

In this section, we use a simple example to show how Thompson Sampling (TS) works.

Example. Beta-Bernoulli Bandit

The scenario is again the Multi-arm Bandit machine.

The machine has K arms. An action $k \in \{1, 2, 3, \dots, K\}$ means pulling the k^{th} arm. Arm k produces a reward of one with probability θ_k and a reward of zero with probability $1 - \theta_k$. The player knows that Bernoulli is a good model, but $\theta_1, \theta_2, \dots, \theta_K$ are unknown but fixed (unchanged)

In round i , when action $X_i \in [K]$ is applied, a reward $r_i \in \{0, 1\}$ is generated with $P(r_i = 1 | X_i, \theta) = \theta_{X_i}$.

After observing r_i , the player updates its estimation of θ .

Recall that, the conjugate prior of Bernoulli distribution is Beta-distribution. So, we use K Beta distributions to model $P(\theta_i)$.

$\theta_i \sim \text{Beta}(\alpha_i, \beta_i)$ which has pdf.

$$p(\theta_i) = \frac{\Gamma(\alpha_i + \beta_i)}{\Gamma(\alpha_i)\Gamma(\beta_i)} \theta_i^{\alpha_i - 1} (1 - \theta_i)^{\beta_i - 1}$$

If action $X_t \in [K]$ is taken, then

$$\begin{aligned} \alpha_i, \beta_i &\longrightarrow \alpha_i, \beta_i && \text{if } X_t \neq i \\ &\longrightarrow \alpha_i + r_t, \beta_i + (1 - r_t) && \text{if } X_t = i, \end{aligned}$$

OK, then we present Greedy Algorithm and TS to show how they are different.

Greedy Algorithm:

for round $t = 1, 2, \dots$ do

Get the MAP estimation of θ by

$$\hat{\theta}_i = d_i / (d_i + \beta_i)$$

Then the action is chosen by taking

$$X_t = \operatorname{argmax}_i \hat{\theta}_i$$

Apply X_t , then update (d_i, β_i) accordingly.

TS Algorithm.

for round $t = 1, 2, \dots$ do

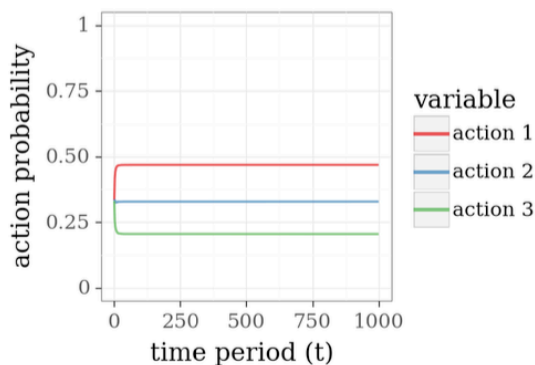
Sample each Beta distribution

y_i is a random sample from $\text{Beta}(d_i, \beta_i)$

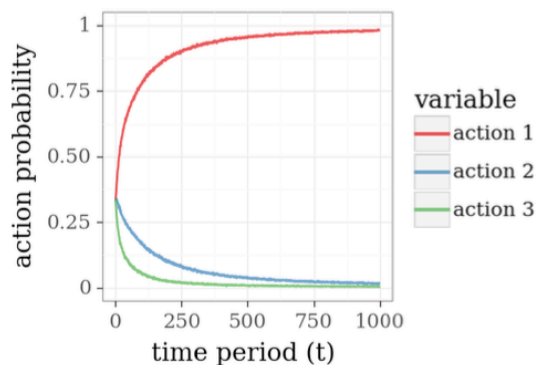
Then the action is chosen by taking

$$X_t = \operatorname{argmax}_i y_i$$

Apply X_t , then update (d_i, β_i) accordingly.



(a) greedy algorithm



(b) Thompson sampling

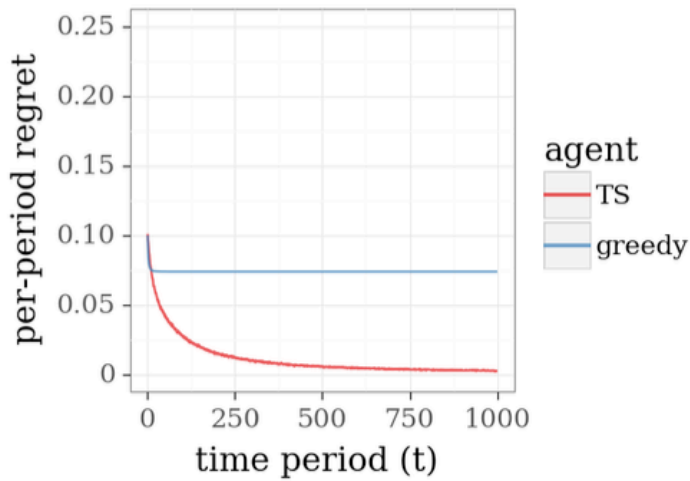
Figure 3.1: Probability that the greedy algorithm and Thompson sampling selects an action.

From the simulation result we can see that Greedy Algorithm may not converge to Action 1, but TS algorithm can.

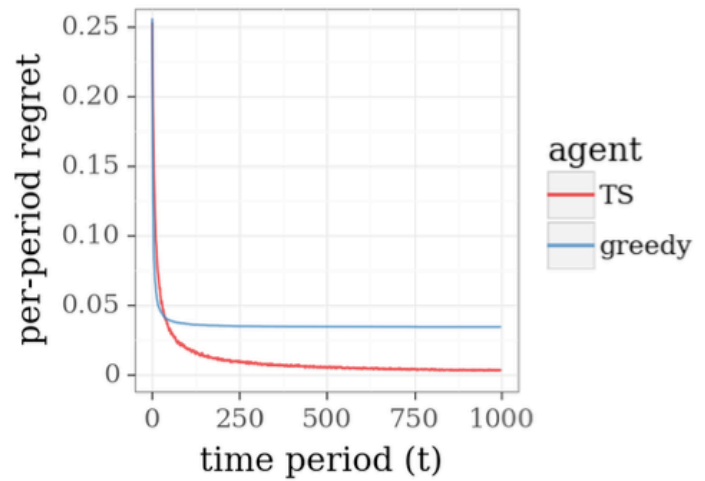
Regret: To evaluate the performance of algorithm, we use

$$\text{regret} \triangleq \sum_t (\max_i r(x_i)) - r_t$$

or $(\max_i r(x_i)) - r_t$ for simultaneous regret.



(a) $\theta = (0.9, 0.8, 0.7)$



(b) average over random θ

Figure 3.2: Regret from applying greedy and Thompson sampling algorithms to the three-armed Bernoulli bandit.