, and then the actual rewards were selected according to a mean  $q_*(a)$  unit variance normal distribution, as suggested by these gray distributions.

# Sutton/Barto figure

- 10 arms
- Each arm has stochastic reward  
 $r \sim N(Q^*(a), 1)$
- Averaged over 2000 bandit problems where each problem starts with  $Q^*(a) \sim N(0,1)$  for all  $a$

$Q^*(a) = \mathbb{E}[r(a)]$   
sample 10 times to get  $Q^*(a_1) \dots Q^*(a_{10})$



# Problems?

never stops exploring  
explores random

# Boltzmann (Softmax) Exploration

$$Q(a) = \frac{1}{n} \sum_{i=1}^n r_i$$

sample average

$$P(a_i) = \frac{\exp(\beta Q(a_i))}{\sum_a \exp(\beta Q(a))}$$

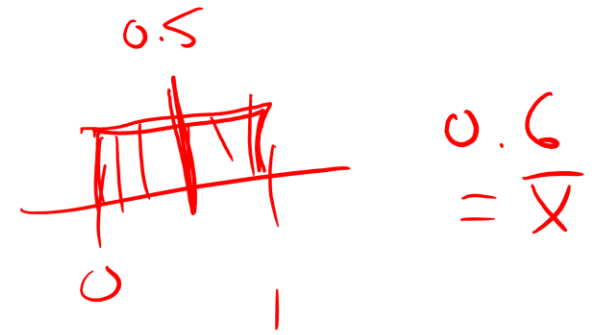
$\beta = \text{inv. temp}$

$\beta = 0$   
 $\Rightarrow$  uniform random

$\beta \rightarrow \infty$   
 $\Rightarrow$  greedy

$\beta = 1, 10$

# Chernoff-Hoeffding Inequality



- Let  $X$  be a random variable in the range  $[0, 1]$  and  $x_1, x_2, \dots, x_n$  be  $n$  independent and identically distributed samples of  $X$ .

- Let  $\bar{X} = \frac{1}{n} \sum_i x_i$  (the empirical average)

- Then we have  $P(\bar{X} \geq \mathbb{E}[X] + c) \leq e^{-2nc^2}$

$$P(\bar{X} \leq \mathbb{E}[X] - c) \leq e^{-2nc^2}$$

$$\Rightarrow P(|\bar{X} - \mathbb{E}[X]| \geq c) \leq 2e^{-2nc^2}$$

$$n \rightarrow \infty$$

$$\frac{1}{e^{2nc^2}} \rightarrow 0$$

$$c \rightarrow \infty$$

$$c = 0.4$$

$$\bar{X} = 0.9$$

# Some fun math

- $P(\bar{X} \geq \mathbb{E}[X] + c) \leq e^{-2nc^2}$
- Typically, we want to pick some kind of high confidence  $1 - \delta$  such that we are very confident about our sample mean being close to the true expectation.
- If we want

$$P(\bar{X} \geq \mathbb{E}[X] + c) \leq \delta$$

What is  $c$  in terms of  $\delta$ ?

$$\delta = e^{-2nc^2} \Rightarrow \log \delta = -2nc^2$$
$$\Rightarrow \sqrt{\frac{\log \delta}{-2n}} = c \Rightarrow$$

$n = 10$   
 $\delta = 0.05$

$-\log \delta = \log(1/\delta)$

$$c = \sqrt{\frac{1}{2n} \log(1/\delta)}$$





# More math

$$C = \sqrt{\frac{1}{2n} \log(1/\delta)}$$

- We can pick  $\delta$  to be whatever we want, so let's pick
- If we select  $\delta = \frac{1}{t^2}$

$$\delta \rightarrow 0 \quad t \rightarrow \infty$$

$$a \log b = \log b^a$$

What is  $c$ ?

$$C = \sqrt{\frac{1}{2n} \log t^2} = \sqrt{\frac{2}{2n} \log t}$$

$$P(\bar{R} \geq \mathbb{E}[R] + c) \leq \frac{1}{t^2} \quad C = \sqrt{\frac{\log t}{n}}$$

# UCB1 (UCB = Upper Confidence Bound)

arms  $a_1 \dots a_K$ , counts  $n_1 \dots n_K$

Key Idea: Optimism in the face of uncertainty

→ • Play each action once to get initial averages of arm values

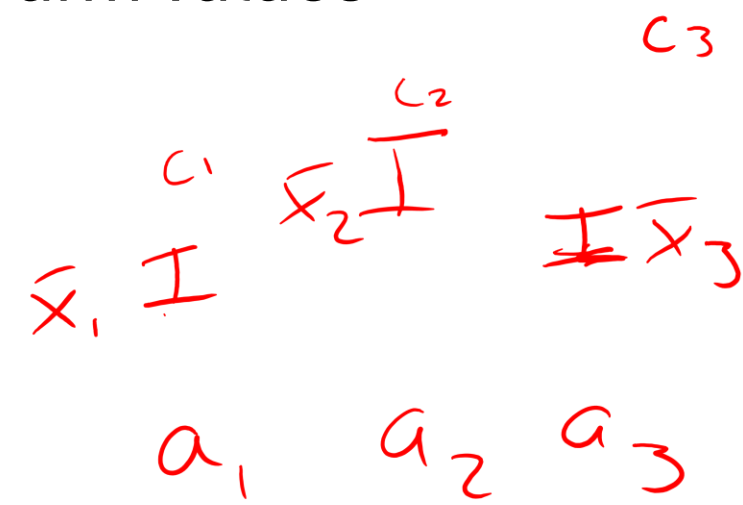
• Keep track of counts of pulls for each arm  $n_i$

• At each step  $t$ , select  $\arg \max_i \bar{X}_i + c(i, t)$

• Where  $c(i, t) = \sqrt{\frac{\log(t)}{n_i}}$

$n_i \rightarrow \infty$   
 $c(i, t) \rightarrow 0$   
 $t \rightarrow \infty$

sample average for arm  $i$



# Regret

$$Q^* = \max_a Q^*(a)$$

- Define  $\mu^*$  as the maximum expected payoff over all  $k$  arms
- $\text{Regret}(T) = T\mu^* - \sum_{t=1}^T r_t$

- Epsilon-Greedy Regret

- $O(T)$

- UCB1 Regret

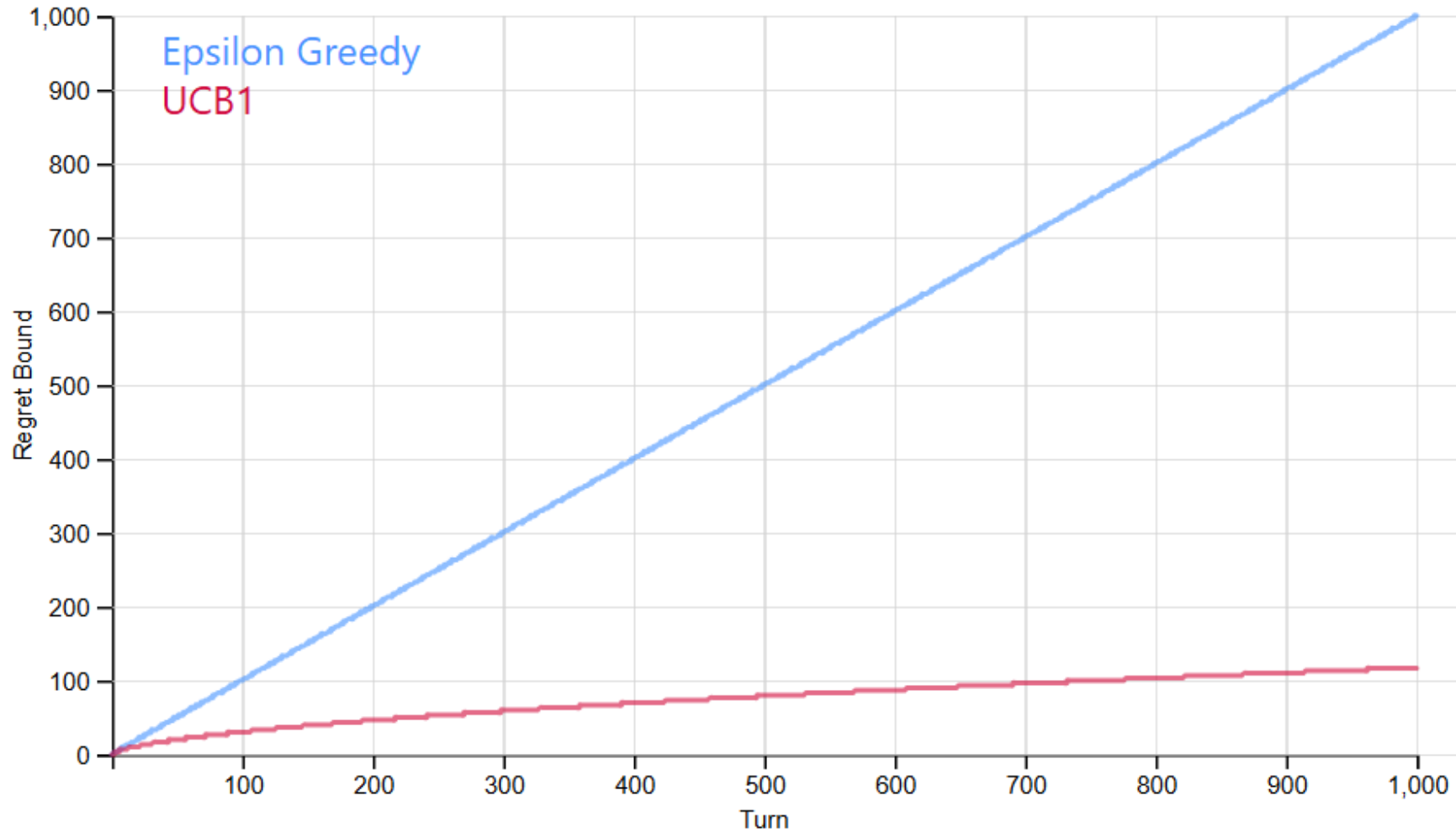
- $O(\sqrt{kT \log(T)})$

↖ # arms

$T = \# \text{ time steps}$

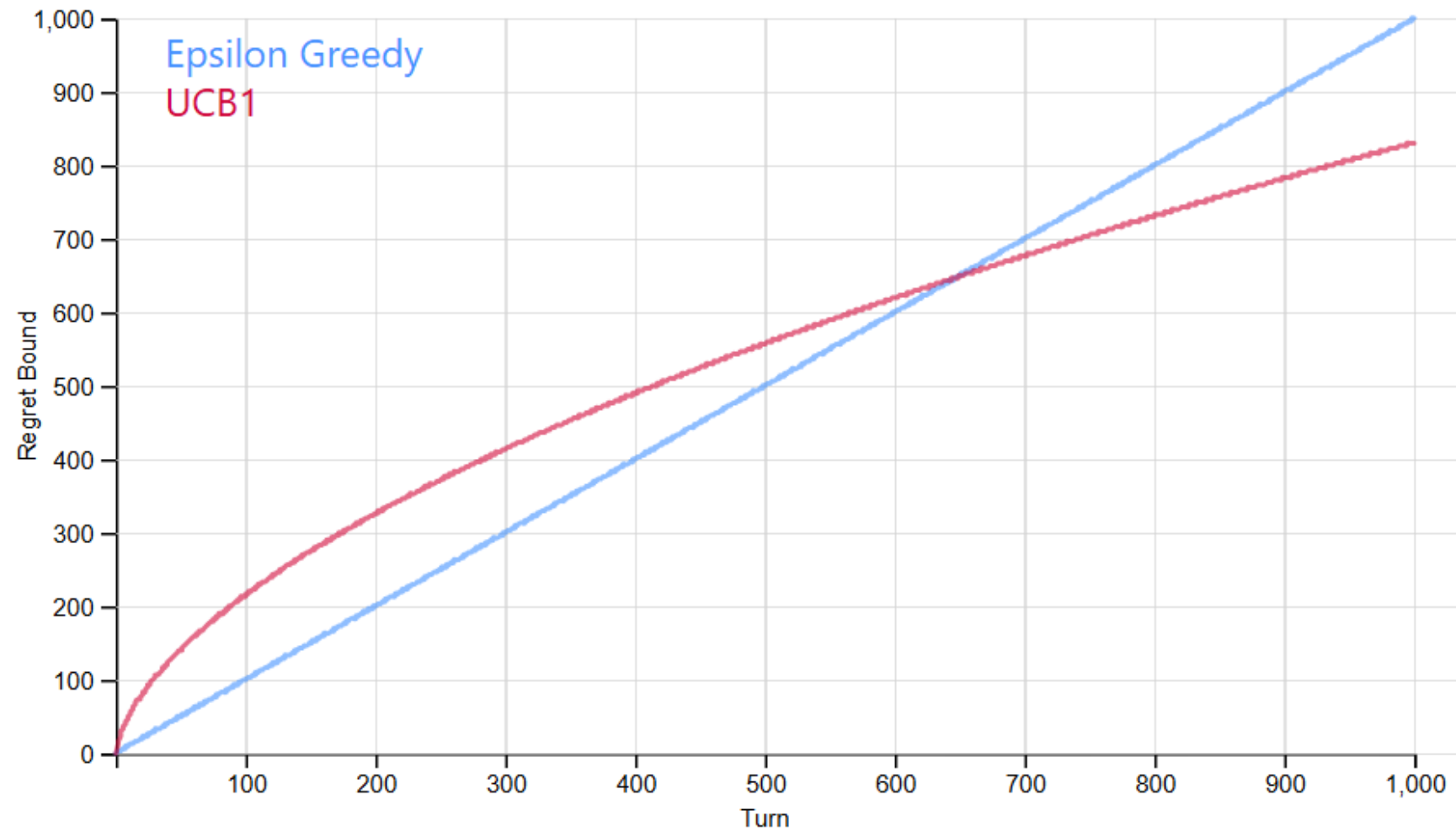
- A **No-Regret** algorithm is such that  $\text{Regret}(T)/T \rightarrow 0$  as  $T \rightarrow \infty$ 
  - Average regret goes to zero

## Regret Bound vs. Turn



$k$  (number of arms):   $T$  (number of steps):

## Regret Bound vs. Turn

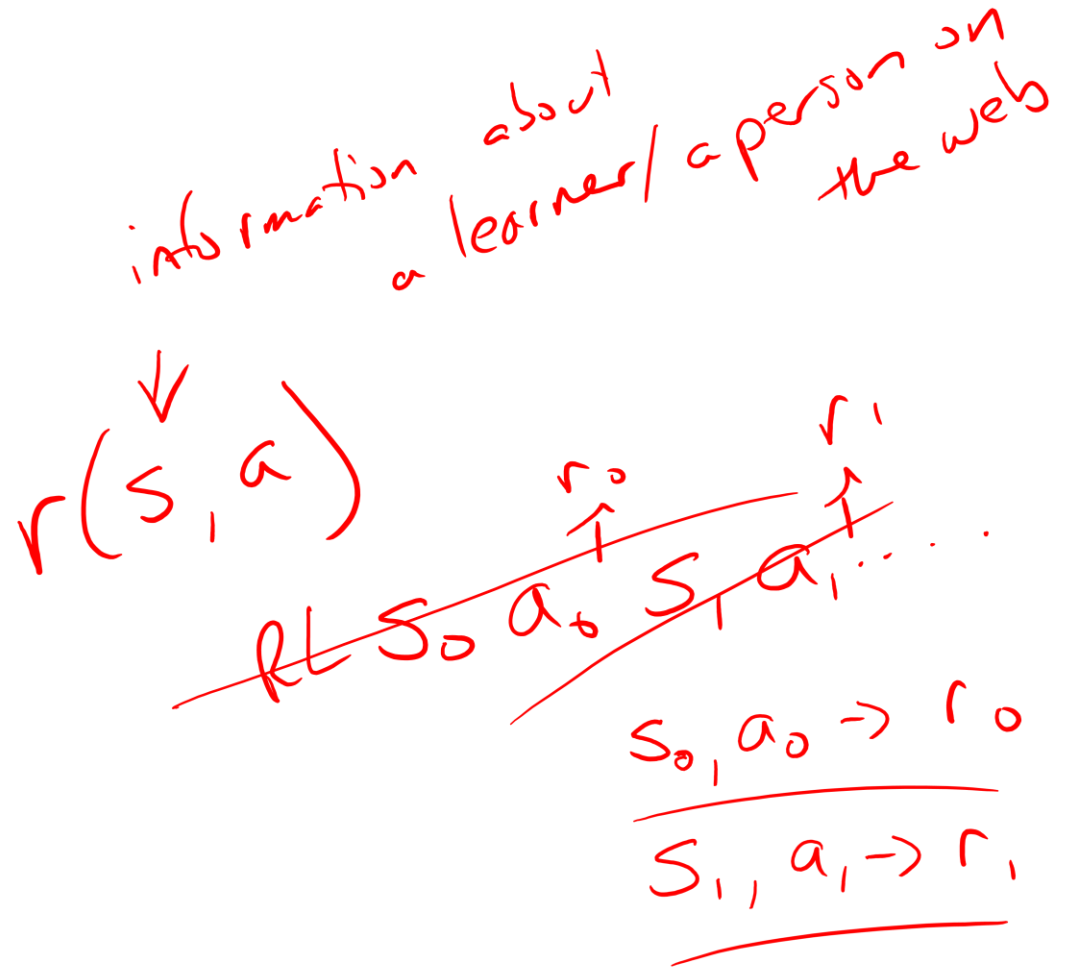


$k$  (number of arms):   $T$  (number of steps):

# Other Bandit Topics

- Thompson Sampling
- Best Arm Identification
- Adversarial Bandits
- Contextual Bandits
  - State information,  $s_t$
  - Reward depends on state, and action
- Linear Bandits
  - Type of contextual bandit
  - Reward is a linear combination of state features.

$$r(s, a) = \theta^T \phi(s, a) = \sum_{i=1}^n \theta_i \phi_i(s, a)$$



known