

# **Stat 260/CS 294-102. Learning in Sequential Decision Problems.**

**Peter Bartlett**

## 1. Thompson sampling

- Bernoulli strategy
- Regret bounds
- Extensions—the flexibility of Bayesian strategies

## Bayesian bandit strategies

Thompson sampling:

- Simple strategy: only requires ability to sample a parameter from the posterior and maximize the expected reward under the sampled parameter.
- Applicable for more complex problems (dependent priors, complex actions, dependent rewards)
- Performs well in practice.
- Strong (frequentist) regret bounds.
- Gittins index has shortcomings:  
only applicable to infinite horizon and independent arm priors;  
complex to compute.

## Bayesian Bernoulli bandits

At time  $t$ , arm  $j$  gives  $X_{j,t} \sim \text{Bernoulli}(\mu_j)$  reward.

Suppose we have fixed, unknown parameters  $\mu_j$ , and we wish to choose a sequence of arms  $I_1, I_2, \dots$ , so as to minimize regret

$$\bar{R}_n = \sum_{t=1}^n (\mu_{j^*} - \mu_{I_t}).$$

## Thompson's idea

Frequentist setting, Bayesian strategy:

- Uniform prior on  $p_j$ .
- Pick action  $j$  with probability that increases with the probability (under posterior distributions) that  $j$  is optimal.

Typically, 'Thompson sampling' refers to *probability matching*. That is, the probability of choosing  $j$  is set to the probability that  $j$  is optimal:

$$\Pr(I_t = j) = \Pr(p_j \text{ is the max}).$$

## Thompson sampling

Given prior  $\pi_{j,0} = \pi_0$ :

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

1. Draw  $p_{1,t}, \dots, p_{k,t}$  from posteriors  $\pi_{j,t-1}$ .
2. Play the arm  $j$  with maximum  $p_{j,t}$ .
3. Observe reward  $X_{j,t}$ .
4. Update posterior:

$$\pi_{j,t}(p) \propto \underbrace{p^{X_{j,t}}(1-p)^{1-X_{j,t}}}_{\text{likelihood}} \underbrace{\pi_{j,t-1}(p)}_{\text{prior}}.$$

Note that  $\Pr(I_t = j) = \Pr(p_{j,t} \text{ is max})$ .

## Thompson sampling

For Beta(1,1) prior:

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

1. Draw  $p_{j,t} \sim \text{Beta}(S_{j,t} + 1, F_{j,t} + 1)$  for  $j = 1, \dots, k$ .
2. Play  $I_t = j$  for  $j$  with maximum  $p_{j,t}$ .
3. Observe reward  $X_{I_t,t}$ .
4. Update posterior:  
Set  $S_{I_t,t+1} = S_{I_t,t} + X_{I_t,t}$ .  
Set  $F_{I_t,t+1} = F_{I_t,t} + 1 - X_{I_t,t}$ .

## Regret of Thompson sampling

**Theorem:** [Agrawal and Goyal]

For every  $\mu_1, \dots, \mu_k$ , there is a constant  $C$  such that for all  $\epsilon > 0$ ,

$$\bar{R}_n \leq (1 + \epsilon) \sum_{j: \Delta_j > 0} \frac{\Delta_j \log n}{d(\mu_j, \mu^*)} + \frac{Ck}{\epsilon^2},$$

where  $d(\mu_j, \mu^*)$  is the KL-divergence between Bernoulli distributions.

## Regret of Thompson sampling: Proof

As always,  $\bar{R}_n = \sum_{j=1}^k \mathbb{E}T_j(n)\Delta_j$ , where  $T_j(n) = \sum_{t=1}^n 1[I_t = j]$ .

For UCB strategies, we were able to bound regret by showing that the upper bounds are (whp) valid, and that, after enough mistakes, they are sufficiently tight that subsequent mistakes are unlikely. Here, we don't have bounds, we're just sampling from distributions that get more concentrated as we sample more.



## Regret of Thompson sampling: Proof

Fix a suboptimal arm  $j \neq j^*$ . Split according to  $\hat{\mu}_j(t-1)$  and  $p_{j,t}$ :

$$\begin{aligned}\mathbb{E}T_j(n) &= \sum_{t=1}^n \Pr(I_t = j) \\ &= \sum_{t=1}^n [\Pr(I_t = j, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} \leq y_j) \\ &\quad + \Pr(I_t = j, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} > y_j) \\ &\quad + \Pr(I_t = j, \hat{\mu}_j(t-1) > x_j)],\end{aligned}$$

where we choose  $\mu_j < x_j < y_j < \mu^*$  as follows:

## Regret of Thompson sampling: Proof

$$\begin{aligned} \mu_j < x_j < \mu^* & \quad \text{s.t.} & \quad d(x_j, \mu^*) &= \frac{d(\mu_j, \mu^*)}{1 + \epsilon} \\ x_j < y_j < \mu^* & \quad \text{s.t.} & \quad d(x_j, y_j) &= \frac{d(x_j, \mu^*)}{1 + \epsilon} \\ & & &= \frac{d(\mu_j, \mu^*)}{(1 + \epsilon)^2}. \end{aligned}$$

This ensures that the relevant divergences are only a constant factor different from  $d(\mu_j, \mu^*)$ .

## Regret of Thompson sampling: Proof

Consider those three probabilities:

1.  $\Pr(I_t = j, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} \leq y_j)$   
counts the times when both the empirical mean  $\hat{\mu}_j(t-1)$  and the sampled  $p_{j,t}$  are not too far above their expectations. For this to be small, we need to be sure that  $j^*$  is pulled a lot.
2.  $\Pr(I_t = j, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} > y_j)$   
counts the times when  $p_{j,t}$  is far above its expectation. As  $T_j$  grows, the distribution of  $p_{j,t}$  gets more concentrated, so this sum is  $O(\log n/d(x_j, y_j))$ .
3.  $\Pr(I_t = j, \hat{\mu}_j(t-1) > x_j)$  counts the times when the empirical mean  $\hat{\mu}_j(t-1)$  is far above its expectation. Its sum is  $O(1)$ .

## Regret of Thompson sampling: Proof

For the 3rd probability: if  $\tau_{j,k}$  is the  $k$ th time when  $I_t = j$ , then

$$\begin{aligned} \sum_{t=1}^n \Pr(I_t = j, \hat{\mu}_j(t-1) > x_j) &\leq 1 + \mathbb{E} \sum_{k=1}^{n-1} \mathbb{1}[\hat{\mu}_j(\tau_{j,k}) > x_j] \\ &\leq 1 + \sum_{k=1}^{n-1} \exp(-kd(x_j, \mu_j)) \\ &\leq 1 + \frac{1}{d(x_j, \mu_j)} = O(1). \end{aligned}$$

## Regret of Thompson sampling: Proof

For the 2nd probability:

$$\begin{aligned} & \sum_{t=1}^n \Pr (I_t = j, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} > y_j) \\ & \leq \sum_{t=1}^n \Pr (I_t = j, p_{j,t} > y_j \mid \hat{\mu}_j(t-1) \leq x_j) \\ & \leq m + \mathbb{E} \sum_{t=\tau_{j,m}+1}^n \Pr (I_t = j, p_{j,t} > y_j \mid \hat{\mu}_j(t-1) \leq x_j, \mathcal{F}_{t-1}), \end{aligned}$$

for any choice of  $m$ . (Recall that  $\tau_{j,k}$  is the  $k$ th time when  $I_t = j$  and  $\mathcal{F}_{t-1}$  is everything up to  $t-1$ .)

## Regret of Thompson sampling: Proof

$$\begin{aligned} & \sum_{t=1}^n \Pr (I_t = j, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} > y_j) \\ & \leq m + \mathbb{E} \sum_{t=\tau_{j,m}+1}^n \exp(-T_j(t-1)d(x_j, y_j)) \\ & \leq \frac{\log n}{d(x_j, y_j)} + \mathbb{E} \sum_{t=\tau_{j,m}+1}^n \frac{1}{n} \\ & \leq \frac{\log n}{d(x_j, y_j)} + 1, \end{aligned}$$

where we choose  $m = \log n / d(x_j, y_j)$  so that

$$\exp(-T_j(t-1)d(x_j, y_j)) \leq \frac{1}{n}.$$

## Regret of Thompson sampling: Proof

For the 1st probability, define  $\alpha_{j,t} = \Pr(p_{j^*,t} > y_j | \mathcal{F}_{t-1})$ . Then

$$\begin{aligned}
 & \sum_{t=1}^n \Pr(I_t = j, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} \leq y_j) \\
 (*) & \leq \sum_{t=1}^n \mathbb{E} \left[ \frac{1 - \alpha_{j,t}}{\alpha_{j,t}} \mathbb{1}[I_t = j^*, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} \leq y_j] \right] \\
 & \leq \sum_{k=0}^{n-1} \mathbb{E} \left[ \frac{1 - \alpha_{j, \tau_{j^*,k}+1}}{\alpha_{j, \tau_{j^*,k}+1}} \sum_{t=\tau_{j^*,k}+1}^{\tau_{j^*,k+1}} \mathbb{1}[I_t = j^*] \right] \\
 & \leq \sum_{k=0}^{n-1} \mathbb{E} \left[ \frac{1}{\alpha_{j, \tau_{j^*,k}+1}} - 1 \right] = \dots = O(1).
 \end{aligned}$$

(Recall  $\tau_{j^*,k}$  is the  $k$ th time when  $I_t = j^*$ .)

## Regret of Thompson sampling: Proof

(\*) Crucial fact: For a bad arm, when its sample average and sampled parameter aren't too large, the probability it's chosen is bounded linearly by the probability that the best arm is chosen.

**Lemma:** For suboptimal  $j$ ,

$$\begin{aligned} & \Pr (I_t = j, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} \leq y_j | \mathcal{F}_{t-1}) \\ & \leq \frac{1 - \alpha_{j,t}}{\alpha_{j,t}} \Pr (I_t = j^*, \hat{\mu}_j(t-1) \leq x_j, p_{j,t} \leq y_j | \mathcal{F}_{t-1}), \end{aligned}$$

where  $\alpha_{j,t} = \Pr (p_{j^*,t} > y_j | \mathcal{F}_{t-1})$ .



## Regret of Thompson sampling: Proof

First, notice that, conditioned on  $\mathcal{F}_{t-1}$ ,  $\hat{\mu}_j(t-1)$  is determined. So assume  $\hat{\mu}_j(t-1) \leq x_j$ .

Define the event  $W_j(t) = \{\forall j' \neq j^*, p_{j,t} \geq p_{j',t}\}$  ( $j$  is the “suboptimal winner”). Then

$$\begin{aligned} & \Pr(I_t = j^* | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \\ & \geq \Pr(I_t = j^*, W_j(t) | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \\ & = \Pr(W_j(t) | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \Pr(I_t = j^* | W_j(t), p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \\ & \geq \Pr(W_j(t) | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \Pr(p_{j^*,t} > y_j | W_j(t), p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \\ & \geq \Pr(W_j(t) | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \Pr(p_{j^*,t} > y_j | \mathcal{F}_{t-1}) \\ & = \Pr(W_j(t) | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \alpha_{j,t}. \end{aligned}$$

## Regret of Thompson sampling: Proof

On the other side:

$$\begin{aligned} & \Pr (I_t = j | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \\ & \leq \Pr (W_j(t), p_{j^*,t} \leq y_j | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \\ & = \Pr (W_j(t) | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) \Pr (p_{j^*,t} \leq y_j | \mathcal{F}_{t-1}) \\ & = \Pr (W_j(t) | p_{j,t} \leq y_j, \mathcal{F}_{t-1}) (1 - \alpha_{j,t}). \end{aligned}$$

## Thompson sampling for bounded rewards

For  $X_{j,t} \in [0, 1]$ , can use an identical approach. But the only analysis applies to Bernoulli-sampled rewards (which introduce a little extra variance):

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

1. Draw  $p_{j,t} \sim \text{Beta}(S_{j,t} + 1, F_{j,t} + 1)$  for  $j = 1, \dots, k$ .
2. Play  $I_t = j$  for  $j$  with maximum  $p_{j,t}$ .
3. Observe reward  $X_{I_t,t}$ .
4. Sample  $R_t \sim \text{Bernoulli}(X_{I_t,t})$ .
5. Update posterior:  
Set  $S_{I_t,t+1} = S_{I_t,t} + R_t$ .  
Set  $F_{I_t,t+1} = F_{I_t,t} + 1 - R_t$ .

## Thompson sampling

One of the advantages of a Bayesian strategy (even in a frequentist setting, where the parameters are fixed but unknown) is that it is easy to generalize it to settings with rewards and (other) outcomes, where the distributions depend on parameters  $\theta$ .

This allows:

- Complex actions (e.g., choose four advertisements on a web page; multi-commodity flow with random capacities).
- Limited observations (e.g., only see rewards for complex actions, not for simple ones).
- Dependent rewards.

## Thompson sampling

Given parameter space  $\Theta$ , action space  $\mathcal{A}$ , outcome space  $\mathcal{Y}$ , prior  $\pi_0$  on  $\Theta$ , likelihood  $\ell(y; a, \theta)$ , and reward distribution:

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

1. Draw  $\theta_t$  from posterior  $\pi_{t-1}$ .
2. Play action  $a_t$  that maximizes  $\mathbb{E}_{\theta_t} R_a$ .
3. Observe outcome  $Y_t$ .
4. Update posterior:

$$\pi_t(\theta) \propto \ell(Y_t; a_t, \theta) \pi_{t-1}(\theta).$$