

PROBABILITY

9.8 MULTI-ARMED BANDITS

ALEKS JOVCIC

SLIDES BY ALEX TSUN

AGENDA

- THE MULTI-ARMED BANDIT (MAB) PROBLEM
- GREEDY/EPSILON-GREEDY
- UPPER CONFIDENCE BOUND (UCB)
- THOMPSON SAMPLING
- MODERN HYPOTHESIS TESTING

MULTI-ARMED BANDIT (MAB) PROBLEM

- K Slot Machines $\{1,2,\dots,K\}$ (aka "Bandits" with "Arms").



MULTI-ARMED BANDIT (MAB) PROBLEM

- K Slot Machines $\{1,2,\dots,K\}$ (aka "Bandits" with "Arms").
- At each time step $t=1,2,\dots,T$: Pull an arm $a_t \in \{1,2,\dots,K\}$ and observe random reward (each arm is independent, and has some reward distribution which doesn't change over time).



MULTI-ARMED BANDIT (MAB) PROBLEM

- K Slot Machines $\{1,2,\dots,K\}$ (aka "Bandits" with "Arms").
- At each time step $t=1,2,\dots,T$: Pull an arm $a_t \in \{1,2,\dots,K\}$ and observe random reward (each arm is independent, and has some reward distribution which doesn't change over time).
- **Goal**: Maximize total (expected) reward after T time steps.



MULTI-ARMED BANDIT (MAB) PROBLEM

- K Slot Machines $\{1,2,\dots,K\}$ (aka "Bandits" with "Arms").
- At each time step $t=1,2,\dots,T$: Pull an arm $a_t \in \{1,2,\dots,K\}$ and observe random reward (each arm is independent, and has some reward distribution which doesn't change over time).
- **Goal**: Maximize total (expected) reward after T time steps.



- **Problem**: At each time step, decide which arm to pull based on past history of rewards.

MULTI-ARMED BANDIT (MAB) PROBLEM

Below has the reward distribution of each of the $K=3$ arms.

What's your strategy to maximize your total (expected) reward?



$Poi(\lambda = 1.36)$



$Bin(n = 10, p = 0.4)$



$\mathcal{N}(\mu = -1, \sigma^2 = 4)$

MULTI-ARMED BANDIT (MAB) PROBLEM

Below has the reward distribution of each of the $K=3$ arms.

What's your strategy to maximize your total (expected) reward?



$Poi(\lambda = 1.36)$



$Bin(n = 10, p = 0.4)$



$\mathcal{N}(\mu = -1, \sigma^2 = 4)$

→
Pull arm 2 every time since it has the highest expected reward!

MULTI-ARMED BANDIT (MAB) PROBLEM

Well actually, we don't know the reward distributions :(.



MULTI-ARMED BANDIT (MAB) PROBLEM

Well actually, we don't know the reward distributions :(.

Have to **estimate** all K expectations



MULTI-ARMED BANDIT (MAB) PROBLEM

Well actually, we don't know the reward distributions :(.

Have to **estimate** all K expectations, WHILE simultaneously maximizing reward!



MULTI-ARMED BANDIT (MAB) PROBLEM

Well actually, we don't know the reward distributions :(.

Have to **estimate** all K expectations, WHILE simultaneously maximizing reward!



This is a hard problem - we know nothing about the K reward distributions!

MULTI-ARMED BANDIT (MAB) PROBLEM

Need to balance the tradeoff between:

Exploitation: Pulling arm(s) we know to be "good".

Exploration: Pulling other arms in the hopes they are also "good" or even better.



BERNOULLI BANDITS

We will handle the case of Bernoulli-bandits. That is, reward of arm $a \in \{1,2,\dots,K\}$ is $\text{Ber}(p_a)$.



$\text{Ber}(p_1)$



$\text{Ber}(p_2)$



$\text{Ber}(p_3)$

BERNOULLI BANDITS

We will handle the case of Bernoulli-bandits. That is, reward of arm $a \in \{1,2,\dots,K\}$ is $\text{Ber}(p_a)$.



$\text{Ber}(p_1)$



$\text{Ber}(p_2)$



$\text{Ber}(p_3)$

We don't know these!

BERNOULLI BANDITS

We will handle the case of Bernoulli-bandits. That is, reward of arm $a \in \{1,2,\dots,K\}$ is $\text{Ber}(p_a)$.

Observe: The expected reward of arm a is just p_a (expectation of Bernoulli).



$\text{Ber}(p_1)$



$\text{Ber}(p_2)$



$\text{Ber}(p_3)$

We don't know these!

REGRET

Regret is the difference between:

- The best possible expected reward (always pull the best arm)
- The actual reward you got



REGRET

Regret is the difference between:

- The best possible expected reward (always pull the best arm)
- The actual reward you got

Let $p^* = \max_{i \in \{1, 2, \dots, K\}} p_i$ denote the highest expected reward from one of the K arms.



REGRET

Regret is the difference between:

- The best possible expected reward (always pull the best arm)
- The actual reward you got



Let $p^* = \max_{i \in \{1, 2, \dots, K\}} p_i$ denote the highest expected reward from one of the K arms.

$$\text{Regret}(T) = Tp^* - \text{Reward}(T)$$

At some "time" T (after T arm pulls)

REGRET

Regret is the difference between:

- The best possible expected reward (always pull the best arm)
- The actual reward you got

Let $p^* = \max_{i \in \{1, 2, \dots, K\}} p_i$ denote the highest expected reward from one of the K arms.

$$\text{Regret}(T) = Tp^* - \text{Reward}(T)$$

At some "time" T (after T arm pulls)



REGRET

Regret is the difference between:

- The best possible expected reward (always pull the best arm)
- The actual reward you got

Let $p^* = \max_{i \in \{1, 2, \dots, K\}} p_i$ denote the highest expected reward from one of the K arms.

$$\text{Regret}(T) = Tp^* - \text{Reward}(T)$$

At some "time" T (after T arm pulls)



REGRET

Regret is the difference between:

- The best possible expected reward (always pull the best arm)
- The actual reward you got



Let $p^* = \max_{i \in \{1, 2, \dots, K\}} p_i$ denote the highest expected reward from one of the K arms.

$$\text{Regret}(T) = Tp^* - \text{Reward}(T)$$

$$\text{Avg-Regret}(T) = p^* - \frac{\text{Reward}(T)}{T}$$

REGRET

Regret is the difference between:

- The best possible expected reward (always pull the best arm)
- The actual reward you got



Let $p^* = \max_{i \in \{1, 2, \dots, K\}} p_i$ denote the highest expected reward from one of the K arms.

$$\text{Regret}(T) = Tp^* - \text{Reward}(T)$$

$$\text{Avg-Regret}(T) = p^* - \frac{\text{Reward}(T)}{T}$$

Want $\text{Avg-Regret}(T) \rightarrow 0$ as $T \rightarrow \infty$. **Minimizing Regret = Maximizing Reward.**

(BERNOULLI) BANDIT FRAMEWORK



How do we choose an arm at each time step (depending on past history), to maximize our total reward?

Algorithm 1 (Bernoulli) Bandit Framework

1: Have K arms, where pulling arm $i \in \{1, \dots, K\}$ gives $Ber(p_i)$ reward \triangleright p_i 's all unknown.

(BERNOULLI) BANDIT FRAMEWORK



How do we choose an arm at each time step (depending on past history), to maximize our total reward?

Algorithm 1 (Bernoulli) Bandit Framework

- 1: Have K arms, where pulling arm $i \in \{1, \dots, K\}$ gives $Ber(p_i)$ reward $\triangleright p_i$'s all unknown.
- 2: **for** $t = 1, \dots, T$ **do**
- 3: At time t , pull arm $a_t \in \{1, \dots, K\}$. \triangleright How do we do decide which arm?

(BERNOULLI) BANDIT FRAMEWORK



How do we choose an arm at each time step (depending on past history), to maximize our total reward?

Algorithm 1 (Bernoulli) Bandit Framework

- 1: Have K arms, where pulling arm $i \in \{1, \dots, K\}$ gives $Ber(p_i)$ reward $\triangleright p_i$'s all unknown.
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: At time t , pull arm $a_t \in \{1, \dots, K\}$. \triangleright How do we decide which arm?
 - 4: Receive reward $r_t \sim Ber(p_{a_t})$. \triangleright Reward is either 1 or 0.
-

(BERNOULLI) BANDIT FRAMEWORK



How do we choose an arm at each time step (depending on past history), to maximize our total reward?

Algorithm 1 (Bernoulli) Bandit Framework

- 1: Have K arms, where pulling arm $i \in \{1, \dots, K\}$ gives $Ber(p_i)$ reward $\triangleright p_i$'s all unknown.
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: At time t , pull arm $a_t \in \{1, \dots, K\}$. \triangleright How do we decide which arm?
 - 4: Receive reward $r_t \sim Ber(p_{a_t})$. \triangleright Reward is either 1 or 0.
-

This is the focus of the rest of this lecture!

MOTIVATION: CLINICAL TRIALS

$K = 4$ Arms (Treatments)



MOTIVATION: CLINICAL TRIALS

$K = 4$ Arms (Treatments)

For patient t , prescribe treatment $a_t \in \{1, 2, 3, 4\}$.



MOTIVATION: CLINICAL TRIALS

$K = 4$ Arms (Treatments)

For patient t , prescribe treatment $a_t \in \{1, 2, 3, 4\}$.

Observe reward $r_t \in \{0, 1\}$. (1 if healed, 0 if not)



MOTIVATION: CLINICAL TRIALS



$K = 4$ Arms (Treatments)

For patient t , prescribe treatment $a_t \in \{1, 2, 3, 4\}$.

Observe reward $r_t \in \{0, 1\}$. (1 if healed, 0 if not)

Maximize: Total number of patients healed.

MOTIVATION: RECOMMENDING MOVIES

K Movies

For visitor t , recommend movie $a_t \in \{1, 2, \dots, K\}$.



MOTIVATION: RECOMMENDING MOVIES



K Movies

For visitor t , recommend movie $a_t \in \{1, 2, \dots, K\}$.

Observe reward $r_t \in \{1, 2, 3, 4, 5\}$. (rating)

Maximize: Total/average rating of recommendations.

MOTIVATION: REAL LIFE?? (FOOD)



K Cuisines/Dishes (a ton)

For meal t , eat dish $a_t \in \{1, 2, \dots, K\}$.



MOTIVATION: REAL LIFE?? (FOOD)



K Cuisines/Dishes (a ton)

For meal t , eat dish $a_t \in \{1, 2, \dots, K\}$.

Observe reward $r_t \in \{1, 2, 3, 4, 5\}$. (happiness rating)

Maximize: Total/average happiness :)



MOTIVATION: REAL LIFE?? (FOOD)



K Cuisines/Dishes (a ton)

For meal t , eat dish $a_t \in \{1, 2, \dots, K\}$.

Observe reward $r_t \in \{1, 2, 3, 4, 5\}$. (happiness rating)

Maximize: Total/average happiness :)

The Question of the Day: Explore or Exploit????



MOTIVATION: REAL LIFE?? (ACTIVITIES)



K Activities

On day t , do activity $a_t \in \{1,2,\dots,K\}$.

Observe reward $r_t \in \{1,2,3,4,5\}$. (happiness rating)

Maximize: Total/average happiness :)

The Question of the Day: Explore or Exploit????



ANY IDEAS ON WHAT STRATEGY WE CAN USE???



GREEDY (NAIVE) ALGORITHM



Algorithm 2 Greedy (Naive) Strategy for Bernoulli Bandits

1: Choose a number of times M to pull each arm initially, with $KM \leq T$.

GREEDY (NAIVE) ALGORITHM



Algorithm 2 Greedy (Naive) Strategy for Bernoulli Bandits

- 1: Choose a number of times M to pull each arm initially, with $KM \leq T$.
- 2: **for** $i = 1, 2, \dots, K$ **do**
- 3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.
- 4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.

GREEDY (NAIVE) ALGORITHM



Algorithm 2 Greedy (Naive) Strategy for Bernoulli Bandits

1: Choose a number of times M to pull each arm initially, with $KM \leq T$.

2: **for** $i = 1, 2, \dots, K$ **do**

3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.

4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.

5: Determine best (empirical) arm $a^* = \arg \max_{i \in \{1, 2, \dots, K\}} \hat{p}_i$.

► We could be wrong...

GREEDY (NAIVE) ALGORITHM



Algorithm 2 Greedy (Naive) Strategy for Bernoulli Bandits

- 1: Choose a number of times M to pull each arm initially, with $KM \leq T$.
 - 2: **for** $i = 1, 2, \dots, K$ **do**
 - 3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.
 - 4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.
 - 5: Determine best (empirical) arm $a^* = \arg \max_{i \in \{1, 2, \dots, K\}} \hat{p}_i$. ▶ We could be wrong...
 - 6: **for** $t = KM + 1, KM + 2, \dots, T$ **do**:
 - 7: Pull arm $a_t = a^*$. ▶ Pull the same arm for the rest of time.
 - 8: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
-

GREEDY (NAIVE) ALGORITHM



Algorithm 2 Greedy (Naive) Strategy for Bernoulli Bandits

- 1: Choose a number of times M to pull each arm initially, with $KM \leq T$.
 - 2: **for** $i = 1, 2, \dots, K$ **do**
 - 3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.
 - 4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.
 - 5: Determine best (empirical) arm $a^* = \arg \max_{i \in \{1, 2, \dots, K\}} \hat{p}_i$. ▶ We could be wrong...
 - 6: **for** $t = KM + 1, KM + 2, \dots, T$ **do**:
 - 7: Pull arm $a_t = a^*$. ▶ Pull the same arm for the rest of time.
 - 8: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
-

If we make a mistake, we will regret our decision for the rest of time....

Can we not do all of our exploration at the beginning?

EPSILON-GREEDY ALGORITHM

ϵ



Explore with probability epsilon!

Algorithm 3 ϵ -Greedy Strategy for Bernoulli Bandits

- 1: Choose a number of times M to pull each arm initially, with $KM \leq T$.
- 2: **for** $i = 1, 2, \dots, K$ **do**
- 3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.
- 4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.

EPSILON-GREEDY ALGORITHM

ϵ



Explore with probability epsilon!

Algorithm 3 ϵ -Greedy Strategy for Bernoulli Bandits

- 1: Choose a number of times M to pull each arm initially, with $KM \leq T$.
- 2: **for** $i = 1, 2, \dots, K$ **do**
- 3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.
- 4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.
- 5: **for** $t = KM + 1, KM + 2, \dots, T$ **do**:
- 6: **if** $\text{Ber}(\epsilon) == 1$: **then**
 - With probability ϵ , explore.
- 7: Pull arm $a_t = \text{Unif}(1, K)$ (discrete).
 - Choose a uniformly random arm.

EPSILON-GREEDY ALGORITHM

 ϵ 

Explore with probability epsilon!

Algorithm 3 ϵ -Greedy Strategy for Bernoulli Bandits

- 1: Choose a number of times M to pull each arm initially, with $KM \leq T$.
- 2: **for** $i = 1, 2, \dots, K$ **do**
- 3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.
- 4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.
- 5: **for** $t = KM + 1, KM + 2, \dots, T$ **do**:
- 6: **if** $\text{Ber}(\epsilon) == 1$: **then**
 - With probability ϵ , explore.
- 7: Pull arm $a_t = \text{Unif}(1, K)$ (discrete).
 - Choose a uniformly random arm.
- 8: **else**
 - With probability $1 - \epsilon$, exploit.
- 9: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \hat{p}_i$.
 - Choose arm with highest estimated reward.

EPSILON-GREEDY ALGORITHM

ϵ



Explore with probability epsilon!

Algorithm 3 ϵ -Greedy Strategy for Bernoulli Bandits

- 1: Choose a number of times M to pull each arm initially, with $KM \leq T$.
 - 2: **for** $i = 1, 2, \dots, K$ **do**
 - 3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.
 - 4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.
 - 5: **for** $t = KM + 1, KM + 2, \dots, T$ **do**:
 - 6: **if** $\text{Ber}(\epsilon) == 1$: **then**
 - ▶ With probability ϵ , explore.
 - 7: Pull arm $a_t = \text{Unif}(1, K)$ (discrete).
 - ▶ Choose a uniformly random arm.
 - 8: **else**
 - ▶ With probability $1 - \epsilon$, exploit.
 - 9: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \hat{p}_i$.
 - ▶ Choose arm with highest estimated reward.
 - 10: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
 - 11: Update \hat{p}_{a_t} (using newly observed reward r_t).
-

EPSILON-GREEDY ALGORITHM

ϵ



Explore with probability epsilon!

Algorithm 3 ϵ -Greedy Strategy for Bernoulli Bandits

- 1: Choose a number of times M to pull each arm initially, with $KM \leq T$.
 - 2: **for** $i = 1, 2, \dots, K$ **do**
 - 3: Pull arm i M times, observing iid rewards $r_{i1}, \dots, r_{iM} \sim \text{Ber}(p_i)$.
 - 4: Estimate $\hat{p}_i = \frac{\sum_{j=1}^M r_{ij}}{M}$.
 - 5: **for** $t = KM + 1, KM + 2, \dots, T$ **do**:
 - 6: **if** $\text{Ber}(\epsilon) == 1$: **then** ▷ With probability ϵ , explore.
 - 7: Pull arm $a_t = \text{Unif}(1, K)$ (discrete). ▷ Choose a uniformly random arm.
 - 8: **else** ▷ With probability $1 - \epsilon$, exploit.
 - 9: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \hat{p}_i$. ▷ Choose arm with highest estimated reward.
 - 10: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
 - 11: Update \hat{p}_{a_t} (using newly observed reward r_t).
-

Can we explore more "naturally"?

UPPER CONFIDENCE BOUND (UCB) ALGORITHM

This algorithm constructs confidence intervals for the estimates of each arm, and chooses the arm with the highest **upper** confidence bound (if the confidence interval is $[a,b]$, we compare only the value of b)



UPPER CONFIDENCE BOUND (UCB) ALGORITHM



This algorithm constructs confidence intervals for the estimates of each arm, and chooses the arm with the highest **upper** confidence bound (if the confidence interval is $[a,b]$, we compare only the value of b)

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i/1$. ▶ Each estimate \hat{p}_i will initially either be 1 or 0.

UPPER CONFIDENCE BOUND (UCB) ALGORITHM



This algorithm constructs confidence intervals for the estimates of each arm, and chooses the arm with the highest **upper** confidence bound (if the confidence interval is $[a,b]$, we compare only the value of b)

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i/1$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .

UPPER CONFIDENCE BOUND (UCB) ALGORITHM



This algorithm constructs confidence intervals for the estimates of each arm, and chooses the arm with the highest **upper** confidence bound (if the confidence interval is $[a,b]$, we compare only the value of b)

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i/1$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .

Point estimate/
Max-likelihood estimate





UPPER CONFIDENCE BOUND (UCB) ALGORITHM

This algorithm constructs confidence intervals for the estimates of each arm, and chooses the arm with the highest **upper** confidence bound (if the confidence interval is $[a,b]$, we compare only the value of b)

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i/1$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .

Point estimate/
Max-likelihood estimate

Takes the upper part of of
the confidence interval.

UPPER CONFIDENCE BOUND (UCB) ALGORITHM



This algorithm constructs confidence intervals for the estimates of each arm, and chooses the arm with the highest **upper** confidence bound (if the confidence interval is $[a,b]$, we compare only the value of b)

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i/1$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.



UPPER CONFIDENCE BOUND (UCB) ALGORITHM

This algorithm constructs confidence intervals for the estimates of each arm, and chooses the arm with the highest **upper** confidence bound (if the confidence interval is $[a,b]$, we compare only the value of b)

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
 - 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
 - 3: Estimate $\hat{p}_i = r_i/1$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
 - 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
 - 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
 - 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
 - 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).
-

UPPER CONFIDENCE BOUND (UCB) ALGORITHM



This algorithm constructs confidence intervals for the estimates of each arm, and chooses the arm with the highest **upper** confidence bound (if the confidence interval is $[a,b]$, we compare only the values of b).

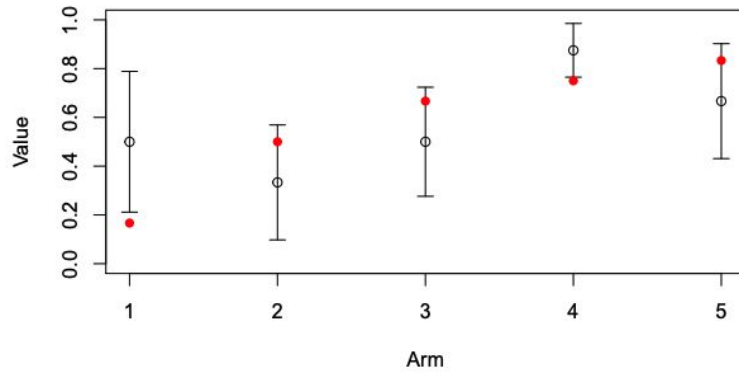
Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
 - 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
 - 3: Estimate $\hat{p}_i = r_i/1$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
 - 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
 - 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
 - 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
 - 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).
-

Exploration is “baked in”: the frequently pulled arms will have narrow confidence intervals (and hence a lower **upper** bound), and the less-frequently pulled arms will have wide intervals (and hence a **higher** upper bound).

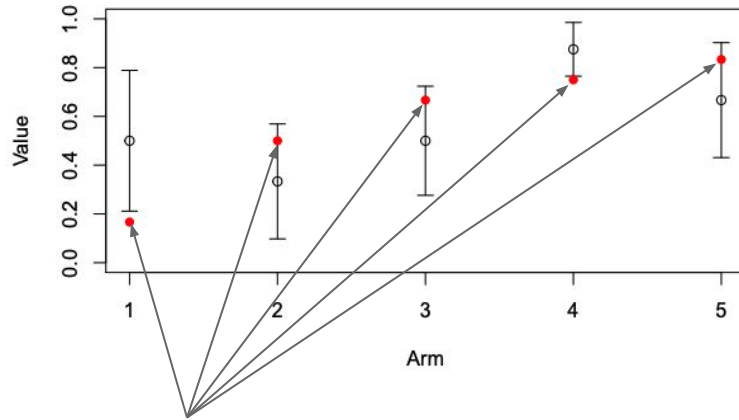
UCB: CONFIDENCE INTERVALS OVER TIME

Confidence Intervals for Mean of Each Arm: $t=10$



UCB: CONFIDENCE INTERVALS OVER TIME

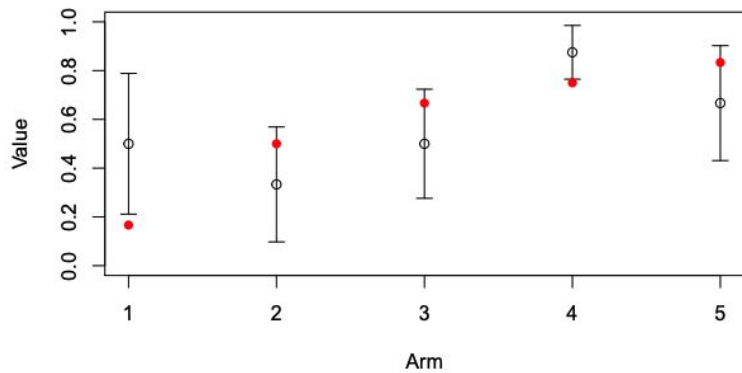
Confidence Intervals for Mean of Each Arm: $t=10$



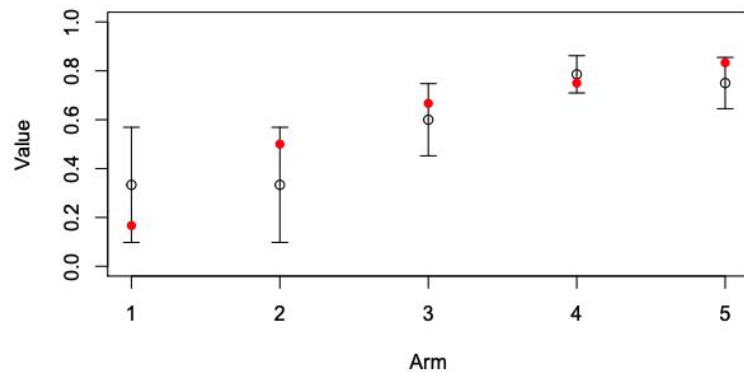
Red Dots: True Means

UCB: CONFIDENCE INTERVALS OVER TIME

Confidence Intervals for Mean of Each Arm: t=10

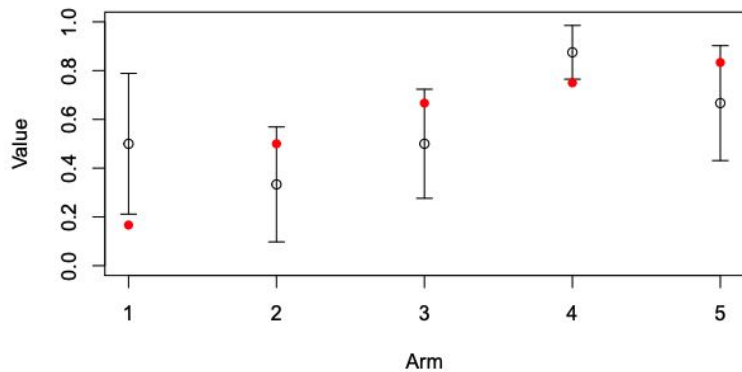


Confidence Intervals for Mean of Each Arm: t=50

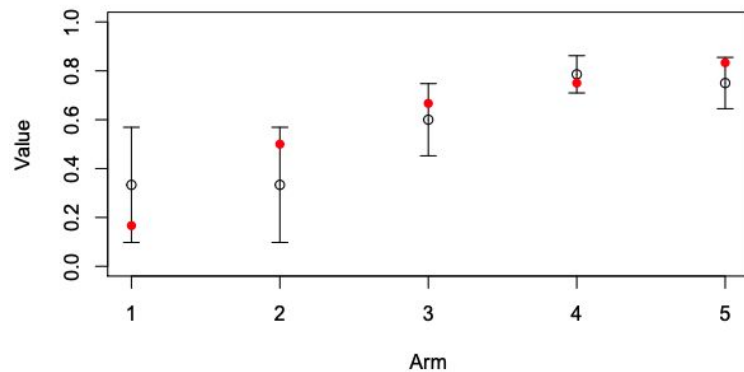


UCB: CONFIDENCE INTERVALS OVER TIME

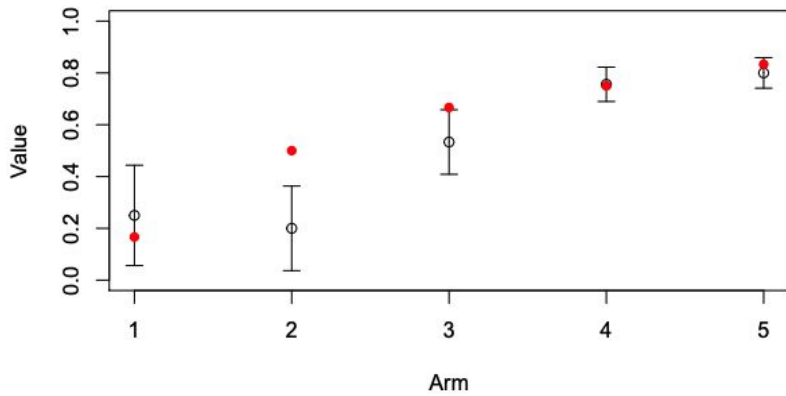
Confidence Intervals for Mean of Each Arm: t=10



Confidence Intervals for Mean of Each Arm: t=50

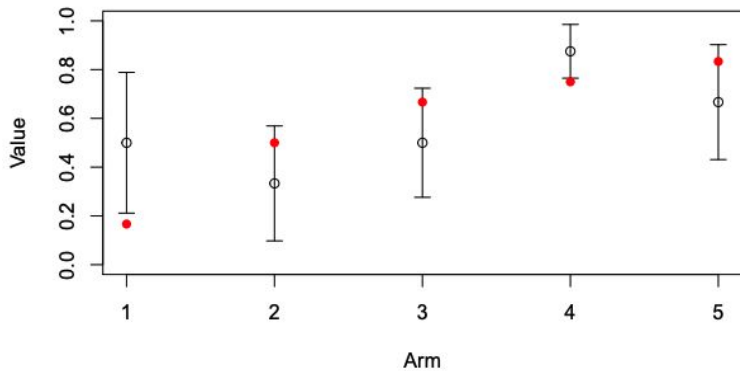


Confidence Intervals for Mean of Each Arm: t=100

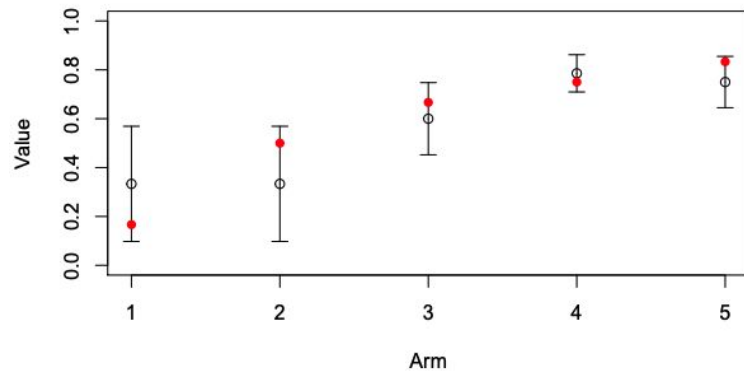


UCB: CONFIDENCE INTERVALS OVER TIME

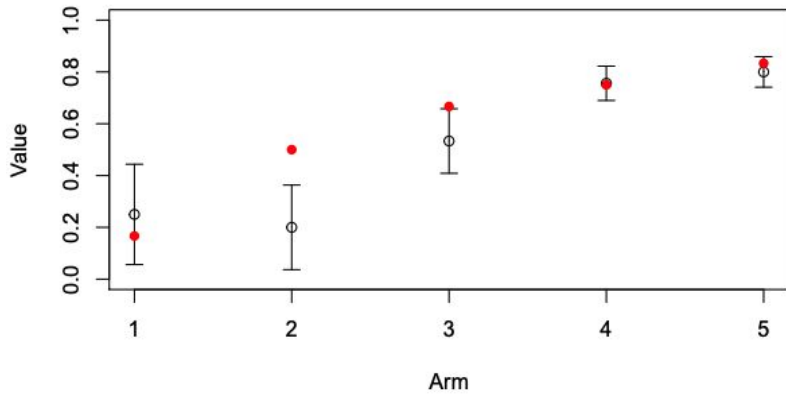
Confidence Intervals for Mean of Each Arm: t=10



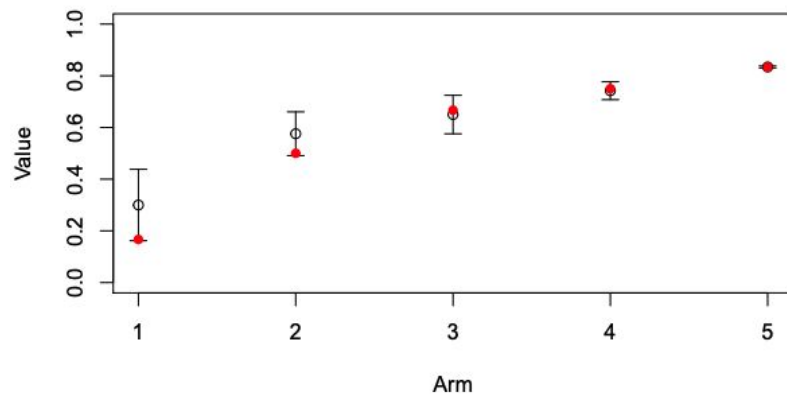
Confidence Intervals for Mean of Each Arm: t=50



Confidence Intervals for Mean of Each Arm: t=100



Confidence Intervals for Mean of Each Arm: t=10000



UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▶ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5				
2	0.2				
3	0.9				

Time (t)	Arm Pulled (a_t)	Reward (r_t)

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▶ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5				
2	0.2				
3	0.9				

Time (t)	Arm Pulled (a_t)	Reward (r_t)

We don't actually know these...

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5				
2	0.2				
3	0.9				

Time (t)	Arm Pulled (a_t)	Reward (r_t)
1	1	0

At time 1, we pull arm 1, and observe either a 1 (with probability 0.5) or a 0 (with probability 1-0.5). We happen to observe a 0.

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▶ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	
2	0.2				
3	0.9				

Time (t)	Arm Pulled (a_t)	Reward (r_t)
1	1	0

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	
2	0.2				
3	0.9				

Time (t)	Arm Pulled (a_t)	Reward (r_t)
2	2	0

At time 2, we pull arm 2, and observe either a 1 (with probability 0.2) or a 0 (with probability 1-0.2). We happen to observe a 0.

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▶ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	
2	0.2	1	0	0/1	
3	0.9				

Time (t)	Arm Pulled (a_t)	Reward (r_t)
2	2	0

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	
2	0.2	1	0	0/1	
3	0.9				

Time (t)	Arm Pulled (a_t)	Reward (r_t)
3	3	1

At time 3, we pull arm 3, and observe either a 1 (with probability 0.9) or a 0 (with probability 1-0.9). We happen to observe a 1.

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	
2	0.2	1	0	0/1	
3	0.9	1	1	1/1	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
3	3	1

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	
2	0.2	1	0	0/1	
3	0.9	1	1	1/1	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4		

At time 4, we must compute all our upper confidence bounds, and choose the best one.

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.665
2	0.2	1	0	0/1	
3	0.9	1	1	1/1	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4		

$$0 + \sqrt{\frac{2 \ln(4)}{1}} \approx 1.665$$

At time 4, we must compute all our upper confidence bounds, and choose the best one.

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▶ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.665
2	0.2	1	0	0/1	1.665
3	0.9	1	1	1/1	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4		

$$0 + \sqrt{\frac{2 \ln(4)}{1}} \approx 1.665$$

At time 4, we must compute all our upper confidence bounds, and choose the best one.

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.665
2	0.2	1	0	0/1	1.665
3	0.9	1	1	1/1	2.665

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4		

$$1 + \sqrt{\frac{2 \ln(4)}{1}} \approx 2.665$$

At time 4, we must compute all our upper confidence bounds, and choose the best one.

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.665
2	0.2	1	0	0/1	1.665
3	0.9	1	1	1/1	2.665

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4	3	

At time 4, arm 3 has the highest UCB so we pull it.

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.665
2	0.2	1	0	0/1	1.665
3	0.9	1	1	1/1	2.665

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4	3	0

At time 4, arm 3 has the highest UCB so we pull it. We observe a reward of 0.

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	
2	0.2	1	0	0/1	
3	0.9	2	1	1/2	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4	3	0

At time 4, arm 3 has the highest UCB so we pull it. We observe a reward of 0. Then we update our estimate for p_3 .

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	
2	0.2	1	0	0/1	
3	0.9	2	1	1/2	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
5		

At time 5, we must compute all our upper confidence bounds, and choose the best one.

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▶ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.794
2	0.2	1	0	0/1	
3	0.9	2	1	1/2	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
5		

$$0 + \sqrt{\frac{2 \ln(5)}{1}} \approx 1.794$$

At time 5, we must compute all our upper confidence bounds, and choose the best one.

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.794
2	0.2	1	0	0/1	1.794
3	0.9	2	1	1/2	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
5		

$$0 + \sqrt{\frac{2 \ln(5)}{1}} \approx 1.794$$

At time 5, we must compute all our upper confidence bounds, and choose the best one.

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.794
2	0.2	1	0	0/1	1.794
3	0.9	2	1	1/2	1.769

Time (t)	Arm Pulled (a_t)	Reward (r_t)
5		

$$\frac{1}{2} + \sqrt{\frac{2 \ln(5)}{2}} \approx 1.769$$

At time 5, we must compute all our upper confidence bounds, and choose the best one.

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.794
2	0.2	1	0	0/1	1.794
3	0.9	2	1	1/2	1.769

Time (t)	Arm Pulled (a_t)	Reward (r_t)
5	1	

At time 5, arms 1 and 2 have the highest UCB so we pull one of them (let's break ties by choosing the smaller index arm). So we pull arm 1.

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	1	0	0/1	1.794
2	0.2	1	0	0/1	1.794
3	0.9	2	1	1/2	1.769

Time (t)	Arm Pulled (a_t)	Reward (r_t)
5	1	0

At time 5, arms 1 and 2 have the highest UCB so we pull one of them (let's break ties by choosing the smaller index arm). So we pull arm 1. We observe a reward of 0.

Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

UCB EXAMPLE



- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	2	0	0/2	
2	0.2	1	0	0/1	
3	0.9	2	1	1/2	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
5	1	0

At time 5, arms 1 and 2 have the highest UCB so we pull one of them (let's break ties by choosing the smaller index arm). So we pull arm 1.

We observe a reward of 0. Then we update our estimate for p_1 .

UCB EXAMPLE



Algorithm 4 UCB1 Algorithm (Upper Confidence Bound) for Bernoulli Bandits

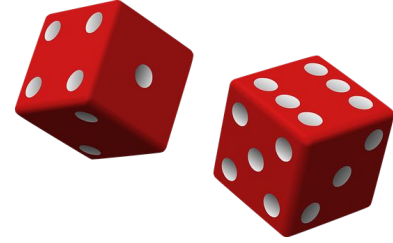
- 1: **for** $i = 1, 2, \dots, K$ **do**
- 2: Pull arm i once, observing $r_i \sim \text{Ber}(p_i)$.
- 3: Estimate $\hat{p}_i = r_i$. ▷ Each estimate \hat{p}_i will initially either be 1 or 0.
- 4: **for** $t = K + 1, K + 2, \dots, T$ **do**:
- 5: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} \left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$, where $N_t(i)$ is the number of times arm i was pulled before time t .
- 6: Receive reward $r_t \sim \text{Ber}(p_{a_t})$.
- 7: Update $N_t(a_t)$ and \hat{p}_{a_t} (using newly observed reward r_t).

Arm (i)	True p_i	# Times Pulled	Total Reward	\hat{p}_i	UCB $\left(\hat{p}_i + \sqrt{\frac{2 \ln(t)}{N_t(i)}} \right)$
1	0.5	2	0	0/2	
2	0.2	1	0	0/1	
3	0.9	2	1	1/2	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
6		

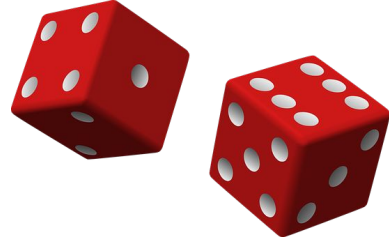
And so on!!! Notice how we started exploring since the confidence bound grows with t for even the unexplored arms!

THOMPSON SAMPLING ALGORITHM



Use MAP: Assume a Beta(1,1) (Uniform) prior on each unknown probability of reward.

THOMPSON SAMPLING ALGORITHM

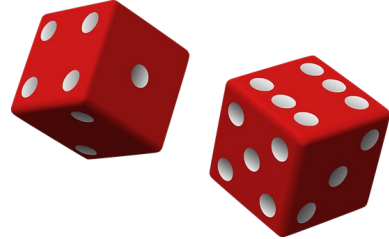


Use MAP: Assume a $\text{Beta}(1,1)$ (Uniform) prior on each unknown probability of reward.

Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. \triangleright Set $\text{Beta}(\alpha_i, \beta_i)$ prior for each p_i .

THOMPSON SAMPLING ALGORITHM



Use MAP: Assume a $Beta(1,1)$ (Uniform) prior on each unknown probability of reward.

Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.

THOMPSON SAMPLING ALGORITHM

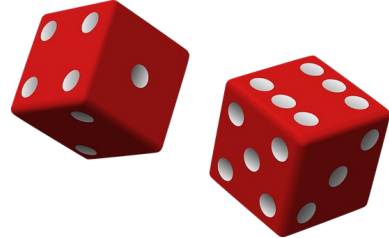


Use MAP: Assume a $Beta(1,1)$ (Uniform) prior on each unknown probability of reward.

Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!

THOMPSON SAMPLING ALGORITHM

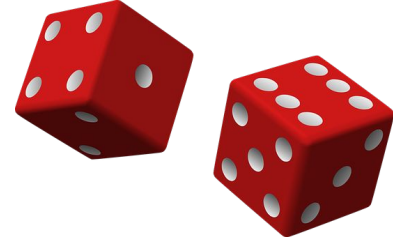


Use MAP: Assume a $Beta(1,1)$ (Uniform) prior on each unknown probability of reward.

Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.

THOMPSON SAMPLING ALGORITHM



Use MAP: Assume a $Beta(1,1)$ (Uniform) prior on each unknown probability of reward.

Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
 - 2: **for** $t = 1, 2, \dots, T$ **do**:
 - 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
 - 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
 - 5: Receive reward $r_t \sim Ber(p_{a_t})$.
 - 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
 - 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.
-

THOMPSON SAMPLING ALGORITHM



Use MAP: Assume a $Beta(1,1)$ (Uniform) prior on each unknown probability of reward.

Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
 - 2: **for** $t = 1, 2, \dots, T$ **do**:
 - 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
 - 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
 - 5: Receive reward $r_t \sim Ber(p_{a_t})$.
 - 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
 - 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.
-

The exploration comes in since we sample from each Beta distribution, rather than just choosing the one with largest expectation or mode (greedy).

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5			
2	0.2			
3	0.9			

Time (t)	Arm Pulled (a_t)	Reward (r_t)

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$.
 - Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$.
 - Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$.
 - This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$.
 - Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$.
 - Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	
2	0.2	1	1	
3	0.9	1	1	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
1		

THOMPSON EXAMPLE

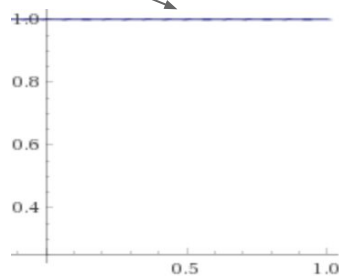


Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.43
2	0.2	1	1	
3	0.9	1	1	

Sample from Beta(1,1) density →



Time (t)	Arm Pulled (a_t)	Reward (r_t)
1		

THOMPSON EXAMPLE

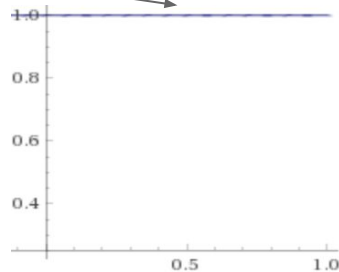


Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.43
2	0.2	1	1	0.75
3	0.9	1	1	

Sample from $Beta(1,1)$ density \rightarrow



Time (t)	Arm Pulled (a_t)	Reward (r_t)
1		

THOMPSON EXAMPLE

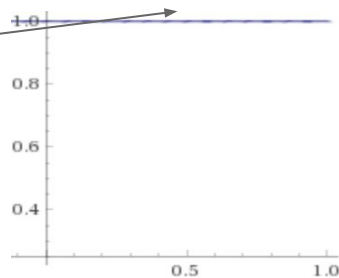


Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.43
2	0.2	1	1	0.75
3	0.9	1	1	0.11

Sample from $Beta(1,1)$ density \rightarrow



Time (t)	Arm Pulled (a_t)	Reward (r_t)
1		

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$s_{i,t}$
1	0.5	1	1	0.43
2	0.2	1	1	0.75
3	0.9	1	1	0.11

Time (t)	Arm Pulled (a_t)	Reward (r_t)
1	2	

Choose arm with highest sample!

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.43
2	0.2	1	1	0.75
3	0.9	1	1	0.11

Time (t)	Arm Pulled (a_t)	Reward (r_t)
1	2	0

Observe reward 1 with probability 0.2
and 0 with probability 0.8.

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	
2	0.2	1	2	
3	0.9	1	1	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
1	2	0

Add a count of 1 to the failures :(.

THOMPSON EXAMPLE

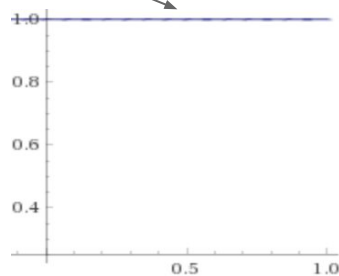


Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.52
2	0.2	1	2	
3	0.9	1	1	

Sample from $Beta(1,1)$ density →



Time (t)	Arm Pulled (a_t)	Reward (r_t)
2		

THOMPSON EXAMPLE

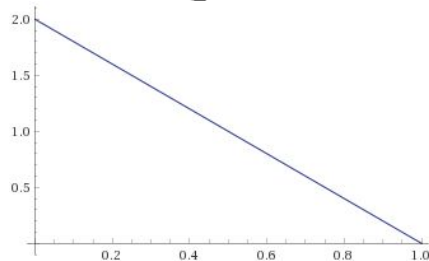


Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.52
2	0.2	1	2	0.05
3	0.9	1	1	

Sample from Beta(1,2) density →



Time (t)	Arm Pulled (a_t)	Reward (r_t)
2		

THOMPSON EXAMPLE

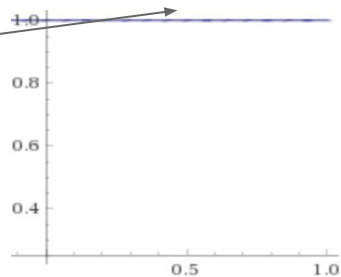


Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.52
2	0.2	1	2	0.05
3	0.9	1	1	0.67

Sample from $Beta(1,1)$ density →



Time (t)	Arm Pulled (a_t)	Reward (r_t)
2		

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.52
2	0.2	1	2	0.05
3	0.9	1	1	0.67

Time (t)	Arm Pulled (a_t)	Reward (r_t)
2	3	

Choose arm with highest sample!

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.52
2	0.2	1	2	0.05
3	0.9	1	1	0.67

Time (t)	Arm Pulled (a_t)	Reward (r_t)
2	3	1

Observe reward 1 with probability 0.9
and 0 with probability 0.1.

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	
2	0.2	1	2	
3	0.9	2	1	

Time (t)	Arm Pulled (a_t)	Reward (r_t)
2	3	1

Add a count of 1 to the successes :).

THOMPSON EXAMPLE



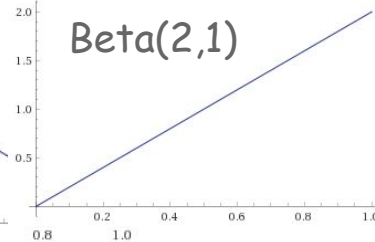
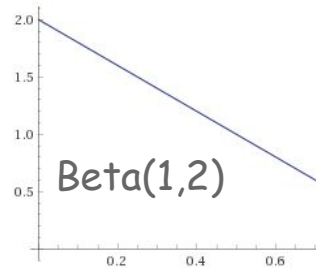
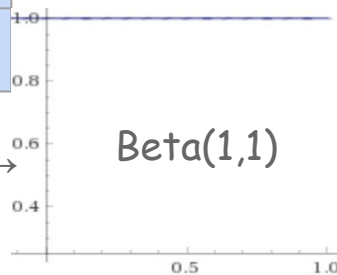
Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.44
2	0.2	1	2	0.27
3	0.9	2	1	0.86

Time (t)	Arm Pulled (a_t)	Reward (r_t)
3		

Sample from each arm's Beta distribution →



THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.44
2	0.2	1	2	0.27
3	0.9	2	1	0.86

Time (t)	Arm Pulled (a_t)	Reward (r_t)
3	3	0

Observe reward 1 with probability 0.9
and 0 with probability 0.1.

THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.44
2	0.2	1	2	0.27
3	0.9	2	2	0.86

Time (t)	Arm Pulled (a_t)	Reward (r_t)
3	3	0

Add a count of 1 to the failures :(.

THOMPSON EXAMPLE



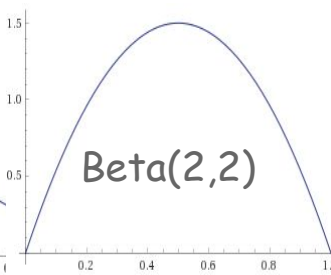
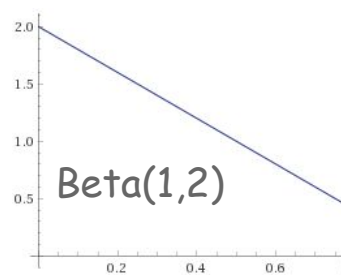
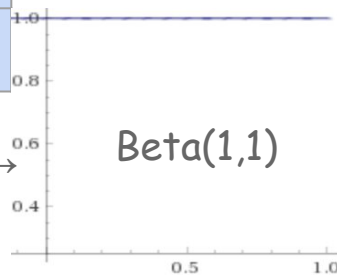
Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.63
2	0.2	1	2	0.15
3	0.9	2	2	0.44

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4		

Sample from each arm's Beta distribution →



THOMPSON EXAMPLE



Algorithm 5 Thompson Sampling Algorithm for Beta-Bernoulli Bandits

- 1: For each arm $i \in \{1, \dots, K\}$, initialize $\alpha_i = \beta_i = 1$. ▷ Set $Beta(\alpha_i, \beta_i)$ prior for each p_i .
- 2: **for** $t = 1, 2, \dots, T$ **do**:
- 3: For each arm i , get sample $s_{i,t} \sim Beta(\alpha_i, \beta_i)$. ▷ Each is a float in $[0, 1]$.
- 4: Pull arm $a_t = \arg \max_{i \in \{1, 2, \dots, K\}} s_{i,t}$. ▷ This “bakes in” exploration!
- 5: Receive reward $r_t \sim Ber(p_{a_t})$.
- 6: **if** $r_t == 1$ **then** $\alpha_{a_t} \leftarrow \alpha_{a_t} + 1$. ▷ Increment number of “successes”.
- 7: **else if** $r_t == 0$ **then** $\beta_{a_t} \leftarrow \beta_{a_t} + 1$. ▷ Increment number of “failures”.

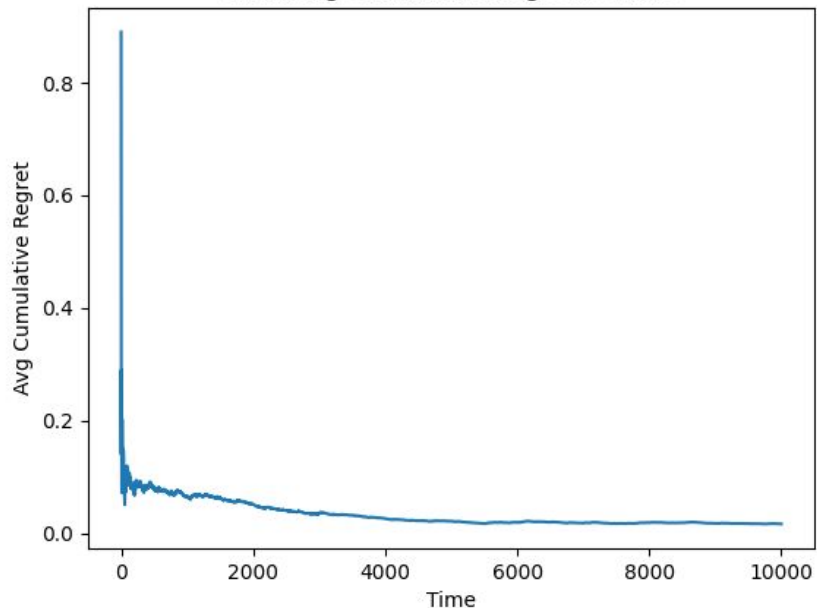
Arm (i)	True p_i	α_i	β_i	$S_{i,t}$
1	0.5	1	1	0.63
2	0.2	1	2	0.15
3	0.9	2	2	0.44

Time (t)	Arm Pulled (a_t)	Reward (r_t)
4		

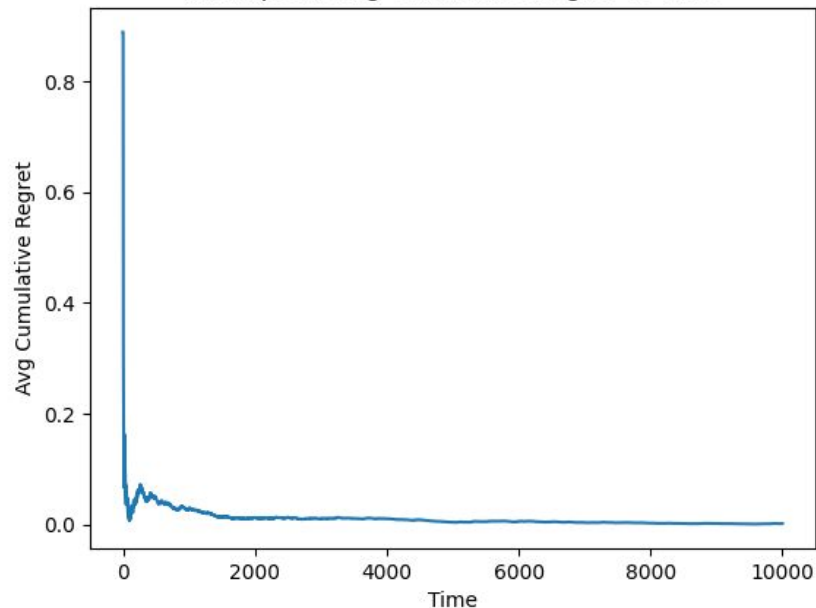
And so on!!! Notice how we explore because there's some chance the “best” arm will have a lower sample occasionally and let other arms win!

UCB VS THOMPSON SAMPLING: AVG REGRET OVER TIME

(ucb) Avg Cumulative Regret vs Time

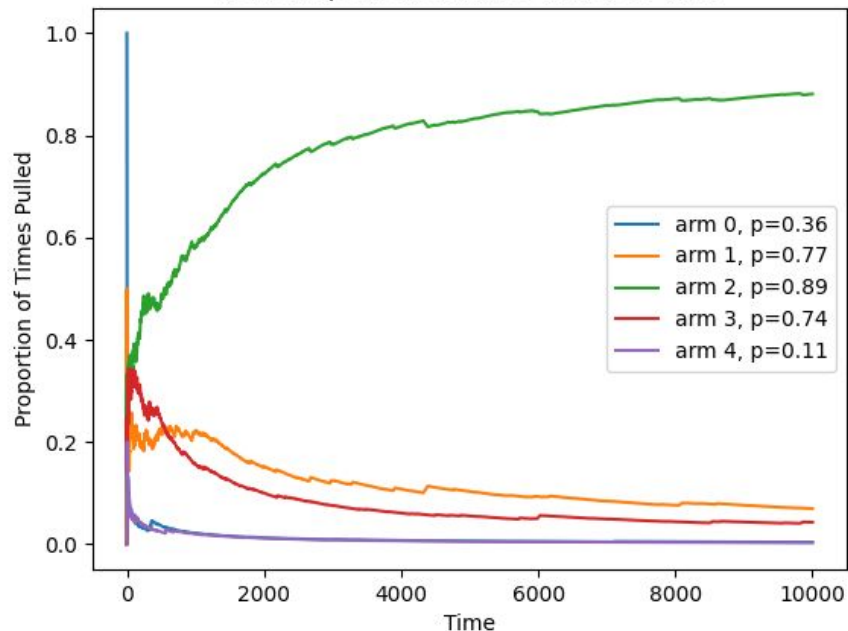


(thompson) Avg Cumulative Regret vs Time

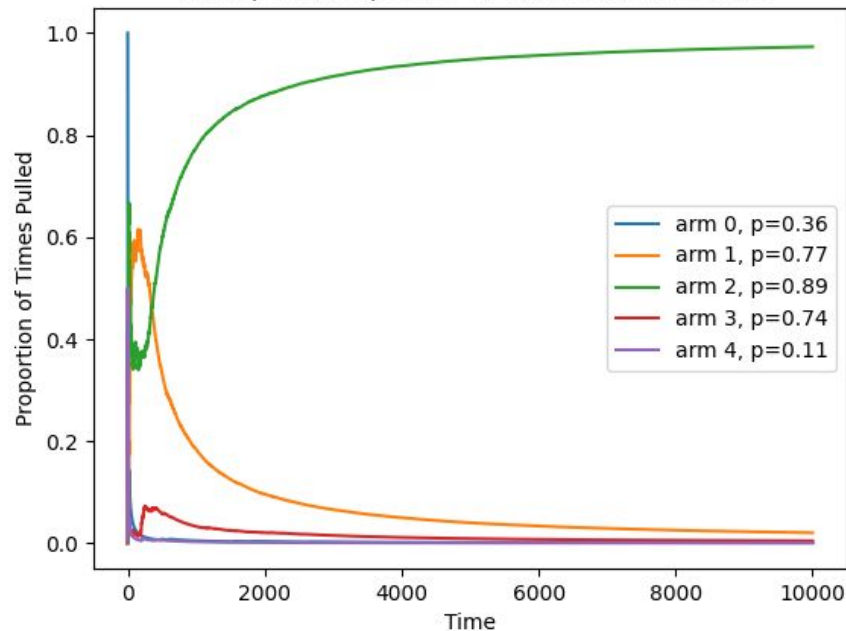


UCB VS THOMPSON SAMPLING: PROPORTION OF TIMES PULLED

(ucb) Proportion of Times Pulled vs Time



(thompson) Proportion of Times Pulled vs Time



TRADITIONAL A/B (HYPOTHESIS) TESTING



A large company wants to experiment releasing a new feature/modification.

Assign

- 99% of population to control group (current feature)
- 1% to experimental group (new feature).

TRADITIONAL A/B (HYPOTHESIS) TESTING



A large company wants to experiment releasing a new feature/modification.

Assign

- 99% of population to control group (current feature)
- 1% to experimental group (new feature).

This has the following consequences:

- If the new feature is "bad", very small percentage of the population sees it, so company protects itself.

TRADITIONAL A/B (HYPOTHESIS) TESTING



A large company wants to experiment releasing a new feature/modification.

Assign

- 99% of population to control group (current feature)
- 1% to experimental group (new feature).

This has the following consequences:

- If the new feature is "bad", very small percentage of the population sees it, so company protects itself.
- If the new feature is "good", very small percentage of the population sees it, so company may lose revenue.

TRADITIONAL A/B (HYPOTHESIS) TESTING



A large company wants to experiment releasing a new feature/modification.

Assign

- 99% of population to control group (current feature)
- 1% to experimental group (new feature).

This has the following consequences:

- If the new feature is "bad", very small percentage of the population sees it, so company protects itself.
- If the new feature is "good", very small percentage of the population sees it, so company may lose revenue.

Can we do better? Can we adaptively assign subjects to each group based on how each is performing rather than deciding at the beginning?

TRADITIONAL A/B (HYPOTHESIS) TESTING



A large company

Assign

- 99% of po
- 1% to exp

This has the fo

- If the new protects it
- If the new may lose r

Can we do better performing rather



company

company

is

MODERN A/B (HYPOTHESIS) TESTING



A large company wants to experiment releasing a new feature/modification.

Assign

- Arm 1: Current Feature
- Arm 2: New feature

MODERN A/B (HYPOTHESIS) TESTING



A large company wants to experiment releasing a new feature/modification.

Assign

- Arm 1: Current Feature
- Arm 2: New feature

When feature is requested by some user, use Multi-Armed Bandit algorithm to decide which feature to show!

MODERN A/B (HYPOTHESIS) TESTING



A large company wants to experiment releasing a new feature/modification.

Assign

- Arm 1: Current Feature
- Arm 2: New feature

When feature is requested by some user, use Multi-Armed Bandit algorithm to decide which feature to show!

(Can have any number of features/arms)

MODERN A/B (HYPOTHESIS) TESTING

When to use Traditional A/B Testing:

- Need to collect data for critical business decisions.
- Need statistical confidence in all your results and impact. Want to learn even about treatments that didn't perform well.
- The reward is not immediate (e.g., if drug testing, don't have time to wait for each patient to finish before experimenting with next patient).
- Optimize/measure multiple metrics, not just one.



MODERN A/B (HYPOTHESIS) TESTING



When to use Traditional A/B Testing:

- Need to collect data for critical business decisions.
- Need statistical confidence in all your results and impact. Want to learn even about treatments that didn't perform well.
- The reward is not immediate (e.g., if drug testing, don't have time to wait for each patient to finish before experimenting with next patient).
- Optimize/measure multiple metrics, not just one.

When to use Multi-Armed Bandits:

- No need for interpreting results, just maximize reward (typically revenue/engagement)
- The opportunity cost is high (if advertising a car, losing a conversion is $\geq \$20,000$)
- Can add/remove arms in the middle of an experiment! Cannot do with A/B tests.

MODERN A/B (HYPOTHESIS) TESTING



When to use Traditional A/B Testing:

- Need to collect data for critical business decisions.
- Need statistical confidence in all your results and impact. Want to learn even about treatments that didn't perform well.
- The reward is not immediate (e.g., if drug testing, don't have time to wait for each patient to finish before experimenting with next patient).
- Optimize/measure multiple metrics, not just one.

When to use Multi-Armed Bandits:

- No need for interpreting results, just maximize reward (typically revenue/engagement)
- The opportunity cost is high (if advertising a car, losing a conversion is $\geq \$20,000$)
- Can add/remove arms in the middle of an experiment! Cannot do with A/B tests.

The study of Multi-Armed Bandits can be categorized as:

- Statistics
- Optimization
- "Reinforcement Learning" (subfield of Machine Learning)