

## Thompson Sampling

Source: Bandit Algorithms, T. Lattimore and C. Szepesvari, Cambridge Univ. Press, 2019 (to appear); Chapter 36.

Thompson sampling is a Bayesian bandit algorithm for a  $k$ -armed bandit.

Setting:  $(\mathcal{E}, \mathcal{G}, \mathcal{Q}, P)$

↑ reward  
parameter space;  
typically a  $k$ -length mean vector  
prior on  $v \in \mathcal{E}$

Example:  $k$ -armed bandit, with Gaussian rewards and Gaussian prior.  $\mathcal{E} = \mathbb{R}^k$ ,  $\mathcal{G} = \sigma(\mathbb{R}^k) =$  Borel- $\sigma$ -algebra over  $\mathbb{R}^k$ ,  $\mathcal{Q}$  a prior over  $v \in \mathcal{E}$ :  $v = (v_1, v_2, \dots, v_k)^T$ .

The interpretation is:  $v_k$  is the (unknown) mean reward for arm  $k$ . This has a generative model:  $v \sim \mathcal{E}_N^k(\cdot)$ , i.e.,  $\mathcal{Q} = N(0, I)$ ;  $P = (P_{kv}, k=1, 2, \dots, k, v \in \mathcal{E})$ ,  
 $\hookrightarrow$  prior is '0' mean, unit variance for all arms

i.e.,  $P_{kv} = N(v_k, 1)$ . An arm pull of the  $k^{th}$  arm generates a reward r.v.  $X$ , where  $X$  has a conditional dist. (given  $v$ ) of  $X \sim P_{kv}$ .

## Thompson Sampling Algorithm:

- ① Given samples  $x_1, x_2, \dots, x_{t-1}$ , update the prior dist. on the environment, i.e., maintain posterior dist:  $Q(\cdot | x_1, \dots, x_{t-1})$ , using Baye's rule starting from the prior dist.  $Q(\cdot)$ .
- ② Draw a sample from this posterior dist: Sample an instance  $v_t = (v_{1,t}, v_{2,t}, \dots, v_{k,t})^\top \sim Q(\cdot | x_1, \dots, x_{t-1})$
- ③ Play arm using the sample  $v_t = (v_{1,t}, \dots, v_{k,t})^\top$ :

$$A_t = \underset{1 \leq j \leq k}{\operatorname{argmax}} v_{j,t}$$

and break ties as the genie policy (which knows the environment instance).

-----

Note 1: The textbook uses a more general notation. Specifically, the environment  $v \in \mathcal{E}$  could include additional parameters, and the mean of the  $j^{\text{th}}$  arm under environment  $v$  is denoted by  $\mu_j(v)$ . In these notes, we are instead simply using  $v_j$  or  $v_{j,t}$  if we need to index both.

an arm  $j$  and a time  $t \in \{1, 2, \dots, n\}$ .

Note 2: We have a doubly stochastic process here. Specifically,  
 $v_t \sim Q(\cdot | x_1, \dots, x_{t-1})$  is a  $k$ -dim. random vector,  
and the reward  $x_t$  is a r.v.; given that  
 $\{A_t = j\} \cap \{v_t = v_t\}, x_t | (j, v_t) \sim N(v_{j,t}, 1)$ .

the r.v.  $\downarrow$  corresp. to the environment at time  $t$

Note 3: A key property of Thompson sampling is that  
the conditional prob. of playing arm  $j$  under Thompson  
Sampling is the same as that of the optimal (genie)  
policy. Formally,

$$\begin{aligned} & P(A_t^* = j | A_1, x_1, \dots, A_{t-1}, x_{t-1}) \xrightarrow{\text{genie}} \textcircled{1} \\ & = P(A_t = j | A_1, x_1, \dots, A_{t-1}, x_{t-1}) \end{aligned}$$

\* This does not mean  $x_{\{A_t^* = j\}} = x_{\{A_t = j\}}$ !

\* Observe that the above does not have any conditioning  
w.r.t. the environment sample  $\triangleright$  that was chosen

at  $t=0$ . This environment  $v$  remains invariant for any  $t=1, 2, \dots, n$  along this sample path. The genie knows this, and thus  $\{A_t^* = j\}$  is a function of  $v$  (specifically,  $A_t^* = \operatorname{argmax}_{i \in [k]} v_i$ ).

Thus to compute  $P(A_t^* = j | A_1, X_1, \dots, A_{t-1}, X_{t-1})$ , we further need to condition and uncondition on the environment. Also finally note that the above property holds only for the marginal dist. of the arm pull.

PF of ①:  $S_{t-1} = \sigma \{A_1, X_1, \dots, A_{t-1}, X_{t-1}\}$ .

$$P(A_t^* = j | S_{t-1}) = \int_{v \in \Sigma} P(A_t^* = j | S_{t-1}, v) dP(v | S_{t-1})$$

↳ conditioning on the true environment

$$\text{But } P(v \in E | S_{t-1}) = Q(v \in E | S_{t-1})$$

This follows because the true environment  $v \sim Q(\cdot)$ , and the update given  $S_{t-1}$  is the posterior as above.

(A)

$$\Rightarrow P(A_t^* = j | S_{t-1}) = \int_{v \in \Sigma} P(A_t^* = j | S_{t-1}, v) \cdot dQ(v | S_{t-1})$$

(for simplicity, we have assumed unique max.; else use tie-breaking rule)

$$S \{ v_j = \max_{i \in [k]} v_i \}.$$

$$P(A_t=j | S_{t-1}) = \int_{v_t \in \mathcal{E}} P(A_t=j | S_{t-1}, v_t) dP(v_t | S_{t-1})$$

↳ conditioning on  
 the sample drawn in Thompson  
 sampling algorithm.

But by defn of Thompson sampling,

$$v_t \sim Q(\cdot | S_{t-1})$$

$$\Rightarrow P(v_t | S_{t-1}) = Q(v_t | S_{t-1})$$

(B)

$$\therefore P(A_t=j | S_{t-1}) = \int_{v_t \in \mathcal{E}} P(A_t=j | S_{t-1}, v_t) \cdot dQ(v_t | S_{t-1})$$

$$\delta \{ v_{j,t} = \max_{i \in \mathcal{S}(k)} v_{i,t} \}$$

Now interpreting  $v$  and  $v_t$  as dummy variables in (A), (B),  
 we observe these expressions are the same.

(T)

### Analyzing Bayesian Regret for Thompson Sampling

$$BR_n(\pi, Q) = E_{v \sim Q, P_\pi, R} \left[ \sum_{t=1}^n (X_{\pi^*} - X_{\pi_t}) \right]$$

environment ↪ ↪ arm reward conditioned on  $v$   
 randomness introduced by algorithm

$$= E_{\nu \sim Q, R} \left[ n v_{A^*} - \sum_{t=1}^n v_{A_t} \right]$$

$$= E_{\nu \sim Q} \left[ \sum_{t=1}^n E_R \left[ v_{A^*} - v_{A_t} \mid \mathcal{F}_{t-1} \right] \right]$$

Theorem (36.1 in text):  $(\mathcal{E}, \mathcal{G}, Q, P)$  a Bayesian bandit, with  $\{P_{v,j}, j=1, 2, \dots, k; v \in \mathcal{E}\}$  uniformly 1-subGaussian, with  $v_i \in [0, 1] \forall i, v$ . Then, Thompson sampling satisfies:

$$BR_n(\pi, Q) \leq C \sqrt{k n \ln(n)}, \quad C > 0.$$

Proof: For  $j = 1, 2, \dots, k$ , and  $t = 1, 2, \dots, n$ , let

$$V_t(j) = \begin{cases} \hat{v}_j(t-1) + \sqrt{\frac{2 \ln(1/\delta)}{\max\{1, T_j(t-1)\}}} & \text{if } T_j(t-1) > 0 \\ 0 & \text{otherwise} \end{cases}$$

(Note:  $\hat{v}_j(\cdot) = 0$  if  $T_j(\cdot) = 0$ ).

↳ clip and project into  $[0, 1]$ .

$\hat{v}_j(t-1)$  = empirical estimate of  $v_j$  = true value of mean in environment  $v \in \mathcal{E}$ .

$$\mathcal{B} = \left\{ \begin{array}{l} \text{if } j \in \{1, \dots, k\} \text{ and } t \in \{1, 2, \dots, n\}, \\ |\hat{v}_j(t-1) - v_j| < \sqrt{\frac{2 \ln(1/\delta)}{\max\{1, T_j(t-1)\}}} \end{array} \right\}$$

Now, recall  $P_{ij} \in \Delta\text{-sub}G$ . Using an union bound over number of samples  $T_j(t-1) \in \{1, 2, \dots, n\}$  and union bound over arms  $j=1, 2, \dots, k$  (very similar to the corresponding step in the UCB algorithm), we have

$$P(B^c) \leq nkS \quad \longrightarrow \textcircled{2}$$

Next, recall:  $\xrightarrow{\textcircled{3}}$

$$BR_n(\pi, Q) = \mathbb{E}_{\pi, Q} \left[ \sum_{t=1}^n \mathbb{E}_R \left[ v_{A^*} - v_{A_t} \mid \mathcal{F}_{t-1} \right] \right]$$

environment  $\longleftarrow$  algo.  $\dashrightarrow$

Thus,  $\xrightarrow{\textcircled{3}} = \sigma \{A_1, X_1, \dots, A_{t-1}, X_{t-1}\}$ .

$$\mathbb{E}_R \left[ v_{A^*} - v_{A_t} \mid \mathcal{F}_{t-1} \right] = \mathbb{E}_R \left[ v_{A^*} - V_t(A_t) + V_t(A_t) \mid \mathcal{F}_{t-1} \right] + v_{A_t}$$

$$\text{Now } \mathbb{E}_R \left[ V_t(A_t) \mid \mathcal{F}_{t-1} \right] = \mathbb{E}_R \left[ V_t(A^*) \mid \mathcal{F}_{t-1} \right].$$

This follows because  $P(A_t=j \mid \mathcal{F}_{t-1}) = P(A^*=j \mid \mathcal{F}_{t-1})$ , which was discussed above in Note 3,  $\textcircled{1}$ .

$$\therefore \mathbb{E}_R \left[ v_{A^*} - v_{A_t} \mid \mathcal{F}_{t-1} \right] = \mathbb{E}_R \left[ v_{A^*} - V_t(A^*) + V_t(A_t) - v_{A_t} \mid \mathcal{F}_{t-1} \right]$$

$\hookrightarrow \textcircled{4}$ .

Substituting ④ in ③, we have:

$$BR_n(\pi, Q) = \mathbb{E}_{\pi, Q} \left[ \sum_{t=1}^n \mathbb{E}_R \left[ v_{A^*} - U_t(A^*) + U_t(A_t) - v_{A_t} \mid \mathcal{F}_{t-1} \right] \right]$$

iterated expectations

$$= \mathbb{E}_{\pi, Q, R} \left[ \sum_{t=1}^n (v_{A^*} - U_t(A^*)) + \sum_{t=1}^n (U_t(A_t) - v_{A_t}) \right]$$

$$= \mathbb{E}_{\pi, Q, R} \left[ \begin{array}{l} \sum_{t=1}^n (v_{A^*} - U_t(A^*) + U_t(A_t) - v_{A_t}) \chi_{B^c} \\ + \left( \sum_{t=1}^n (v_{A^*} - U_t(A^*)) \right) \chi_B + \\ \left( \sum_{t=1}^n (U_t(A_t) - v_{A_t}) \right) \chi_B \end{array} \right]$$

$$\textcircled{I} \leq 2n \quad (\text{because all terms } \in [0, 1]).$$

$$\textcircled{II} \leq 0 \quad \text{because } B \Rightarrow v_{A^*} \leq \hat{v}_{A^*} + \sqrt{\frac{2 \ln(1/\delta)}{\max\{1, T_j(n)\}}}$$

$$\text{Now } U_t(A^*) = \left[ \hat{v}_{A^*} + \sqrt{\frac{2 \ln(1/\delta)}{\max\{1, T_j(n)\}}} \right]_0^1$$

If term inside  $> 1$ , immediately  
true that  $\textcircled{II} \leq 0$ .

$$\text{If term inside } < 0, U_t(A^*) \geq \hat{v}_{A^*} + \sqrt{\frac{2 \ln(1/\delta)}{\max\{1, T_j(n)\}}}$$

$$\text{else, } U_t(A^*) = \hat{v}_{A^*} + \sqrt{\frac{2 \ln(1/\delta)}{\max\{1, T_j(n)\}}}$$

In either case,  $B \Rightarrow v_{A^*} - U_t(A^*) \leq 0$ .

$$\text{III} = \mathbb{E}_{\nu_{nQ,R}} \left[ \sum_{t=1}^n (U_t(A_t) - v_{A_t}) X_B \right]$$

$$= \mathbb{E}_{\nu_{nQ,R}} \left[ \sum_{t=1}^n \sum_{j=1}^k \chi_{\{A_t=j\}} \mathbb{E} (U_t(j) - v_j) X_B \right]$$

$$\leq \mathbb{E}_{\nu_{nQ,R}} \left[ \sum_{j=1}^k \sum_{t=1}^n \chi_{\{A_t=j\}} \sqrt{\frac{2 \ln(1/\delta)}{\max\{1, T_j(n)\}}} X_B \right]$$

$$\leq \mathbb{E}_{\nu_{nQ,R}} \left[ \sum_{j=1}^k \sum_{t=1}^n \chi_{\{A_t=j\}} \sqrt{\frac{2 \ln(1/\delta)}{\max\{1, T_j(n)\}}} \right]$$

$$\leq \mathbb{E}_{\nu_{nQ,R}} \left[ \sum_{j=1}^k \int_0^{T_j(n)} \sqrt{\frac{2 \ln(1/\delta)}{s}} ds \right]$$

$$= \mathbb{E}_{\nu_{nQ,R}} \left[ \sum_{j=1}^k \sqrt{4 \ln(1/\delta) T_j(n)} \right]$$

$\sum_j T_j(n) = n$

$$\leq \mathbb{E}_{\nu_{nQ,R}} \left[ \sqrt{4 \ln(1/\delta) nk} \right]$$

deterministic  $\Rightarrow$  drop  $E[\cdot]$ .

Choose  $\delta = 1/n^2$ ,  
and result follows.



Frequentist (environment dependent) analysis of Thompson Sampling:

Algorithm:  $k$ -arm setting

Initial:  $\{F_1(1), F_2(1), \dots, F_k(1)\}$  prior CDF for the  $k$  arms.

time = 1

arm 1

arm = k

for each  $t=1, 2, \dots, n$ :

① Sample possible environment: Choose  $\theta_i(t) \sim F_i(t)$ ,  $i=1, 2, \dots, k$

② Action:  $A_t = \arg \max_i \theta_i(t)$ .

③ Observe / Update: Observe reward  $X_t$  (corresponding to arm  $A_t$ ).

Update  $F_{A_t}(t+1) = \text{Update}(F_{A_t}(t), A_t, X_t)$

$F_i(t+1)(x)$

$= P(X_i \leq x | \mathcal{E}_t)$   $F_i(t+1) = F_i(t) \quad i \neq A_t$

Bayesian posterior update for TS

Notation:  $F_{i, \text{arm}}(\cdot) = \text{Posterior CDF of arm } i \text{ using }$

# of samples

S samples. (We will use this in an appropriate union bound as in UCB proof to switch from  $\{T_i(t)\}, t=1, 2, \dots, n\}$  to samples  $\{S=1, 2, \dots, n\}$ .)

Theorem (3b.2 in text): for arm  $i$ , and  $\varepsilon \in \mathbb{R}$ ,

$$E[T_i(n)] \leq 1 + E\left[\sum_{s=0}^{n-1} \left(\frac{1}{G_{is}} - 1\right)\right]$$

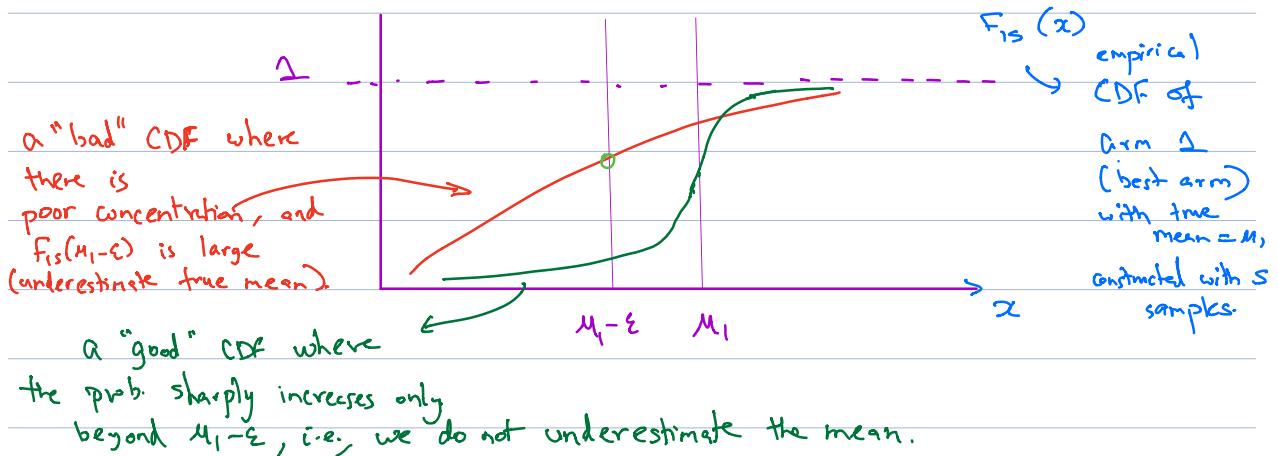
expectation only over randomness of arm reward and algorithm; NOT over environment

$$+ E\left[\sum_{s=0}^{n-1} \chi_{\{G_{is} \geq 1/n\}}\right]$$

where  $G_{is} = 1 - F_{is}(M_i - \varepsilon)$ .

Intuition: (assume  $\varepsilon > 0$ , a small true constant)

$$\textcircled{1} \quad \left(\frac{1}{G_{is}} - 1\right) = \frac{1 - G_{is}}{G_{is}} = \left(\frac{F_{is}(M_i - \varepsilon)}{1 - F_{is}(M_i - \varepsilon)}\right)$$



The first term  $E \left[ \sum_{s=0}^{n-1} \left( \frac{1}{G_{is}} - 1 \right) \right]$  determines the regret due to underestimation of the true mean  $\mu_i$ . This can cause sub-optimal arms to be played too frequently thus leading to larger regret.

We will later show that in a Gaussian setting, and with appropriate priors s.t.  $F_i(t) \sim N(\hat{\mu}_i(t), \frac{1}{t})$ , and  $\mu_1 > \mu_2 \geq \dots \geq \mu_K$ , we have:

$$E \left[ \sum_{s=0}^{n-1} \left( \frac{1}{G_{is}} - 1 \right) \right] \leq \frac{C}{\varepsilon^2} \ln \left( \frac{1}{\varepsilon} \right) \quad (\text{problem 36.6 @ in text})$$

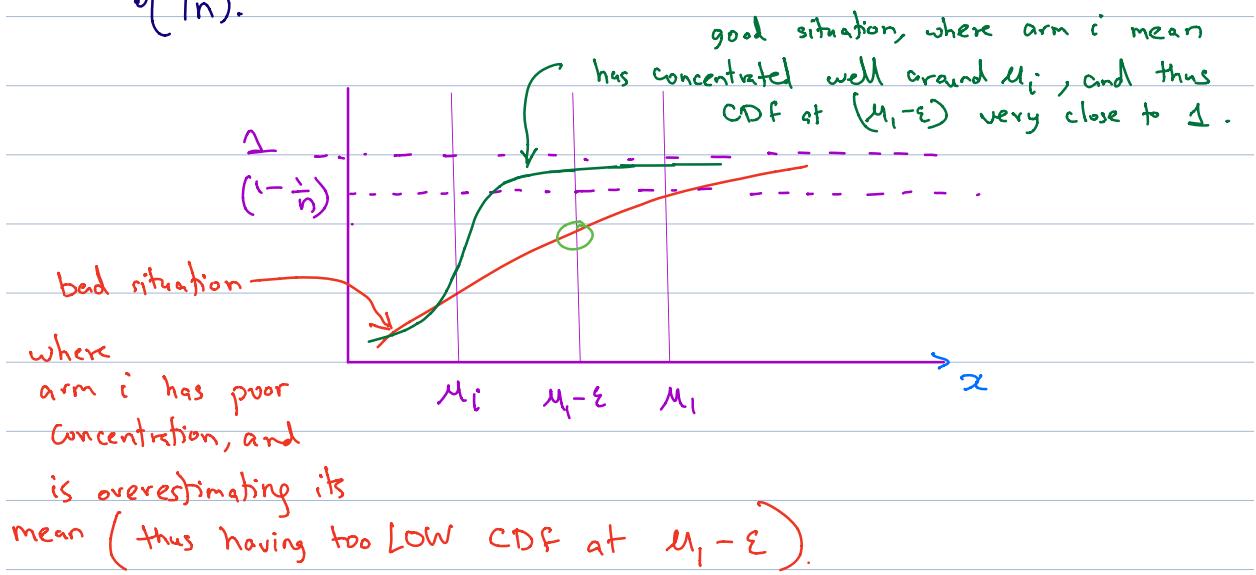
$$\textcircled{2} \quad \chi_{\{G_{is} > \frac{1}{n}\}} = \chi_{\{1 - F_{is}(\mu_i - \varepsilon) > \frac{1}{n}\}}$$

$$= \chi_{\{F_{is}(\mu_i - \varepsilon) < 1 - \frac{1}{n}\}}$$

$\uparrow \downarrow$   
note difference. i vs. 1.

This controls the overestimation of the  $i^{\text{th}}$  arm (which is suboptimal), when compared to the mean of the best arm  $\mu_1$ . Ideally,  $F_{is}(\mu_i + \tilde{\varepsilon}) \approx 1$  (i.e., the  $i^{\text{th}}$  arm's  $\uparrow$  some small  $\tilde{\varepsilon}$ , unrelated to  $\varepsilon$ )

mean is estimated correctly). With  $\mu_i < \mu_i + \tilde{\epsilon} < \mu_i - \epsilon < \mu_i$ , we thus want to be in a regime where we have  $F_{is}(\mu_i - \epsilon) \approx 1$ , where the approx. is extremely small, i.e.,  $\mathcal{O}(\frac{1}{n})$ .



In the same Gaussian setting as ① above, we will later show that:

(problem 36.6 ⑥ in text)

$$E \left[ \sum_{s=0}^{n-1} \chi_{\{G_{is} > \frac{1}{n}\}} \right] \leq 1 + \frac{2}{(\Delta_i - \epsilon)^2} \left( \ln(n) + \sqrt{\pi \ln(n)} + 1 \right).$$

Note 1: Explicit Expressions in a Gaussian setting.

Suppose that we are in a Gaussian prior - Gaussian rewards

setting, i.e.;  $Q_i \sim N(\mu_p, \sigma_p^2)$  and samples for arm  $i$  follow  $X_i \sim N(\mu_i, 1)$ , where  $\mu_i \sim Q_i$ , i.e., the mean is sampled from the prior. Further suppose

$$\mu_p, \sigma_p^2 \rightarrow \infty \text{ s.t. } \sigma_p^2 \rightarrow \infty \text{ and } \frac{\mu_p}{\sigma_p^2} \rightarrow 0.$$

Thus, we initially draw from the environment  $\mu_i \sim Q_i$ .

Then, the posterior of arm  $i$  at time  $t$  is given by:

$$F_{is}(x) = N(\hat{\mu}_{is}, \frac{1}{s}), \text{ i.e.,}$$

function valued r.v. but with finite dimensional parameterization

$$f_{is}(x) = \sqrt{\frac{s}{2\pi}} e^{-\frac{(x-\hat{\mu}_{is})^2}{2s}}$$

empirical estimate of mean at time  $t$

pdf  $\rightarrow$  integrate over  $(-\infty, z] \rightarrow$  get  $\begin{cases} F_i(t)(z) \\ F_{is}(z) \end{cases}$   $x \in \mathbb{R}$ ,  $t = 1, 2, \dots, n$ .

Further,

$$F_{is}(x) = \sqrt{\frac{s}{2\pi}} e^{-\frac{(x-\hat{\mu}_{is})^2}{2s}}$$

empirical estimate of mean with  $s$  samples

arm index  $i$      $s$  samples

CDF at  $x \in (-\infty, \infty)$

$$\text{Finally, } G_{is} = 1 - F_{is}(\mu_i - \varepsilon)$$

$\hookrightarrow$  r.v. that is the value of the random function at  $(\mu_i - \varepsilon)$ .

$$\text{Thus, } G_{is} = \int_{\mu_i - \varepsilon}^{\infty} \sqrt{\frac{s}{2\pi}} e^{-(x - \hat{\mu}_{is})^2/(2s)} dx$$

$\Rightarrow G_{is}(\hat{\mu}_{is})$

Note: This uses the facts that  $X_i$  are iid,  $X_i \sim N(\mu_i, 1)$   $\rightarrow \hat{\mu}_{is} \sim N(\mu_i, \frac{1}{s})$

or equivalently,

$$F_{is}(\mu_i - \varepsilon) = \int_{-\infty}^{\mu_i - \varepsilon} \sqrt{\frac{s}{2\pi}} e^{-(x - \hat{\mu}_{is})^2/(2s)} dx.$$

To compute  $E\left[\left(\frac{1}{G_{is}} - 1\right)\right]$ , we can use this form.

Defining:

$$H(z) = \int_{-\infty}^z \sqrt{\frac{s}{2\pi}} e^{-y^2/(2s)} dy \quad \begin{pmatrix} \text{CDF of} \\ N(0, \frac{1}{s}) \text{ r.v.} \end{pmatrix}$$

We observe by a (sample-pathwise) change of variable:

$y = x - \hat{\mu}_{is}$ , we have:

$$F_{is}(\mu_i - \varepsilon) = \int_{-\infty}^{\mu_i - \hat{\mu}_{is} - \varepsilon} \sqrt{\frac{s}{2\pi}} e^{-y^2/(2s)} dy$$

$\rightarrow (\mu_i - \mu_i)$ .

$$= H(\mu_i - \hat{\mu}_{is} - \varepsilon) = H(\Delta_i - \varepsilon + R)$$

$\hookrightarrow N(\mu_i, \frac{1}{s}) \quad \hookrightarrow N(0, \frac{1}{s}).$

Note: This uses the facts that  $X_i$  are iid,  $X_i \sim N(\mu_i, 1)$   $\rightarrow$

$$\text{Thus, } E\left[\left(\frac{1}{G_{1s}} - 1\right)\right] = E\left[\frac{1 - G_{1s}}{G_{1s}}\right] = E\left[\frac{F_{1s}(\mu_1 - \varepsilon)}{1 - F_{1s}(\mu_1 - \varepsilon)}\right]$$

$$= E_R\left[\frac{H(R - \varepsilon)}{1 - H(R - \varepsilon)}\right] \quad (\because \Delta_1 = 0)$$

$$= \int_{-\infty}^{\infty} h(y) \cdot \frac{H(y - \varepsilon)}{1 - H(y - \varepsilon)} dy$$

$$= \int_{-\infty}^{\infty} h(y + \varepsilon) \frac{H(y)}{1 - H(y)} dy, \quad h(y) = \sqrt{\frac{s}{2\pi}} e^{-y^2/(2s)}$$

Note: This proof below is general, i.e.,  
 Back to Theorem: does not require posterior Sampling, or any  
assumptions on prior or dist. of samples.  
 Algorithm model on Page 10 (i.e., sample from  $F_i(t)$  and update it based on samples).

Theorem (3b.2 in text): for arm  $i$ , and  $\varepsilon \in \mathbb{R}$ ,

$$E[T_i(n)] \leq 1 + E\left[\sum_{s=0}^{n-1} \left(\frac{1}{G_{is}} - 1\right)\right]$$

↓  
 expectation only over randomness of arm reward and algorithm; NOT over environment

$$+ E\left[\sum_{s=0}^{n-1} \chi_{\{G_{is} > 1/n\}}\right]$$

where  $G_{is} = 1 - F_{is}(\mu_i - \varepsilon)$ .

Proof:  $\mathcal{S}_t = \sigma\{A_1, X_1, \dots, A_t, X_t\}$  = information until (and including) time  $t$ .

$$W_j(t) = \{\theta_j(t) \leq \mu_j - \varepsilon\}, j=1, 2, \dots, k$$

where  $\theta_j(t) \sim F_j(t)$  is the sample drawn from

↳ simulation sample

the posterior  $F_j(t) = \tilde{F}_{j, T_j(t-1)}$

Note: We would typically expect (if all estimates are good) that  $\theta_j(t) \leq \mu_j - \varepsilon$ , for arm  $j$  at the end of time  $(t-1)$ .

and indeed, with large enough number of samples,  $\theta_j(t) \approx \mu_j$ .

$$\begin{aligned} \text{Further } P(\theta_1(t) \geq \mu_1 - \varepsilon \mid \mathcal{S}_{t-1}) &= 1 - F_{1, T_1(t-1)}(\mu_1 - \varepsilon) \\ &= G_{1, T_1(t-1)} \quad \text{a.s.} \end{aligned}$$

$$\text{Now, } E[T_j(n)] = E \left[ \sum_{t=1}^n \chi_{\{A_t=j\}} \right]$$

$$= E \left[ \sum_{t=1}^n \chi_{\{A_t=j \wedge W_j(t)\}} \right] \quad \uparrow \{ \theta_j(t) \leq \mu_j - \varepsilon \}$$

(I)

$$+ E \left[ \sum_{t=1}^n \chi_{\{A_t=j \wedge W_j^c(t)\}} \right]$$

(II)

(I): Here the  $j^{th}$  arm sample is below the true mean of the best arm ( $\mu_1$ ). In this case, we will play arm  $j$  (and thus incur regret) if the sample from arm 1,  $\theta_1(t)$ , is atypical and below its true value  $\mu_1$ .

Below, we upper bound (I):  $E \left[ \sum_{t=1}^n X_{\{A_t=j \wedge W_j(t)\}} \right]$

Let  $A'_t = \arg \max_{j \neq 1} \theta_j(t)$ . Then, we have:

$$P(A_t=j, W_j(t) | \mathcal{F}_{t-1}) \geq P(A'_t=j, W_j(t), \theta_1(t) > \mu_1 - \varepsilon | \mathcal{F}_{t-1})$$

$\{\theta_l(t) \leq \mu_1 - \varepsilon\}$   $\hookrightarrow$  here  $\theta_1(t) > \theta_j(t) \geq \theta_{l+1}(t)$

$\Rightarrow$  this is one way

for arm 1 to be the best.

$\{\theta_1(t), \dots, \theta_k(t)\}$  conditionally indep. given  $\mathcal{F}_{t-1}$ .

$$= P(\theta_1(t) > \mu_1 - \varepsilon | \mathcal{F}_{t-1}) \cdot P(A'_t=j, W_j(t) | \mathcal{F}_{t-1})$$

$= G_{1, T_1(t-1)}$

$\hookrightarrow$  function of  $\{\theta_r(t) \text{ } \forall r \neq 1\}$

Now observe:

①

$$P(A_t=j, W_j(t) | \mathcal{F}_{t-1}) = P(A'_t=j, W_j(t), \underbrace{\theta_1(t) < \theta_j(t)}_{\substack{\text{best arm skipping arm 1 is } j, \\ \text{arm 1 below arm } j}} | \mathcal{F}_{t-1})$$

$\hookrightarrow$  best arm =  $j$  and  $\theta_j(t) \leq \mu_1 - \varepsilon$

$$\leq P(A_t' = j, W_j(t), \Theta_1(t) \leq \mu_1 - \varepsilon \mid \mathcal{F}_{t-1})$$

$$= P(A_t' = j, W_j(t) \mid \mathcal{F}_{t-1}) P(\Theta_1(t) \leq \mu_1 - \varepsilon \mid \mathcal{F}_{t-1})$$

$\underbrace{\quad \quad \quad}_{\begin{array}{l} \{ \Theta_1(t), \dots, \Theta_k(t) \} \text{ conditionally indep.} \\ \text{given } \mathcal{F}_{t-1}. \end{array}}$

$1 - G_{1, T_1(t-1)}$

$$\Rightarrow P(A_t' = j, W_j(t) \mid \mathcal{F}_{t-1}) \geq \underbrace{P(A_t = j, W_j(t) \mid \mathcal{F}_{t-1})}_{1 - G_{1, T_1(t-1)}}$$

$\hookrightarrow (2)$ .

(1) + (2)  $\Rightarrow$

$$\begin{aligned} P(A_t = j, W_j(t) \mid \mathcal{F}_{t-1}) &\leq \left( \frac{1}{G_{1, T_1(t-1)}} - 1 \right) P(A_t = 1, W_j(t) \mid \mathcal{F}_{t-1}) \\ &\leq \left( \frac{1}{G_{1, T_1(t-1)}} - 1 \right) P(A_t = 1 \mid \mathcal{F}_{t-1}). \end{aligned}$$

$\hookrightarrow (3)$ .

$$\therefore E \left[ \sum_{t=1}^n \chi_{\{A_t = j \wedge W_j(t)\}} \right]$$

iterated  
expectations

$$= E \left[ \sum_{t=1}^n E \left[ \chi_{\{A_t = j, W_j(t)\}} \mid \mathcal{F}_{t-1} \right] \right]$$

$$= E \left[ \sum_{t=1}^n P(A_t=j, W_j(t) | S_{t-1}) \right]$$

from ③

$$\leq E \left[ \sum_{t=1}^n \left( \frac{1}{G_{1,T_1(t-1)}} - 1 \right) P(A_t=1 | S_{t-1}) \right]$$

$E[X_{\{A_t=1\}} | S_{t-1}]$

$$\leq E \left[ \sum_{t=1}^n \left( \frac{1}{G_{1,T_1(t-1)}} - 1 \right) X_{\{A_t=1\}} \right]$$

Switching from time t  
to samples s

$$\leq E \left[ \sum_{s=1}^n \left( \frac{1}{G_{1,s}} - 1 \right) \right]$$

The last step is similar to the trick in the UCB proof.

It follows because every time  $A_t=1$ ,  $T_1(t-1)$  increments by 1, and thus a term  $G_{1,T_1(t-1)}$  appears at most once in the summation for any fixed value of  $r=T_1(\cdot)$ .

(II): Upper Bounding  $E \left[ \sum_{t=1}^n X_{\{A_t=j \cap W_j^c(t)\}} \right]$

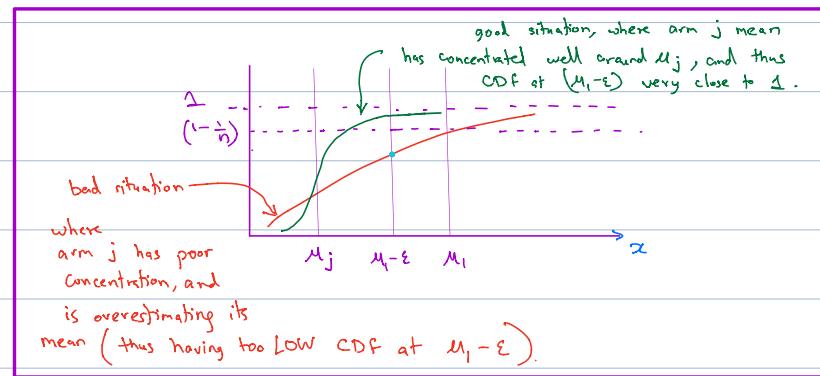
$\{\theta_j(t) > \mu_j - \varepsilon\}$

This second term controls the number of times that the sub optimal arm  $j$  is played when sample from arm  $j$  is atypically high.

$$\text{Let } T = \left\{ t \in \{1, 2, \dots, n\} : 1 - F_{j, T_j(t-1)}(\mu_j - \varepsilon) > \frac{1}{n} \right\}$$

$$= \left\{ t \in \{1, 2, \dots, n\} : G_{j, T_j(t-1)} > \frac{1}{n} \right\}$$

These correspond to time slots s.t. the empirical dist. of arm  $j$  has an abnormally large amount of mass above  $(\mu_j - \varepsilon)$  which is well above  $\mu_j$ , the true mean for arm  $j$ .



Corresponds to the red empirical CDF in the figure.

Aside: In the proof of TS, we need to control for randomness at two levels: (i) due to the empirical dist. being atypical, and (ii) the sample  $\theta_j(t) \sim F_j(t)$  being atypical. The time slots  $T$  control for (i).

$$\textcircled{II}: E \left[ \sum_{t=1}^n \chi_{\{A_t=j, w_j^c(t)\}} \right]$$

$$\leq E \left[ \sum_{t \in T} \chi_{\{A_t=j\}} \right] \xrightarrow{\textcolor{red}{\longrightarrow}} \textcircled{4}.$$

$$\textcircled{III} + E \left[ \sum_{t \notin T} \chi_{\{w_j^c(t)\}} \right]$$

\textcircled{IV}

Now, suppose  $A_t=j$ , and  $T_j(t-1) = (s-1)$ ,  $T_j(t) = s$  for some  $t, s$ . Then, observe that  $G_{j, T_j(t-1)} > \frac{1}{n}$ , when  $t \in T$ . (Ref: Solutions to Bandit Algorithms, 36.5)

$$\textcircled{III} : E \left[ \sum_{j \in T} \chi_{\{A_t=j\}} \right]$$

$$\leq E \left[ \sum_{t=1}^n \sum_{s=1}^n \chi_{\{T_j(t-1) = s-1, T_j(t) = s, G_{j, T_j(t-1)} > \frac{1}{n}\}} \right]$$

replace by  $(s-1)$

$$= E \left[ \sum_{s=1}^n \chi_{\{G_{j, (s-1)} > \frac{1}{n}\}} \underbrace{\sum_{t=1}^n \chi_{\{T_j(t-1) = s-1, T_j(t) = s\}}} \right]$$

$\leq 1$  (for any fixed  $s$ ,

this can occur at most once)

$$\leq E \left[ \sum_{s=1}^n \chi_{\{G_{j, (s-1)} > \frac{1}{n}\}} \right]$$

replace  $(s-1)$  by  $s$  and bound continues to hold.

$$\text{IV} : E \left[ \sum_{t \in T} \underbrace{\chi_{\{w_j^c(t)\}}}_{\substack{\{t : G_j \tau_j(t-1) \leq \frac{1}{n}\}}} \right] \quad \begin{array}{l} \text{Ref: Solutions,} \\ \# 36.5; \text{ see .} \\ \text{textbook website} \end{array}$$

$\theta_j(t) > \mu_i - \varepsilon$

$\{t : G_j \tau_j(t-1) \leq \frac{1}{n}\}$

$$= \sum_{t=1}^n E \left[ E \left[ \underbrace{\chi_{\{G_j \tau_j(t-1) \leq \frac{1}{n}, w_j^c(t)\}}}_{S_{t-1} \text{ measurable}} \mid S_{t-1} \right] \right] = \chi_{S_{t-1}} \chi_{S_{t-1}}$$

Now recall:

$$E \left[ \chi_{\{w_j^c(t)\}} \mid S_{t-1} \right] = P(w_j^c(t) \mid S_{t-1})$$

$$= P(\theta_j(t) > \mu_i - \varepsilon \mid S_{t-1}) = G_j \tau_j(t-1).$$

$$= \sum_{t=1}^n E \left[ \chi_{\{G_j \tau_j(t-1) \leq \frac{1}{n}\}} \cdot \underbrace{G_j \tau_j(t-1)}_{\leq \frac{1}{n}} \right]$$

$$\leq 1$$

The result now follows by putting these together.  $\square$

Main Theorem: Thompson Sampling has logarithmic regret  
on the instance dependent frequentist setting. Specifically,  
(Thm 36.3 in text) in a Gaussian-Gaussian setting,

suppose  $F_{js} \sim N(\hat{\mu}_j(s), \frac{1}{s})$ ,  $s=1, 2, \dots, n$ ,

$j=1, 2, \dots, k$

(see Note 1 on pp. 13).

Then,  $\lim_{n \rightarrow \infty} \frac{R_n(\pi, v)}{\ln(n)} \leq \sum_{\{j : \Delta_j > 0\}} \left( \frac{v}{\Delta_j} \right)$ , where

$v \in \Sigma_N^k(i)$ .

Proof: (Adapted from Solutions to Bandit Algorithms, #3 6.b)

Recall from prev Thm that:  $j \in \{1, 2, \dots, k\}$ , we have

$$E[T_j(n)] \leq 1 + E\left[ \sum_{s=0}^{n-1} \left( \frac{1}{G_{js}} - 1 \right) \right]$$

$$+ E\left[ \sum_{s=0}^{n-1} \chi_{\{G_{js} \geq 1/n\}} \right]$$

Further, regret decomposition implies

$$R_n(\pi, v) = \sum_{\{j : \Delta_j > 0\}} \Delta_j E[T_j(n)]$$

Below, we explicitly derive bounds for the two

terms in  $E[\tau_j(n)]$  in the Gaussian setting.

(I): From Note 1 on pp. 13/14, we have:

$$E\left[\left(\frac{1}{G_{1s}} - 1\right)\right]$$

This step  $X_i \text{ iid}, X_i \sim N(\mu_i, 1)$   
from Note 1

$$= \int_{-\infty}^{\infty} h(y+\varepsilon) \frac{H(y)}{1-H(y)} dy, \quad h(y) = \sqrt{\frac{s}{2\pi}} e^{-y^2/(2s)}$$

$$H(y) = \int_{-\infty}^y h(z) dz.$$

The rest is (intricate) Gaussian calculations.

$$H'(y) \geq \frac{e^{-sy^2/2}}{y\sqrt{s} + \sqrt{sy^2 + 4}}$$

Fact:  $\phi(z)$  CDF of  $N(0, 1)$

$$1 - \phi(z) \geq \frac{e^{-z^2/2}}{z + \sqrt{z^2/4}}$$

Now

$$\int_{-\infty}^{\infty} h(y+\varepsilon) \frac{H(y)}{1-H(y)} dy = \int_0^{\infty} h(y+\varepsilon) \underbrace{\frac{H(y)}{1-H(y)}}_{\leq 1} dy + \int_{-\infty}^0 h(y+\varepsilon) \underbrace{\frac{H(y)}{1-H(y)}}_{\geq 1} dy$$

$$\leq \int_0^{\infty} \frac{h(y+\varepsilon)}{1-H(y)} dy + 2 \int_{-\infty}^0 h(y+\varepsilon) H(y) dy$$

$$\leq \int_0^\infty \frac{h(y+\varepsilon) e^{-sy^2/2}}{\sqrt{yN\varepsilon + \sqrt{sy^2+4}}} dy + 2 \int_{-\infty}^\infty h(y+\varepsilon) H(y) dy$$

$$\leq 2e^{-s\varepsilon^2/4} \int_0^\infty e^{-sy^2} (y\sqrt{N\varepsilon+1}) \sqrt{\frac{s}{2\pi}} dy + 2e^{-s\varepsilon^2}$$

$$\leq \frac{2 \left(1 + \sqrt{N\varepsilon}\right)}{\varepsilon^2 \sqrt{2\pi}} e^{-s\varepsilon^2/2} + 2e^{-s\varepsilon^2}$$

$$\therefore E \left[ \sum_{s=1}^n \left( \frac{1}{G_{js}} - 1 \right) \right] \leq \sum_{s=1}^n \left( \dots \right) \leq \sum_{s=1}^\infty \left( \dots \right)$$

$$\leq \frac{C}{\varepsilon^2} \ln \left( \frac{1}{\varepsilon} \right) \quad \begin{matrix} \text{(see solutions)} \\ \text{manuel 36.6 a)} \end{matrix}$$

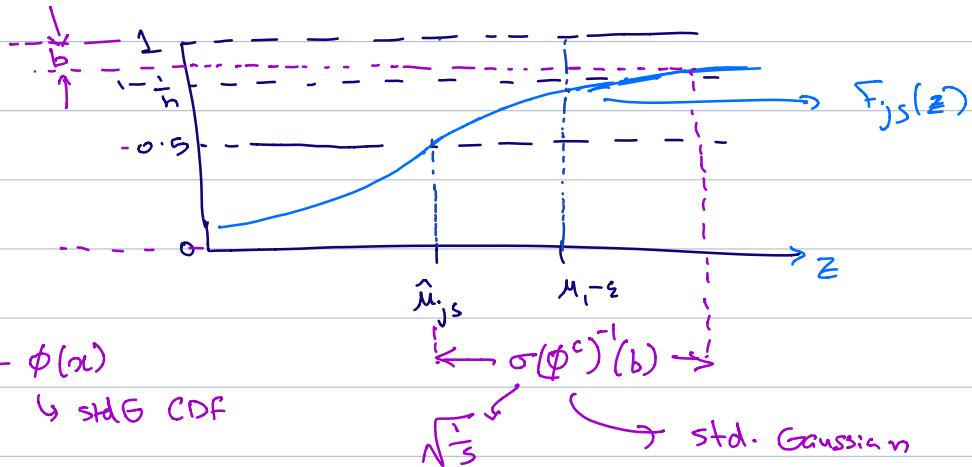
$$\textcircled{II}: E \left[ \sum_{s=0}^{n-1} \chi_{\{G_{js} > \frac{1}{n}\}} \right] = \sum_{s=1}^{n-1} P \left( G_{js} > \frac{1}{n} \right)$$

Now, suppose

$$G_{js} > \frac{1}{n} \iff 1 - F_{js}(m_1 - \varepsilon) > \frac{1}{n} \quad \xrightarrow{N(\bar{\mu}_{js}, \frac{1}{s})}$$

$$\iff F_{js}(\mu_1 - \varepsilon) \leq \left(1 - \frac{1}{n}\right).$$

(i.e.,  $\hat{\mu}_{js}$  is "close to"  $\mu_1 - \varepsilon$ )  
 $\hookrightarrow$   $\text{ort}(\cdot)$



$$\text{Now: } X \sim N(0,1) \Rightarrow P(X > x) \leq \frac{e^{-x^2/2}}{\sqrt{2\pi x^2}} \leq e^{-x^2/2} \xrightarrow{x > \frac{1}{\sqrt{2\pi}}}.$$

$\therefore$  for  $x = \sqrt{2 \ln(n)}$ ,  $n \geq 3$ , we have:

$$F_{js}(\mu_1 - \varepsilon) \leq 1 - \frac{1}{n} \implies \hat{\mu}_{js} + \sqrt{\frac{2 \ln(n)}{s}} > (\mu_1 - \varepsilon).$$

$$\therefore P(G_{js} > \frac{1}{n}) \leq P\left(\hat{\mu}_{js} + \sqrt{\frac{2 \ln(n)}{s}} > \mu_1 - \varepsilon\right)$$

$\downarrow$   
 $\mu_j + \Delta_j$

$$\leq e^{-\frac{s}{2} \left( \Delta_j - \varepsilon - \sqrt{\frac{2\ln(n)}{s}} \right)^2}$$

$$\therefore \sum_{s=1}^n P(G_{j,s} > \frac{1}{n}) \leq 1 + \frac{2}{(\Delta_j - \varepsilon)^2} \left( \ln(n) + \sqrt{\pi \ln(n)} + 1 \right)$$

(see 36.2(b), Solutions for more details).

The result now follows by putting the bounds together in the regret decomposition.  $\square$

Note: The main theorem provides the regret bound as long as the algorithm uses the sampling distribution:

$$F_{is}(x) \sim N\left(\hat{\mu}_{is}, \frac{1}{s}\right),$$

$$i = 1, 2, \dots, k$$

$$s = 1, 2, \dots, n.$$

and the rewards  $X_i \sim N(\mu_i, 1)$ ,  $i = 1, 2, \dots, k$ .

- (a) The regret bound holds without any assumption on the prior, and whether  $F_{is}(\cdot)$  above is the correct posterior.
- (b) The regret bound matches that of KL-UCB, and thus is a strict improvement over UCB.
- (c) In the case when the prior is **uninformative**, i.e.,  $Q \sim N(\mu_p, \sigma_p^2)$  with  $\sigma_p^2 \rightarrow \infty$ ,  
 $\frac{\mu_p}{\sigma_p^2} \rightarrow 0$ , then this bound matches the regret lower bound.