

Imperial College London  
Department of Mathematics

# Mixture Model Bandit

Jiaxu Li

CID: 01378821

Supervised by Ciara Pike-Burke

12th December 2021

Submitted in partial fulfilment of the requirements for the MSc in Statistics of  
Imperial College London

The work contained in this thesis is my own work unless otherwise stated.

Signed: Jiaxu Li

Date: 20th August 2021

# Abstract

This main topic of this project is mixture model bandits. This project is mainly a simulation study that aims to compare the performances of different algorithms under different conditions with the mixture model assumptions. All the bandit algorithms seek to maximize the cumulative rewards in each term or play. Some algorithms are UCB-based, meaning they exploit the upper confidence bounds as a criterion of choosing arms. Other algorithms include Thompson sampling, which places a prior on the parameters of the bandit environment, and also epsilon-greedy algorithm, which uses explore-and-then-commit principle.

A main tool used here is the EM algorithm, which is used to take into account the mixture model structure and calculate estimates of the mixture distributions. A main question here is whether the implementation of EM algorithm in bandit algorithms would lead to a significant improvement in the performance. EM algorithm takes more computational efforts but brings us with a more accurate estimation of the mixture model parameters.

In this project, simulations are firstly done for UCB and UCB-V with and without EM, Thompson sampling, UCB-Normal algorithms under the Gaussian setting, i.e., the rewards of arms are mixtures of Gaussians. Then, Beta reward distributions are also studied to compare UCB-V with and without EM.

# Acknowledgements

I would like to thank Dr Ciara Pike-Burke for her expert advice and encouragement during the completion of this project.

# Table of Contents

<b>1. Introduction</b>	<b>7</b>
1.1. Backgrounds . . . . .	7
1.2. The Language of Bandits . . . . .	7
1.3. Stochastic Bandits, Regrets, and Mixture Model Assumptions . . . . .	8
1.3.1. Mixture Model Assumptions . . . . .	9
<b>2. Known Variance</b>	<b>11</b>
2.1. The Upper Confidence Bound Algorithm . . . . .	11
2.1.1. Assumptions and Theoretical Basis for UCB . . . . .	11
2.1.2. The Optimism Principle . . . . .	13
2.2. Method 1 - without EM* . . . . .	14
2.3. Method 2 - Expectation Maximization Algorithm* . . . . .	15
2.3.1. Introduction to EM Algorithm . . . . .	15
2.3.2. EM for Gaussian Mixture Model . . . . .	16
2.3.3. Using EM in the UCB Algorithm* . . . . .	17
2.4. Empirical Results for Known Variance Case* . . . . .	19
<b>3. Unknown Variance Case</b>	<b>24</b>
3.1. UCB-V (Empirical Bernstein UCB) . . . . .	24
3.1.1. Assumptions for Using UCB-V . . . . .	24
3.1.2. The UCB-V Algorithm . . . . .	24
3.1.3. Natural Choices for the Hyperparameters in the UCB-V Algorithm* . . . . .	25
3.1.4. The UCB-V Algorithm for Mixture Model Bandits* . . . . .	25
3.1.5. The UCB-V Algorithm for Mixture Model Bandits with EM* . . . . .	26
3.2. The UCB1-Normal Algorithm for Mixture Model Bandits* . . . . .	27
3.3. Empirical Results and Regret Comparison for Unknown Variance . . . . .	29
<b>4. Thompson Sampling</b>	<b>34</b>
4.1. The Bayesian Bandit Environment . . . . .	34
4.2. Posterior Distributions in Bayesian Bandits . . . . .	35
4.3. The Bayesian Regret . . . . .	36
4.4. Introduction to Thompson Sampling . . . . .	36
4.5. Thompson Sampling for Mixture Model Bandits* . . . . .	37
4.5.1. Introduction to Variational Bayesian Methods . . . . .	37
4.5.2. Python Implementation of the Thompson Sampling* . . . . .	38

---

The sections marked with (\*) are mostly my original work.

4.6. Empirical Results for Thompson Sampling*	39
<b>5. Mixture Model Bandits with Beta as Component Distributions</b>	<b>44</b>
5.1. Introduction to the E-MM Algorithm	45
5.1.1. The Method of Moments for Beta Distribution	45
5.1.2. The Algorithm	45
5.1.3. Empirical Results for E-MM Algorithm*	47
5.2. Implementation of the E-MM Algorithm in the UCB-V*	47
5.3. Empirical Results for Beta Mixture Model Bandits*	48
<b>6. Conclusion*</b>	<b>51</b>
<b>Bibliography</b>	<b>52</b>
<b>A. Python Codes</b>	<b>53</b>
A.1. bandit_class.py	53
A.2. bandit_class_beta.py	55
A.3. UCB_known_sd.py	56
A.4. UCB_Normal.py	59
A.5. UCB_V.py	61
A.6. Thompson_Sampling.py	62
A.7. UCB_V_Beta.py	64
A.8. regret_comparison.py	70

# 1. Introduction

## 1.1. Backgrounds

In real life, a two-armed bandit is a usually a (gambling) machine of which humans could decide to pull either the left or the right arm, each giving a reward. However, the player has generally no information at all about the distribution of the reward from each different machine. Intuitively, the human player wants to maximize his cumulative rewards throughout his plays. It was named a two-armed bandit in paying tribute to the one-armed bandit, an archaic name for a lever-operated coin machine.

Bandit problems were initially introduced by William R. Thompson. He was interested in medical trials and the humanity in running them. The name can be traced back to 1950s, when Frederick Mosteller and Robert Bush did experiments on animal learning through trials on mice and subsequently on humans. In the experiment, the mice were placed in at the bottom of a T-shaped box and had to choose which way to go to find food. The same experiment was then repeated on the mice for multiple times to study the learning process of the mice and see whether the mice could maximize the cumulative times of finding food at the end of the passage of their choices.

The importance of studying bandit problems is becoming increasingly important these days. All human have to make decisions with uncertainty, and bandit problems can provide an applicable model for this problem.

In a clinical trial design, researchers have used bandit models to generate ideas for their works for nearly 80 years. Adaptive trial design is becoming more and more popular and is recommended by the US Food and Drug Administration because not doing so may result in effective drugs being withheld until a beneficial effect has long been demonstrated. Large technology companies have used bandit algorithms in applications such as news recommendations, dynamic pricing, and ad placement are used. A Bandit algorithm had also been implemented in Monte Carlo Tree Search, an algorithm made famous by the recent feat of Alpha-Go (Lattimore and Szepesvári, 2020).

## 1.2. The Language of Bandits

A bandit problem is a sequential game between a learner and an environment. The game is played over  $n$  rounds, where  $n$  is a positive natural number called the horizon.

In each round  $t \in [n]$ , the learner first chooses an action  $A_t$  from a given set  $\mathcal{A}$ , and the environment then reveals a reward  $X_t \in \mathbb{R}$ .

The learner cannot predict the future when deciding their actions, which means that  $A_t$  should only be dependent on the history  $H_{t-1} = (A_1, X_1, \dots, A_{t-1}, X_{t-1})$ . A policy is a mapping from histories to actions: A learner accept a policy to have interactions with an environment. An environment is a mapping from history sequences ending in actions to rewards. Both the learner and the environment may randomize their decisions. The most common goal of the learner is to choose actions that lead to the largest possible cumulative reward over all  $n$  rounds,  $\sum_{t=1}^n X_t$ .

The main challenge in bandit problems is that the environment is unknown to the learner. All the information that a learner has access to is that the true environment lies in some set  $\mathcal{E}$  called the environment class.

There are several measures to evaluate a learner. Most of the efforts in designing a successful policy are dedicated to understand the regret. There are multiple ways to define this quantity. The most general definition is firstly given here.

### 1.3. Stochastic Bandits, Regrets, and Mixture Model Assumptions

A stochastic bandit is a specific kind of bandit problem, which contains a collection of distributions  $\nu = (P_a : a \in \mathcal{A})$ , where  $\mathcal{A}$  is the set of available actions. The learner and environment interact sequentially over  $n$  rounds. In each round  $t \in \{1, \dots, n\}$ , the learner chooses an action  $A_t \in \mathcal{A}$ , which is fed to the environment. The environment then samples a reward  $X_t \in \mathbb{R}$  from distribution  $P_{A_t}$  and reveals  $X_t$  to the learner. The interaction between the learner (or policy) and environment induces a probability measure on the sequence of outcomes  $A_1, X_1, A_2, X_2, \dots, A_n, X_n$ . In most cases, the horizon  $n$  is finite, but sometimes it's allowed that the interaction continues infinitely ( $n = \infty$ ). The sequence of actions and rewards should satisfy the following assumptions (Lattimore and Szepesvári, 2020):

- (a) The conditional distribution of reward  $X_t$  given  $A_1, X_1, \dots, A_{t-1}, X_{t-1}, A_t$  is  $P_{A_t}$ , which captures the intuition that the environment samples  $X_t$  from  $P_{A_t}$  in round  $t$ .
- (b) The conditional law of action  $A_t$  given  $A_1, X_1, \dots, A_{t-1}, X_{t-1}$  is  $\pi_t(\cdot \mid A_1, X_1, \dots, A_{t-1}, X_{t-1})$ , where  $\pi_1, \pi_2, \dots$  is a sequence of probability kernels that characterize the learner. The most important assumption is that the learner cannot use the future observations in current decisions.

The learner's goal is to maximize the total reward  $S_n = \sum_{t=1}^n X_t$ , which is a random variable that depends on both the actions of the learner and the rewards sampled from



the environment.

In the previous section, we have defined the regret in a general term as being the deficiency suffered by the learner relative to the optimal policy. Let  $\nu = (P_a : a \in \mathcal{A})$  be a stochastic bandit and define

$$\mu_a(\nu) = \int_{-\infty}^{\infty} x dP_a(x)$$

Then let  $\mu^*(\nu) = \max_{a \in \mathcal{A}} \mu_a(\nu)$  be the largest mean of all the arms.

**Definition 1.** (Lattimore and Szepesvári, 2020) *The regret of policy  $\pi$  on bandit instance  $\nu$  is*

$$R_n(\pi, \nu) = n\mu^*(\nu) - \mathbb{E} \left[ \sum_{t=1}^n X_t \right],$$

where the expectation is taken with respect to the probability measure on outcomes induced by the interaction of  $\pi$  and  $\nu$ . Minimizing the regret is at the same time maximizing the expectation of  $S_n$ .

In a bandit setting, when measuring the risk, we can either have the random regret or pseudo-regret defined by

$$\begin{aligned} \hat{R}_n &= n\mu^* - \sum_{t=1}^n X_t, & (\text{random regret}) \\ \bar{R}_n &= n\mu^* - \sum_{t=1}^n \mu_{A_t}. & (\text{pseudo-regret}) \end{aligned} \tag{1.1}$$

In this project and all the plots presented afterwards, pseudo-regrets are used because while  $\hat{R}_n$  is influenced by the noise  $X_t - \mu_{A_t}$  in the rewards, the pseudo-regret sieves this out by expectation, thus making it a better measure of the ‘skill’ of a policy.

### 1.3.1. Mixture Model Assumptions

In this project, we focus on a particular bandit environment where the reward of each arm follows a mixture of independent distributions. Suppose  $X_t$  is the reward of an arm of the stochastic bandits defined above, then we can write the distribution of  $X_t$  as:

$$f_X(x) = \sum_{i=1}^m \pi_i g_i(x),$$

where  $g_i(x)$  for  $i = 1, \dots, m$  such that  $m \geq 2$  are independent component distributions and  $\alpha_i$  are the individual probabilities assigned to the corresponding component distributions such that  $\sum_{i=1}^m \pi_i = 1$ .

It's assumed that all the arms in mixture model bandit environments follow the distribution defined above. What's flexible are the parameters of the component distributions and the categorical probabilities assigned to these components.

## 2. Known Variance

To start with, let's first assume that the variances of each arm are known to the learner and also that the component distributions are Gaussian. This means for each arm the reward obtained follows the distribution

$$f_X(x) = \sum_{i=1}^m \pi_i g_i(x), \quad \text{s.t. } g_i(x) = \frac{1}{\sigma_i \sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \left( \frac{x - \mu_i}{\sigma_i} \right)^2 \right\} \text{ for } i = 1 \rightarrow m,$$

where  $\sigma_1, \dots, \sigma_K$  are known to the learner throughout the plays.

### 2.1. The Upper Confidence Bound Algorithm

One very commonly used and well-studied algorithm for stationary stochastic bandits is the UCB algorithm, namely the upper confidence bound algorithm. Before adapting it to account for the mixture model assumption, this section aims to establish the basic assumptions of the UCB algorithm and its theoretical basis.

#### 2.1.1. Assumptions and Theoretical Basis for UCB

One assumption that must be met in order to use UCB algorithm is sub-gaussianity.

**Definition 2.** (Lattimore and Szepesvári, 2020) *A random variable  $X$  is  $\sigma$ -subgaussian if for all  $\lambda \in \mathbb{R}$ , it holds that  $\mathbb{E}[\exp(\lambda X)] \leq \exp(\lambda^2 \sigma^2 / 2)$ .*

Gaussian variables are clearly sub-gaussian themselves by the definition and the moment generating function of the gaussian. Theorem 4 explains the origin of the term ‘subgaussian’. The tails of a  $\sigma$ -subgaussian random variable decay approximately as fast as that of a Gaussian with zero mean and the same variance.

Note that for random variables that are not centred ( $\mathbb{E}[X] \neq 0$ ), we abuse notation by saying that  $X$  is  $\sigma$ -subgaussian if the noise  $X - \mathbb{E}[X]$  is  $\sigma$ -subgaussian. A distribution is called  $\sigma$ -subgaussian if a random variable drawn from that distribution is  $\sigma$ -subgaussian. Subgaussianity is in fact a property to be possibly satisfied by both a random variable and the measure on the space on which it is defined (Lattimore and Szepesvári, 2020).

**Lemma 3.** *For any random variable  $X$  and  $\varepsilon > 0$ , the following holds:*

(a) (Markov):  $\mathbb{P}(|X| \geq \varepsilon) \leq \frac{\mathbb{E}[|X|]}{\varepsilon}$ .

(b) (Chebyshev):  $\mathbb{P}(|X - \mathbb{E}[X]| \geq \varepsilon) \leq \frac{\mathbb{V}[X]}{\varepsilon^2}$ .

We need the above lemmas when proving the theorem 4, which provides a bound for the tail behaviour of sub-gaussian random variables.

**Theorem 4.** (Lattimore and Szepesvári, 2020) If  $X$  is  $\sigma$ -subgaussian, then for any  $\varepsilon \geq 0$ ,

$$\mathbb{P}(X \geq \varepsilon) \leq \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right)$$

*Proof.* We take a generic approach called the Cramér-Chernoff method. Let  $\lambda > 0$  be some constant to be tuned later. Then

$$\begin{aligned} \mathbb{P}(X \geq \varepsilon) &= \mathbb{P}(\exp(\lambda X) \geq \exp(\lambda \varepsilon)) \\ &\leq \mathbb{E}[\exp(\lambda X)] \exp(-\lambda \varepsilon) \quad (\text{Markov's inequality}) \\ &\leq \exp\left(\frac{\lambda^2 \sigma^2}{2} - \lambda \varepsilon\right) \quad (\text{Def. of subgaussianity}) \end{aligned}$$

Choosing  $\lambda = \varepsilon/\sigma^2$  completes the proof.  $\square$

**Lemma 5.** (Lattimore and Szepesvári, 2020) Suppose that  $X$  is  $\sigma$ -subgaussian and  $X_1$  and  $X_2$  are independent and  $\sigma_1$  and  $\sigma_2$ -subgaussian, respectively, then :

- (a)  $\mathbb{E}[X] = 0$  and  $\mathbb{V}[X] \leq \sigma^2$ .
- (b)  $cX$  is  $|c|\sigma$ -subgaussian for all  $c \in \mathbb{R}$ .
- (c)  $X_1 + X_2$  is  $\sqrt{\sigma_1^2 + \sigma_2^2}$ -subgaussian.

**Corollary 6.** (Lattimore and Szepesvári, 2020) Assume that  $X_i - \mu$  are independent,  $\sigma$ -subgaussian random variables. Then for any  $\varepsilon \geq 0$ ,

$$\mathbb{P}(\hat{\mu} \geq \mu + \varepsilon) \leq \exp\left(-\frac{n\varepsilon^2}{2\sigma^2}\right) \quad \text{and} \quad \mathbb{P}(\hat{\mu} \leq \mu - \varepsilon) \leq \exp\left(-\frac{n\varepsilon^2}{2\sigma^2}\right)$$

where  $\hat{\mu} = \frac{1}{n} \sum_{t=1}^n X_t$ .

*Proof.* By Lemma 5, it holds that  $\hat{\mu} - \mu = \sum_{i=1}^n (X_i - \mu) / n$  is  $\sigma/\sqrt{n}$ -subgaussian. Then apply Theorem 4.  $\square$

The above results is equivalent to the statement that under the assumptions of the result, for any  $\delta \in [0, 1]$ , with probability at least  $1 - \delta$  (Lattimore and Szepesvári, 2020):

$$\mu \leq \hat{\mu} + \sqrt{\frac{2\sigma^2 \log(1/\delta)}{n}} \quad (2.1)$$

### 2.1.2. The Optimism Principle

The UCB algorithm is built upon *the principle of optimism in the face of uncertainty*, which states that a player should choose actions as if the environment is as nice as plausibly possible. The principle is useful not only for the finite-armed stochastic bandit problem (Lattimore and Szepesvári, 2020).

For bandits, the optimism principle translates to using the data (the rewards from each arm) observed so far to label each arm with a value, called the upper confidence bound that with overwhelming probability is an overestimate of the unknown mean. Note that the upper confidence bounds are updated after each play. The intuitive reason why this algorithm works is not hard to see. If the upper confidence bound assigned to the optimal arm is indeed an overestimate, then another arm can only be played if its upper confidence bound is greater than that of the optimal arm, which is itself larger than the mean of the optimal arm. And yet this cannot happen too frequently since data obtained by playing a suboptimal arm will provide more information about that sub-optimal arm. The upper confidence bound for this arm will eventually fall below that of the optimal arm.

In order to make the argument rigorous, it's necessary to define the upper confidence bound. Let  $(X_t)_{t=1}^n$  be a sequence of independent 1-subgaussian random variables with mean  $\mu$  and  $\hat{\mu} = \frac{1}{n} \sum_{t=1}^n X_t$ . By Eq. 2.1,

$$\mathbb{P} \left( \mu \geq \hat{\mu} + \sqrt{\frac{2 \log(1/\delta)}{n}} \right) \leq \delta \quad \text{for all } \delta \in (0, 1). \quad (2.2)$$

When deciding its possible actions in play  $t$ , the learner has observed  $T_i(t-1)$  samples from arm  $i$  and received rewards from that arm with an empirical mean of  $\hat{\mu}_i(t-1)$ . Then a reasonable candidate ‘as large as plausibly possible’ for the unknown mean of the  $i$  th arm can be defined as:

$$\text{UCB}_i(t-1, \delta) = \begin{cases} \infty & \text{if } T_i(t-1) = 0 \\ \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}} & \text{otherwise} \end{cases} \quad (2.3)$$

To sum up, the UCB algorithm can therefore be presented as follows:

---

**Algorithm 1:** Upper Confidence Bound Algorithm (Lattimore and Szepesvári, 2020)

---

**Data:**  $k$  and  $\delta$

For  $t = 1$ , sample each arm once ;

**for**  $t \in 2, \dots, n$  **do**

    Choose action  $A_t = \operatorname{argmax}_i \text{UCB}_i(t-1)$ ;

    Observe reward  $X_{A_t}$  and update upper confidence bounds;

**end**

---

## 2.2. Method 1 - without EM\*

The first method for tackling this problem is to deliberately ignore the mixture-model assumption. This is based on the fact that mixture distributions may inherit some properties from their component distributions, for example if both component distributions are bounded then the mixture distribution will also be bounded. In this case, a mixture of subgaussian random variables is also subgaussian by the following derivation. This establishes that UCB (Upper Confidence Bound) algorithm can be also applied here with assumptions properly satisfied.

For now, the assumption would be that the variances of  $m$  component distributions are known. Suppose that an arm has a mixture model reward  $X$  with  $m$  components  $X_1, \dots, X_m$  with given variances  $\sigma_1^2, \dots, \sigma_m^2$  but unknown means  $\mu_i$  and weights  $\pi_i$ . Then, we can derive an upper bound for the sub-gaussian parameter of  $X$  with the definition of sub-gaussian as in Def. 2:

$$\begin{aligned}
 \mathbb{E} \left[ e^{\lambda X} \right] &= \int e^{\lambda x} \cdot f_x(x) dx \\
 &= \int e^{\lambda x} \cdot \sum_{i=1}^m \pi_i f_{X_i}(x) dx \\
 &= \sum_{i=1}^m \pi_i \int e^{\lambda x} f_{X_i}(x) dx \\
 &= \sum_{i=1}^m \pi_i e^{\lambda^2 \sigma_i^2 / 2} \\
 &\leq e^{\lambda^2 \max_i \{ \sigma_i^2 \} / 2}
 \end{aligned} \tag{2.4}$$

This means the reward here is at most  $\max\{\sigma_1, \dots, \sigma_m\}$ -subgaussian by Def. 2 and also note that this value is likely an overestimate of the sub-gaussian parameter of the mixture model reward.

From Eq. 2.2, given  $(X_t)_{t=1}^n$ , a sequence of  $\max\{\sigma_1, \dots, \sigma_m\}$ -subgaussian random variables, the inequality can be established as:

$$\mathbb{P} \left( \mu \geq \hat{\mu} + \max_i \{\sigma_i\} \sqrt{\frac{2 \log(1/\delta)}{n}} \right) \leq \delta \quad \text{for all } \delta \in (0, 1)$$

From Eq. 2.3, this can be rewritten in this form to account for the times that each arm  $i$  has been sampled at time  $t$ ,  $T_i(t-1)$ :

$$\text{UCB}_i(t-1, \delta) = \hat{\mu}_i(t-1) + \max_i \{\sigma_i\} \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}} \quad (2.5)$$

For asymptotically optimal UCB, a better regret bound can be achieved by replacing the confidence level term  $\delta$  with a function of  $t$  (Lattimore and Szepesvári, 2020).

$$\text{UCB}_i(t-1) = \hat{\mu}_i(t-1) + \max_i \{\sigma_i\} \sqrt{\frac{2 \log f(t)}{T_i(t-1)}}, \text{ where } f(t) = 1 + t \log^2(t) \quad (2.6)$$

This derived bound can then be implemented in the aforementioned UCB algorithm, which hopefully will give us a low cumulative regret.

## 2.3. Method 2 - Expectation Maximization Algorithm\*

The other method would be to consider the mixture model assumption and apply the EM (expectation-maximization algorithm) within a bandit algorithm to fit the mixture reward distribution of each arm.

### 2.3.1. Introduction to EM Algorithm

In statistics, an expectation-maximization (EM) algorithm is an iterative method to find (local) maximum likelihood or maximum a posteriori (MAP) estimates of parameters in statistical models, where the model depends on unobserved or so-called latent variables. The EM iterates between a expectation (E) step, which creates a function for the expectation of the log-likelihood evaluated using the current estimate for the parameters, and a maximization (M) step, which computes parameters maximizing the expected log-likelihood found on the E step. These parameter-estimates are then used to determine the distribution of the latent variables in the next E step. A formal description of the procedure firstly introduced by Dempster, Laird, and Rubin (1977) is given below:

Given the statistical model which generates a set  $\mathbf{X}$  of observed data, a set of unobserved latent data or missing values  $\mathbf{Z}$ , and a vector of unknown parameters  $\boldsymbol{\theta}$ , along with a likelihood  $L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$ , the maximum likelihood estimate (MLE) of the unknown parameters is determined by maximizing the marginal likelihood of the observed data:

$$L(\boldsymbol{\theta}; \mathbf{X}) = p(\mathbf{X} | \boldsymbol{\theta}) = \int p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) d\mathbf{Z} = \int p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}) p(\mathbf{X} | \boldsymbol{\theta}) d\mathbf{Z}$$

However, this quantity is almost always intractable since  $\mathbf{Z}$  is unobserved and the distribution of  $\mathbf{Z}$  is unknown before attaining  $\boldsymbol{\theta}$

The EM algorithm seeks to find the MLE of the marginal likelihood by iteratively applying these two steps (Dempster, Laird, and Rubin, 1977):

- (a) Expectation step (E step): Define  $Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)})$  as the expected value of the log likelihood function of  $\boldsymbol{\theta}$ , with respect to the current conditional distribution of  $\mathbf{Z}$  given  $\mathbf{X}$  and the current estimates of the parameters  $\boldsymbol{\theta}^{(t)}$  :

$$Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}) = \mathbb{E}_{\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\theta}^{(t)}} [\log L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})]$$

- (b) Maximization step (M step): Find the parameters that maximize this quantity:

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)})$$

### 2.3.2. EM for Gaussian Mixture Model

It's assumed in this project that the reward from each arm follows the Gaussian Mixture Model, abbreviated as GMM. In this section, the component variances of each arm are known to the player (but not the weights  $\pi_i$  and variances  $\sigma_i^2$ ).

There are many existing methods for estimating the parameters of a GMM, but EM based on the maximum likelihood estimation is perhaps the most appropriate one to be applied here.

For instance, consider the case where a dataset consists of many points that happen to be generated by two different distributions. The samples for each process follow a Gaussian distribution, but the data is mixed and thus forming a mixture. The distributions are overlapped and close enough that it is not apparent to see which distribution a given sample may originate from.

The stochasticity generating the data point from different component distributions is represented by a latent random variable. It influences the data but is not observable. As such, the EM algorithm is an appropriate approach to use to estimate the parameters of the distributions.

In the EM algorithm, the estimation-step try to find a value of the latent variable for each sample, and the maximization step would try to search for the optimal parameters of the probability distributions to try to best explain the data. The process is iterated until a set threshold has been met. The algorithm can be summarized as:

- (a) E-Step. Estimate the expected value for each latent variable. This can be shown to be equivalent to calculating the probability that sample  $i$  belongs to the cluster



$k$ , which can be denoted as  $w_{ik}$ , which is calculated using Bayesian theorem as illustrated below. Here,  $z_{ik}$  stands for the latent variable indicating that the sample  $i$  comes from the cluster  $k$ ,  $\Theta$  the parameters of the model  $(\mu, \Sigma)$ ,  $X$  the Gaussian mixture model random variable of interest.

$$\begin{aligned} w_{ik} &= p(z_{ik} = 1 | \mathbf{x}_i, \Theta) = \frac{p(X = \mathbf{x}_i | z_{ik} = 1, \Theta) p(z_{ik} = 1)}{p(X = \mathbf{x}_i | \Theta)} \\ &= \frac{p_k(\mathbf{x}_i | z_k, \mu_k, \Sigma_k) \cdot \alpha_k}{\sum_{m=1}^K p_m(\mathbf{x}_i | z_m, \mu_m, \Sigma_m) \cdot \alpha_m} \end{aligned}$$

- (b) M-Step. Optimize the parameters of the distribution using maximum likelihood. In the known variance case, only the means and weights  $\pi_i$  are being iteratively updated (the first two terms). In a more common unknown variance case, update the parameters by:

$$\pi_k = \frac{\sum_{n=1}^N w_{nk}}{N}, \quad \mu_k^* = \frac{\sum_{n=1}^N w_{nk} \mathbf{x}_n}{\sum_{n=1}^N w_{nk}}, \quad \left( \Sigma_k^* = \frac{\sum_{n=1}^N w_{nk} (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T}{\sum_{n=1}^N w_{nk}} \right)$$

### 2.3.3. Using EM in the UCB Algorithm\*

One can apply EM within the UCB bandit algorithm to fit each arm. Suppose we have  $m$  components of Gaussian. If we can use EM algorithm to estimate the weights  $\hat{\pi}_i$ , then an accurate estimate for variance can be established. Suppose the we have a mixture model reward  $X$  with  $m$  components  $X_1, \dots, X_m$  with weights  $\pi_1, \dots, \pi_m$  such that  $\sum_{i=1}^m \pi_i = 1$ . Suppose the components are all centered and variances are denoted as  $\sigma_1^2, \dots, \sigma_m^2$ . Then,

$$\begin{aligned} \mathbb{E}[X] &= \mu = \sum_{i=1}^m \pi_i \mu_i = 0, \\ \mathbb{E}[(X - \mu)^2] &= \sigma^2 = \mathbb{E}[X^2] - \mu^2 = \sum_{i=1}^m \pi_i \mathbb{E}[X_i^2] = \sum_{i=1}^m \pi_i \sigma_i^2. \end{aligned} \tag{2.7}$$

From this, we know that the mixture model reward here is at least  $\sqrt{\sum_{i=1}^m \pi_i \sigma_i^2}$ -subgaussian. This is apparently a better estimate for the sub-gaussian parameter, compared to what we have derived before ignoring the mixture model assumption:  $\max\{\sigma_1, \dots, \sigma_m\}$ -subgaussian.

Therefore, this can be used to give a better upper confidence bound with estimated weights  $\hat{\pi}_i$  in a similar fashion as in Eq. 2.3

$$\text{UCB}_i(t-1, \delta) = \begin{cases} \infty & \text{if } T_i(t-1) = 0 \\ \hat{\mu}_i(t-1) + \sum_{i=1}^m \hat{\pi}_i \sigma_i^2 \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}} & \text{otherwise} \end{cases} \tag{2.8}$$

Similarly, this can be adapted to the asymptotically optimal UCB as in Eq. 2.6:

$$\text{UCB}_i(t-1) = \hat{\mu}_i(t-1) + \sum_{i=1}^m \hat{\pi}_i \sigma_i^2 \sqrt{\frac{2 \log f(t)}{T_i(t-1)}}, \text{ where } f(t) = 1 + t \log^2(t) \quad (2.9)$$

Then, the algorithm can be summarized in the following page. Note that as we have known variance in this case, the only thing that EM algorithm is improving the estimation of is the weights  $\pi_i$ . The parameters that are being updated in the EM algorithm are  $\pi_i$  and  $\mu_i$  for  $i \in 1, \dots, m$ . However, the component means are not being used here as the estimates for the means of rewards because the simple sample mean  $\frac{1}{N} \sum_{i=1}^N X_i$  is already an unbiased estimator for the mean. Therefore, there is no need to use a more sophisticated expression  $\sum_{i=1}^m \pi_i \mu_i$  that also involves the output from EM algorithm as the estimate for the mean of the rewards.

---

**Algorithm 2:** UCB using EM (known variance)

---

**Input :**  $x, \dots, x_N$ , acc (accuracy),  $n$  (the horizon),  $k$  (the number of arms),  $m$  (the number of mixture components per arm),  $\sigma_{ij}^2$  (the variances) for  $i = 1 \rightarrow k$  and  $j = 1 \rightarrow m$

**for**  $t = 1, \dots, k$  **do** // Initialization  
    Observe from arm  $t$  for  $m$  times and observe  $X_{t,1}, \dots, X_{t,m}$ .  
    Initialize  $\bar{X}_{t,m}$  as  $\frac{1}{m} \sum_{i=1}^m X_{t,i}$ .  
    Initialize  $\pi_{t,j}$  randomly for  $j = 1 \rightarrow m$ .  
**end**

**for**  $t = km + 1, \dots, n$  **do** // Standard UCB-EM  
    Choose the action  $\mathcal{A}_t$  that maximizes the UCB bound as in Eq. 2.9.  
    Write  $\mathcal{A}_t = i$  as the action, and observe  $X_t$  from the arm  $i$ .  
    Update  $\hat{\pi}_{ij}$  for  $j = 1 \rightarrow m$  using the EM algorithm given  $X_{i,1}, \dots, X_{i,T_i(t)}$  and  $\sigma_{ij}^2$  for  $j = 1 \rightarrow m$  (here  $T_i(t)$  stands for the number of previous samples from the arm  $i$ ).  
    Update  $\hat{\mu}_i(t)$  as  $\frac{1}{T_i(t)} \sum_{j=1}^{T_i(t)} X_{i,j}$  (samples from the arm  $i$  until time  $t$ ).  
**end**

---

Note that the initialization stage is necessary because we need at least some samples to start with in order to use the EM algorithm. If the number of samples to start with is too small, we have two problems. Firstly is that the program would not allow it (the package `sklearn.mixture.GaussianMixture` is used for the case when the variance is **unknown**, in the **known** variance case the EM is coded up by myself to account for the known variance feature). Secondly is that even when the algorithm converges, the  $\pi_i$  is likely to be stuck in the local minimum 0 or 1, which is undesirable.

## 2.4. Empirical Results for Known Variance Case\*

In this section, the empirical results about the performances of the two algorithms aforementioned are given: the one ignoring the mixture model structure, and the one using EM. Note that in all the plots afterwards, psuedo-regrets are used as a means of comparing the performances of algorithms (policies) as described in Eq. 1.1.

$$\bar{R}_n = n\mu^* - \sum_{t=1}^n \mu_{A_t} \quad (\text{pseudo-regret})$$

For all the plots afterwards, the parameters information is written in the titles:

- **k**: the number of arms.
- **n**: the horizon (number of plays).
- **m**: the number of component distributions (in the initial case of mixture distribution,  $m = 2$ )
- **r**: the number of repeated experiments under the same bandit environment.
- **$a < \mu < b$** : this indicates that the means of all the component distributions of this bandit environment are sampled from the uniform distribution over  $[a, b]$ .
- **$a < \sigma < b$** : this indicates that the standard deviations of all the component distributions of this bandit environment are sampled from the uniform distribution over  $[a, b]$ .

For each bandit environment, the experiment is repeated for 30 times because the regret is generated stochastically and so is the reward from each arm. A 95% confidence interval is generated while the middle line stands for the mean regret over 30 experiments.

The main observation here is that UCB-EM performs uniformly better and beats UCB in almost all the bandit environments. Each page is afterwards is denoted to a fixed number of components. For example, for figure 2.1 to figure 2.6, the reward is a mixture of 2 gaussians, i.e.,  $f_X(x) = \alpha f_1(x) + (1 - \alpha)f_2(x)$ .

In the 6 plots drawn each page, the left 3 ones is devoted to a more challenging bandit environment with means of gaussians initialized from the range  $[0, 1]$ , while for the right 3 plots the means are generated from the range  $[0, 5]$ . The reason why the these particular bandit environments are deemed easy or difficult is not hard to see: imagine a 2-arm bandit problem such that the first arm keeps giving rewards between 0 and 1 while the second arm only gives very high reward randomly in the range 100 and 101. This is arguably a very easy bandit problem if we assume that the rewards are random uniform variates. In other words, the bandit problem is deemed more difficult to minimize regret if the mean rewards of arms are more clustered together and close to each other.

The UCB-EM algorithm approximately takes 10 times more time to run compared to UCB alone. Even though the computational efforts are definitely bigger with the implementation of EM algorithm, the efforts pay off in the regret and have somehow improved the performance of the UCB policy in almost all bandit environments in which experiments have done.

One other observation is that in a less challenging bandit environments (the right 3 plots on each page), the gap between the regrets of these algorithms are visibly smaller. That's because in an easier bandit environment, as explained before, both of the algorithms can learn the optimal arm in relatively less time, thus overshadowing the benefits brought by a more accurate estimate of the variance term by EM algorithm. Thus, one can choose to use either UCB or UCB-EM in a more careless way if the bandit problem is deemed 'easy'.

There can be some way of detecting whether a bandit problem is easy or difficult. If a UCB-based method is implemented, the relative ratios among the number of plays for different machine can be used as a criterion. If we want to use something flashy, **the Gini coefficient** is a possible choice here as originally a measure for the statistical dispersion. A player can firstly do a number of plays, and then calculate the Gini coefficient to see whether all the rewards are more or less coming from a single arm or are actually alternating among different arms, indicating that the algorithm finds it difficult to locate the optimal arm. A subsequent decision on which algorithm to use can then be made. As the rewards and actions are recorded as history, there would be no big problems stemming from changing an algorithm to another during the plays as long as the player is trying to minimize the cumulative regret.

Another observable trend is that when the number of components  $m$  increases, so is the difference in performances of the UCB and UCB-EM algorithms under a more challenging bandit environment. The EM algorithm really has done more efforts in estimating a more accurate variance when the number of components increases. However, as  $m$  increases, the difference in performance becomes less distinguishable when the bandit environment is easy.

The following figures compare the UCB and UCB-EM algorithm with  $m = 2$  (2 components).

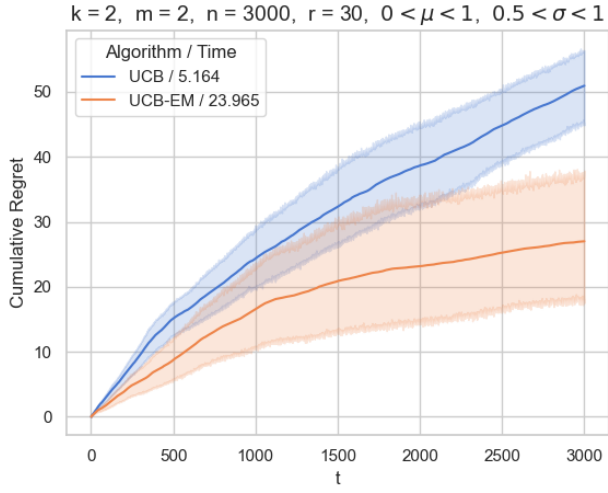


Figure 2.1

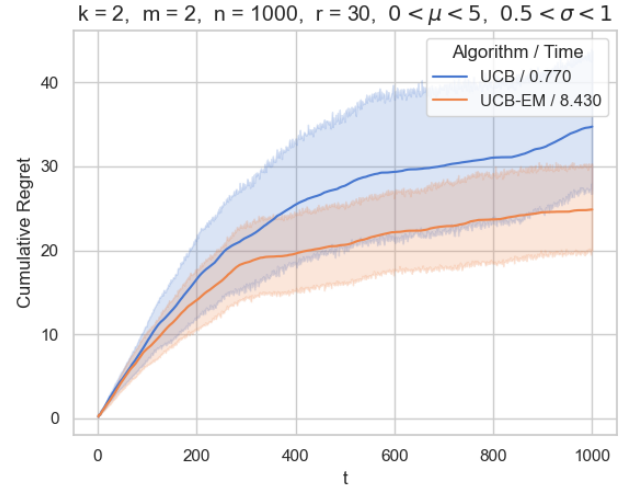


Figure 2.2

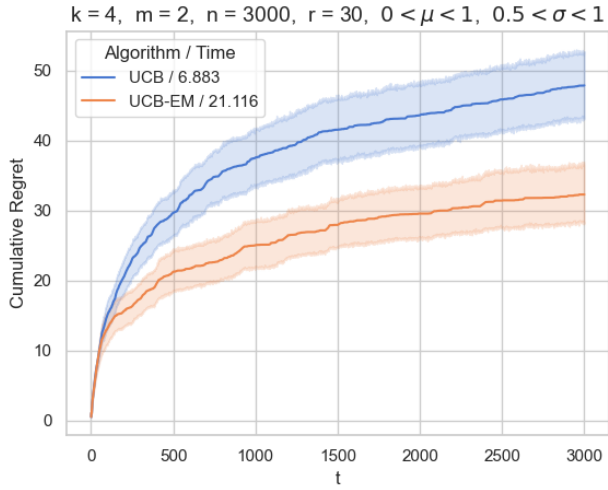


Figure 2.3

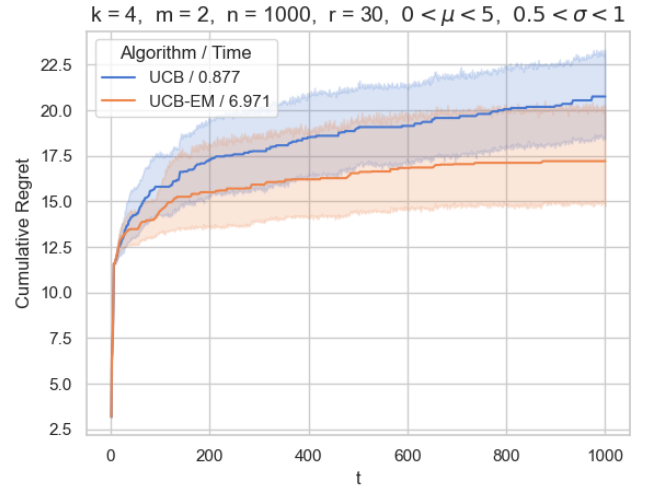


Figure 2.4

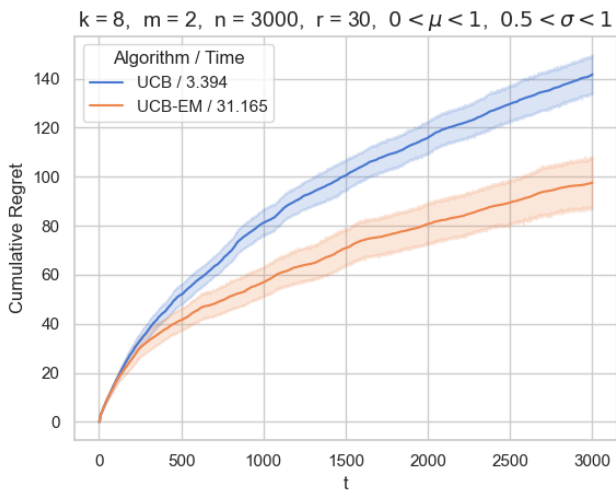


Figure 2.5

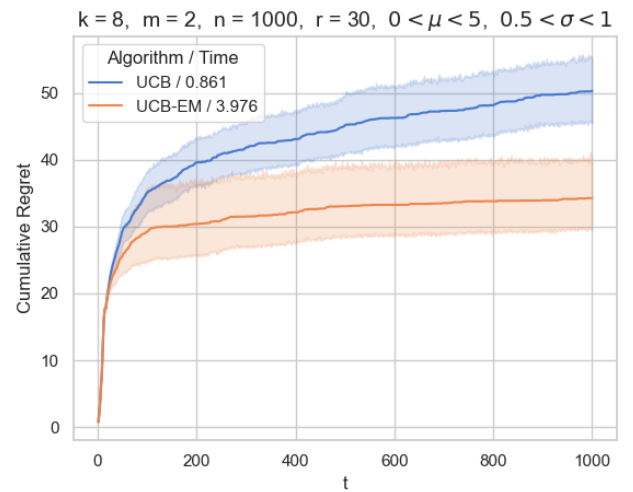


Figure 2.6

The following figures compare the UCB and UCB-EM algorithm with  $m = 3$  (3 components).

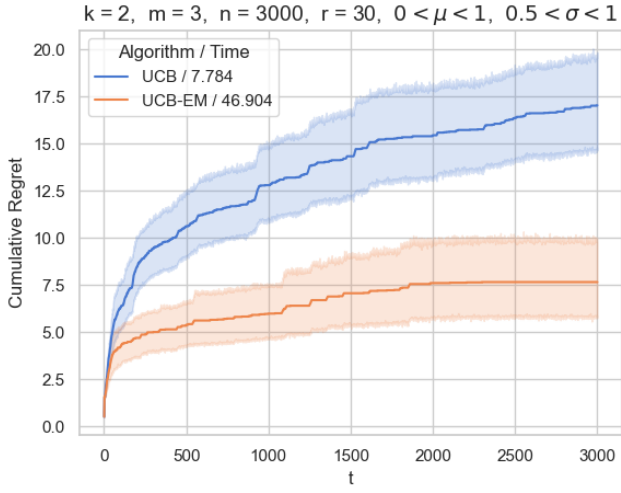


Figure 2.7

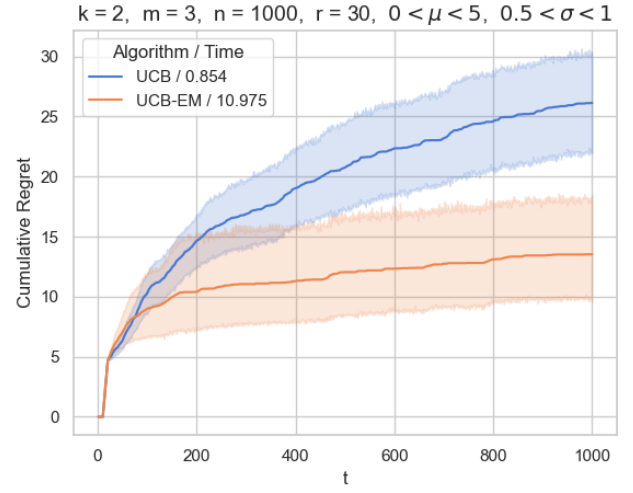


Figure 2.8

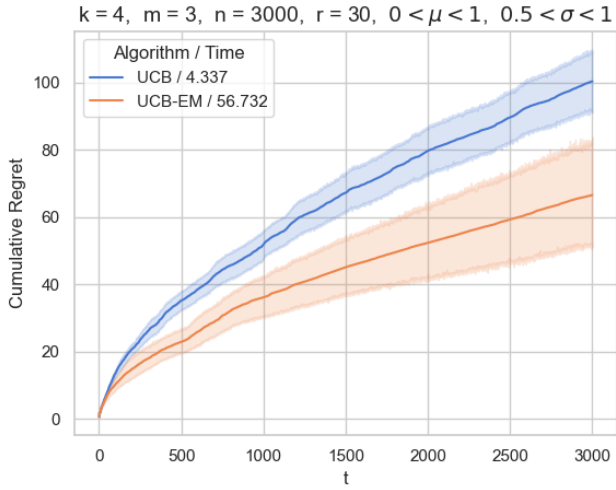


Figure 2.9

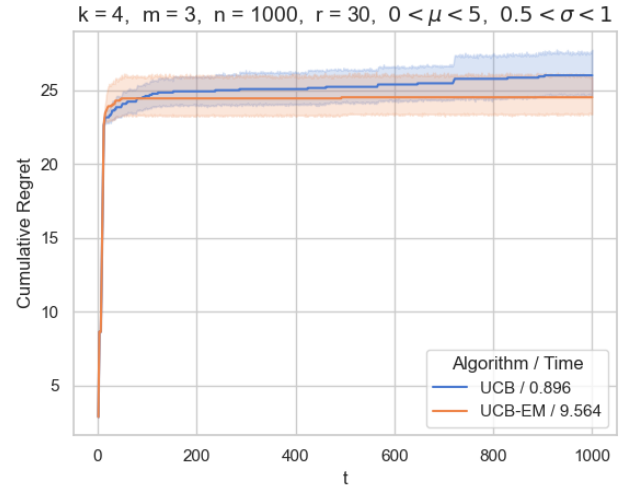


Figure 2.10

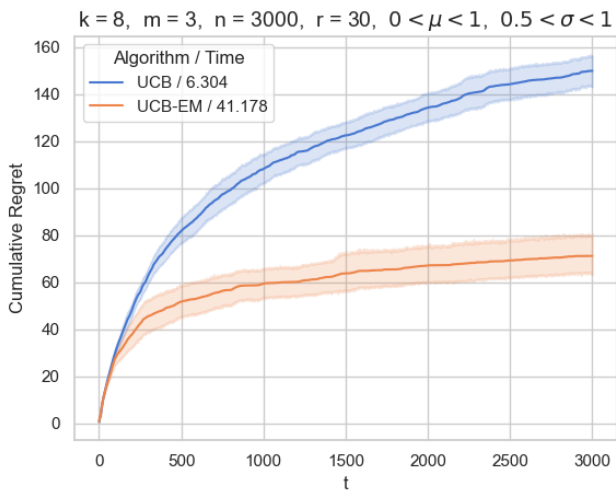


Figure 2.11

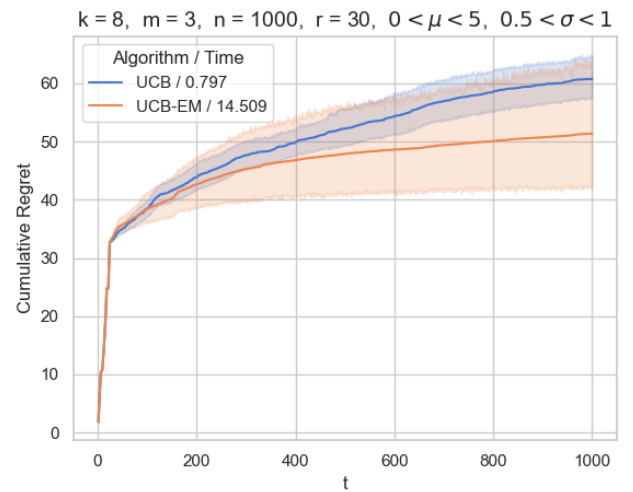


Figure 2.12

The following figures compare the UCB and UCB-EM algorithm with  $m = 4$  (4 components).

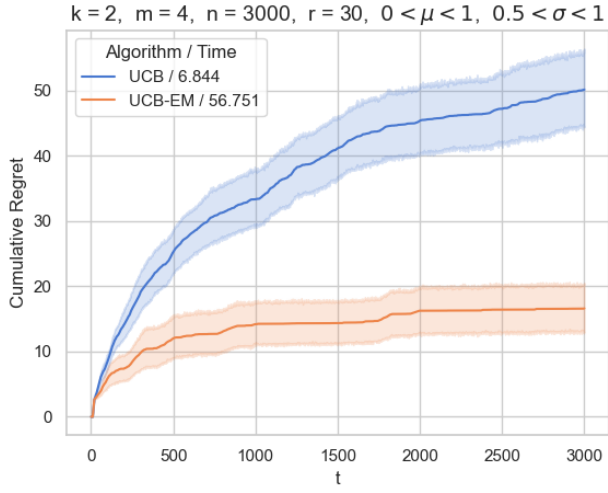


Figure 2.13

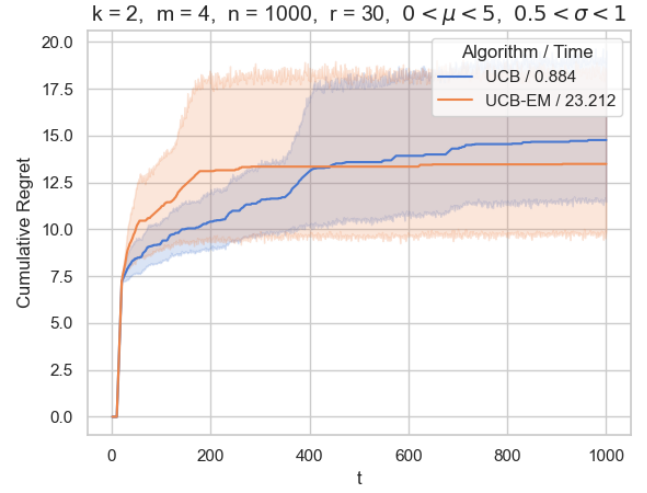


Figure 2.14

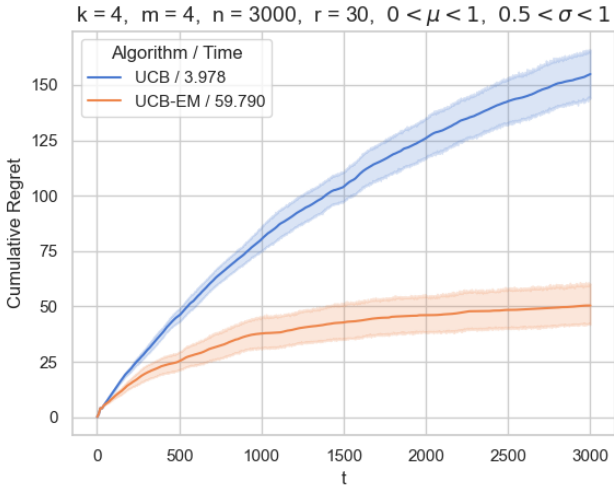


Figure 2.15

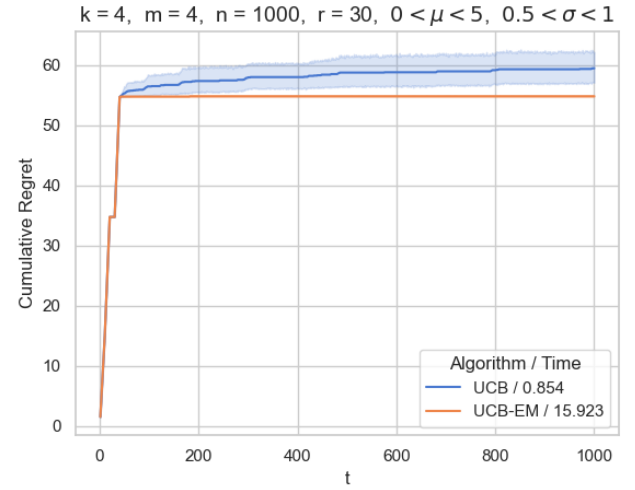


Figure 2.16

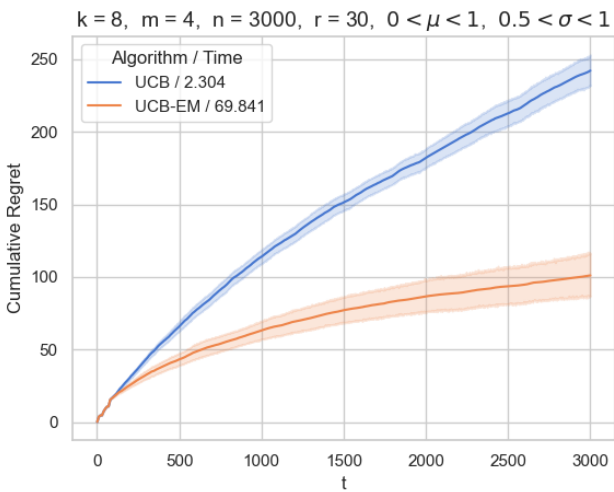


Figure 2.17

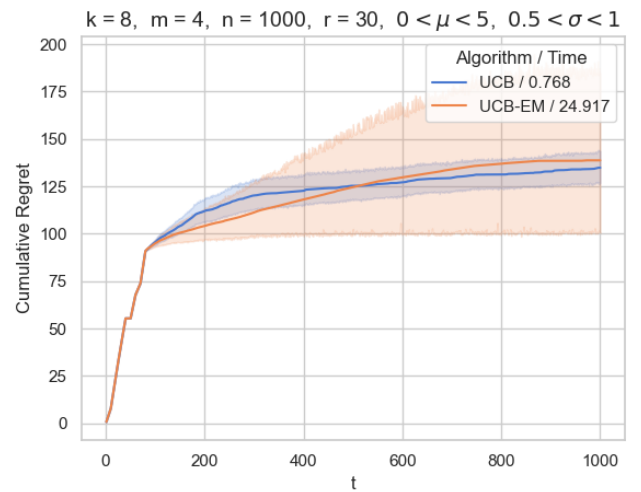


Figure 2.18

### 3. Unknown Variance Case

In this section, it's assumed that no prior information about the parameters of the arms is given to the player. Aside from the mixture model and Gaussian assumptions formulated previously, the means, variances, and weights of the component distributions are all unknown. For each arm, the reward  $X$  obtained follows a mixture model distribution

$$f_X(x) = \sum_{i=1}^m \pi_i g_i(x), \quad \text{s.t. } g_i(x) = \frac{1}{\sigma_i \sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \left( \frac{x - \mu_i}{\sigma_i} \right)^2 \right\} \text{ for } i = 1 \rightarrow m.$$

#### 3.1. UCB-V (Empirical Bernstein UCB)

UCB-V algorithm is introduced in the paper by Audibert, Munos, and Szepesvári (2007). It is also commonly referred to as the empirical Bernstein UCB because it has used the Bennett's and Bernstein's inequalities in deriving the upper confidence bound and the regret bound. It is particularly useful here because the algorithm assumes no information is given about the variance of the rewards. Empirical estimates for variances and means are used here. EM algorithm can also be adapted to provide a better estimate for sample variance.

##### 3.1.1. Assumptions for Using UCB-V

The assumptions of the algorithm are: Let  $K > 2, \nu_1, \dots, \nu_K$  distributions over reals with support  $[0, b]$ . For  $1 \leq k \leq K$ , let  $\{X_{k,t}\} \sim \nu_k$  be an i.i.d. sequence of random variables specifying the rewards for arm  $k$ . Assume that the rewards of different arms are independent of each other, i.e., for any  $k, k', 1 \leq k < k' \leq K, t \in \mathbb{N}^+$ , the collection of random variables,  $(X_{k,1}, \dots, X_{k,t})$  and  $(X_{k',1}, \dots, X_{k',t})$ , are independent of each other (Audibert, Munos, and Szepesvári, 2007).

##### 3.1.2. The UCB-V Algorithm

For any  $k \in \{1, \dots, K\}$  and  $t \in \mathbb{N}$ , let  $\bar{X}_{k,t}$  and  $V_{k,t}$  be the empirical estimates of the mean payoff and variance of arm  $k$  :

$$\bar{X}_{k,t} \triangleq \frac{1}{t} \sum_{i=1}^t X_{k,i} \quad \text{and} \quad V_{k,t} \triangleq \frac{1}{t} \sum_{i=1}^t (X_{k,i} - \bar{X}_{k,t})^2 \quad (3.1)$$



Let  $c \geq 0$ . Let  $\mathcal{E} = (\mathcal{E}_{s,t})_{s \geq 0, t \geq 0}$  be non-negative real numbers such that for any  $s \geq 0$ , the function  $t \mapsto \mathcal{E}_{s,t}$  is non-decreasing.  $\mathcal{E}$  (a function of  $(s, t)$ ) is called the exploration function. For any arm  $k$  and any non-negative integers  $s, t$ , introduce

$$B_{k,s,t} \triangleq \bar{X}_{k,s} + \sqrt{\frac{2V_{k,s}\mathcal{E}_{s,t}}{s}} + c\frac{3b\mathcal{E}_{s,t}}{s} \quad (3.2)$$

with the convention  $1/0 = +\infty$ .

**UCB-V algorithm:** At any time  $t$ , play an arm maximizing  $B_{k, T_k(t-1), t}$  (Audibert, Munos, and Szepesvári, 2007).

### 3.1.3. Natural Choices for the Hyperparameters in the UCB-V Algorithm\*

As we can see from the formula in Eq. 3.2, it involves unspecified exploration function  $\epsilon$  and  $c$ . In the paper by Audibert, Munos, and Szepesvári (2007), it's been pointed out that these are hyperparameters to be tuned under different environment. However, the choice can't be chosen without any constraint.

It's been established and proved by Audibert, Munos, and Szepesvári (2007) that in order to achieve sub-linear and logarithmic regret in particular, natural values for the constants appearing in the bound are:

$$B_{k,s,t} \triangleq \bar{X}_{k,s} + \sqrt{\frac{2V_{k,s} \log t}{s}} + \frac{b \log t}{2s} \quad (3.3)$$

This choice corresponds to the critical exploration function  $\mathcal{E}_t = \log t$  and to  $c = 1/6$ , which is also the minimal associated value of  $c$  to ensure that the UCB-V algorithm doesn't suffer a polynomial loss. The proofs that lead to this result are beyond the scope of this project.

### 3.1.4. The UCB-V Algorithm for Mixture Model Bandits\*

Here, UCB-V algorithm can be directly applied to the mixture model bandit problem with a minor breach of the assumptions. The empirical variances can be directly estimated by the sample variance as defined in Eq. 3.1.

We require a value of  $b$  in the Eq. 3.1, an upper bound for the range of the rewards of **all** arms, which doesn't allow us to use the UCB-V algorithm here strictly speaking. However, with the sub-gaussian property of the mixture gaussian random variables as discussed before, we can assume that the tail behaviors of the mixture of gaussians

allows us to believe most of the samples from this mixture distribution would stay in an approximated interval  $[\hat{b}_1, \hat{b}_2]$ . Then, we can use this estimated  $b_2$  to replace the  $b$  in Eq. 3.3. In practice, I choose to estimate the  $b$  by the maximum value of the rewards we get from playing any arms, i.e.,  $\hat{b} = \max_{k,s} \{X_{k,s} \mid \forall k, \forall s\}$ . As a result, we can expect  $\hat{b}$  to be reasonably close to what we need here eventually for a horizon that is sufficiently large.

Note that there is **no** theoretical guarantee or regret bound obtained for this revised algorithm, but empirically this works well as shown in the cumulative regret plots given in the following sections.

---

**Algorithm 3:** UCB-V for Mixture Model Bandits without EM

---

**Input :**  $x, \dots, x_N$ ,  $n$  (the horizon),  $k$  (the number of arms),  $m$  (the number of mixture components per arm)

**for**  $t = 1, \dots, n$  **do**

    Choose the action  $\mathcal{A}_t$  that maximizes the UCB bound as defined in Eq. 3.3.

    Denote  $i = \mathcal{A}_t$  as the action, and observe  $X_{i,T_i(t-1)}$ .

    Update  $\bar{X}_{i,T_i(t-1)}$  and  $\bar{V}_{i,T_i(t-1)}$  as defined in Eq. 3.1

    Update  $\hat{b} = \max_{k,s} \{X_{k,s} \mid \forall k, \forall s\}$

**end**

---

### 3.1.5. The UCB-V Algorithm for Mixture Model Bandits with EM\*

The next algorithm to be considered here is an ‘advancement’ for what we have done in the previous section. We want to consider the mixture model assumption here. Therefore, EM algorithm can be used to provide an estimate for the empirical variance term  $\bar{V}_{k,s}$  with assumption about the mixture model reward for each arm.

Suppose the means of each components are  $\mu_1, \dots, \mu_m$  and variances are  $\sigma_1^2, \dots, \sigma_m^2$ , then we know:

$$\begin{aligned}
\mathbb{E}[X] &= \mu = \sum_{i=1}^m \pi_i \mu_i \\
\mathbb{E}[(X - \hat{\mu})^2] &= \sigma^2 = \text{Var}[X] \\
&= \mathbb{E}[X^2] - \mu^2 \\
&= \sum_{i=1}^m \pi_i \cdot \mathbb{E}[X_i^2] - \mu^2 \\
&= \sum_{i=1}^m \pi_i (\sigma_i^2 + \mu_i^2) - \mu^2
\end{aligned} \tag{3.4}$$

The empirical variance can therefore be estimated as:

$$\bar{V}_{k,t} \triangleq \sum_{i=1}^m \hat{\pi}_i (\hat{\sigma}_i^2 + \hat{\mu}_i^2) - \hat{\mu}^2, \tag{3.5}$$

where  $\hat{\pi}_i$  are weights,  $\hat{\sigma}_i^2$  and  $\hat{\mu}_i$  are component variances and means, which can be estimated from the EM algorithm with samples  $X_{k,1}, \dots, X_{k,t}$ .

By a similar way of implementing EM within the UCB-V algorithm, we call this algorithm UCB-V-EM algorithm, which is summarized as below:

---

**Algorithm 4:** UCB-V for Mixture Model Bandits with EM

---

**Input :**  $x, \dots, x_N$ ,  $n$  (the horizon),  $k$  (the number of arms),  $m$  (the number of mixture components per arm)

**for**  $t = 1, \dots, k$  **do** // Initialization

Observe from arm  $t$  for  $m$  times and observe  $X_{t,1}, \dots, X_{t,m}$ .

Initialize  $\bar{X}_{t,m}$  as  $\frac{1}{m} \sum_{i=1}^m X_{t,i}$  and  $\bar{V}_{t,m}$  as 1.

Initialize  $\pi_{t,j}$  randomly for  $j = 1 \rightarrow m$ .

**end**

**for**  $t = km, \dots, n$  **do** // Standard UCB-V-EM Steps

Choose the action  $\mathcal{A}_t$  that maximizes the UCB bound as defined in Eq. 3.3.

Denote  $\mathcal{A}_t = i$  as the action, and observe  $X_{i,T_i(t-1)}$ .

Update  $\bar{X}_{i,T_i(t-1)}$  as defined in Eq. 3.1.

Update  $\bar{V}_{i,T_i(t-1)}$  as defined in Eq. 3.5 through EM algorithm.

Update  $\hat{b} = \max_{k,s} \{X_{k,s} \mid \forall k, \forall s\}$ .

**end**

---

### 3.2. The UCB1-Normal Algorithm for Mixture Model Bandits\*

Another policy considered here is the UCB1-NORMAL algorithm introduced by Auer, Cesa-Bianchi, and Fischer (2002). It's designed for the multi-armed bandit problem un-

der the assumptions that each machine has a reward that follows a normal distribution with unknown means and variances.

The algorithm is especially concerned with the finite-armed bandit problem with Gaussian rewards. It's been shown by Auer, Cesa-Bianchi, and Fischer (2002) that UCB1-NORMAL has logarithmic regret. The proof has used bounds on the tails of the  $\chi^2$  and the Student distribution that has only been verified numerically.

As our assumption is that the reward follows a mixture of Gaussian, it's therefore compelling to use UCB1-NORMAL algorithm here because we can reasonably expect that the tail behaviours of mixture of Gaussian are similar to those of a single normal distribution. Therefore, the UCB1-NORMAL algorithm is directly used here to assess its performance in the mixture model bandit problem.

Again, the interest here is whether considering the mixture model assumption of the reward would lead to an improvement of the performance of the different algorithms, and also about the extra computational cost induced by the EM algorithm. Therefore, the sample variance term in the Eq. 3.7 in the UCB1-NORMAL algorithm,  $\frac{q_j - n_j \bar{x}_j^2}{n_j - 1}$ , can be replaced by the one estimated using EM algorithm as we have derived before in Eq. 3.5:

$$\bar{V}_{k,t} \triangleq \sum_{i=1}^m \hat{\pi}_i \left( \hat{\sigma}_i^2 + \hat{\mu}_i^2 \right) - \hat{\mu}^2, \quad (3.6)$$

---

**Algorithm 5:** UCB1-NORMAL (Auer, Cesa-Bianchi, and Fischer, 2002)

---

**for**  $n = 1, 2, \dots$  **do**

If there is an arm that has been played less than  $\lceil 8 \log n \rceil$  times, play this arm.

Otherwise, play arm  $j$  that maximizes

$$\bar{x}_j + \sqrt{16 \cdot \frac{q_j - n_j \bar{x}_j^2}{n_j - 1} \cdot \frac{\ln(n-1)}{n_j}} \quad (3.7)$$

where  $\bar{x}_j$  is the average reward obtained from machine  $j$ ,  $q_j$  is the sum of squared rewards obtained from machine  $j$ , and  $n_j$  is the number of times machine  $j$  has been played so far.

Update  $\bar{x}_j$  and  $q_j$  with the obtained reward  $x_j$

**end**

---

---

**Algorithm 6:** Altered UCB1-NORMAL with EM

---

**for**  $n = 1, 2, \dots$  **do**

    If there is an arm that has been played less than  $\lceil 8 \log n \rceil$  times, play this arm.

    Otherwise, play arm  $j$  that maximizes

$$\bar{x}_j + \sqrt{16 \cdot \bar{V}_{j,n_j} \cdot \frac{\ln(n-1)}{n_j}}$$

    where  $\bar{x}_j$  is the average reward obtained from machine  $j$ ,  $\hat{\sigma}_j^2$  is the estimated sample variance from machine  $j$ ,  $n_j$  is the number of times machine  $j$  has been played so far, and  $\bar{V}_{j,n_j}$  is the empirical variance estimated from using the  $n_j$  samples from arm  $j$  as written in Eq. 3.5.

    Update  $\bar{x}_j$  and  $\hat{\sigma}_j^2$  with the obtained reward  $x_j$  using EM algorithm.

**end**

---

### 3.3. Empirical Results and Regret Comparison for Unknown Variance

This section presents the empirical results by comparing the cumulative regrets of these policies that assume unknown variance: UCB-V, UCB-V-EM, UCB1-NORMAL, and UCB1-NORMAL-EM.

The regret used here is same as previous sections, which is the pseudo-regret as discussed in Eq. 1.1:  $\bar{R}_n = n\mu^* - \sum_{t=1}^n \mu_{A_t}$ . The notations in the titles of the figures have the same meanings as those in the previous figures (see page 19 for details) .

From observations of the plots, the main conclusion here is to **never** use UCB1-NORMAL or UCB1-NORMAL-EM under mixture model assumptions. In all the plots, UCB1-NORMAL, no matter whether implemented with EM or not, can't achieve the regret as low as the UCB-V algorithm. This is probably caused by the first step of the UCB1-NORMAL algorithm, which selects an arm to play whenever it's been played less than  $\lceil 8 \log n \rceil$  times. Even though the algorithm achieves logarithmic regrets under Gaussian setting, this step might result in large regret because it chooses an arm not based on its upper confidence bound. The ruggedness of the cumulative regret of the UCB1-NORMAL algorithm is also caused by this.

However, the EM algorithm did somehow play a role in improving the regret of the UCB1-NORMAL algorithm, even though the UCB1-NORMAL algorithm is not fit to be used in the mixture model setting. This can be observed most clearly in the figure 3.4, 3.9, 3.11, 3.13.

Unfortunately, EM didn't do a great job in reducing the regret of the UCB-V algorithm. The mean and the confidence interval of the cumulative regrets don't seem to be affected by the introduction of the EM. The expectation is that with the mixture model assumption taking into account, a better estimate for the variance can be attained. However, this result is also expected because we know that the sample variance is already itself an unbiased estimator for the overall variance of an arm. With the mixture model assumptions, only one piece of new information is introduced which is the number of components  $m$ . This is arguably quite limited if we assume that the variances are unknown.

The time taken by the UCB-V-EM algorithm to run is on average 10 times more than that of the single UCB-V algorithm. So, **the plain UCB-V is the optimal one to use** so far among the four aforementioned algorithms for mixture model bandits because it has the best regret and takes the smallest computational cost.

UCB-V, UCB-V-EM, UCB1-NORMAL, UCB1-NORMAL-EM, 2 arms

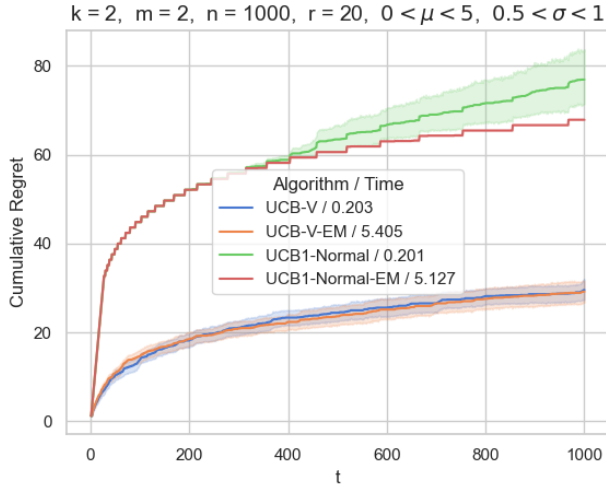


Figure 3.1

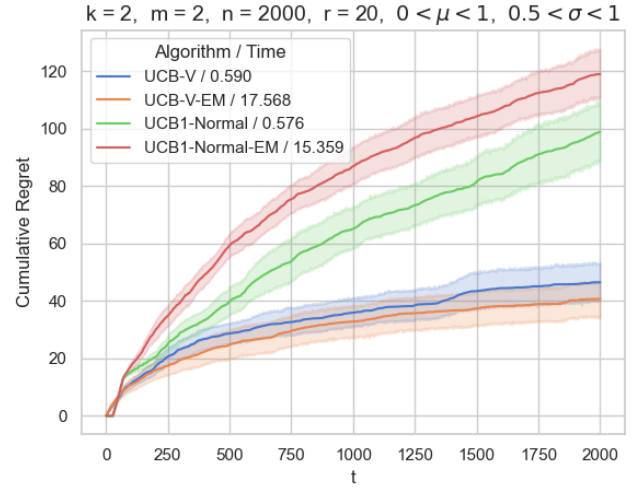


Figure 3.2

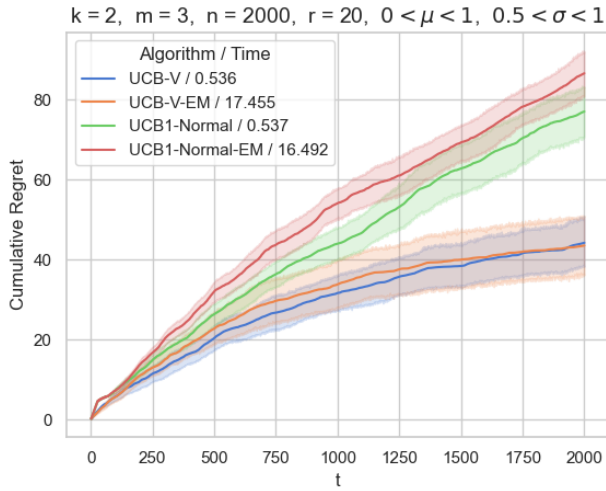


Figure 3.3

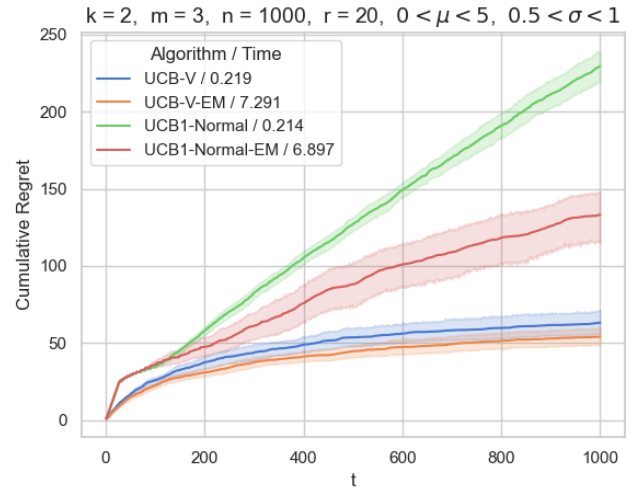


Figure 3.4

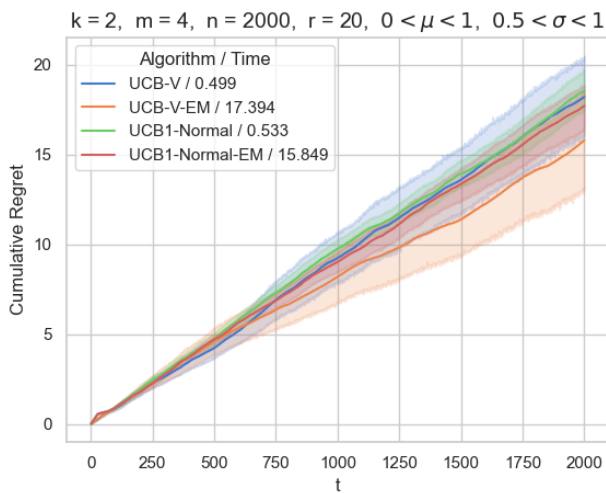


Figure 3.5

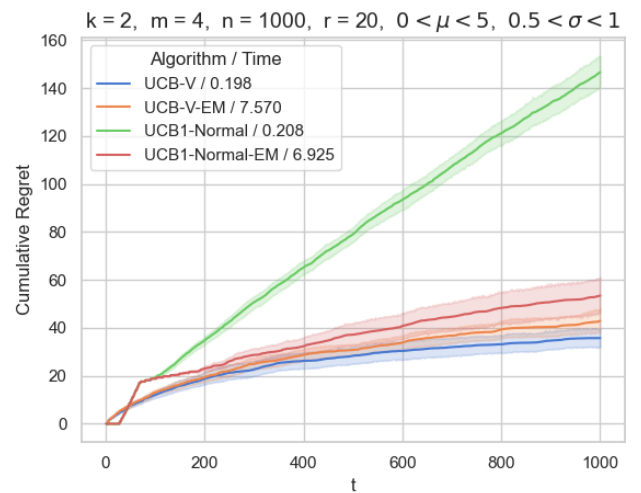


Figure 3.6

UCB-V, UCB-V-EM, UCB1-NORMAL, UCB1-NORMAL-EM, 4 arms

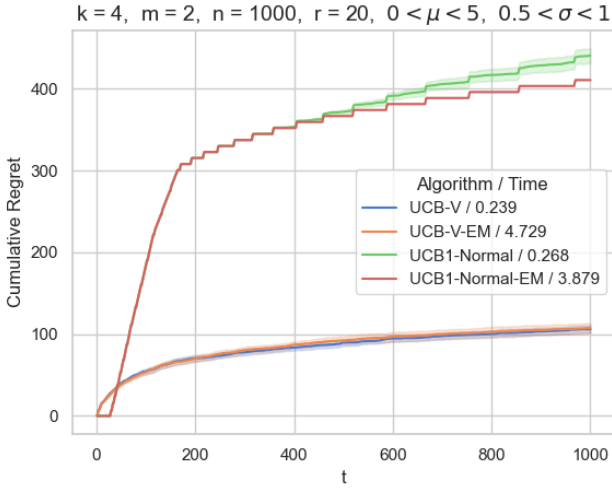


Figure 3.7

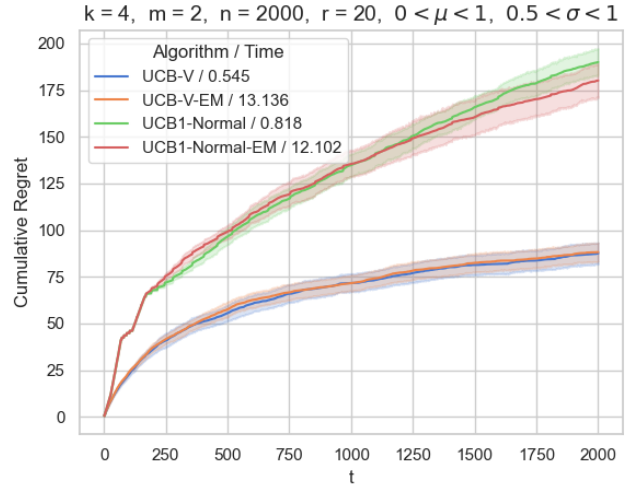


Figure 3.8

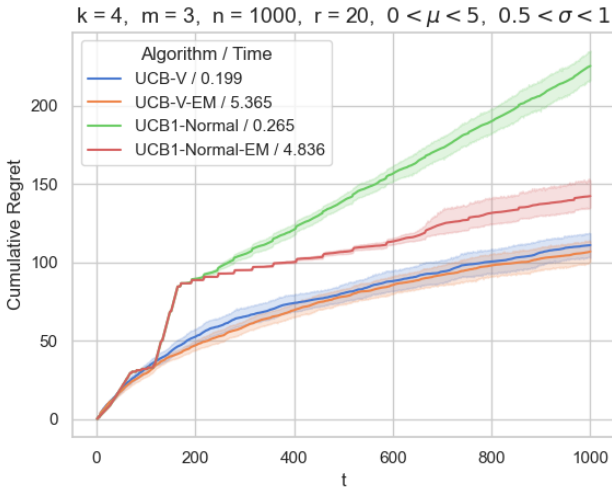


Figure 3.9

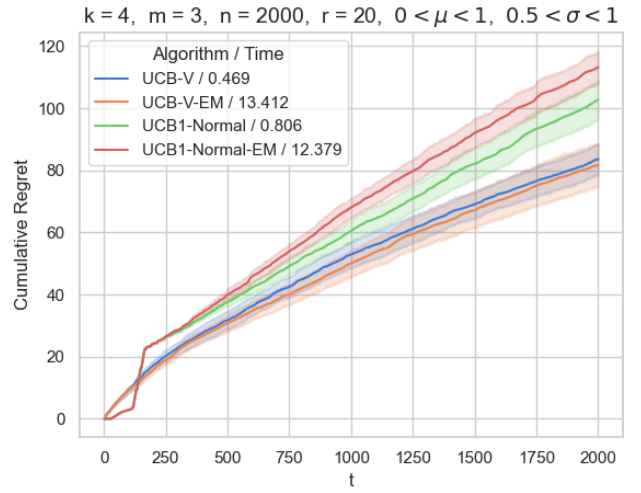


Figure 3.10

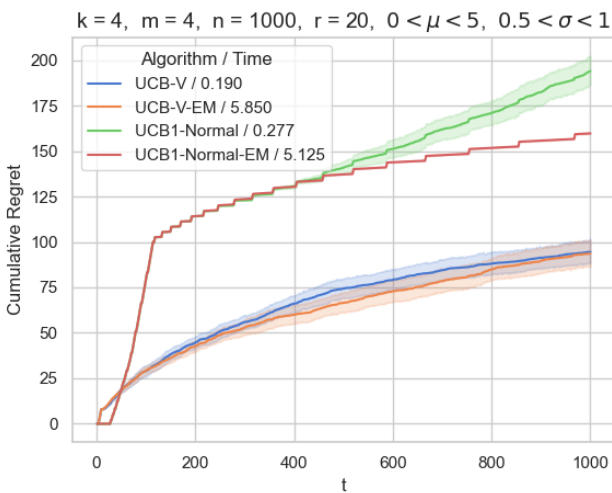


Figure 3.11

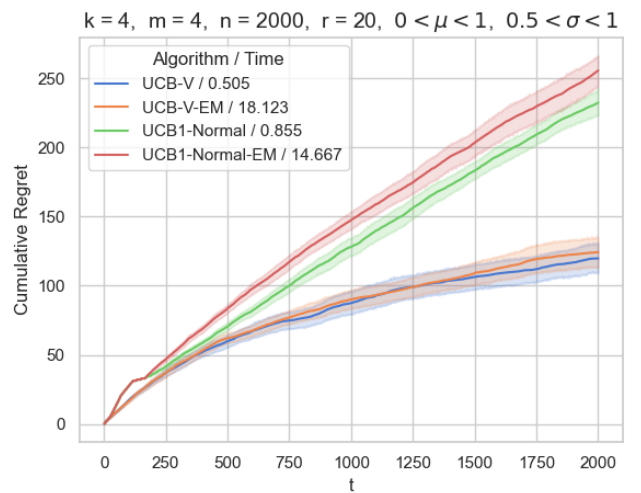


Figure 3.12



UCB-V, UCB-V-EM, UCB1-NORMAL, UCB1-NORMAL-EM, 8 arms

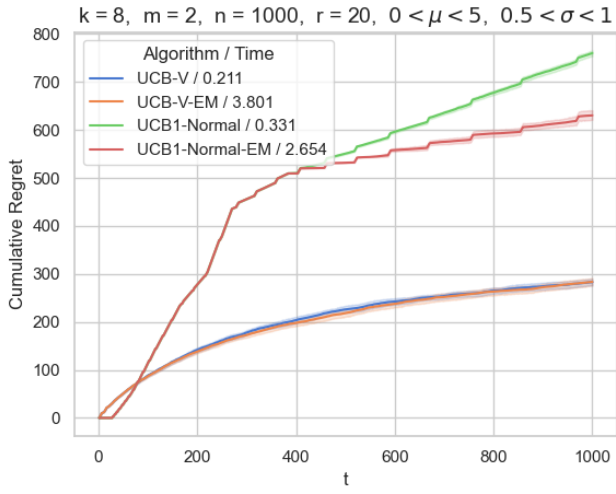


Figure 3.13

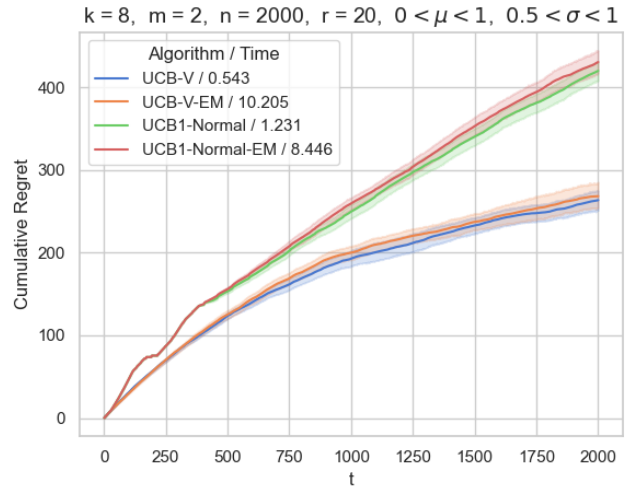


Figure 3.14

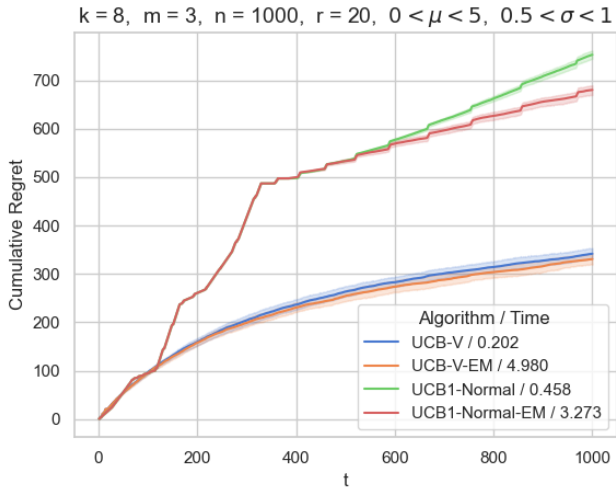


Figure 3.15

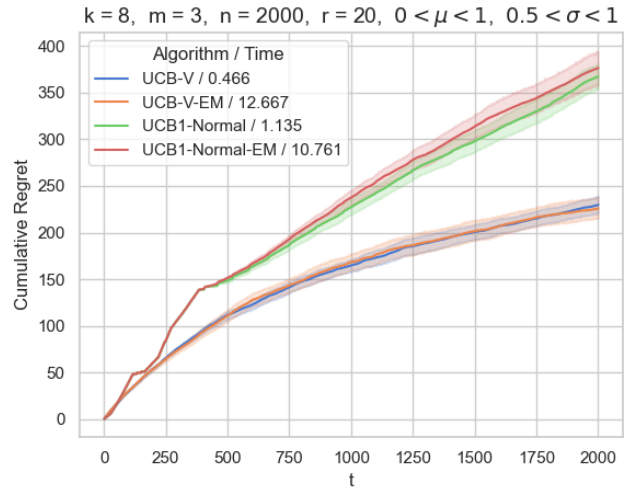


Figure 3.16

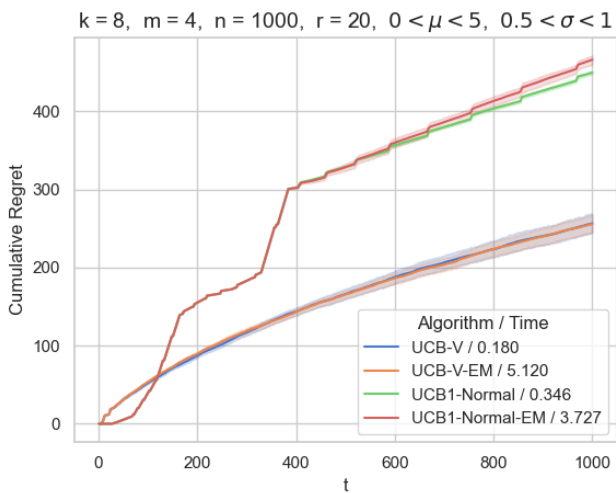


Figure 3.17

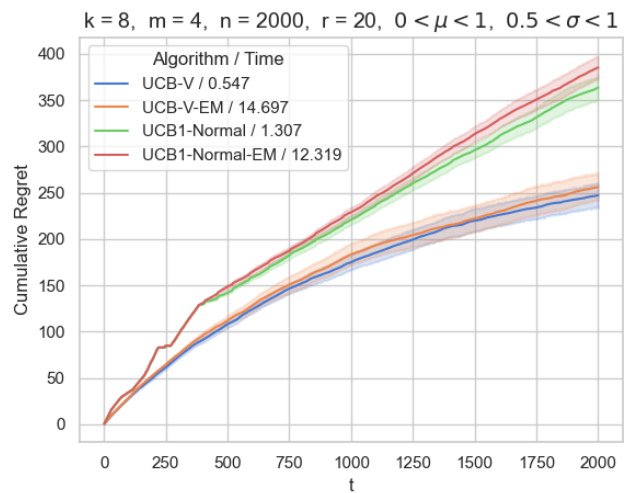


Figure 3.18

## 4. Thompson Sampling

In this section, it's assumed that the bandit environment is in a Bayesian setting. More specifically, the Bayesian bandit environment requires that the parameters of each arm are sampled from prior distributions. A more detailed definition is given in the following section.

The goal is to try to explore the performance of Thompson sampling used in the mixture model bandits and compare it with the UCB-V algorithm.

Note that it doesn't make sense here to 'ignore' the mixture model assumptions here because in Bayesian bandits a prior is placed on the original bandit environment which contains information about the component means, variances, and weights. It would be impractical to ignore this and place a prior on overall means and variances, etc., because that would be a violation against the assumption about the mixture model reward distribution.

In this section, it's also assumed that the component variances, means, weights are unknown and that the reward follows a mixture of independent Gaussian distributions.

### 4.1. The Bayesian Bandit Environment

The Bayesian bandit model is different in its own way that at the start of the bandit game, an environment is randomly sampled from the prior. Similarly, the sampled environment is not accessible to the learner. In this project, it's assumed we have only a finite number of arms.

**Definition 7** (Bayesian bandit environment, Lattimore and Szepesvári (2020)). *A  $k$ -armed Bayesian bandit environment is a tuple  $(\mathcal{E}, \mathcal{G}, Q, P)$ , where  $(\mathcal{E}, \mathcal{G})$  is a measurable space and  $Q$  is a probability measure on  $(\mathcal{E}, \mathcal{G})$  called the prior. The last element  $P = (P_{\nu i} : \nu \in \mathcal{E}, i \in [k])$  is a probability kernel from  $\mathcal{E} \times [k]$  to  $(\mathbb{R}, \mathfrak{B}(\mathbb{R}))$ , where  $P_{\nu i}$  is the reward distribution associated with the  $i$ th arm in bandit  $\nu$ . A Bayesian bandit environment and policy  $\pi = (\pi_t)_{t=1}^n$  interact to produce a collection of random variables,  $\nu \in \mathcal{E}$ ,  $(A_t)_{t=1}^n$  and  $(X_t)_{t=1}^n$  with  $A_t \in [k]$  and  $X_t \in \mathbb{R}$  that satisfy:*

$$(a) \quad \mathbb{P}(\nu \in \cdot) = Q(\cdot)$$

- (b) the conditional distribution of action  $A_t$  given  $\nu, A_1, X_1, \dots, A_{t-1}, X_{t-1}$  is  $\pi_t(\cdot \mid A_1, X_1, \dots, A_{t-1}, X_{t-1})$  almost surely; and
- (c) the conditional distribution of the reward  $X_t$  given  $\nu, A_1, X_1, \dots, A_t$  is  $P_{\nu A_t}$  almost surely.

*Example 8.* (Lattimore and Szepesvári, 2020) A  $k$ -armed Bayesian Bernoulli bandit environment could be defined by letting  $\mathcal{E} = [0, 1]^k$ ,  $\mathcal{G} = \mathfrak{B}(\mathcal{E})$  and  $P_{\nu i} = \mathcal{B}(\nu_i)$ . A natural prior in this case would be a product of  $\text{Beta}(\alpha, \beta)$  distributions (Lattimore and Szepesvári, 2020):

$$Q(A) = \int_A \prod_{i=1}^k q_i(x_i) dx$$

where  $q_i(x) = x^{\alpha-1}(1-x)^{\beta-1}\Gamma(\alpha+\beta)/(\Gamma(\alpha)\Gamma(\beta))$ .

## 4.2. Posterior Distributions in Bayesian Bandits

Let  $(\mathcal{E}, \mathcal{G}, Q, P)$  be a  $k$ -armed Bayesian bandit environment. Assuming that  $(\mathcal{E}, \mathcal{G})$  is a Borel space, The following theorem guarantees the existence of the posterior: a probability kernel  $Q(\cdot \mid \cdot)$  from the space of histories to  $(\mathcal{E}, \mathcal{G})$  with the form

$$Q(A \mid a_1, x_1, \dots, a_t, x_t) \tag{4.1}$$

which can also be written as  $\mathbb{E}[\mathbb{I}_A(\nu) \mid A_1, X_1, \dots, A_t, X_t]$  (Lattimore and Szepesvári, 2020).

**Theorem 9.** (Lattimore and Szepesvári, 2020) *If  $(\Theta, \mathcal{G})$  is a Borel space, then there exists a probability kernel  $Q : \mathcal{X} \times \mathcal{G} \rightarrow [0, 1]$  such that  $Q(A \mid X) = \mathbb{P}(\theta \in A \mid X)$  simultaneously for all  $A \in \mathcal{G}$  outside some  $P$ -null set. Furthermore, for any two probability kernels  $Q, Q'$  satisfying this condition,  $Q(\cdot \mid x) = Q'(\cdot \mid x)$  for all  $x$  in some set of  $\mathbb{P}_X$ -probability one.*

*Proof.* The proof of the above result is beyond the scope of this project.  $\square$

*Example 10.* (Lattimore and Szepesvári, 2020) The posterior for the Bayesian bandit in Example 8 in terms of its density with respect to the Lebesgue measure is

$$q(\theta \mid \underbrace{a_1, x_1, \dots, a_t, x_t}_{h_t}) \propto \prod_{i=1}^k \theta_i^{\alpha+s_i(h_t)-1} (1-\theta_i)^{\beta+t_i(h_t)-s_i(h_t)-1}$$

where  $s_i(h_t) = \sum_{u=1}^t x_u \mathbb{I}\{a_u = i\}$  and  $t_i(h_t) = \sum_{u=1}^t \mathbb{I}\{a_u = i\}$ . This means the posterior is also the product of Beta distributions, each updated according to the observations from the relevant arm.

### 4.3. The Bayesian Regret

We have defined before that the regret of policy  $\pi$  in  $k$ -armed bandit environment  $\nu$  over  $n$  rounds is

$$R_n(\pi, \nu) = n\mu^* - \mathbb{E} \left[ \sum_{t=1}^n X_t \right] \quad (4.2)$$

where  $\mu^* = \max_{i \in [k]} \mu_i$  and  $\mu_i$  is the mean of  $P_{\nu i}$ .

**Definition 11** (Bayesian Regret, Lattimore and Szepesvári (2020)). *Given a  $k$ -armed Bayesian bandit environment  $(\mathcal{E}, \mathcal{G}, Q, P)$  and a policy  $\pi$ , the Bayesian regret is*

$$\text{BR}_n(\pi, Q) = \int_{\mathcal{E}} R_n(\pi, \nu) dQ(\nu) \quad (4.3)$$

The Bayesian optimal regret is  $\text{BR}_n^*(Q) = \inf_{\pi} \text{BR}_n(\pi, Q)$ , and the optimal policy is

$$\pi^* = \text{argmin}_{\pi} \text{BR}_n(\pi, Q) \quad (4.4)$$

### 4.4. Introduction to Thompson Sampling

Thompson sampling, named after William R. Thompson, is an algorithm for choosing actions that addresses the exploration-exploitation dilemma in the multi-armed bandit problem. Intuitively, it states that the player should select the action that optimizes the expected reward with a randomly drawn belief based on previous actions and rewards history. The heuristic was firstly introduced by Thompson (1933).

Formally, let  $(\mathcal{E}, \mathfrak{B}(\mathcal{E}), Q, P)$  be a  $k$ -armed Bayesian bandit environment as defined in the previous sections. The player then chooses actions  $(A_t)_{t=1}^n$  and obtains rewards  $(X_t)_{t=1}^n$ , and the posterior given  $t$  observations is a probability kernel  $Q(\cdot \mid \cdot)$  from  $([k] \times \mathbb{R})^t$  to  $(\mathcal{E}, \mathfrak{B}(\mathcal{E}))$ . Denote the mean of the  $i$  th arm in bandit  $\nu \in \mathcal{E}$  by  $\mu_i(\nu) = \int_{\mathbb{R}} x dP_{\nu i}(x)$ . In round  $t$ , Thompson sampling algorithm samples a bandit environment  $\nu_t$  from the posterior of  $Q$  given  $A_1, X_1, \dots, A_{t-1}, X_{t-1}$  and then chooses the arm with the largest mean.

---

**Algorithm 7:** Thompson Sampling (Lattimore and Szepesvári, 2020)

---

**Input :** Bayesian bandit environment  $(\mathcal{E}, \mathfrak{B}(\mathcal{E}), Q, P)$

**for**  $n = 1, 2, \dots, n$  **do**

Sample  $\nu_t \sim Q(\cdot \mid A_1, X_1, \dots, A_{t-1}, X_{t-1})$   
 Choose  $A_t = \text{argmax}_{i \in [k]} \mu_i(\nu_t)$

**end**

---

The validity and soundness of the Thompson sampling algorithm can also be reasoned as the existence of a bound for the Bayesian regret, similar to the UCB-based algorithms.

**Theorem 12** (Bayesian Regret of Thompson Sampling, Lattimore and Szepesvári (2020)). *Let  $(\mathcal{E}, \mathfrak{B}(\mathcal{E}), Q, P)$  be a  $k$ -armed Bayesian bandit environment such that for all  $\nu \in \mathcal{E}$*

and  $i \in [k]$ , the distribution  $P_{\nu_i}$  is 1-subgaussian (after centring) with mean in  $[0, 1]$ . Then the policy  $\pi$  of Thompson sampling satisfies

$$\text{BR}_n(\pi, Q) \leq C \sqrt{kn \log(n)}$$

where  $C > 0$  is a universal constant.

*Proof.* The proof of the above result is beyond the scope of this project.  $\square$

## 4.5. Thompson Sampling for Mixture Model Bandits\*

From the algorithm of Thompson sampling, it's obvious that the main task and challenge here is the calculation of posterior distribution  $Q(\cdot \mid A_1, X_1, \dots, A_{t-1}, X_{t-1})$ . In the non-Bayesian setting discussed in the previous sections, the EM (Expectation and Maximization) algorithm is used to work out the parameters of a mixture distribution. Here, under a Bayesian setting, the problem is more complex because the posterior distribution is in most cases intractable with hidden states, in this case the labels of samples indicating which components it comes from.

The method to work out the posterior of a mixture of Gaussian is the *Variational Bayes* (VB) method, also known as variational inference methods. One python function, namely, `sklearn.mixture.Bayesian.GaussianMixture` has implemented the variational Bayesian method for graphical models introduced and discussed in the paper by Attias et al. (1999), and the VB inference for Dirichlet process mixtures discussed in the paper by Blei and Jordan (2006) for working out the posterior of a mixture of Gaussian, which is exactly what's needed here.

Even though the details about the derivation of this method and the proof of the validity and efficiency of the VB methods are beyond the focus and the scope of this project, a general idea about how it works is summarized in the following section.

### 4.5.1. Introduction to Variational Bayesian Methods

This section is an introduction to the basis of the mathematical derivation of variational Bayesian methods which will be used to estimate the posterior of the mixture of Gaussian in the Thompson sampling algorithm. One of the earliest literature that formulates the variational Bayes problems and introduces the VB algorithm and its generalizations can be found in the book by Beal (2003).

In variational Bayes methods, we seek to approximate the posterior distribution over a set of unobserved variables  $\mathbf{Z} = \{Z_1 \dots Z_n\}$  given data  $\mathbf{X}$  by a variational distribution,  $Q(\mathbf{Z})$ :

$$P(\mathbf{Z} \mid \mathbf{X}) \approx Q(\mathbf{Z})$$

Therefore, the task is transformed to minimize the dissimilarity between these two distributions. It's common to use the Kullback-Leibler divergence (KL-divergence) of  $P$  from  $Q$  as the measure for dissimilarity.

$$D_{\text{KL}}(Q\|P) \triangleq \sum_{\mathbf{Z}} Q(\mathbf{Z}) \log \frac{Q(\mathbf{Z})}{P(\mathbf{Z} | \mathbf{X})}.$$

By Bayes theorem, we know that

$$P(\mathbf{Z} | \mathbf{X}) = \frac{P(\mathbf{X} | \mathbf{Z})P(\mathbf{Z})}{P(\mathbf{X})} = \frac{P(\mathbf{X} | \mathbf{Z})P(\mathbf{Z})}{\int_{\mathbf{Z}} P(\mathbf{X}, \mathbf{Z}') d\mathbf{Z}'}.$$

The difficulty in this problem stems from the fact the integral over  $\mathbf{Z}$  to calculate the marginal distribution  $P(\mathbf{X})$  in the denominator is often intractable (both computationally and analytically). Therefore, we seek an approximation, using  $Q(\mathbf{Z}) \approx P(\mathbf{Z} | \mathbf{X})$ .

Given that  $P(\mathbf{Z} | \mathbf{X}) = \frac{P(\mathbf{X}, \mathbf{Z})}{P(\mathbf{X})}$ , the KL-divergence above can also be written as

$$\begin{aligned} D_{\text{KL}}(Q\|P) &= \sum_{\mathbf{Z}} Q(\mathbf{Z}) \left[ \log \frac{Q(\mathbf{Z})}{P(\mathbf{Z}, \mathbf{X})} + \log P(\mathbf{X}) \right] \\ &= \sum_{\mathbf{Z}} Q(\mathbf{Z}) [\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})] + \sum_{\mathbf{Z}} Q(\mathbf{Z}) [\log P(\mathbf{X})] \\ &= \sum_{\mathbf{Z}} Q(\mathbf{Z}) [\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})] + \log P(\mathbf{X}) \\ &= \mathbb{E}_{\mathbf{Q}} [\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})] + \log P(\mathbf{X}), \end{aligned} \tag{4.5}$$

which can be rearranged to become

$$\log P(\mathbf{X}) = D_{\text{KL}}(Q\|P) - \mathbb{E}_{\mathbf{Q}} [\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})] = D_{\text{KL}}(Q\|P) + \mathcal{L}(Q).$$

As  $\log P(\mathbf{X})$  is constant with respect to  $Q$ , to minimize the  $D_{\text{KL}}(Q\|P)$ , it's equivalent to maximize  $\mathcal{L}(Q)$ . By appropriate choice of  $Q$ ,  $\mathcal{L}(Q)$  would be tractable enough to maximize easily. The variational distribution  $Q(\mathbf{Z})$  is usually assumed to factorize over some partition of the latent variables, i.e. for some partition of the latent variables  $\mathbf{Z}$  into  $\mathbf{Z}_1 \dots \mathbf{Z}_M$  (Attias et al., 1999; Beal, 2003):

$$Q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i | \mathbf{X}).$$

#### 4.5.2. Python Implementation of the Thompson Sampling\*

Information is given in this section about the python function `sklearn.mixture.BayesianGaussianMixture` which is implemented in the Thompson sampling algorithm.

The basic assumption here is that for each arm we have a mixture of  $m$  independent gaussians as the reward. Then, the prior for means of component distributions

are themselves Gaussian  $\mu_i \sim \mathcal{N}(\hat{\mu}_i, \bar{\sigma}^2)$ . Here  $\hat{\mu}$  can be worked out from k-means algorithm as the mean for the prior (this is done automatically in the program). Here  $\bar{\sigma}$  is the prior knowledge we have about the variance of the mean distribution, which is denoted `mean_precision_prior` as the inverse of the variance and as a parameter in the `BayesianGaussianMixture` function.

The prior for covariance is the Wishart distribution as the default. As we are only dealing with one dimensional reward, Wishart distribution is equivalent to a Gamma distribution as the prior. A specific prior with customized parameters can be set with the option `covariance_prior`. Here, it's omitted so that it can be initialized automatically as the empirical covariance of the given samples.

`weight_concentration_prior_type` specifies the type of the prior for the mixture probabilities. Mixture models can be either finite or infinite. If it's infinite, a Dirichlet process is set as the prior. In this case, as only  $m$  components are considered. Therefore, what's needed is only a finite version of the Dirichlet process, the Dirichlet distribution. Dirichlet distribution is the conjugate prior for the categorical distribution characteristic of the weights. Here, the prior of  $\pi_i$  is set as the default  $\text{Dir}(1/m, \dots, 1/m)$  because it's assumed there are  $m$  components and the prior doesn't have a preference for one component over any of the other components.

After generating an instance from the `BayesianGaussianMixture` and fitting it with data, 3 attributes are used here to compute the posterior means, which are `means_`, `mean_precision_`, and `dirichlet_concentration`.

As a result,  $\pi_i$  is sampled from the posterior Dirichlet distribution with concentration parameters `dirichlet_concentration`.  $\mu_i$  is sampled from the posterior of the means, which is a gaussian distribution with parameters `means_` and `mean_precision_`. Subsequently, the posterior mean of a given bandit reward with  $m$  components is equal to  $\sum_{i=1}^m \pi_i \mu_i$ .

Note that the posterior covariances of the components are not needed because the calculation of the mean of the Bayesian bandits doesn't involve the covariances.

*(More details about the usage of `sklearn.mixture.BayesianGaussianMixture` can be found at [this link](#).)*

## 4.6. Empirical Results for Thompson Sampling\*

This section contains the empirical results comparing Thompson sampling with UCB-V-EM. Only algorithms that assume the unknown variances and use EM are compared with Thompson sampling because it's desirable to keep the assumptions of algorithms as consistent as possible (even though not exactly the same). UCB-V-EM use EM, which is

a non-Bayesian version of the variational Bayes method used in the Thompson Sampling.

UCB1-NORMAL algorithms are not taken into consideration because they perform uniformly worse than the UCB-V under the mixture model assumptions.

The regret used here is same as previous sections, which is the pseudo-regret as discussed in Eq. 1.1:  $\bar{R}_n = n\mu^* - \sum_{t=1}^n \mu_{A_t}$ . The notations in the titles of the figures have the same meanings as those in the previous figures (see page 19 for more details) .

The experiments have done for a large variety of bandit environments. Experiments are done for 2, 4, and 6 arms cases. The number of components ranges from 2 to 4. For a more challenging bandit problem, the means of each component are sampled from a uniform distribution with support  $[0, 1]$ , compare to  $[0, 5]$  for an easier counterpart. The possible standard deviations are all sampled uniformly from  $[0.5, 1]$  for the all the plots below.

The main result is that Thompson sampling overall has **much better** performance than the UCB-V-EM. The improvement in the performance is most evident in less challenging bandit environments (the left 3 plots on each page). In more difficult bandit problems, the reduction in cumulative regret is less significant. In addition, it seems like Thompson sampling has slightly larger confidence intervals in the cumulative regret plots for some cases (more noticeably in 2-armed cases).

Thompson sampling in general increases the running time more than by 20-fold, which means it requires significantly **more computational efforts** than UCB-V.



## 2 arms case, regret comparison for Thompson Sampling and UCB-V-EM

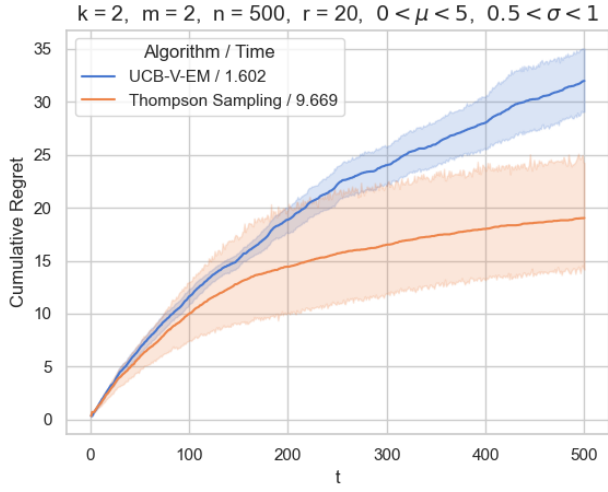


Figure 4.1

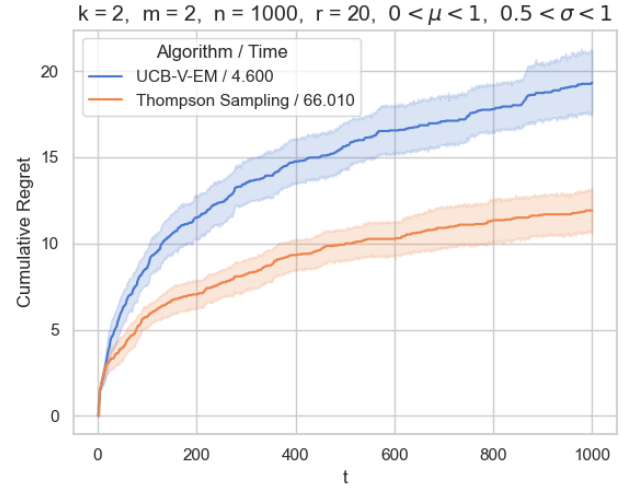


Figure 4.2

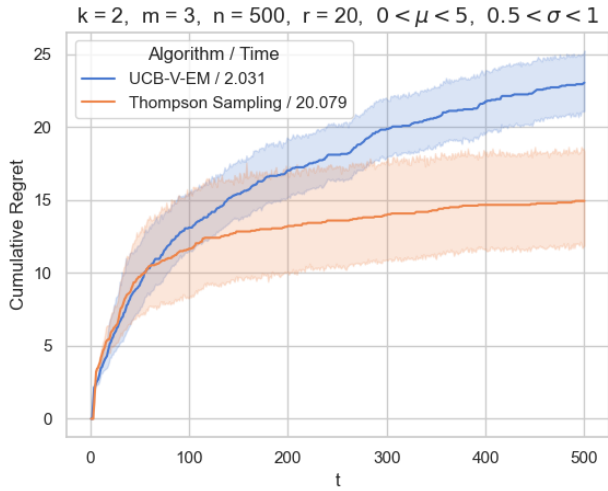


Figure 4.3

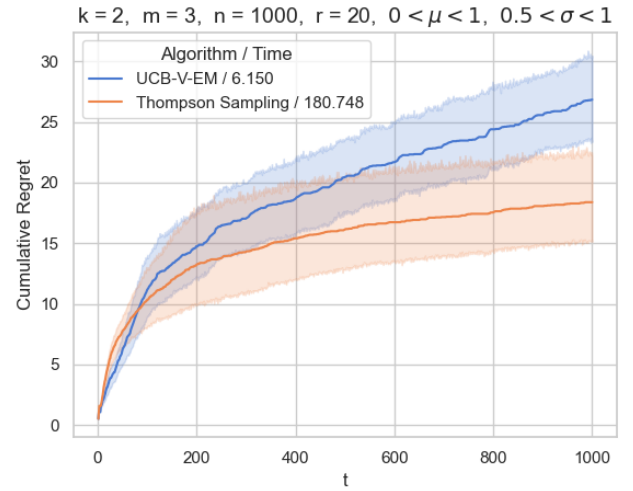


Figure 4.4

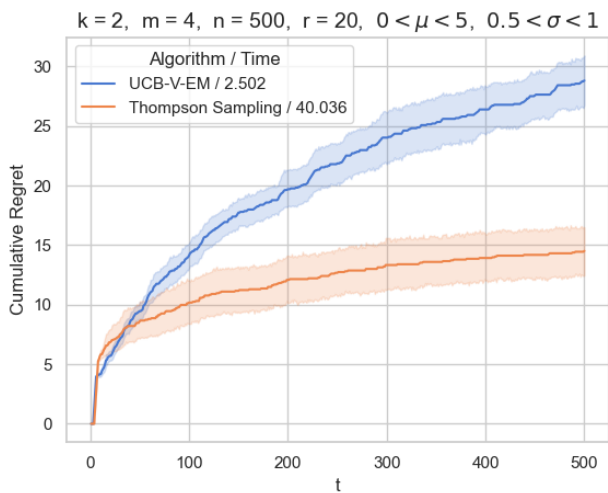


Figure 4.5

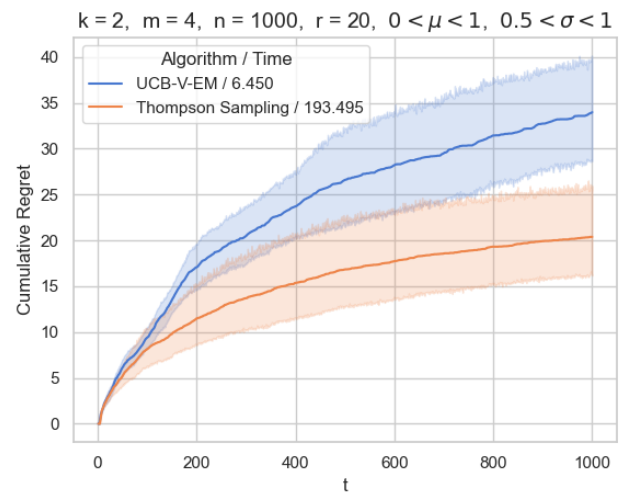


Figure 4.6

#### 4 arms case, regret comparison for Thompson Sampling and UCB-V-EM

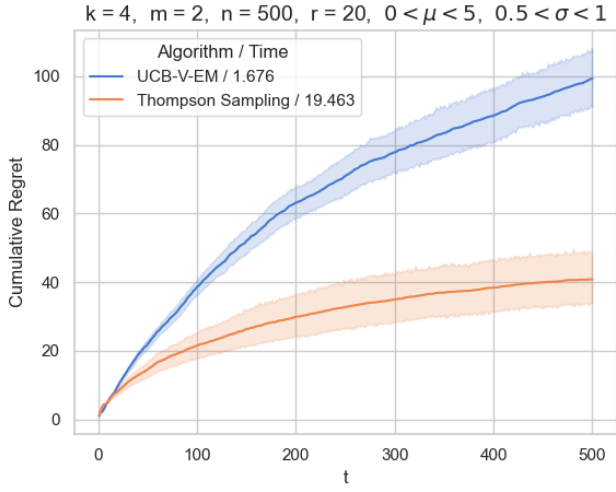


Figure 4.7

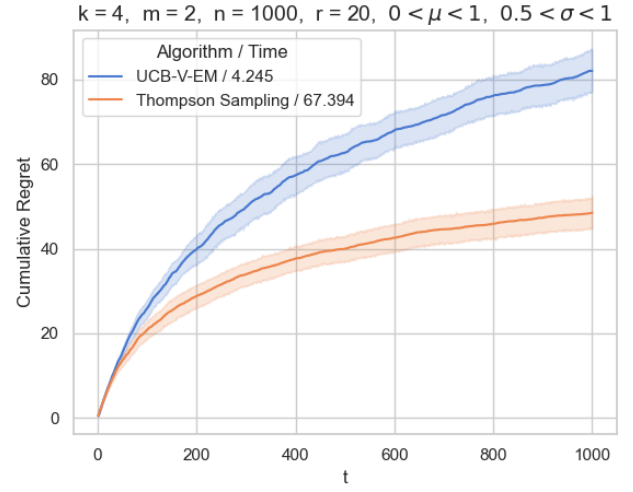


Figure 4.8

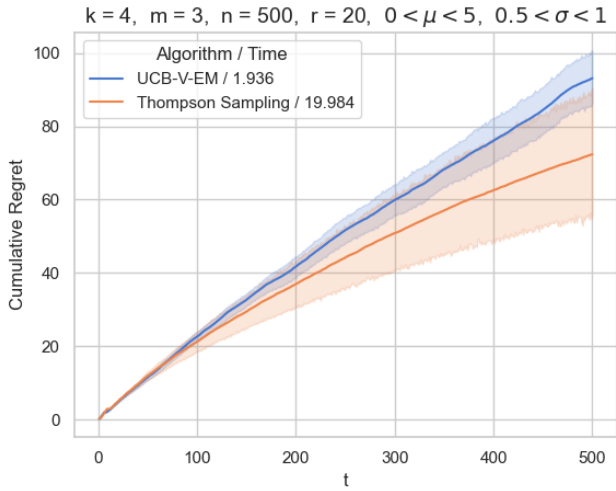


Figure 4.9

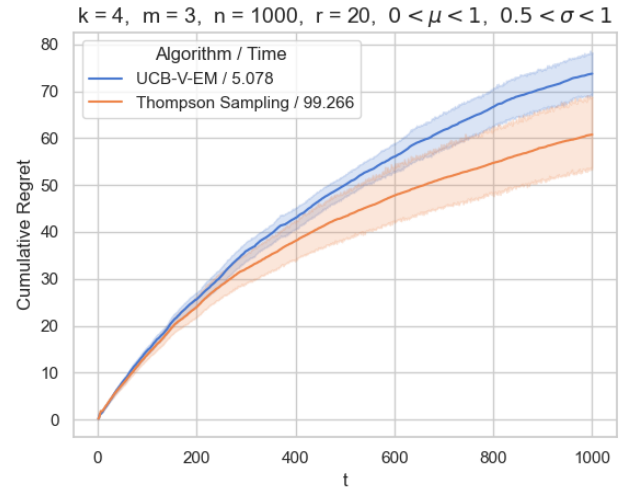


Figure 4.10

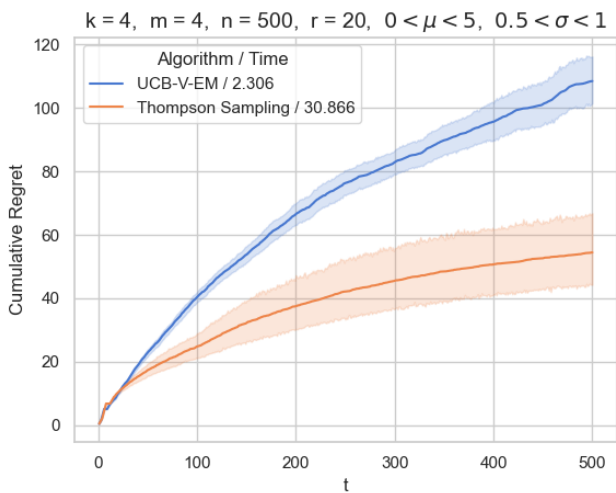


Figure 4.11

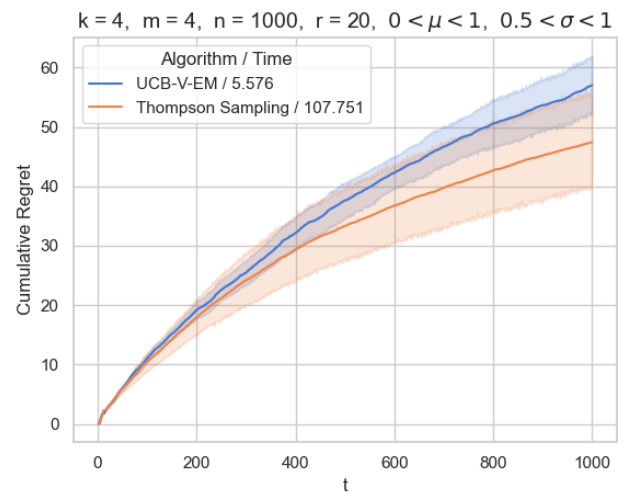


Figure 4.12

**8 arms case, regret comparison for Thompson Sampling and UCB-V-EM**

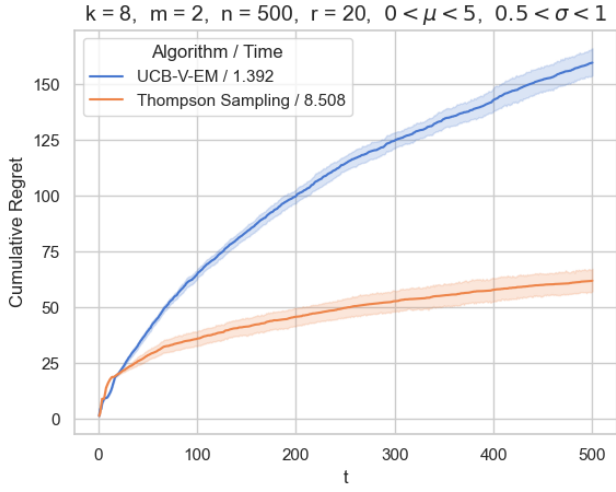


Figure 4.13

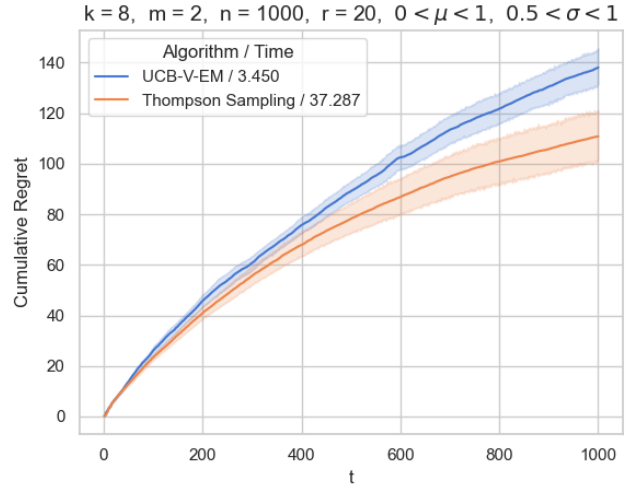


Figure 4.14

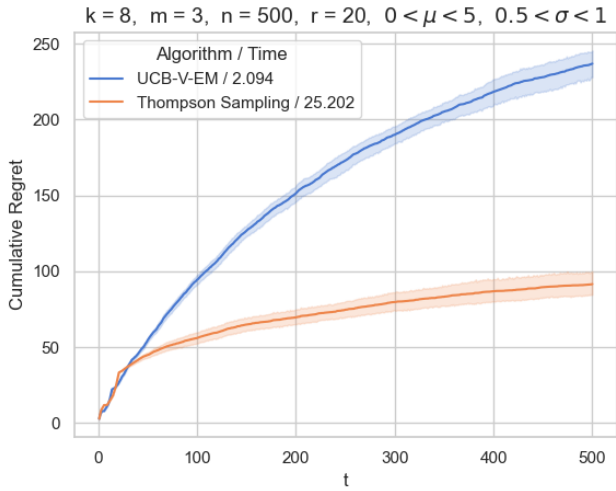


Figure 4.15

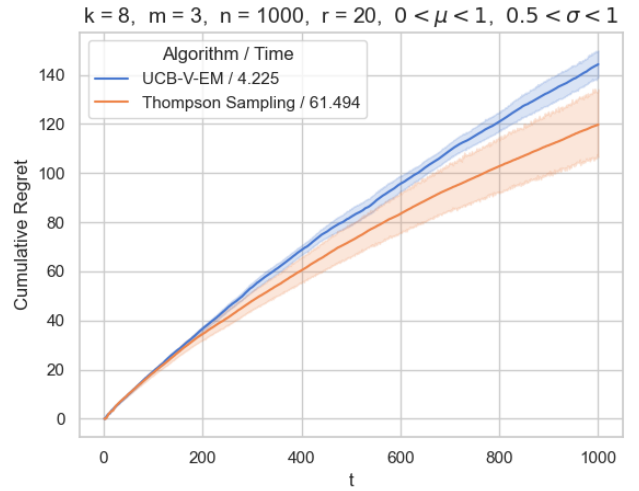


Figure 4.16

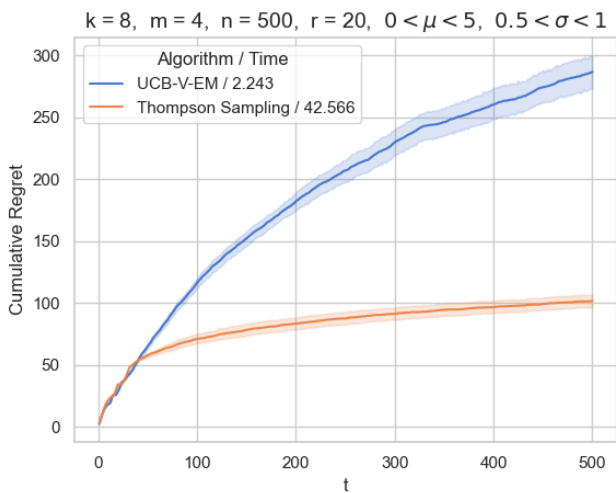


Figure 4.17

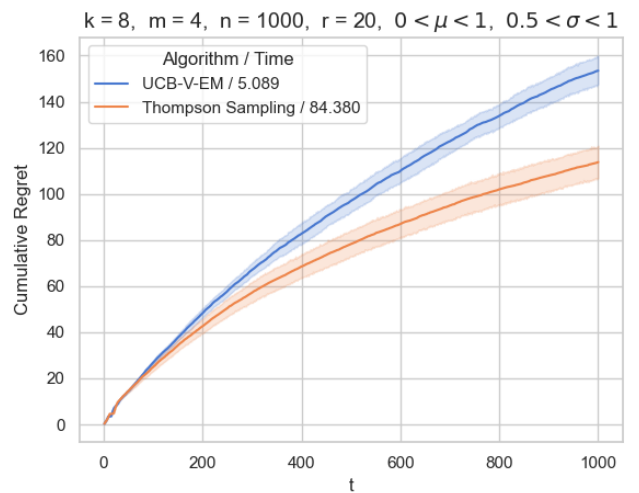


Figure 4.18

## 5. Mixture Model Bandits with Beta as Component Distributions

Beta distribution is also studied as another kind of reward distribution in this project. Here, the bandit environment is consisted of  $k$  arms with reward that follows a mixture of independent Beta distributions. In other words, suppose the distribution of the reward of an arm is  $f_X(x)$ , then:

$$f_X(x) = \sum_{j=1}^m \pi_j \cdot b_{\alpha_j, \beta_j}(x), \quad b_{\alpha, \beta}(x) = \frac{1}{B(\alpha, \beta)} \cdot x^{\alpha-1} \cdot (1-x)^{\beta-1},$$

where  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$ ,  $x \in [0, 1]$ .

The UCB-V algorithm is a main tool here because it assumes distributions over reals with bounded support  $[0, b]$ . The aim of this section is to see whether the performance of UCB-V algorithm can be improved with the implementation of a hybrid algorithm, which is also an iterative algorithm similar to the EM, proposed by Schröder and Rahmann (2017), to specifically deal with the problem of parameter estimations for the Beta mixture model.

The reasons why it's not recommended to directly use EM are given here. One of the main reasons why the EM algorithm is extensively used in practice for parameter estimation for mixture models is the existence of the log-likelihood. By Jensen's inequality, the likelihood increases in each EM iteration, and when it eventually stops increasing, it means the algorithm has found a stationary point. However, this solution might only be locally optimal. Global optimal solution can be found by repeating the run for several times and then comparing the log-likelihood of these stationary points.

In beta mixtures, several obstacles exist with the EM algorithm. First, the weights  $W_{i,j}$  are not properly defined for  $x_i = 0$  or  $x_i = 1$  for some beta distribution (specifically, if one of  $\alpha$  and  $\beta$  is less or equal to 1). Then, the maximization step cannot be carried out if the data contains any 0 or 1. Also, even if all  $x_i \in (0, 1)$ , the resulting mixtures are highly sensitive to perturbation of the data. Finally, for the EM algorithm for Beta mixtures, each M-step itself requires an iterative optimization procedure, the computational efforts for several EM runs are a great burden. E-MM algorithm, a hybrid algorithm that combines EM and method of moment proposed by Schröder and Rahmann (2017), is implemented in the bandit algorithm to work out the mixture model parameters in

this project.

## 5.1. Introduction to the E-MM Algorithm

E-MM algorithm is named so because it is a hybrid of the EM algorithm and the method of moments.

The algorithm proceeds iteratively as in the EM framework and alternates between EM's E step and Pearson's method of moments until the algorithm reaches convergence within a set threshold (Pearson, 1894).

To estimate  $Q$  free parameters, the method of moments calculates  $Q$  moments of the target distribution by expressing them in terms of the parameters, and equate them to the  $Q$  sample moments. This boils down to find the solution to a system of  $Q$  non-linear equations.

Thus, to avert the aforementioned problems specific to Beta mixtures, the algorithm substitute the M-step in EM by a method of moments step (MM-step). It therefore blends the idea of latent weights from EM algorithm with a method of moments estimation. It also avoids the exorbitant computational efforts from using pure moment based estimation for mixture models.

### 5.1.1. The Method of Moments for Beta Distribution

The beta distribution with parameters  $\alpha > 0$  and  $\beta > 0$  is a continuous probability distribution on the unit interval  $[0, 1]$ .

If  $X$  is a random variable with a beta distribution, then its expected value  $\mu$  and variance  $\sigma^2$  are:

$$\mu := \mathbb{E}[X] = \frac{\alpha}{\alpha + \beta}, \quad \sigma^2 := \text{Var}[X] = \frac{\mu(1 - \mu)}{\alpha + \beta + 1} = \frac{\mu(1 - \mu)}{1 + \phi}.$$

To use the method of moments, conversely, the parameters  $\alpha$  and  $\beta$  may be expressed in  $\mu$  and  $\sigma^2$ :

$$\phi = \frac{\mu(1 - \mu)}{\sigma^2} - 1; \quad \text{then} \quad \alpha = \mu\phi, \quad \beta = (1 - \mu)\phi. \quad (5.1)$$

### 5.1.2. The Algorithm

This section contains the details of the E-MM algorithm proposed by Schröder and Rahmann (2017).

### Initialization

The algorithm adapts the  $D^2$ -weighted initialization stage originally used in the famous k-means algorithm (Johnson, Wichern, et al., 2014). Denote  $X \subset [0, 1]$  to be the set of samples. Let  $Y \subset X$  be the set of initial component centres, such that  $Y = \{\}$  in the beginning. Let  $D_Y(x) := \min_{y \in Y} |x - y|$ . Then the initialization stage of the E-MM algorithm consists of:

- Sample the first point  $y$  randomly and uniformly from  $X$  to form  $Y := \{y\}$ .
- Repeat until  $|Y| = c$ : Choose  $y \in X \setminus Y$  with a categorical probability proportional to  $D_Y(y)^2$ ; then set  $Y := Y \cup \{y\}$ .
- Sort  $Y$  such that  $y_1 < \dots < y_c$
- Expectation and variance of component  $j = 1, \dots, c$  are initially estimated from the corresponding sample moments of all data points in the interval  $[y_j - 0.5, y_j + 0.5]$

### E-step

The E-step contains a minor modification from the expectation stage in EM in that the algorithm set weights of data points  $x_i = 0$  and  $x_i = 1$  in way that would not result in singularities.

Let  $j_0$  be the component index  $j$  with the smallest  $\alpha_j$ . If there is more than one, choose the one with the largest  $\beta_j$ . Then,  $W_{i,j_0} = 1$  and  $W_{i,j} = 0$  for  $j \neq j_0$ . Similarly, let  $j_1$  be the component index  $j$  with the smallest  $\beta_j$ . For all  $i$  with  $x_i = 1$ , set  $W_{i,j_1} = 1$  and  $W_{i,j} = 0$  for  $j \neq j_1$ .

### MM-step

The MM-step estimates mean and variance of each component  $j$  by responsibility-weighted sample moments,

$$\mu_j = \frac{\sum_{i=1}^n W_{ij} \cdot x_i}{\sum_{i=1}^n W_{ij}} = \frac{\sum_{i=1}^n W_{ij} \cdot x_i}{n \cdot \pi_j}, \quad \sigma_j^2 = \frac{\sum_{i=1}^n W_{ij} \cdot (x_i - \mu_j)^2}{n \cdot \pi_j} \quad (5.2)$$

Then  $\alpha_j$  and  $\beta_j$  are by the method of moments as specified in Eq. 5.1 and new weights through Bayes rules.

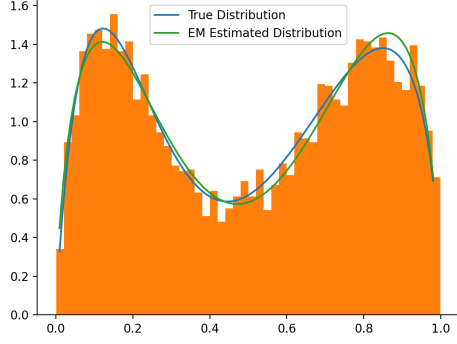
$$W_{i,j} = \frac{\pi_j b_{\alpha_j, \beta_j}(x_i)}{\sum_k \pi_k b_{\alpha_k, \beta_k}(x_i)} \quad \text{and} \quad \pi_j^+ = \frac{1}{n} \sum_{i=1}^n W_{i,j}. \quad (5.3)$$

### Termination

Let  $\theta_q$  be any real-valued parameter to be estimated and  $T_q$  a given threshold for  $\theta_q$ . After each MM-step, we compare  $\theta_q$  (old value) and  $\theta_q^+$  (updated value) by the relative change  $\kappa_q := |\theta_q^+ - \theta_q| / \max(|\theta_q^+|, |\theta_q|)$ . (If  $\theta_q^+ = \theta_q = 0$ , we set  $\kappa_q := 0$ .) We say that  $\theta_q$  is stationary if  $\kappa_q < T_q$ . The algorithm terminates when all parameters are stationary.

### 5.1.3. Empirical Results for E-MM Algorithm\*

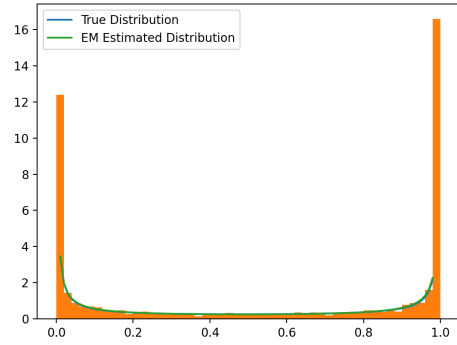
Below are the empirical results that show the algorithm's effectiveness at approximating Beta mixtures.



**Figure 5.1:** True distribution vs. The estimated distribution from E-MM algorithm

	component 1	component 2
$\alpha$	1.92	4.23
$\alpha^*$	1.80	4.59
$\beta$	7.68	1.57
$\beta^*$	6.82	2.43
$\pi$	0.42	0.58
$\pi^*$	0.50	0.55

**Table 5.1.:** Theoretical and Estimated Parameters of Mixture Beta Distribution



**Figure 5.2:** True distribution vs. The estimated distribution from E-MM algorithm

	component 1	component 2
$\alpha$	1.68	0.15
$\alpha^*$	1.81	0.09
$\beta$	1.43	0.13
$\beta^*$	1.68	0.14
$\pi$	0.03	0.97
$\pi^*$	0.11	0.86

**Table 5.2.:** Theoretical and Estimated Parameters of Mixture Beta Distribution

## 5.2. Implementation of the E-MM Algorithm in the UCB-V\*

UCB-V algorithm (Audibert, Munos, and Szepesvári, 2007) is implemented to study whether the proposed E-MM algorithm is able to improve the performance.

Again, the goal here is to use E-MM algorithm to obtain a hopefully better estimate for the empirical variance. For any  $k \in \{1, \dots, K\}$  and  $t \in \mathbb{N}$ , let  $\bar{X}_{k,t}$  and  $V_{k,t}$  be the

empirical estimates of the mean payoff and variance of arm  $k$ , then:

$$\bar{X}_{k,t} \triangleq \frac{1}{t} \sum_{i=1}^t X_{k,i} \quad \text{and} \quad V_{k,t} \triangleq \frac{1}{t} \sum_{i=1}^t (X_{k,i} - \bar{X}_{k,t})^2$$

If we use E-MM Algorithm to estimate the variance term, as derived in Eq. 3.5, it can be estimated as:

$$\bar{V}_{k,t} \triangleq \sum_{i=1}^m \hat{\pi}_i (\hat{\sigma}_i^2 + \hat{\mu}_i^2) - \hat{\mu}^2,$$

If the mixture model assumption is ignored, then the UCB-V algorithm seeks to maximize the bound in Eq. 3.3 in each play:

$$B_{k,s,t} \triangleq \bar{X}_{k,s} + \sqrt{\frac{2V_{k,s} \log t}{s}} + \frac{b \log t}{2s},$$

where  $V_{k,s}$  is the simple sample variance as written above,  $b$  is the supremum of the reward,  $k$  the arm,  $s$  the number of samples from arm  $k$  until time  $t$ . Note that  $b = 1$  in the Beta mixture case and there is no need to estimate it, because the support of the Beta distribution is on  $[0, 1]$ .

If the mixture model structure is considered, then use the E-MM algorithm to estimate  $V_{k,s}$  by the mixture model variance specified in Eq. 3.5.

### 5.3. Empirical Results for Beta Mixture Model Bandits\*

Even though the E-MM algorithm works quite well for estimating parameters of single Beta mixture model, the algorithms doesn't seem to have improved the performance of the UCB-V algorithm significantly as indicated in the regret plots presented afterwards.

The algorithm is arguably more complicated than EM and therefore takes more computational efforts as expected. It increases running time by the scale or nearly 1000 folds but doesn't pay off with a significant improvement in the cumulative regret.

Unfortunately, the bandit algorithm itself also doesn't work quite well here. Note that in most of the plots (especially those in the 4-components case) the algorithm doesn't seem to have achieved logarithmic regret in the given horizon. This could be due to the fact that a mixture of Beta distributions (which can have a big variety of shapes with different values of  $\alpha$  and  $\beta$ ) might be too challenging for the bandit algorithm to learn quickly.



# Regret Plots for UCB-V-EM and UCB-V for Beta Mixture, 2 Components

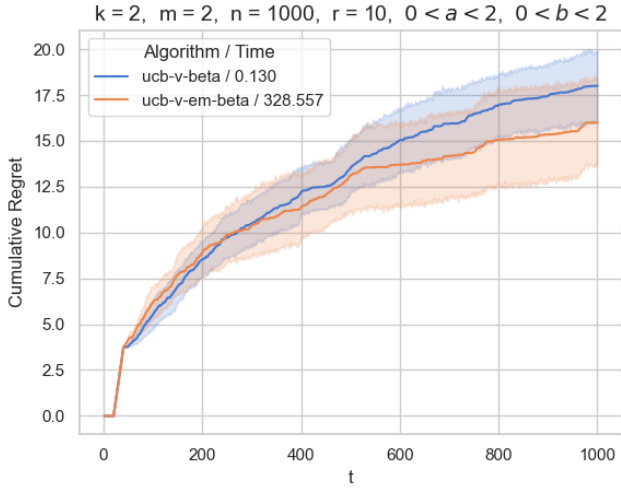


Figure 5.3

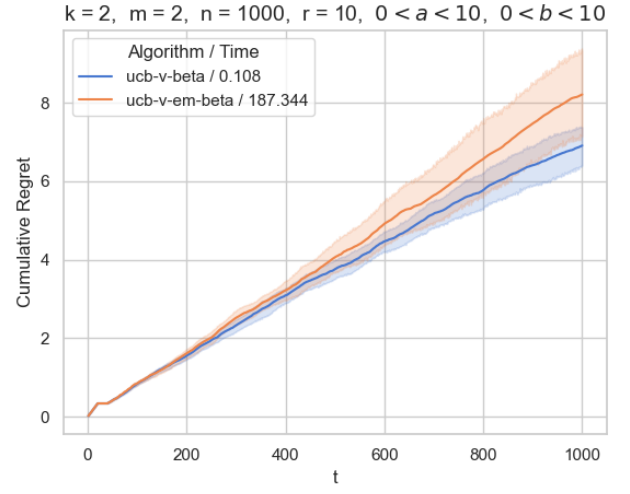


Figure 5.4

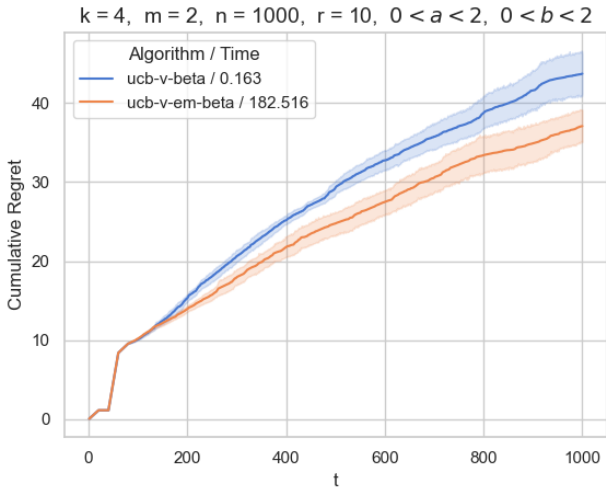


Figure 5.5

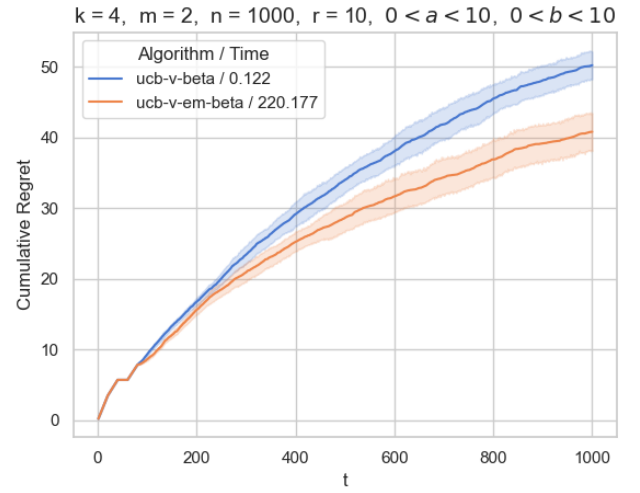


Figure 5.6

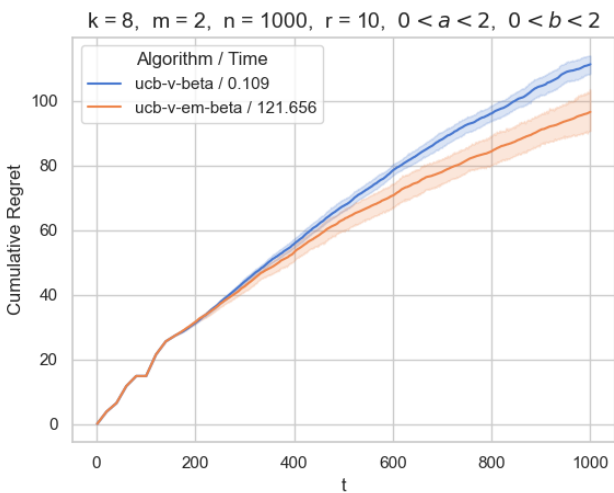


Figure 5.7

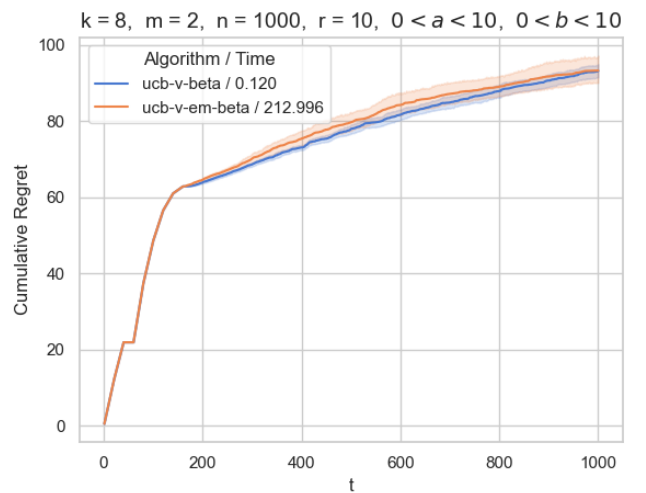


Figure 5.8

# Regret Plots for UCB-V-EM and UCB-V for Beta Mixture, 3 Components

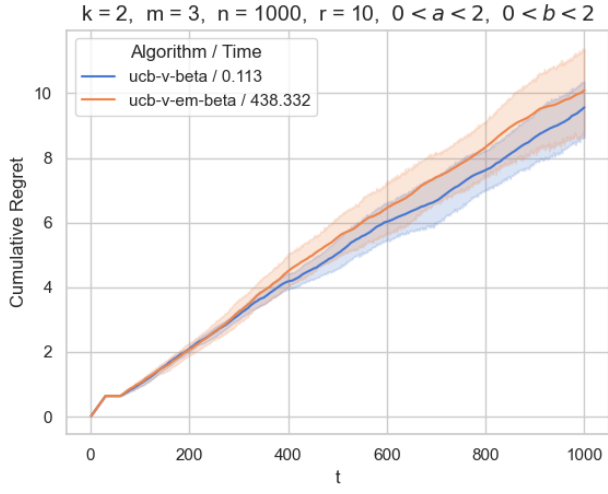


Figure 5.9

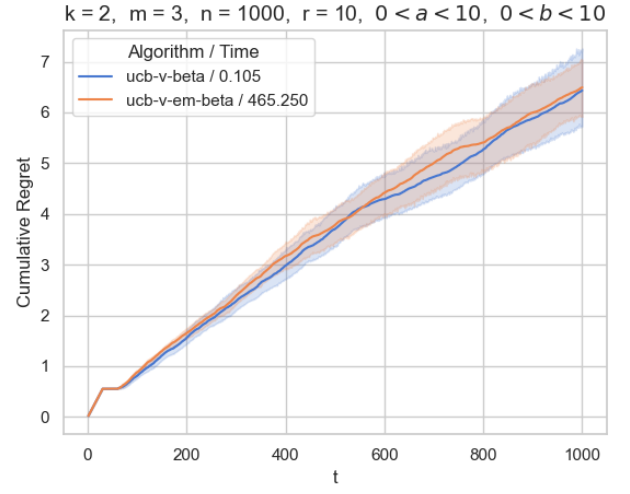


Figure 5.10

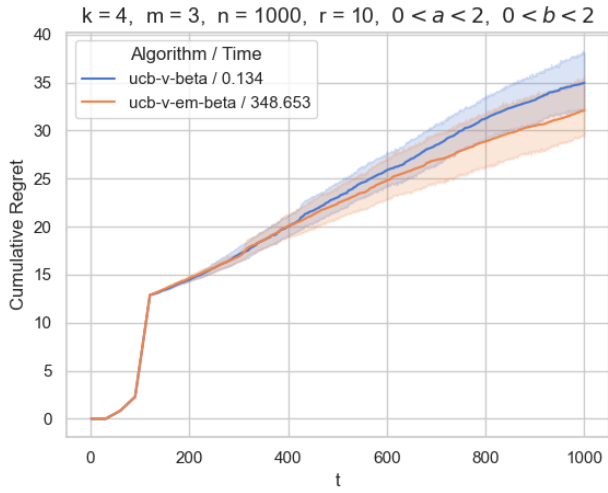


Figure 5.11

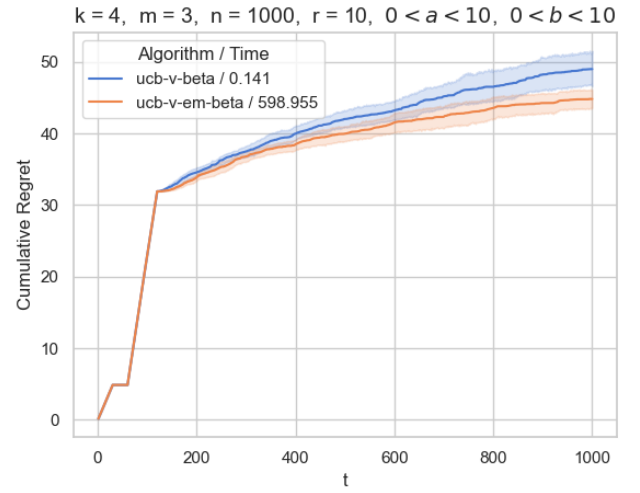


Figure 5.12

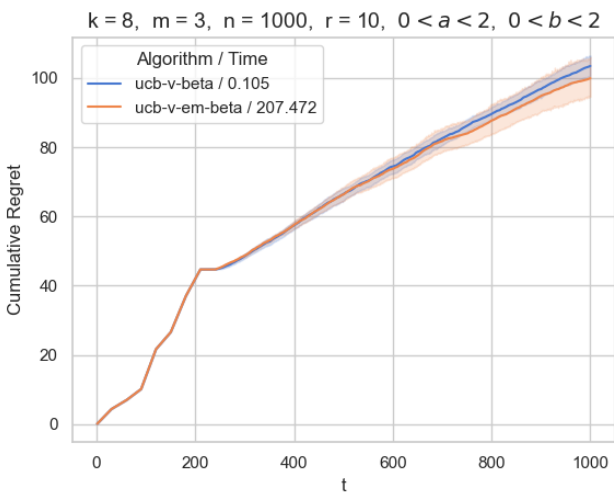


Figure 5.13

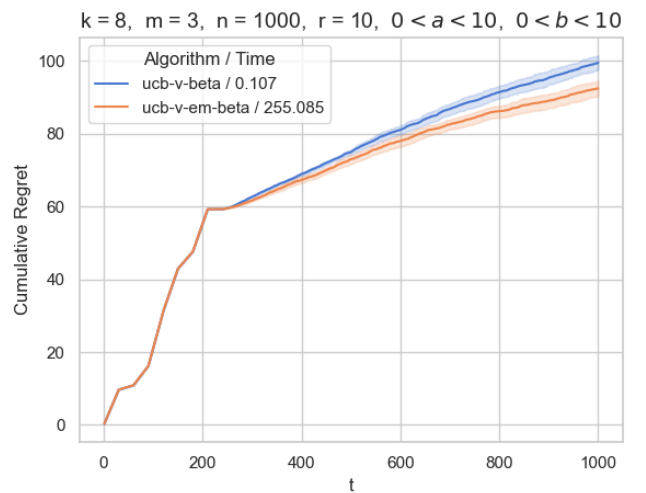


Figure 5.14

## 6. Conclusion\*

In the known variance case, experiments have been done for UCB and UCB-EM algorithm for mixture model bandit problem. The improvement in the performance from EM is apparent from the regret plots. By estimating the weights of each arm, EM has provided us with a more accurate estimate or bound for the overall subgaussianity of a mixture model and therefore improved the performance of the algorithm significantly.

However, EM doesn't bring us with significant benefits when the assumptions change to that the variances are totally unknown to the player. The empirical results are that the performances of the UCB-V and UCB-V-EM are almost entirely the same except that UCB-V-EM is slower because of more computational burden. This is also expected because what EM can really do here is to utilize this only new information, the number of components, to try to provide a better estimation for the overall variances. At the same time, the simple empirical sample variance is already an asymptotically unbiased estimator for the overall variance no matter we assume mixture model assumptions or not.

UCB1-NORMAL completely falls out of discussion because its performance is strictly dominated by that of the UCB-V algorithm under mixture model assumption. This is also a quite reasonable and explainable result because the assumption of UCB1-NORMAL is that the rewards are strictly Gaussian, which is not the case here.

Thompson sampling is arguably the best algorithm here in terms of sheer cumulative regret. However, it takes much more computational efforts and should not be used without considering the practical needs.

In the Beta mixture model case, the improvement by EM is not evident.

# Bibliography

- [Attias et al. 1999] ATTIAS, Hagai et al.: A Variational Bayesian Framework for Graphical Models. In: *NIPS* Bd. 12 Citeseer, 1999
- [Audibert et al. 2007] AUDIBERT, Jean-Yves ; MUNOS, Rémi ; SZEPEŠVÁRI, Csaba: Tuning bandit algorithms in stochastic environments. In: *International conference on algorithmic learning theory* Springer, 2007, S. 150–165
- [Auer et al. 2002] AUER, Peter ; CESA-BIANCHI, Nicolo ; FISCHER, Paul: Finite-time analysis of the multiarmed bandit problem. In: *Machine learning* 47 (2002), Nr. 2, S. 235–256
- [Beal 2003] BEAL, Matthew J.: *Variational algorithms for approximate Bayesian inference*. University of London, University College London (United Kingdom), 2003
- [Blei and Jordan 2006] BLEI, David M. ; JORDAN, Michael I.: Variational inference for Dirichlet process mixtures. In: *Bayesian analysis* 1 (2006), Nr. 1, S. 121–143
- [Dempster et al. 1977] DEMPSTER, Arthur P. ; LAIRD, Nan M. ; RUBIN, Donald B.: Maximum likelihood from incomplete data via the EM algorithm. In: *Journal of the Royal Statistical Society: Series B (Methodological)* 39 (1977), Nr. 1, S. 1–22
- [Johnson et al. 2014] JOHNSON, Richard A. ; WICHERN, Dean W. et al.: *Applied multivariate statistical analysis*. Bd. 6. Pearson London, UK:, 2014
- [Lattimore and Szepesvári 2020] LATTIMORE, Tor ; SZEPEŠVÁRI, Csaba: *Bandit algorithms*. Cambridge University Press, 2020
- [Pearson 1894] PEARSON, Karl: Contributions to the mathematical theory of evolution. In: *Philosophical Transactions of the Royal Society of London. A* 185 (1894), S. 71–110
- [Schröder and Rahmann 2017] SCHRÖDER, Christopher ; RAHMANN, Sven: A hybrid parameter estimation algorithm for beta mixtures and applications to methylation state classification. In: *Algorithms for Molecular Biology* 12 (2017), Nr. 1, S. 1–12
- [Thompson 1933] THOMPSON, William R.: On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. In: *Biometrika* 25 (1933), Nr. 3/4, S. 285–294

# A. Python Codes

## A.1. bandit\_class.py

```
1 import numpy as np
2 from scipy.stats import norm, uniform
3 from tabulate import tabulate
4 from pprint import pprint
5
6
7 class MixtureBandit:
8     # Mixture Model Bandit class
9     def __init__(self, arm_num, components_num=2, mean_range=None,
10                  sd_range=None):
11         # arm_num: k, how many arms
12         # components_num: how many components
13         # mean_range: the range from which the component means are
14             sampled from
15         # sd_range: the range from which the component standard
16             deviations are sampled from
17
18         self.is_beta = False
19
20         self.arm_num = arm_num
21         self.components_num = components_num
22
23         self.mean_range = mean_range
24         self.sd_range = sd_range
25
26         self.alpha = np.random.dirichlet([1]*components_num, size=arm_num)
27
28         self.mean_list = np.random.uniform(low=mean_range[0], high=
29             mean_range[1], size=[arm_num, components_num])
30         self.sd_list = np.random.uniform(low=sd_range[0], high=sd_range
31             [1], size=[arm_num, components_num])
32
33         self.weighted_mean = [np.dot(self.alpha[i], self.mean_list[i])
34             for i in range(self.arm_num)]
35         self.optimal_arm_mean = max(self.weighted_mean)
36         self.optimal_arm = self.weighted_mean.index(self.optimal_arm_mean)
37             + 1
38
39     def print_optimal_arm(self):
40         print(f"\nBest arm is the arm {self.optimal_arm} "
41             f"with mean {self.optimal_arm_mean:1.3f}\n")
```

```

35     def get_para(self):
36         # output the parameters of the environment
37         return {'mean': self.mean_list, 'sd': self.sd_list, 'alpha': self
            .alpha}
38
39     def pull(self, arm, size):
40         # get rewards from an arm
41         arm_mean = self.mean_list[arm-1]
42         arm_sd = self.sd_list[arm-1]
43         my_sample = []
44         for i in range(size):
45             component_chosen = list(np.random.multinomial(1, self.alpha[
                arm-1]))
46             component_chosen = component_chosen.index(max(
                component_chosen))
47             my_sample.append(norm.rvs(loc=arm_mean[component_chosen],
                scale=arm_sd[component_chosen], size=size)[0])
48         return my_sample
49
50     def parameters(self, latex: bool):
51         # return a tabulate that displays the parameters of a bandit
            environment
52         table = list([])
53
54         for component_index in range(self.components_num):
55             mean = ['mean ' + str(component_index+1)] + list(np.round([self
                .mean_list[i][component_index] for i in range(self.arm_num
                )], 2))
56             sd = ['std ' + str(component_index + 1)] + list(np.round([
                self.sd_list[i][component_index] for i in range(self.
                arm_num)], 2))
57             alpha = ['alpha ' + str(component_index + 1)] + list(np.round
                ([self.alpha[i][component_index] for i in range(self.
                arm_num)], 2))
58             table.append(alpha)
59             table.append(mean)
60             table.append(sd)
61             table.append(['-' + ',' * self.arm_num])
62             table.append(['weighted mean'] + list(np.round(self.weighted_mean
                , 2)))
63             headers = ['arm ' + str(x + 1) if x + 1 != self.optimal_arm
64                 else 'arm ' + str(x + 1) + '*' for x in range(self.
                arm_num)]
65
66             if latex:
67                 return tabulate(table, headers=headers, tablefmt='latex')
68             return tabulate(table, headers=headers, tablefmt='simple')
69
70 # my_bandit = MixtureBandit(arm_num=2, components_num=3, mean_range=[0,
71     10], sd_range=[0.5, 1])
72 # print(my_bandit.parameters(latex=True))
73 # my_bandit.print_optimal_arm()
74 # print(np.mean(my_bandit.pull(arm=0, size=1000)))

```

## A.2. bandit\_class\_beta.py

```
1 import numpy as np
2 from scipy.stats import norm, uniform
3 from tabulate import tabulate
4 from pprint import pprint
5
6
7 class MixtureBanditBeta:
8     # the Beta mixture model Bandit class, the structure is
9     # mostly the same as 'MixtureBandit'
10    def __init__(self, arm_num, alpha_range, beta_range, components_num
11                =2):
12        self.is_beta = True
13
14        self.arm_num = arm_num
15        self.components_num = components_num
16
17        self.alpha_range = alpha_range
18        self.beta_range = beta_range
19
20        self.pi = np.random.dirichlet([1]*components_num, size=arm_num)
21        self.alpha_list = np.random.uniform(low=alpha_range[0], high=
22            alpha_range[1], size=[arm_num, components_num])
23        self.beta_list = np.random.uniform(low=beta_range[0], high=
24            beta_range[1], size=[arm_num, components_num])
25
26        self.mean_list = self.alpha_list/(self.alpha_list+self.beta_list)
27        self.var_list = self.alpha_list*self.beta_list/\
28            (self.alpha_list+self.beta_list)**2\
29            /(self.alpha_list+self.beta_list+1)
30
31        self.weighted_mean = [np.dot(self.mean_list[i], self.pi[i]) for i
32                                in range(self.arm_num)]
33        self.optimal_arm_mean = max(self.weighted_mean)
34        self.optimal_arm = self.weighted_mean.index(self.optimal_arm_mean
35            ) + 1
36
37    def print_optimal_arm(self):
38        print(f"\nBest arm is the arm {self.optimal_arm} "
39            f"with mean {self.optimal_arm_mean:1.3f}\n")
40
41    def get_para(self):
42        return {'alpha': self.alpha_list, 'beta': self.beta_list, 'pi':
43            self.pi}
44
45    def pull(self, arm, size):
46        arm_alpha = self.alpha_list[arm-1]
47        arm_beta = self.beta_list[arm-1]
48        my_sample = []
49        for i in range(size):
50            component_chosen = list(np.random.multinomial(1, self.pi[arm
51                -1]))
```

```

45         component_chosen = component_chosen.index(max(
46             component_chosen))
47         my_sample.append(np.random.beta(a=arm_alpha[component_chosen],
48             b=arm_beta[component_chosen], size=size)[0])
49     return my_sample
50
51 def parameters(self, latex: bool):
52     table = list([])
53
54     for component_index in range(self.components_num):
55         alpha = ['alpha ' + str(component_index+1)] + \
56             list(np.round([self.alpha_list[i][component_index]
57                 for i in range(self.arm_num)], 2))
58         beta = ['beta ' + str(component_index + 1)] + \
59             list(np.round([self.beta_list[i][component_index]
60                 for i in range(self.arm_num)], 2))
61         pi = ['pi ' + str(component_index + 1)] + \
62             list(np.round([self.pi[i][component_index]
63                 for i in range(self.arm_num)], 2))
64         table.append(alpha)
65         table.append(beta)
66         table.append(pi)
67         table.append(['-' + '-' * self.arm_num])
68         table.append(['weighted mean'] + list(np.round(self.weighted_mean, 2)))
69         headers = ['arm ' + str(x + 1) if x + 1 != self.optimal_arm
70             else 'arm ' + str(x + 1) + '*' for x in range(self.arm_num)]
71         if latex:
72             return tabulate(table, headers=headers, tablefmt='latex')
73         return tabulate(table, headers=headers, tablefmt='simple')
74
75 # my_bandit = MixtureBanditBeta(arm_num=3, components_num=2, alpha_range
76 #     =[0, 5], beta_range=[0, 5])
77 # print(my_bandit.parameters(latex=False))
78 # my_bandit.print_optimal_arm()
79 # print(my_bandit.pull(arm=1, size=10))
80 #

```

### A.3. UCB\_known\_sd.py

```

1 import numpy as np
2 from scipy.stats import norm, uniform
3 from sklearn.cluster import KMeans
4 from bandit_class import MixtureBandit
5
6
7 def ucb_bound(t, prev_t, mu, sd, alpha_value):
8     # This is the upper confidence bound for unknown variance case
9     # input: parameters for the bound
10    # output: the UCB bound
11    def ft(x): return 1 + x * np.log(x) ** 2

```



```

12     if alpha_value is not None:
13         var_est = np.dot(alpha_value, sd ** 2)
14         return mu + np.sqrt(2 * var_est * np.log(ft(t)) / prev_t)
15     return mu + np.sqrt(2 * np.max(sd ** 2) * np.log(ft(t)) / prev_t)
16
17
18 def em_known_var(x, sd, m, acc, prev_mean, prev_alpha):
19     # EM algorithm for known variance case
20     # INPUT:
21     # x: samples,
22     # sd: standard deviations,
23     # m: components number
24     # acc: accuracy criterion for convergence
25     # prev_mean: component means when the algorithm converged last time
26     # prev_alpha: weights when the algorithm converged last time
27
28     # OUTPUT:
29     # [component means, component weights]
30     x = np.array(x)
31     big_n = len(x)
32     if (np.min(prev_alpha) > 0.01) & (np.min(1 - prev_alpha) > 0.01) & \
33         (np.min(np.abs(prev_alpha - 1/m)) > 0.01) & (big_n > 50):
34         mean = prev_mean
35         alpha = prev_alpha
36     else:
37         k_means = KMeans(n_clusters=m, random_state=0).fit(x.reshape(-1,
38         1))
39         mean = np.squeeze(k_means.cluster_centers_)
40         alpha = uniform.rvs(loc=0.4, scale=0.3, size=m)
41
42     while True:
43         p = np.array([norm.pdf(x, loc=mean[i], scale=sd[i]) for i in
44         range(m)])
45         p_x = np.array([np.dot([p[j][i] for j in range(m)], alpha) for i
46         in range(big_n)])
47         w = np.array([alpha[i] * p[i] / p_x for i in range(m)])
48
49         new_alpha = np.array([np.sum(w[i]) for i in range(m)]) / big_n
50         new_mean = np.array([1 / np.sum(w[i]) * np.dot(w[i], x) for i in
51         range(m)])
52
53         if (np.max(np.abs(mean - new_mean)) < acc) & (np.max(np.abs(
54         new_alpha - alpha)) < acc):
55             break
56         alpha = new_alpha
57         mean = new_mean
58     return [mean, alpha]
59
60 def ucb(bandit: MixtureBandit, n: int):
61     # UCB algorithm for the known variance case
62     # INPUT:
63     # bandit: a bandit instance from MixtureBandit class

```

```

60     # n: horizon
61
62     # OUTPUT:
63     # [[action 1, reward 1],...,[action n, reward n]]
64     k = bandit.arm_num # how many arms?
65     m = bandit.components_num
66     samples = []
67     mu_lst = np.array([0.] * k) # means of each arm
68     sd = np.array(bandit.sd_list)
69     count_lst = [0] * k
70
71     size = np.max([m, 10])
72     for i in range(k):
73         sample = bandit.pull(arm=i + 1, size=size)
74         for j in range(size):
75             samples.append([i + 1, sample[j]])
76             count_lst[i] += size
77             mu_lst[i] = np.mean(sample)
78
79     for t in range(k * size, n):
80         if t % 100 == 0:
81             print('t = ', t)
82             ucb_lst = [ucb_bound(t=t,
83                                 prev_t=count_lst[j],
84                                 mu=mu_lst[j],
85                                 sd=sd[j],
86                                 alpha_value=None) for j in range(k)]
87
88             action = ucb_lst.index(max(ucb_lst)) + 1
89             count_lst[action - 1] += 1
90             new_x = bandit.pull(arm=action, size=1)
91             samples.append([action, new_x])
92             mu_lst[action - 1] = np.mean([x[1] for x in samples if x[0] ==
93                                           action])
94     print(count_lst)
95     return samples
96
97 def ucb_em(bandit: MixtureBandit, n: int):
98     # UCB-EM Algorithm for the known variance case
99     # INPUT
100     # bandit: a bandit instance from the MixtureBandit class
101     # n: horizon
102
103     # OUTPUT
104     # [[action 1, reward 1],...,[action n, reward n]]
105     print(bandit.parameters(latex=False))
106     k = bandit.arm_num
107     m = bandit.components_num
108     samples = []
109     mu_lst = np.array([0.] * k)
110     sd = bandit.sd_list
111

```

```

112     prev_means = np.array([[0.] * m] * k)
113     prev_alphas = np.random.dirichlet([1]*m, size=k)
114     count_lst = [0] * k
115
116     size = np.max([10, m])
117     for i in range(k):
118         sample = bandit.pull(arm=i + 1, size=size)
119         count_lst[i] += size
120         for j in range(size):
121             samples.append([i + 1, sample[j]])
122         mu_lst[i] = np.mean(sample)
123         prev_means[i] = [np.mean(sample)]*m
124
125     action = 0
126     ucb_lst_em = 0
127     for t in range(k * size, n): # do the UCB algorithm for remaining of
128         t from k+1 to n
129         if t % 100 == 0:
130             print(f"t = {t}, chosen arm: {action}, ucb_em: {ucb_lst_em},
131                   count_lst: {count_lst}")
132
133         ucb_lst_em = [ucb_bound(t=t, prev_t=count_lst[j], mu=mu_lst[j],
134                                sd=sd[j],
135                                alpha_value=prev_alphas[j]) for j in
136                        range(k)]
137
138         action = ucb_lst_em.index(max(ucb_lst_em)) + 1
139         new_x = bandit.pull(arm=action, size=1)[0]
140         samples.append([action, new_x])
141         count_lst[action-1] += 1
142
143         prev_samples = [x[1] for x in samples if x[0] == action]
144         [means, alpha] = em_known_var(x=prev_samples,
145                                       sd=sd[action - 1],
146                                       m=m,
147                                       acc=1e-3,
148                                       prev_mean=prev_means[action - 1],
149                                       prev_alpha=prev_alphas[action - 1])
150
151         prev_means[action - 1] = means
152         prev_alphas[action - 1] = alpha
153         mu_lst[action - 1] = np.mean([x[1] for x in samples if x[0] ==
154                                     action])
155     print(count_lst)
156     return samples

```

## A.4. UCB\_Normal.py

```

1 import numpy as np
2 from UCB_known_sd import *
3 from UCB_V import gm_sol
4
5
6 def ucb_normal(bandit: MixtureBandit, n: int, em=False):

```

```

7   # UCB-NORMAL algorithm
8   # INPUT:
9   # bandit: a bandit instance from the MixtureBandit class
10  # n: horizon
11  # em: boolean, use EM or not?
12  # OUTPUT:
13  # [[action 1, reward 1],...,[action n, reward n]]
14  k = bandit.arm_num
15  samples = []
16  count_lst = [0] * k
17  sd_lst = [1] * k
18  bounds = 0
19  for t in range(1, n+1):
20      if t % 200 == 0:
21          print(f"t = {t}, bounds = {bounds}, count = {count_lst}")
22      check_log = count_lst < np.ceil(8 * np.log(t))
23      if t == 1:
24          check_log[0] = True
25      cond = False
26      for i in range(1, k+1):
27          if check_log[i-1]:
28              samples.append([i, bandit.pull(arm=i, size=1)[0]])
29              count_lst[i-1] += 1
30              cond = True
31              break
32      if cond:
33          continue
34
35      bounds = [0]*k
36      for arm in range(1, k+1):
37          prev_samples = np.array([x[1] for x in samples if x[0] == arm
38                                  ])
39          n = len(prev_samples)
40          x_mean = np.mean(prev_samples)
41
42          if em:
43              bounds[arm - 1] = x_mean + np.sqrt(16 * np.square(sd_lst[
44                  arm-1]) * np.log(t - 1) / n)
45          else:
46              q_j = np.sum(np.square(prev_samples))
47              bounds[arm - 1] = x_mean + np.sqrt(16 * q_j/(n - 1) * np.
48                  log(t-1)/n)
49
50      action = bounds.index(max(bounds)) + 1
51      new_x = bandit.pull(arm=action, size=1)[0]
52      if em:
53          prev_samples = np.array([x[1] for x in samples if x[0] ==
54                                  action])
55          solutions = gm_sol(prev_samples, bandit.components_num)
56          sd_lst[action-1] = solutions['sd']
57      samples.append([action, new_x])
58      count_lst[action - 1] += 1
59  print(count_lst)

```

```

56     return samples
57
58
59 def ucb_normal_em(bandit, n):
60     return ucb_normal(bandit, n, em=True)

```

## A.5. UCB\_V.py

```

1  import numpy as np
2  from bandit_class import MixtureBandit
3  from sklearn.mixture import GaussianMixture
4  from UCB_known_sd import *
5
6
7  def gm_sol(x, m):
8      # Function for solving gaussian mixture parameters
9      # INPUT:
10     # x: samples from an arm
11     # m: how many components?
12     # OUTPUT:
13     # sd: estimated overall standard deviation of the arm
14     x = np.array(x)
15     x_bar = np.mean(x)
16     gm = GaussianMixture(n_components=m, random_state=0, warm_start=True)
17         .fit(x.reshape(-1, 1))
18     component_means = gm.means_.reshape(m)
19     component_vars = gm.covariances_.reshape(m)
20     my_cov = np.dot(component_means**2+component_vars, gm.weights_) -
21         x_bar**2
22     return {'sd': np.sqrt(my_cov)}
23
24 def ucb_bern_bound(mu, var, b, t, big_t):
25     # UCB bernstein bound
26     return mu + np.sqrt(2 * var * np.log(t) / big_t) + b * np.log(t) /
27         (2*big_t)
28
29 def ucb_v(bandit: MixtureBandit, n: int, em=False):
30     # UCB-V Algorithm
31     # bandit: a bandit instance from the MixtureBandit class
32     # n: horizon
33     # em: use EM algorithm or not?
34     # OUTPUT:
35     # [[action 1, reward 1],...,[action n, reward n]]
36     k = bandit.arm_num
37     m = bandit.components_num
38     samples = []
39     mu_lst = [0] * k
40     sd_lst = [0] * k
41     count_lst = [0] * k
42     b = 1

```

```

43     for i in range(k):
44         initial_samples = bandit.pull(arm=i+1, size=m-1)
45         if max(initial_samples) > b:
46             b = max(initial_samples)
47         for sample in initial_samples:
48             samples.append([i+1, sample])
49         mu_lst[i] = np.mean(initial_samples)
50         sd_lst[i] = 1
51         count_lst[i] += m-1
52
53     ucb_lst = 0
54     for t in range(k*(m-1), n):
55         if t % 200 == 0:
56             print(f"t = {t}, bernstein bounds = {ucb_lst}, count = {
57                 count_lst}")
58
59         ucb_lst = [ucb_bern_bound(t=t, big_t=count_lst[i],
60                                mu=mu_lst[i],
61                                var=sd_lst[i] ** 2,
62                                b=b) for i in range(k)]
63         action = ucb_lst.index(max(ucb_lst)) + 1
64         new_x = bandit.pull(arm=action, size=1)[0]
65         if new_x > b:
66             b = new_x
67         samples.append([action, new_x])
68         count_lst[action-1] += 1
69
70         prev_samples = [x[1] for x in samples if x[0] == action]
71         mu_lst[action-1] = np.mean(prev_samples)
72         if not em:
73             sd_lst[action-1] = np.std(prev_samples)
74         else:
75             solution = gm_sol(prev_samples, m)
76             sd_lst[action-1] = solution['sd']
77     print(count_lst, ', b = ', b)
78     return samples
79
80 def ucb_v_em(bandit, n):
81     return ucb_v(bandit, n, em=True)

```

## A.6. Thompson\_Sampling.py

```

1 from sklearn.mixture import BayesianGaussianMixture
2 from bandit_class import MixtureBandit
3 from scipy.stats import norm
4 from pprint import pprint
5 import numpy as np
6
7
8 def posterior_para(sample, bandit):
9     # Posterior of the parameters
10    # INPUT:

```

```

11 # sample: samples
12 # bandit: a bandit instance from the bandit class
13 # OUTPUT:
14 # a dictionary containing:
15 # 'dirichlet_concentration': concentration parameters of the
    posterior
16 #
    dirichlet distribution for the weights
17 # 'mean': means of the posterior of component mean
18 # 'mean_sd': standard deviations of the posterior of component mean
19 x = np.squeeze(np.array(sample)).reshape(-1, 1)
20 bgm = BayesianGaussianMixture(n_components=bandit.components_num,
21                               mean_precision_prior=0.1,
22                               covariance_prior=[[np.mean(bandit.
    sd_range)**2]],
23                               max_iter=5000,
24                               weight_concentration_prior_type='
    dirichlet_distribution',
25                               warm_start=True).fit(x)
26
27 return {'dirichlet_concentration': np.squeeze(bgm.
    weight_concentration_),
28         'mean': np.squeeze(bgm.means_),
29         'mean_sd': 1/np.sqrt(bgm.mean_precision_)}
30
31
32 def thompson_sampling(bandit: MixtureBandit, n: int):
33     # Thompson sampling algorithm
34     # INPUT:
35     # bandit: a bandit instance from the bandit class
36     # n: horizon
37     # OUTPUT:
38     # [[action 1, reward 1], ..., [action n, reward n]]
39     k = bandit.arm_num
40     m = bandit.components_num
41     para_list = [0]*k
42     samples = []
43     count_lst = [0]*k
44     for i in range(k):
45         first_samples = bandit.pull(arm=i+1, size=m)
46         para_list[i] = posterior_para(first_samples, bandit)
47         for j in range(m):
48             samples.append([i + 1, first_samples[j]])
49         count_lst[i] += m
50
51     for i in range(m*k, n):
52         if i % 200 == 0:
53             print('t = ', i)
54             # pprint(para_list, indent=2)
55         mean_list = [0]*k
56         for j in range(k):
57             para_arm_j = para_list[j]
58
59             means = para_arm_j['mean']

```

```

60     mean_sds = para_arm_j['mean_sd']
61     dirichlet_concentration = para_arm_j['dirichlet_concentration']
62
63     mean_posterior_list = [float(norm.rvs(loc=means[i], scale=
64         mean_sds[i], size=1)) for i in range(m)]
65     dirichlet_sample = np.squeeze(np.random.dirichlet(
66         dirichlet_concentration, size=1))
67
68     weighted_mean = np.dot(mean_posterior_list, dirichlet_sample)
69     mean_list[j] = weighted_mean
70
71     action = mean_list.index(max(mean_list))+1
72     count_lst[action-1] += 1
73     new_x = bandit.pull(arm=action, size=1)[0]
74     samples.append([action, new_x])
75
76     prev_samples = np.array([x[1] for x in samples if x[0] == action
77         ])
78     para_list[action-1] = posterior_para(prev_samples, bandit)
79
80     print(count_lst)
81     for dic in para_list:
82         for key, value in dic.items():
83             print(f'{key:25}{value}')
84     return samples

```

## A.7. UCB\_V\_Beta.py

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3  from bandit_class_beta import MixtureBanditBeta
4  from UCB_known_sd import *
5  from scipy.stats import beta
6  from scipy.interpolate import *
7  from pprint import pprint
8  from regret_comparison import *
9
10
11 def check_cond(new_para, prev_para, threshold):
12     # function that checks whether the algorithm has converged
13     # new_para: the new set of parameters
14     # prev_para: the previous set of parameters
15     # threshold: the threshold
16
17     # OUTPUT:
18     # Boolean, TRUE or FALSE
19
20     if new_para == prev_para or np.max([np.abs(new_para), np.abs(
21         prev_para)]) == 0:
22         return True
23
24     Kq = np.abs(new_para-prev_para) / np.max([np.abs(new_para), np.abs(
25         prev_para)])

```



```

23     if Kq < threshold:
24         return True
25     else:
26         return False
27
28
29 def em_beta(x, m, acc, cache_para=None):
30     # the EM-M algorithm for solving parameters of a beta mixture
31     # x: samples
32     # m: components number
33     # acc: accuracy for convergence
34     # cache_para: previous converged parameters when the algorithm ran
35         last time
36
37     # OUTPUT:
38     # a dictionary containing:
39     # parameters of a beta mixtures: alpha, beta, weights (pi), mu (
40         component means)
41     # var (component variances), and also
42     # loop count: how many loops have this algorithm taken to finish?
43     N = len(x)
44     if cache_para:
45         alpha_lst = cache_para['alpha']
46         beta_lst = cache_para['beta']
47         pi = cache_para['pi']
48     else:
49         y1 = np.random.uniform(low=0, high=1, size=1)[0]
50         y = np.array([y1])
51         while len(y) < m:
52             candidates = np.array([i for i in x if i not in y])
53             Dy = np.array([np.min(np.abs(y - c)) for c in candidates])**2
54             new_y_index = list(np.random.multinomial(1, Dy/sum(Dy)))
55             new_y_index = new_y_index.index(max(new_y_index))
56             y = np.append(y, candidates[new_y_index])
57         y.sort()
58         alpha_lst = np.random.rand(m)
59         beta_lst = np.random.rand(m)
60         pi = np.array([1/m]*m)
61
62         for component_num in range(m):
63             samples = np.array([sample for sample in x if (sample <= y[
64                 component_num]+0.5) and
65                 (sample >= y[component_num] - 0.5)])
66             # print(samples)
67             mu = np.mean(samples)
68             if len(samples) <= 1:
69                 samples = x
70             sigma2 = np.std(samples)**2
71             phi = mu*(1-mu)/sigma2 - 1
72             alpha_lst[component_num] = mu * phi
73             beta_lst[component_num] = (1 - mu) * phi
74
75     loop_count = 0

```

```

73 start_time = time.perf_counter()
74 while True:
75     loop_count += 1
76     W = np.array([[pi[j] * beta.pdf(x[i], a=alpha_lst[j], b=beta_lst[
77         j]]) /
78         np.dot(pi, [beta.pdf(x[i], a=alpha_lst[k], b=
79             beta_lst[k])
80             for k in range(m)])
81         if not np.isnan(pi[j] * beta.pdf(x[i], a=alpha_lst
82             [j], b=beta_lst[j])) /
83         np.dot(pi, [beta.pdf(x[i], a=alpha_lst[k], b=
84             beta_lst[k])
85             for k in range(m)])]) else np.random.
86         uniform(low=0, high=1, size=1)[0]
87         for j in range(m)] for i in range(N)])
88
89     prev_pi, prev_alpha, prev_beta = pi, alpha_lst, beta_lst
90     pi = np.sum(W, axis=0) / N
91     mu = np.dot(x, W) / (N * pi)
92     sigma2 = np.array([np.dot(W[:, j], (x - mu[j])**2) / (N * pi[j])
93         for j in range(m)])
94     phi_lst = mu * (1 - mu) / sigma2 - 1
95     alpha_lst, beta_lst = mu * phi_lst, (1 - mu) * phi_lst
96
97     my_cond = True
98     for i in range(m):
99         if not (check_cond(alpha_lst[i], prev_alpha[i], acc)
100             & check_cond(prev_beta[i], beta_lst[i], acc)
101             & check_cond(prev_pi[i], pi[i], acc)):
102             my_cond = False
103             break
104     if my_cond | (time.perf_counter()-start_time >= 5.0 * 3600
105         /6/10/1000):
106         break
107
108     return {'alpha': alpha_lst, 'beta': beta_lst, 'pi': pi, 'mu': mu, '
109         var': sigma2,
110         'loop_count': loop_count}
111
112
113
114 def plot_beta_mixture(alphas, betas, pi, label):
115     # plot a beta mixture
116     x = np.arange(0.01, 0.99, 0.01)
117     y = np.array([0.]*len(x))
118     for i in range(len(alphas)):
119         y += pi[i] * beta.pdf(x, a=alphas[i], b=betas[i])
120     plt.plot(x, y, label=label)
121
122
123
124 def plot_est_and_true(beta_bandit, arm, para, samples):
125     # plot a beta mixture given a bandit arm
126     # INPUT:
127     # beta_bandit: a beta_bandit instance from the MixtureBanditBeta

```

```

class
118 # arm: which arm?
119 # para: dictionary containing alphas, betas, and weights (from EM-M
    algorithm)
120 # samples: samples
121
122 # OUTPUT:
123 # a plot containing histogram, a theoretical beta mixture densities,
    an estimated
124 # beta density using parameters worked out by EM-M algorithm
125 ind = arm - 1
126 plot_beta_mixture(alphas=beta_bandit.alpha_list[ind],
127                   betas=beta_bandit.beta_list[ind],
128                   pi=beta_bandit.pi[ind],
129                   label='True Distribution')
130
131 plt.hist(samples, bins=50, density=True)
132
133 plot_beta_mixture(alphas=para['alpha'],
134                   betas=para['beta'],
135                   pi=para['pi'],
136                   label='EM Estimated Distribution')
137 plt.legend()
138 plt.show()
139
140
141 def ucb_v_beta(bandit, n: int, em=False):
142 # UCB-V Algorithm for Beta Mixture Bandits
143 # bandit: a beta mixture bandit instance
144 # n: horizon
145 # em: use EM or not?
146 # OUTPUT:
147 # [[action 1, reward 1], ..., [action n, reward n]]
148 k = bandit.arm_num
149 m = bandit.components_num
150 samples = []
151 mu_lst = np.array([0.] * k)
152 sd_lst = np.array([0.] * k)
153 cache_alphas = np.random.rand(k, m)
154 cache_betas = np.random.rand(k, m)
155 cache_pi = np.random.rand(k, m)
156 count_lst = [0] * k
157 loop_count = ['init']
158
159 each_arm_sample_size = int(np.min([m * 10, np.floor(n / k)]))
160 for i in range(k):
161     initial_samples = bandit.pull(arm=i+1, size=each_arm_sample_size)
162     for sample in initial_samples:
163         samples.append([i+1, sample])
164     mu_lst[i] = np.mean(initial_samples)
165     sd_lst[i] = np.std(initial_samples)
166     count_lst[i] += each_arm_sample_size
167

```

```

168 ucb_lst = 0
169 for t in range(k*each_arm_sample_size, n):
170     if em:
171         if t % 10 == 0:
172             print(f"EM-beta, t = {t}, bernstein bounds = {ucb_lst},
173                   count = {count_lst},"
174                   f" loop count is {loop_count[-1]}")
175         else:
176             if t % 200 == 0:
177                 print(f"Non-EM-beta, t = {t}, bernstein bounds = {ucb_lst
178                       }, count = {count_lst}")
179
180 ucb_lst = [ucb_bern_bound(t=t, big_t=count_lst[i],
181                          mu=mu_lst[i],
182                          var=sd_lst[i] ** 2,
183                          b=1) for i in range(k)]
184 action = ucb_lst.index(max(ucb_lst)) + 1
185 new_x = bandit.pull(arm=action, size=1)[0]
186 samples.append([action, new_x])
187 count_lst[action-1] += 1
188 prev_samples = np.array([x[1] for x in samples if x[0] == action
189 ])
190 mu_lst[action-1] = np.mean(prev_samples)
191 if not em:
192     sd_lst[action-1] = np.std(prev_samples)
193 else:
194     if len(prev_samples) <= 50:
195         para = em_beta(prev_samples, m, acc=0.1)
196     else:
197         para = em_beta(prev_samples, m, acc=0.1,
198                        cache_para={'alpha': cache_alphas[action
199 -1],
200                                'beta': cache_betas[action-1],
201                                'pi': cache_pi[action-1]})
202 var = np.dot(para['pi'], para['var']+para['mu']**2) - mu_lst[
203 action-1]**2
204 sd_lst[action - 1] = np.sqrt(var)
205 cache_alphas[action-1] = para['alpha']
206 cache_betas[action-1] = para['beta']
207 cache_pi[action-1] = para['pi']
208 loop_count.append(para['loop_count'])
209 print(count_lst)
210 return samples
211
212 def ucb_v_em_beta(bandit, n):
213     return ucb_v_beta(bandit, n, em=True)
214
215 # sample = my_beta_mixture.pull(arm=1, size=5000)
216 # print(my_beta_mixture.parameters(latex=True))
217 #
218 # para_est = em_beta(sample, m=my_beta_mixture.components_num, acc=0.01)

```

```

216 # plot_est_and_true(my_beta_mixture, arm=1, para=para_est, samples=sample
    )
217
218
219 def make_plots_beta(k_lst, components_lst, n_lst, alpha_range_lst,
    beta_range_lst,
220                     rep, folder_name, algorithm_lst, plots_num=1):
221     # function for drawing and saving multiple plots and parameters
222     # k_lst: a list containing the number of arms an environment carries
223     # components_lst: a list containing possible numbers of components
224     # n_lst: a list containing horizons
225     # alpha_range_list: list of the form [[a1, b1], ..., [an, bn]]
        denoting the range of alphas
226     # beta_range_list: list of the form [[a1, b1], ..., [an, bn]]
        denoting the range of betas
227     # rep: how many experiments to repeat under a given environment
228     # folder_name: folder name
229     # algorithm_lst: which algorithms to test on? of the form [[algorithm
        , name],...]
230     # plots_num: repeat all for how many plots?
231     for a_plot in range(plots_num):
232         for k in k_lst:
233             for index in range(len(alpha_range_lst)):
234                 for m in components_lst:
235                     print(f"k = {k}, m = {m}, alpha range is {
                        alpha_range_lst[index]}, "
236                           f"beta range is {beta_range_lst[index]}".center
                            (130, '@'))
237                     my_beta_mixture = MixtureBanditBeta(arm_num=k,
238                                                           alpha_range=
239                                                           alpha_range_lst
240                                                           [index],
241                                                           beta_range=
242                                                           beta_range_lst
243                                                           [index],
244                                                           components_num=m)
245                     regret_comparison(bandit=my_beta_mixture,
246                                       algorithm_lst=algorithm_lst,
247                                       n=n_lst[index], rep=rep, show=False
248                                       , folder_name=folder_name)
249
250 make_plots_beta(k_lst=[2],
251                alpha_range_lst=[[0, 2], [0, 10]],
252                beta_range_lst=[[0, 2], [0, 10]],
253                n_lst=[1000, 1000],
254                components_lst=[3],
255                rep=10,
256                folder_name='Beta_mixture_m_3',
257                algorithm_lst=[[ucb_v_beta, 'ucb-v-beta'],
258                               [ucb_v_em_beta, 'ucb-v-em-beta']])

```

## A.8. regret\_comparison.py

```
1 from Thompson_Sampling import *
2 from UCB_Normal import *
3 from UCB_V import *
4 import os
5 import pandas as pd
6 import time
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9
10
11 def regret_comparison(bandit, algorithm_list, n, rep, folder_name, show:
    bool):
12     # function to compare regrets and save or display one plot,
13     # also saving the parameters to a txt file.
14
15     # bandit: a bandit instance, either gaussian or beta
16     # algorithm_list: [[algorithm, algorithm name], ...]
17     # n: horizon
18     # rep: how many repeated experiments for this bandit problem?
19     # folder_name: folder name
20     # show: boolean, show the plot immediately or save it?
21
22     print(bandit.parameters(latex=False))
23     [optimal_mean, optimal_action] = [bandit.optimal_arm_mean, bandit.
        optimal_arm]
24     optimal_samples = [(i + 1) * optimal_mean for i in range(n)]
25     print(f"The optimal action is the arm {optimal_action}. \n\n")
26
27     time_record = []
28     regret_df = pd.DataFrame({})
29     for [algorithm, algorithm_name] in algorithm_list:
30         print(f" {algorithm_name} ".center(120, '#'))
31         print(bandit.parameters(latex=False))
32         total_time = 0
33
34         for i in range(rep):
35             print('{0}, {1:1d}th iteration'.format(algorithm_name, i + 1)
                )
36
37             tic = time.perf_counter()
38             samples = algorithm(bandit, n)
39             toc = time.perf_counter()
40             total_time += toc - tic
41
42             pseudo_means = [bandit.weighted_mean[x[0] - 1] for x in
                samples]
43             regret = np.array(optimal_samples) - np.cumsum(pseudo_means)
44             regret = pd.DataFrame({'Cumulative Regret': regret,
                'Algorithm / Time': algorithm_name,
                't': np.arange(n) + 1})
45             regret_df = pd.concat([regret_df, regret], axis=0)
```

```

48         print(f"Time: {toc - tic:0.3f} secs, approximately remaining:
49             "
50             f" {total_time / (i + 1) * (rep - i - 1) / 60:0.2f}
               minutes. \n")
51
52     time_record.append(total_time / rep)
53
54     for i in range(len(time_record)):
55         print(f"{algorithm_list[i][1]} on average {time_record[i]:0.3f}
               secs per experiment.")
56         name_with_time = f"{algorithm_list[i][1]} / {time_record[i]:0.3f}
               "
57         regret_df.loc[regret_df['Algorithm / Time'] == algorithm_list[i]
               ][1],
58                     'Algorithm / Time'] = name_with_time
59
60     if bandit.is_beta:
61         [alpha1, alpha2] = bandit.alpha_range
62         [beta1, beta2] = bandit.beta_range
63         name_info = f"k = {bandit.arm_num:d}, m = {bandit.components_num
               :d}, n = {n:d}, r = {rep:d}"
64         name_identify = f"{name_info}, [{alpha1}, {alpha2}], [{beta1}, {
               beta2}]"
65         name_without_suffix = name_identify
66     else:
67         [mean1, mean2] = bandit.mean_range
68         [sd1, sd2] = bandit.sd_range
69         name_info = f"k = {bandit.arm_num:d}, m = {bandit.components_num
               :d}, n = {n:d}, r = {rep:d}"
70         name_identify = f"{name_info}, [{mean1}, {mean2}]"
71         name_without_suffix = name_identify
72
73     print('\n\nSaving parameters.....\n\n')
74
75     try:
76         os.mkdir(os.path.join(os.getcwd(), folder_name))
77         os.mkdir(os.path.join(os.getcwd(), folder_name, 'parameters'))
78         os.mkdir(os.path.join(os.getcwd(), folder_name, 'plots'))
79     except FileExistsError:
80         pass
81
82     my_int = 1
83     while os.path.exists(os.path.join(os.getcwd(), folder_name, '
               parameters', name_identify)):
84         name_identify = f"{name_without_suffix} ({my_int})"
85         my_int += 1
86
87     my_directory_para = os.path.join(os.getcwd(), folder_name, '
               parameters', name_identify)
88     my_directory_png = os.path.join(os.getcwd(), folder_name, 'plots',
               name_identify)
89     f = open(my_directory_para, 'w+')

```

```

90     f.write(bandit.parameters(latex=True))
91     f.close()
92
93     print('Saving plots.....\n\n')
94     sns.set(style="whitegrid", palette="muted", color_codes=True)
95     sns.lineplot(data=regret_df, x="t", y='Cumulative Regret', hue='
Algorithm / Time')
96     plt.rcParams.update({'font.size': 100})
97
98     if bandit.is_beta:
99         plt.title(f"{name_info},   $\alpha_1 < a < \alpha_2$ ,   $\beta_1 <
b < \beta_2$ ",
100                  fontsize=14)
101     else:
102         plt.title(f"{name_info},   $\mu_1 < \mu < \mu_2$ ,   $\sigma_1 < \
sigma < \sigma_2$ ",
103                  fontsize=14)
104     plt.savefig(f"{my_directory_png}.png")
105     print(' DONE! '.center(120, '#')+'\n\n')
106     if show:
107         plt.show()
108     plt.clf()
109
110
111 def make_plots(k_lst, mean_range_lst, n_lst, components_num, rep,
folder_name, algorithm_lst, plots_num=1):
112     # function for drawing and saving multiple plots and parameters
113
114     # k_lst: a list containing the number of arms an environment carries
115     # components_lst: a list containing possible numbers of components
116     # n_lst: a list containing horizons
117     # mean_range_lst: list of the form [[a1, b1], ..., [an, bn]]
denoting the range of means
118     # rep: how many experiments to repeat under a given environment
119     # folder_name: folder name
120     # algorithm_lst: which algorithms to test on? list of the form [[
algorithm, name], ...]
121     # plots_num: repeat all for how many plots?
122
123     for a_plot in range(plots_num):
124         for k in k_lst:
125             for index in range(len(mean_range_lst)):
126                 for m in components_num:
127                     print(f"k = {k}, m = {m}, mean range is {
mean_range_lst[index]".center(130, '@'))
128                     bandit_1 = MixtureBandit(arm_num=k,
components_num=m,
129                                                 mean_range=mean_range_lst[
index],
130                                                 sd_range=[0.5, 1])
131
132                     regret_comparison(bandit_1,
n=n_lst[index],
133                                     rep=rep,

```



```

135         folder_name=folder_name,
136         show=False,
137         algorithm_list=algorithm_lst)
138
139
140 def algo_lst(index):
141     algo = [[ucb_v, 'UCB-V'], # 1
142             [ucb_v_em, 'UCB-V-EM'], # 2
143             [ucb_normal, 'UCB1-Normal'], # 3
144             [ucb_normal_em, 'UCB1-Normal-EM'], # 4
145             [thompson_sampling, 'Thompson Sampling'], # 5
146             [ucb, 'UCB'], # 6
147             [ucb_em, 'UCB-EM']] # 7
148     return [algo[i-1] for i in index]
149
150
151 make_plots(k_lst=[2, 4, 8],
152           mean_range_lst=[[0, 1], [0, 5]],
153           n_lst=[1000, 500],
154           components_num=[2, 3, 4],
155           rep=20,
156           folder_name='Thompson_sampling',
157           algorithm_lst=algo_lst([2, 5]))

```