# Day 5. Thompson Sampling

## NPEX Reinforcement Learning

July 30, 2021

Jaeuk Shin, Minkyu Park

CORE
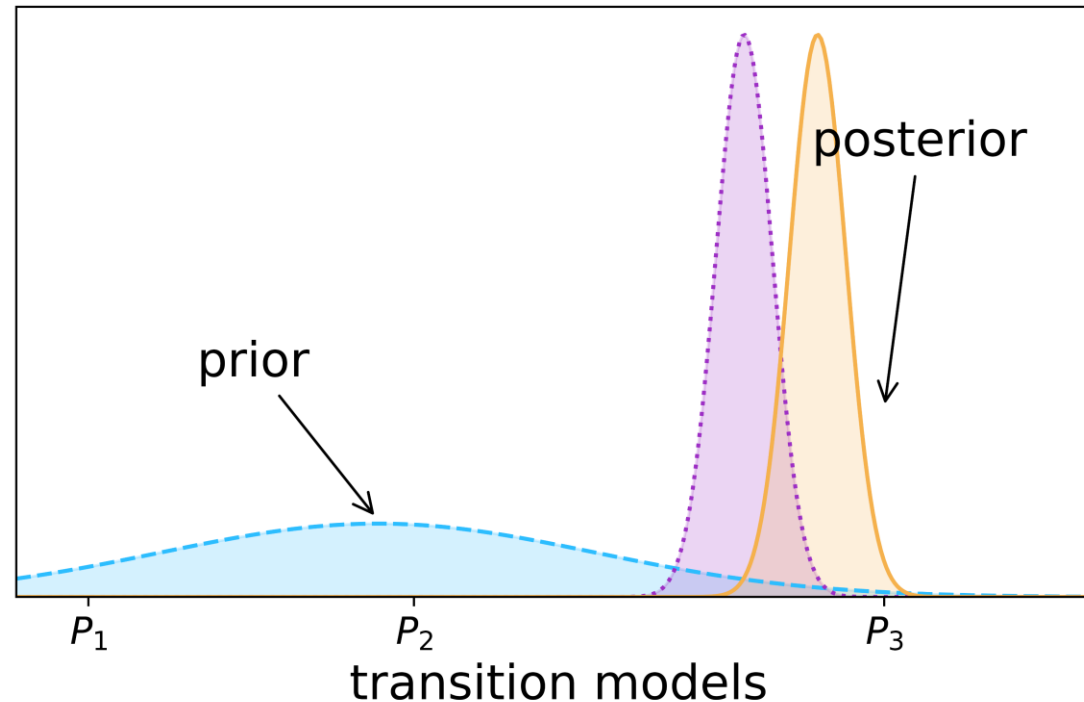Control + Optimization Research Lab

# Thompson Sampling

**Key Idea** : When the underlying dynamics is unknown, model-based approaches always repeat the following steps:

1. Act optimally according to our model of the dynamics.

   - Model Predictive Control
   - **Dynamic Programming** (Finite MDP, Linear Quadratic Control)

2. Observe the result of the action.

3. Update our model based on new information.

   - Least-square Fitting
   - **Bayesian update** $\Longleftarrow$ Spirit of Thompson Sampling

CORE
Control + Optimization Research Lab

# Thompson Sampling

What do we mean by **Bayesian update**?

- Using Thompson sampling, we do not update our model deterministically.

- Instead, whenver new information is given, we update the **distribution of models!**
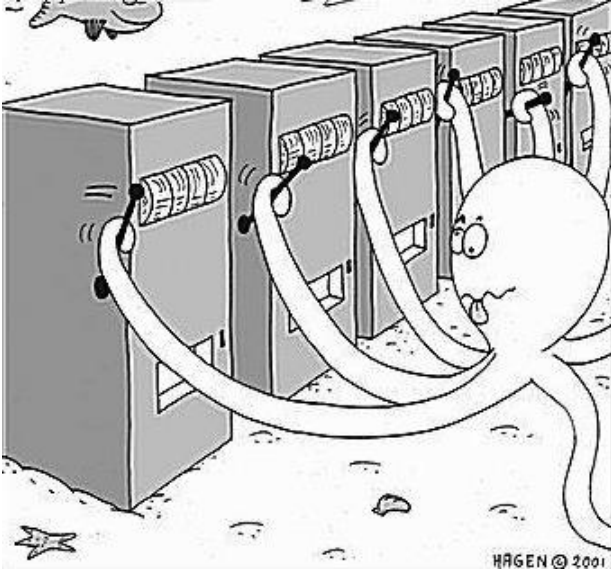
# Thompson Sampling

3 applications of Thompson Sampling in RL:

1.  Multi-armed Bandit Problem

2.  Model-based RL via Thompson Sampling

3.  Learning Linear Quadratic Control vis Thompson Sampling

# Application 1. Multi-armed Bandit Problem

# Multi-Armed Bandit Problem



1-step MDP problem

Framework to handle **exploration-exploitation**

With prior belief $\approx p(r|a)$, build polcy $\pi$

Take action, get reward $(r, a)$ to update $p(r|a)$ via Bayes' rule

Regret : $r(a^*) - r(\pi)$

2 approaches

Use action to maximize $E[r|a]$ - greedy selection

Sample $r \sim p(r|a)$ and choose action - Thompsom sampling

Since we utilize $model \; p(r|a)$, it is model-based RL approach

# Multi-Armed Bandit Problem
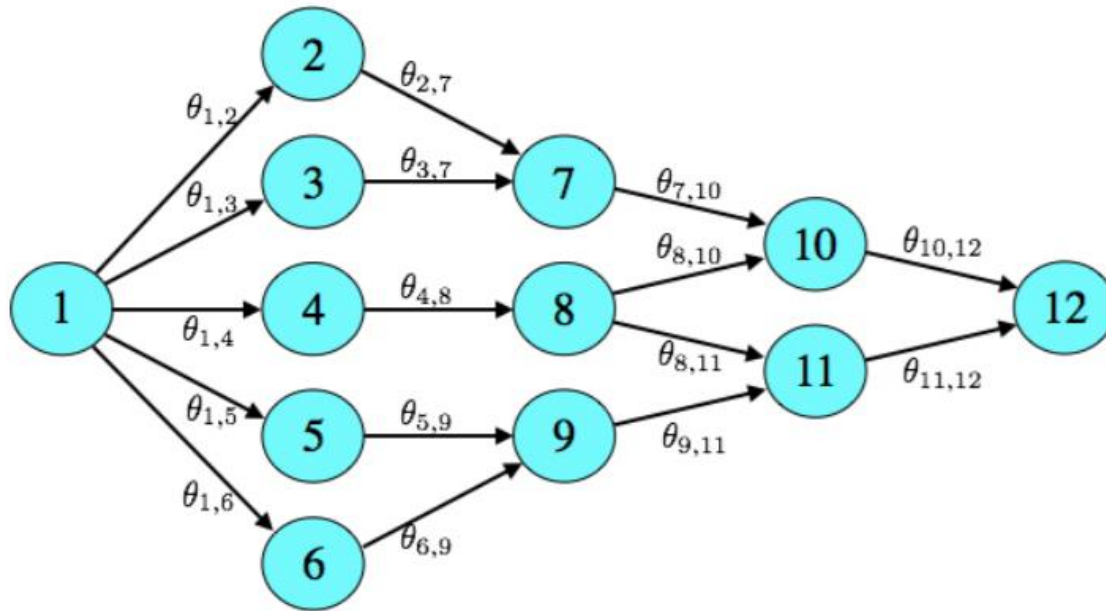
Implementation : travel time optimization



**Figure 1.1:** Shortest path problem.

Table for road number naming

| | | | |
|---|---|---|---|
| 0 | 5 | 10 | 14 |
| 1 | 6 | 11 | 15 |
| 2 | 7 | 12 | |
| 3 | 8 | 13 | |
| 4 | 9 | | |

Time on road $\theta$ follows log-gaussian : $y_t \sim \mathcal{N}(\log(\theta) - \tilde{\sigma}^2/2, \ \tilde{\sigma}^2)$

There are total 16 path (16 actions) available

# Multi-Armed Bandit Problem

Implementation : travel time optimization

$$\theta \sim \mathcal{N}(\mu_e, \sigma_e)$$

$\mu_e$ is $d_e + w$, where $d_e$ is distance of each path and $w$ is random noise

Agent will use prior $\theta \approx \mathcal{N}(d_e, \hat{\sigma}_e)$

When $y_t$ is given, agent can perform Bayes' rule to update its $(\hat{\mu}_e, \hat{\sigma}_e)$

Reward will be given with $-\sum y_t$

# Multi-Armed Bandit Problem

Implementation : travel time optimization

```python
1    def sampling(self):
2        self.theta = []
3        for mu_e, sigma_e in zip(self.mu, self.sigma):
4            theta = np.random.normal(mu_e, sigma_e)
5            self.theta.append(np.exp(theta))
6
7        action_score = []
8        for path in self.path_set:
9            time = 0
10           for e_idx in path:
11               # Thompson sampling samples [y_t] to choose action
12               mean = np.log(self.theta[e_idx]) - 0.5*(sigma_wave**2)
13               std = sigma_wave
14               y_t = np.random.lognormal(mean, std)
15               time += y_t
16           action_score.append(-time)
17
18       return np.argmax(action_score)
```

Each path corresponds to bandit action

Sample $r \sim p(r|a)$ and choose action - Thompsom sampling

CORE
Control + Optimization Research Lab

# Multi-Armed Bandit Problem

Implementation : travel time optimization

```python
1      def greedy(self, eps=0):
2          self.theta = []
3          for mu_e, sigma_e in zip(self.mu, self.sigma):
4              theta = mu_e
5              self.theta.append(np.exp(theta))
6
7          action_score = []
8          for path in self.path_set:
9              time = 0
10             for e_idx in path:
11                 # Greedy poilcy use E[y_t|theta] to choose action
12                 # TODO : calculate E[y_t|theta]
13                 # Hint : for lognormal distribution ( log(m)-0.5*(n**2), n ), m = E[y_t|theta]
14                 # Hint : find y_t at sampling function!
15                 y_t = self.theta[e_idx]
16                 time += y_t
17             action_score.append(-time)
18
19         r = np.random.uniform(0,1)
20         if eps < r:
21             return np.random.randint(0,len(self.path_set))
22         else:
23             return np.argmax(action_score)
```
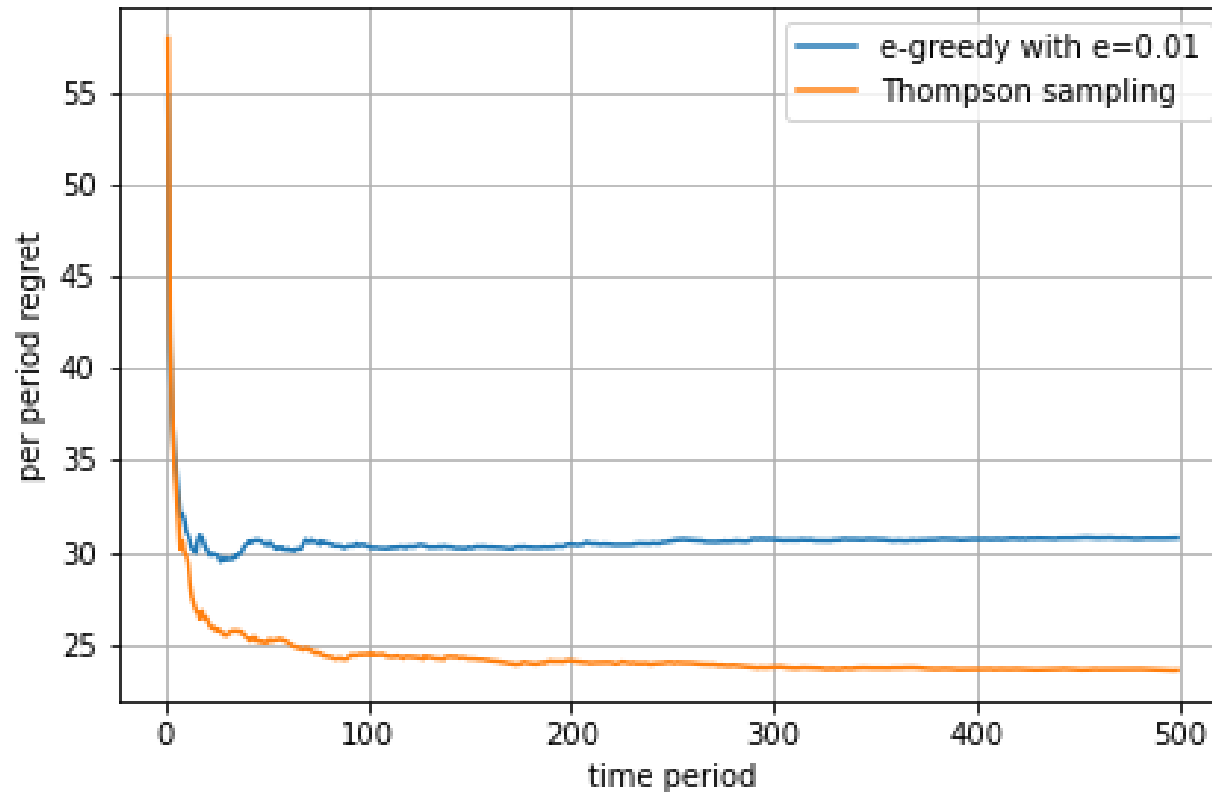
Use mean value of prior distribution, rather than sampling

$\epsilon$-greedy

CORE
Control + Optimization Research Lab

# Multi-Armed Bandit Problem

Using Thompson sampling reduce regret much larger than greedy policy

# Application 2. Model-based Reinforcement Learning via Thompson Sampling

In this application, we assume

- MDP has $n$ states & $m$ actions

- transition probability $P$ : **unknown**

- reward matrix $R$ : known

**Recall** : transition probability : matrix of the form

$$P = \begin{pmatrix} p(s_1|s_1, a_1) & \cdots & p(s_n|s_1, a_1) \\ \vdots & \vdots & \vdots \\ p(s_1|s_n, a_m) & \cdots & p(s_n|s_n, a_m) \end{pmatrix}$$

$\implies$ each row : probability vector of transitions from a single $(s, a)$ pair

# Model-based RL via Thompson Sampling : Implementation

Our strategy is to repeat the followings:

1. Sample $P$ from the prior/posterior distribution (make a guess)

   - Sample each row $(p_1, \cdots, p_n)$ from **Dirichlet distribution**
     - Why? $\implies$ Dirichlet : conjugate prior of Categorical

2. Act optimally according to $P$

   - Use **relative value iteration** (for average reward infinite-horizon MDP)

3. Update distribution using **Bayes rule**

   - Update parameters of each Dirichlet distributions ($n \times m$ distributions in total)

CORE
Control + Optimization Research Lab

# Model-based RL via Thompson Sampling : Implementation

```python
class TSRLAgent:
    def __init__(self, num_states, num_actions, alpha0, R):
        self.num_states = num_states
        self.num_actions = num_actions
        self.alpha = alpha0
        self.R = R
        self.policy = None

    def reset(self):
        P = self.infer_model()
        self.policy = value_iteration(P=P, R=self.R)[1]

    def act(self, state):
        return self.policy[state]

    def update(self, state, action, next_state):
        self.alpha[self.num_states * action + state, next_state] += 1

    def infer_model(self):
        nrows = self.num_states * self.num_actions
        P = np.zeros((nrows, self.num_states))
        for i in range(nrows):
            P[i] = np.random.dirichlet(self.alpha[i])
        return P
```

prior distribution

Sample new model & compute optimal policy via DP

Posterior update : if $(s^t, a^t, s^{t+1})$ is observed, then adjust $\alpha(s^t, a^t)$.

Sample new model!
$p^{(t)}(\cdot|s, a) \sim \mathbf{Dirichlet}\left(\alpha^{(t)}(s, a)\right)$ for each $(s, a)$

CORE
Control + Optimization Research Lab

# Model-based RL via Thompson Sampling : Implementation

```python
def TSRL(env):
    alpha0 = .1 * np.ones((num_states * num_actions, num_states))
    agent = TSRLAgent(num_states=num_states, num_actions=num_actions, alpha0=alpha0, R=env.R_true)
    t = 0
    t_init = 0
    ep_len = 1
    visit_count = np.zeros((num_states, num_actions), dtype=int)
    state = env.reset()
    while True:
        agent.reset()
        visit_count_init = visit_count
        while t <= t_init + ep_len:
            action = agent.act(state)
            next_state, reward = env.step(action)
            visit_count[state, action] += 1
            agent.update(state, action, next_state)
            state = next_state
            if visit_count[state, action] > 2 * visit_count_init[state, action]:
                break
            t += 1
        ep_len = t - t_init
        t_init = t
```

Forget about t_init or ep_len...

Sample new model & compute optimal policy via DP

Act as if the sampled model is correct

Bayesian update

resampling criteria

CORE
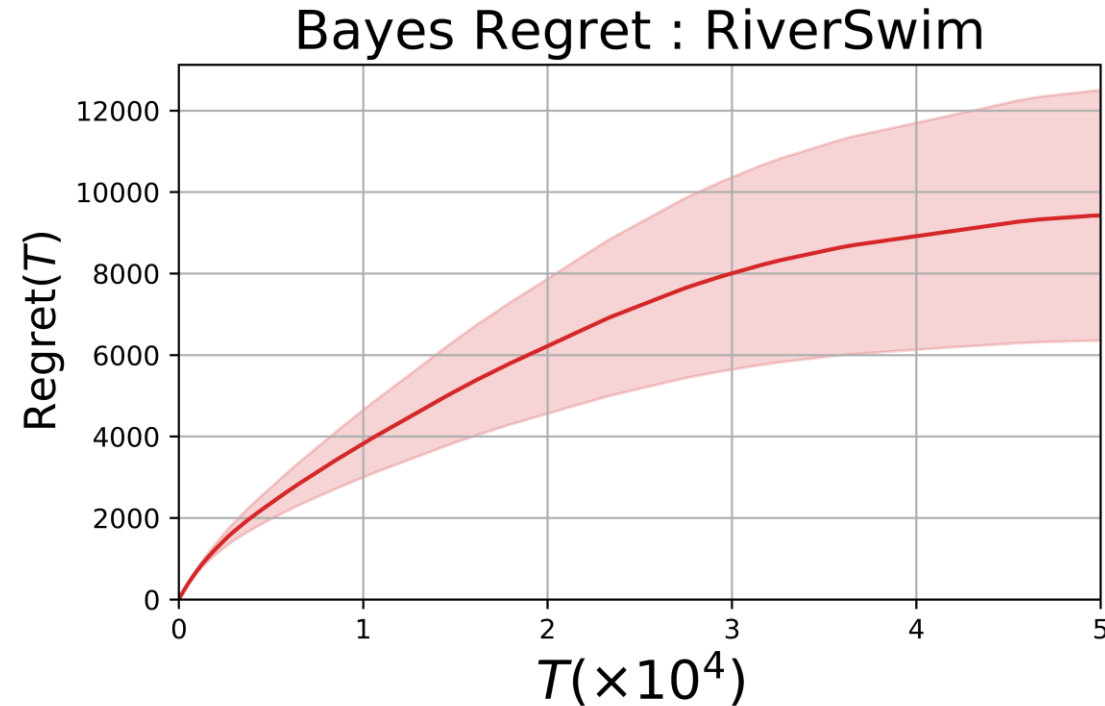Control + Optimization Research Lab

RiverSwim Example

- reward : -0.8 if LEFT at $s_1$, 0 if RIGHT at $s_n$, -1 everywhere else

- LEFT always succeeds, but RIGHT succeeds with probability $p$

# Model-based RL via Thompson Sampling : Implementation

10 states, $p = 0.5$



Bayes Regret : RiverSwim

$$\text{Regret}(T) = O(\sqrt{T})$$

# Application 3. Learning Linear Quadratic Control via Thompson Sampling

Consider the following linear system with stochastic noise:

$$x_{t+1} = Ax_t + Bu_t + w_t, \quad w_t \sim \mathcal{N}(0, I).$$

with a quadratic cost function

$$c(x_t, u_t) = x_t^\top Q x_t + u_t^\top R u_t.$$

Here we assume

- system matrices $A$ & $B$ : **Unknown**

- cost function $c(x, u)$ : known

Why linear system? $\implies$ Many interesting systems can be approximated as a linear one!

# Learning Linear Quadratic Control via Thompson Sampling

The approach is similar to the previous one: we simply repeat

1. Sample $\Theta = \begin{bmatrix} A \\ B \end{bmatrix} \in \mathbb{R}^{(n+m) \times n}$ from the prior/posterior distribution (make a guess)

   - Sample each column from **Gaussian distribution**

2. Act optimally according to $\Theta$

   - apply theory of **Algebraic Riccati equation**

3. Update distribution using **Bayes rule**

   - Update parameters of each Gaussian distributions ($n$ distributions in total)

# Learning Linear Quadratic Control via Thompson Sampling

```python
1   class Distribution:
2       def __init__(self, system: LinearSystem, Q, R, mean_prior, cov_prior):
3           ...
4
5       def sample(self, delta):
6           L = np.linalg.cholesky(self.cov)
7           d, n = self.mean.shape
8           accept = False
9           while not accept:
10              theta = self.mean + L @ np.random.randn(d, n)
11              G = compute_gain(theta, self.Q, self.R)
12              accept = self.system.is_stable(G, delta)
13          return theta
14
15      def update(self, z, x_next):
16          cov_times_z = self.cov @ z
17          denom = 1 + np.dot(z, cov_times_z)
18          self.mean += np.outer(cov_times_z, x_next - z @ self.mean) / denom
19          self.cov -= np.outer(cov_times_z, cov_times_z) / denom
20          return
```

sample $\Theta = \begin{bmatrix} A \\ B \end{bmatrix}$ from Gaussian prior/posterior

Given $(\underbrace{x_t, u_t}_{z_t}, x_{t+1})$, update the posterior

CORE
Control + Optimization Research Lab

# Learning Linear Quadratic Control via Thompson Sampling

```python
def PSRL_LQ(system: LinearSystem, Q, R, mean_prior, cov_prior, delta=0.999):
    distribution = Distribution(system, Q, R, mean_prior, cov_prior)
    t = 0
    t_init = 0
    ep_len = 1
    x = system.reset()
    while True:
        theta = distribution.sample(delta=delta)
        G = compute_gain(theta, Q, R)
        init_cov_sz = distribution.cov_sz
        current_cov_sz = init_cov_sz
        while (t <= t_init + ep_len) and (current_cov_sz >= 0.5 * init_cov_sz):
            u = G @ x
            state_arr.append(x)
            cost += np.dot(x, Q @ x) + np.dot(u, R @ u)
            cost_arr.append(cost)
            z = np.concatenate([x, u])
            x_next = system.step(u)
            distribution.update(z=z, x_next=x_next)
            x = x_next
            t += 1
        ep_len = t - t_init
        t_init = t
```

sample $\Theta = \begin{bmatrix} A \\ B \end{bmatrix}$ from Gaussian prior/posterior

obtain optimal policy $u_t = G(\Theta)x_t$ by solving Algebraic Ricatti eqn.

resampling criteria

Given $(\underbrace{x_t, u_t}_{z_t}, x_{t+1})$, update the posterior

CORE
Control + Optimization Research Lab

# Learning Linear Quadratic Control via Thompson Sampling

Path tracking example

Goal : control the car to follow the laser point of constant velocity

linearized dynamics with sampling time $h = 0.1\text{sec}$:

$$A = \begin{pmatrix} 1 & 0 & -h^2 \\ 0 & 1 & h^2 \\ 0 & 0 & 1 \end{pmatrix}, \quad B = \begin{pmatrix} h/\sqrt{2} & 0 \\ h/\sqrt{2} & 0 \\ 0 & 1 \end{pmatrix}$$

with state vector $x_t = x_{\text{car}} - x_{\text{laser}}$ where $x_{\text{car}} = (x_{\text{car},t}, y_{\text{car},t}, \theta_{\text{car},t})$

CORE
Control + Optimization Research Lab

# Learning Linear Quadratic Control via Thompson Sampling