

max Total Reward
min Ave Regret

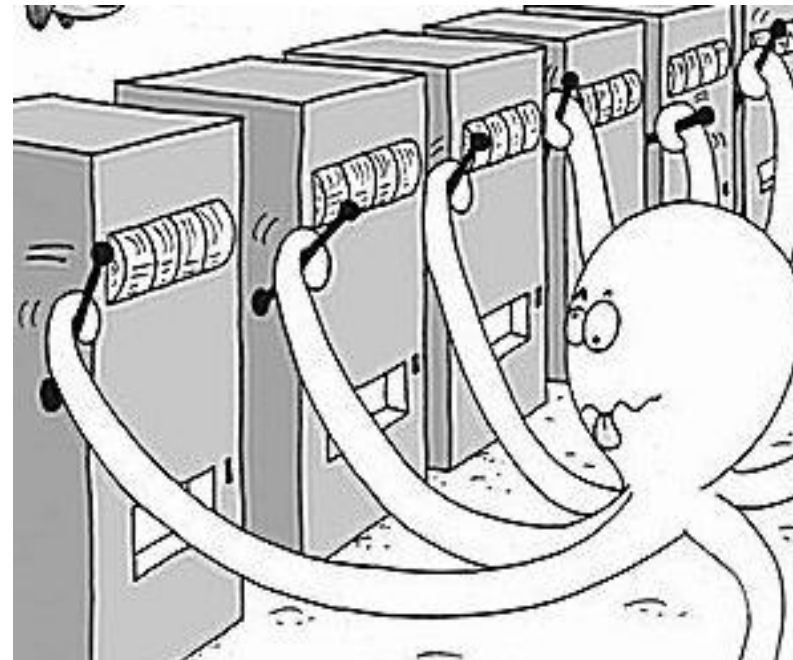
$$A = \{1, \dots, n\}$$

$v_a^* \sim \text{Distribution}$

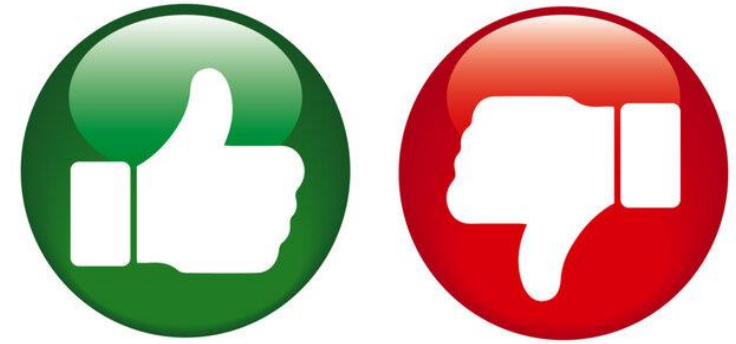
Multi-Armed Bandits

~~$D = \{(s, a) \dots (s', a)\}$~~

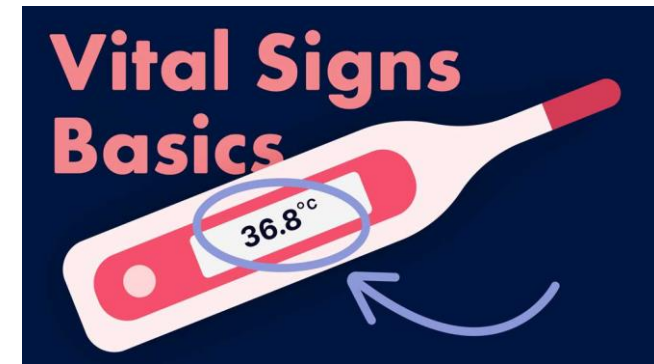
Daniel Brown



Evaluative feedback



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

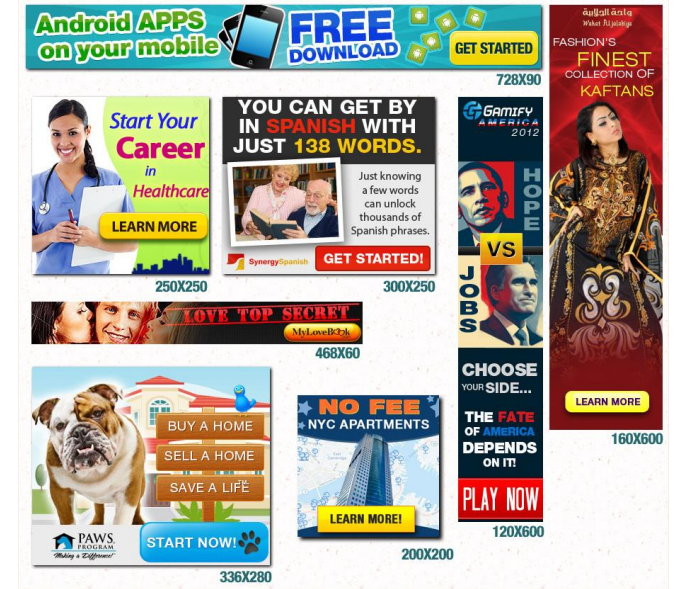


Applications

$A^?$ $R^?$
 $\pi \rightarrow A$

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

$MAB \rightarrow$ Contextual Bandit

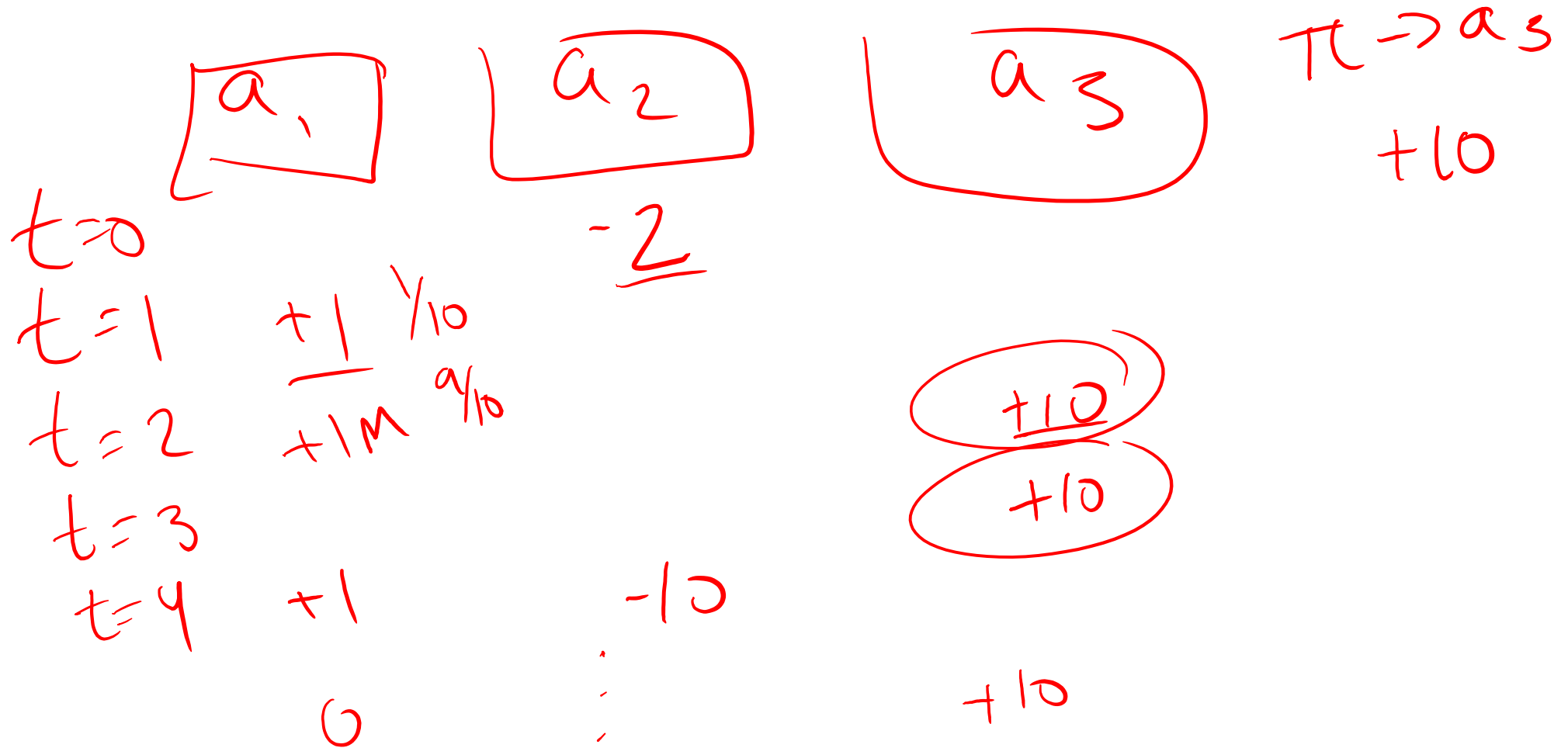


Problem formalism

- ^{Actions} Arms $\mathcal{A} = \{a_1, \dots, a_k\}$
 - Each arm is associated with an unknown reward distribution
- Rewards $r_t(a_i)$ - $\mathbb{E}[r(a_i)] = \mu_{a_i}$
- 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

\hookrightarrow if visited $\epsilon \triangleq 0.1$

How should we solve this problem?

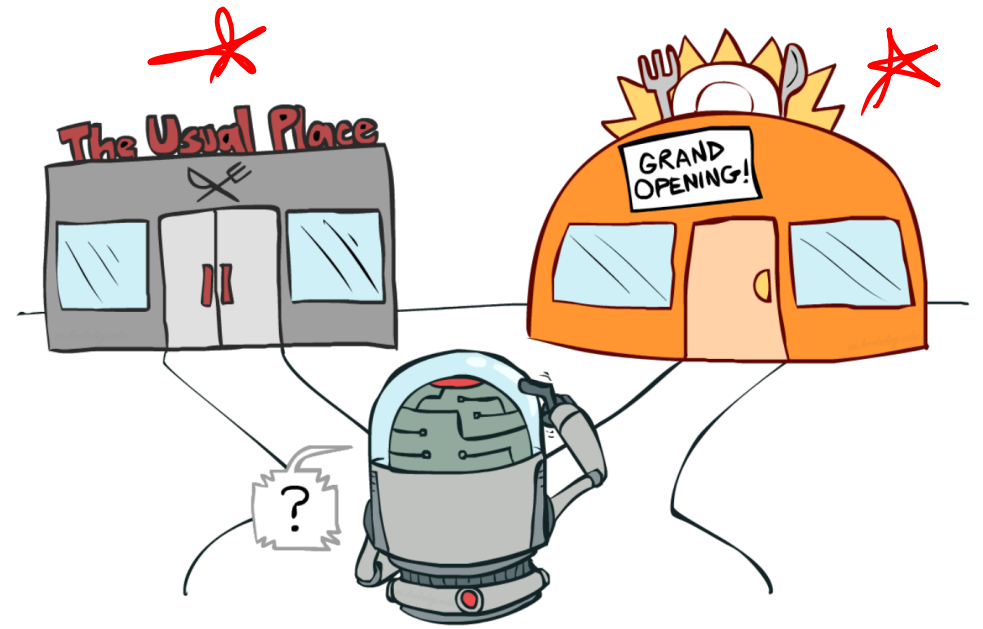


Random

Big Regret

Greedy

Exploration



ϵ -Greedy

$\epsilon \in (0, 1)$

e.g. 0.1, 0.05, 0.01
* anneal/reduce/decay
 ϵ over time

loop: rand $\in [0, 1]$

random (if rand $< \epsilon$:
explore action at random

greedy (else:
act greedy

2.3 The 10-armed Testbed

To roughly assess the relative effectiveness of the greedy and ε -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, $q_*(a)$, $a = 1, \dots, 10$,

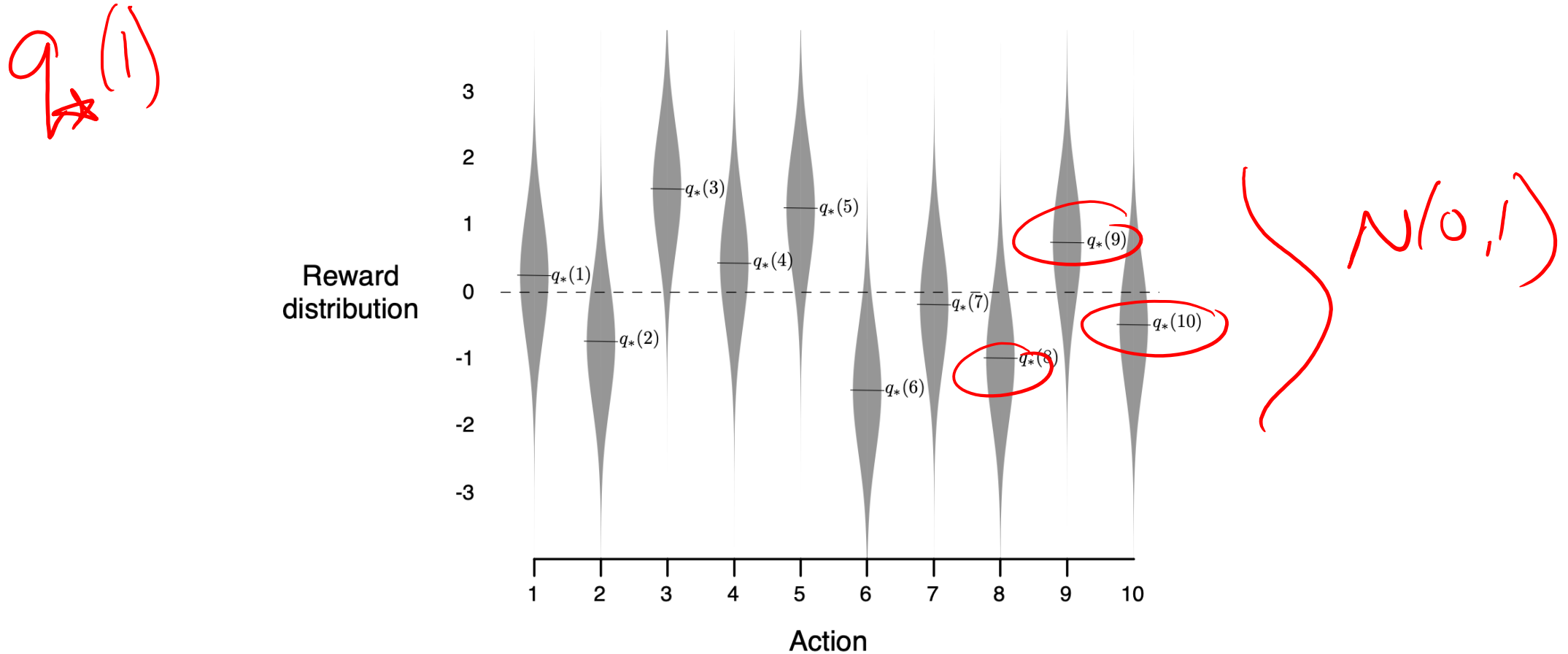


Figure 2.1: An example bandit problem from the 10-armed testbed. The true value $q_*(a)$ of each of the ten actions was selected according to a normal distribution with mean zero and unit variance, and then the actual rewards were selected according to a mean $q_*(a)$ unit variance normal distribution,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

✓ ↓ ↓

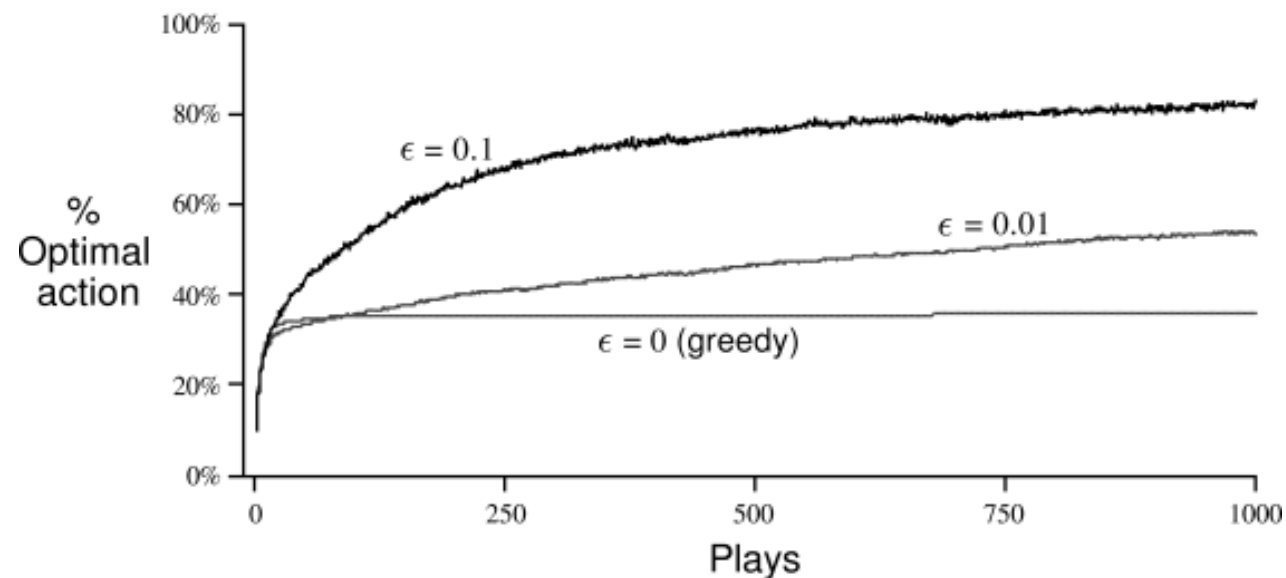
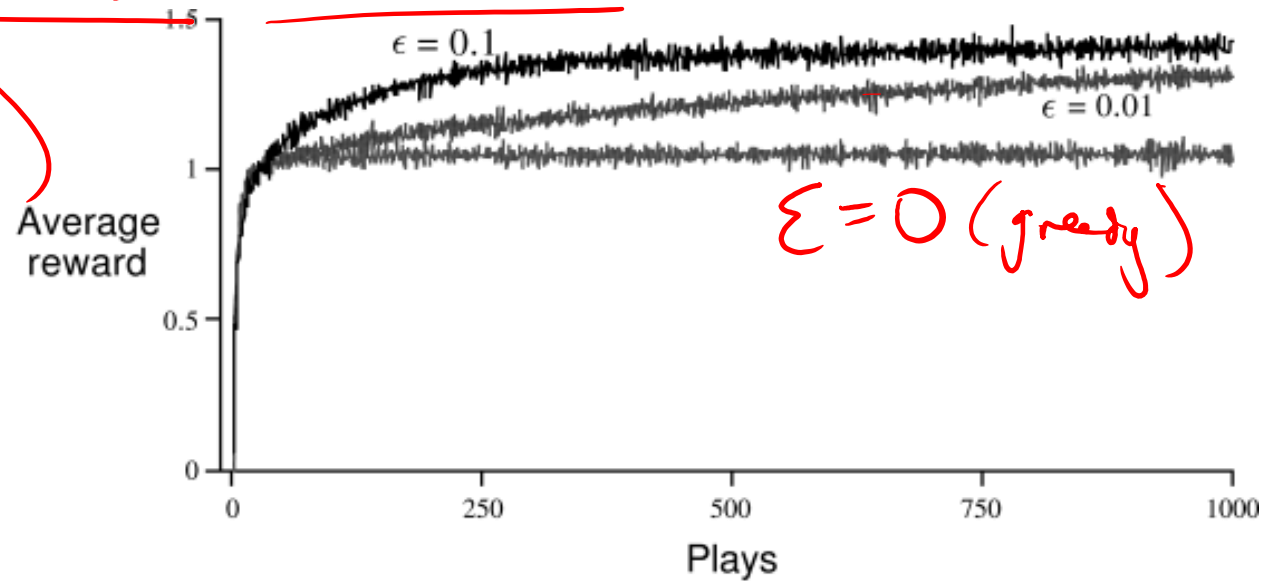
$N(0.1, 1)$ $N(-0.2, 1)$ $N(1.5, 1)$

$$s \sim N(a, b)$$

$$r \sim N(0, 1)$$

$$s = a + br$$

random, random(
 $\sim N(0, 1)$



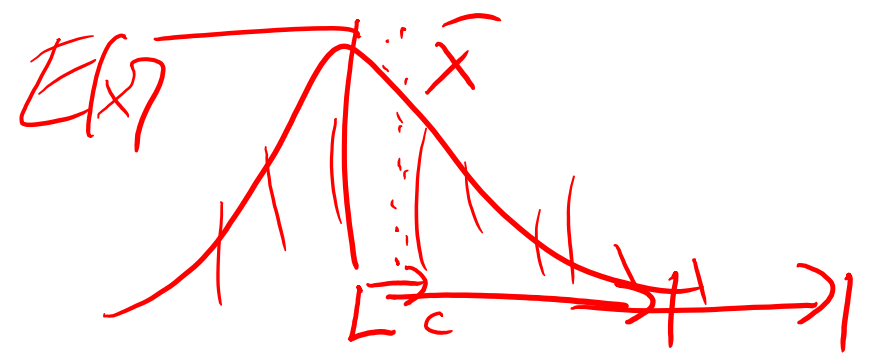
Problems?

$\epsilon?$

Boltzmann (Softmax) Exploration

$$\Pr(a) = \frac{e^{\beta \hat{Q}(a)}}{\sum_{a' \in A} e^{\beta \hat{Q}(a')}}$$

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 . - $r \in$
- Let $\bar{X} = \frac{1}{n} \sum_i x_i$ (the empirical average) $\hat{Q}(a)$
- Then we have $P(\bar{X} \geq \mathbb{E}[X] + c) \leq e^{-2nc^2}$

Some fun math!

Handwritten notes in red:

$\mathbb{E}[X] \xrightarrow{c} \bar{X}$ 99% 95% $\delta = 0.05$

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

Handwritten notes in red:

$\delta = e^{-2nc^2}$ 0.01

$$c = \sqrt{\frac{\log \delta}{-2n}}$$

More math

- We can pick δ to be whatever we want, so let's pick
- If we select $\delta = t^{-4}$

What is c ?

Handwritten notes and calculations:

$\leftarrow c \rightarrow$
 $E[X] \quad \frac{1}{X}$
 $t \rightarrow \infty \quad \frac{1}{t^4}$
 $\delta \rightarrow 0$

$$C = \sqrt{\frac{\log \delta}{-2n}} = \sqrt{\frac{\log t^{-4}}{-2n}} = \sqrt{\frac{-4 \log t}{-2n}}$$

Result boxed:

$$= \sqrt{\frac{2 \log t}{n}}$$

UCB1 (UCB = Upper Confidence Bound)

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 $c(i, t) = \sqrt{\frac{2 \cdot \log(t)}{n_i}}$

$t =$ total # actions

$t = \infty$



Regret

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

- Epsilon-Greedy Regret
 - $O(T)$

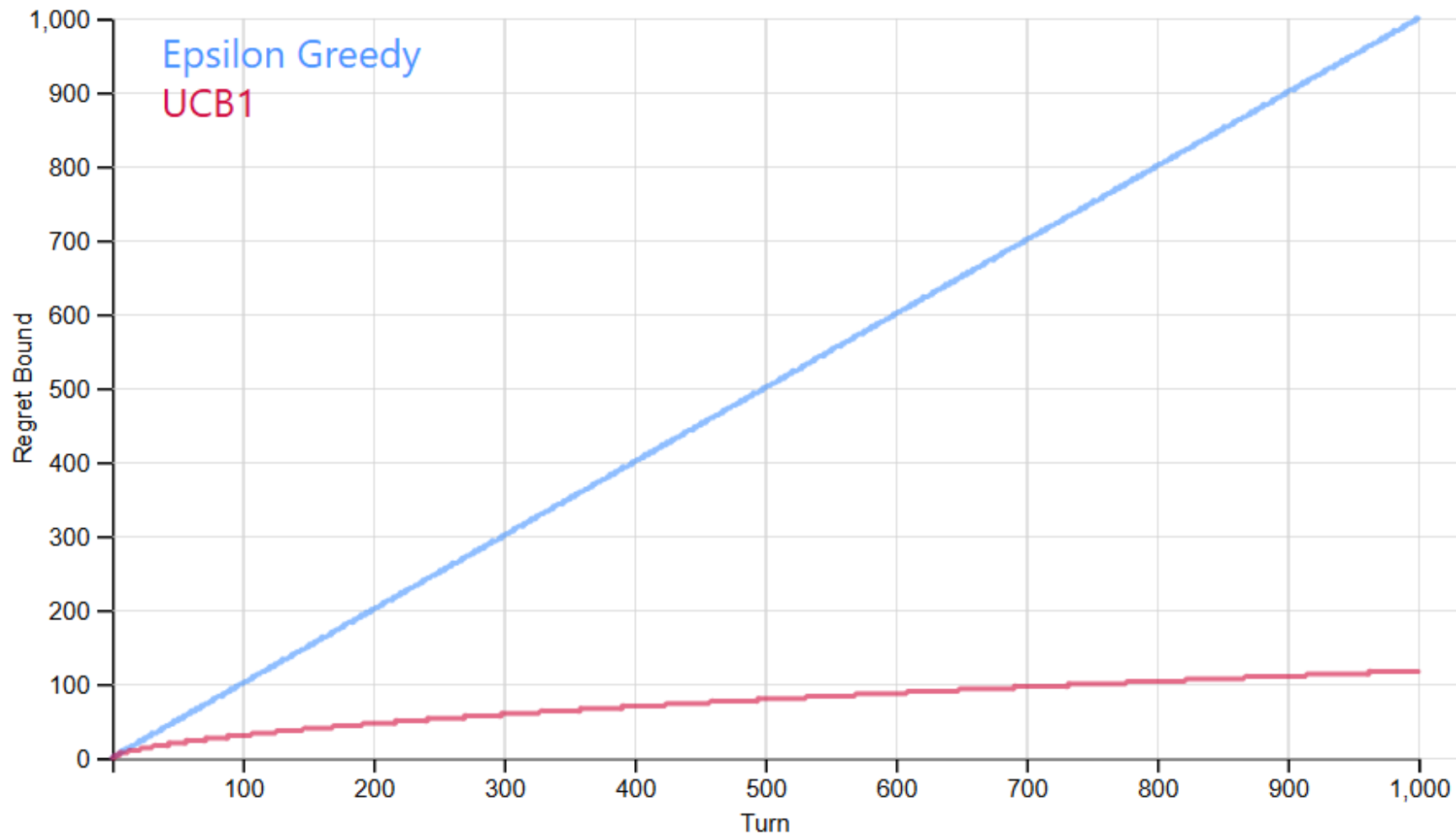
fixed ϵ

$K = \# \text{ arms}$
 $K \ll T$

- UCB1 Regret
 - $O(\sqrt{kT \log(T)})$

- 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):

Notes

- The version we derived is for rewards in range $[0,1]$
- What happens if rewards are in range $[a,b]$?
 - Just scale the upper confidence value

$$c(i, t) = (b - a) \sqrt{\frac{2 \cdot \log(t)}{n_i}}$$

In practice the scaling term is often just treated as a hyperparameter that controls exploration vs. exploitation.

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

What problems do we have with vanilla MAB?

Bad/Strong Assumptions

- Stationarity
- no state/context/obs



Contextual Bandits

- **High-level definition**

At each round t :

- Observe **context** x_t
- Choose action $a_t \in \{1, \dots, K\}$
- Observe reward $r_t(a_t, x_t)$

MAB
no context
 $r_t(a_t)$

Each arm has a context-dependent reward function:

$$\mathbb{E}[r \mid x, a] = f_a(x)$$

LinUCB

$$\begin{array}{l} x_0, a_0, r_0(a_0, x_0) \\ \cancel{x_1, a_1, r_1(a_1, x_1)} \\ x_2, a_0, r_2(a_0, x_2) \end{array}$$

$$\begin{array}{l} x \rightarrow y \\ x \rightarrow r \end{array}$$

- **Assumption**

- For each arm a :
- Reward is linear in features
- Separate parameter vector per arm

$$\mathbb{E}[r \mid x, a] = x^\top \theta_a$$

$$\mathcal{D}_{a_1} = \left\{ (x, r), \dots \right\}$$
$$\theta_{a_1}^\top x \rightarrow r_{a_1}$$

- How should we choose arms (actions) given a context vector x ?

- We can maintain one linear regression model per arm
- Plus uncertainty (e.g. upper confidence bound)

$$\theta_{a_2}^\top x \rightarrow r_{a_2}$$

$$\vdots$$
$$\theta_{a_n}^\top x \rightarrow r_{a_n}$$

Aside: Linear Regression

$$x \in \mathbb{R}^d$$

- Model assumption

$$y = x^\top \theta + \varepsilon$$

- Closed form solution

$$\hat{\theta} = (X^\top X)^{-1} X^\top \vec{y}$$

$$\begin{pmatrix} x_0, y_0 \\ x_1, y_1 \\ \vdots \\ x_n, y_n \end{pmatrix}$$

$$\vec{y} = \begin{bmatrix} y_0 \\ \vdots \\ y_n \end{bmatrix}$$

$$X = \begin{bmatrix} -x_1 - \\ -x_2 - \\ \vdots \\ -x_n - \end{bmatrix}$$

$$X^\top \vec{y} = \underbrace{X^\top}_{d \times n} \underbrace{X}_{n \times d} \hat{\theta}$$
$$(X^\top X)^{-1} X^\top \vec{y} = \hat{\theta}$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}_{N \times d}$$

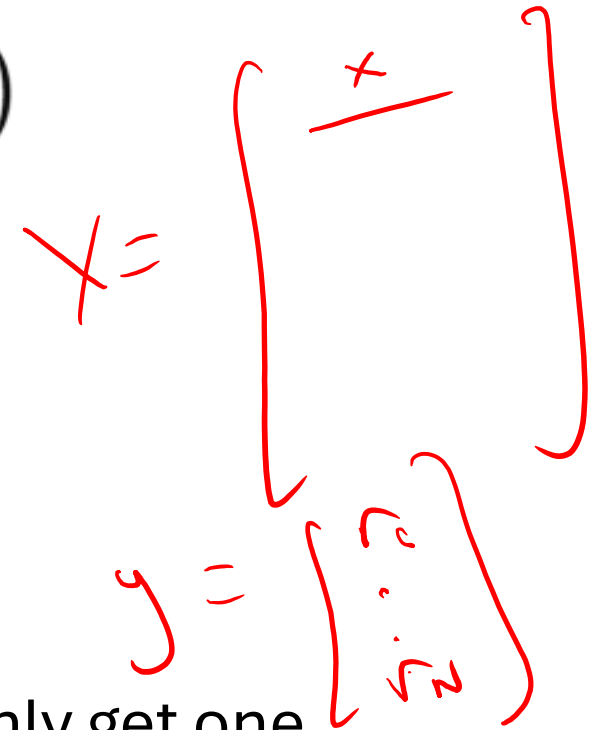
Ridge Regression (adds L2 regularization)

$$\hat{\theta} = \arg \min_{\theta} (\|X\theta - y\|^2 + \lambda \|\theta\|^2)$$

- Also has a closed form solution:

$$\hat{\theta}_a = (X_a^\top X_a + \lambda I)^{-1} X_a^\top y_a$$

- How would you solve this in an online way? E.g., you only get one sample (x_t, y_t) at each timestep and you want to iteratively update θ



We want something like this but that can be updated as we get new data

$$\hat{\theta} = \underbrace{(X^T X + \lambda I)}_A \underbrace{X^T y}_b$$

We can write $\hat{\theta} = A^{-1}b$

At each time t you observe (x_t, y_t)

What is A_t and b_t such that we can estimate θ given the data so far?

$$X = \begin{bmatrix} x_1 & y_1 \\ \vdots & \vdots \\ x_t & y_t \end{bmatrix}$$

$$A_t = \lambda I + \sum_{s=1}^t x_s x_s^T + x_{t+1} x_{t+1}^T$$

$$b_t = \sum_{s=1}^t y_s x_s + y_{t+1} x_{t+1}$$

$$\hat{\theta} = A_t^{-1} b_t$$

LinUCB



- Optimism in the face of uncertainty in ***function space***.

LinUCB = online ridge regression + optimism via confidence ellipsoids

LinUCB Algorithm

Assumption

For each arm $a \in \{1, \dots, K\}$:

$$\mathbb{E}[r \mid x, a] = x^\top \theta_a$$

For rounds $t = 1, 2, \dots$

Observe context

$$x_t \in \mathbb{R}^d$$

For each arm a :

$$\hat{\theta}_a = A_a^{-1} b_a$$

$$\text{UCB}_a(x_t) = x_t^\top \hat{\theta}_a + \alpha \sqrt{x_t^\top A_a^{-1} x_t}$$

Select arm

$$a_t = \arg \max_a \text{UCB}_a(x_t)$$

Initialization

For each arm a :

Handwritten red annotations: a lambda symbol and a circle.

$$A_a \leftarrow \lambda I_d \quad b_a \leftarrow 0_d$$

Observe reward

$$r_t \in \mathbb{R}$$

Update (chosen arm only)

$$A_{a_t} \leftarrow A_{a_t} + x_t x_t^\top$$

$$b_{a_t} \leftarrow b_{a_t} + r_t x_t$$

LinUCB

- LinUCB puts confidence ellipsoids on models!
- Works really well in practice!
- Very data efficient.
- Has been used heavily in production systems (ads, recommendations, etc)

$$(x, a, r)$$

Modern, Deep Learning Approaches

- Neural Contextual Bandits

- Replace linear model with neural network
- Train online from feedback
- Figuring out good exploration/uncertainty is challenging

$$\theta_a^T x$$

$$r \approx f_{\theta}(x, a)$$

$$f_{\theta_2}(x, a)$$

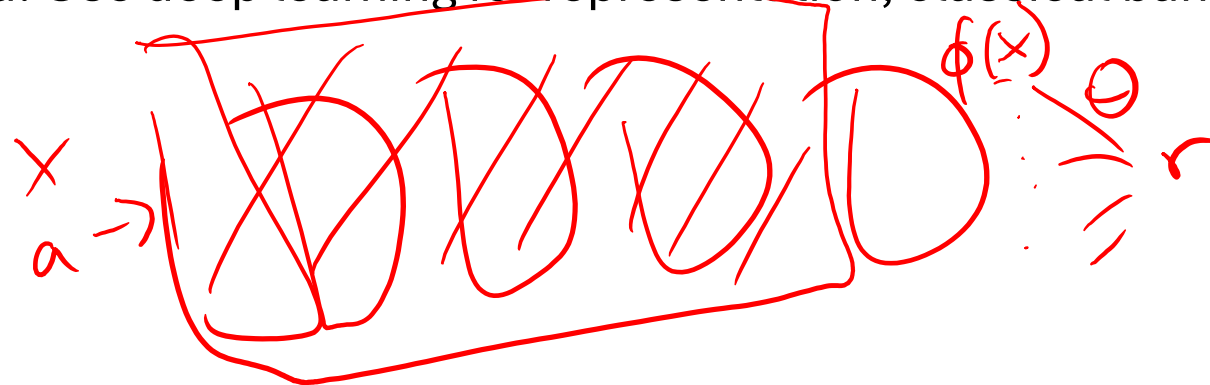
$$f_{\theta_3}(x, a)$$

$$x, a, a_w$$

- Uncertainty via approximations

- Ensembles
- NeuralUCB (linear UCB on last-layer features)
 - Key idea: Use deep learning for representation, classical bandits for exploration.

$$\arg \max_x \bar{f}_{\theta}(x, a) + \text{var}$$



Transitioning to Reinforcement Learning

- At what point does a contextual bandit become full reinforcement learning?

$$R_t = \sum_{t=0}^T r_t$$

$$a \rightarrow x \xrightarrow{a} x$$