# 1   Proof of theorem 1 from lecture 6 continued

We adopt all notations from the previous lecture. Recall that we ended with bounding the expected number of times that a suboptimal arm $a$ is played:

$$\mathbb{E}\left[N_T(a)\right] = \sum_{t=1}^{T} \underbrace{\mathbb{P}\left(A_t = a, E_{a,t}^{\mu}, E_{a,t}^{\theta}\right)}_{S_1} + \underbrace{\mathbb{P}\left(A_t = a, E_{a,t}^{\mu}, \overline{E_{a,t}^{\theta}}\right)}_{S_2} + \underbrace{\mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\mu}}\right)}_{S_3}$$

## 1.1   Bounding $S_1$

Last lecture gave a bound for $S_1$. Summarizing, we had

$$S_1 \leq \mathbb{E}\left[\sum_{k=0}^{T-1} \frac{1 - p_{a,\tau_{k+1}}}{p_{a,\tau_{k+1}}}\right]$$

and the following lemma

**Lemma 1.** *Let $\tau_k$ be the time at which the arm $a$ is pulled for the $k^{th}$ time. Then*

$$\mathbb{E}\left[\frac{1 - p_{a,\tau_{k+1}}}{p_{a,\tau_{k+1}}}\right] \leq \begin{cases} \frac{3}{\Delta_a'} & \text{if } k < \frac{8}{\Delta_a'} \\ \Theta\left(e^{-k\Delta_a'^2/2} + \frac{1}{(k+1)\Delta_a'^2}e^{-kD_a} + \frac{1}{e^{k\Delta_a'^2/4}-1}\right) & \text{otherwise} \end{cases}$$

*where $\Delta_a' = \mu_0 - y_a$ and $D_a = d(y_a, \mu_0)$.*

Since we are not proving lemma 1, we will give some heuristics as to why it is true. Recall that $p_{a,\tau_{k+1}} = \mathbb{P}\left(\theta_{0,\tau_{k+1}} > y_a \mid \mathcal{H}_t\right)$. If $\theta_{0,\tau_{k+1}}$ were to be replaced by the sample mean $\hat{\mu}_{0,k}$, then Chernoff bound would imply that $p_{a,\tau_{k+1}} \leq 1 - e^{-kC}$ (nevermind the fact that $\tau_{k+1}$ is random). Furthermore, $\frac{e^{-kC}}{1-e^{-kC}} \approx e^{-kC}$ for $k \gg 0$. This gives some heuristics as to why the terms inside the $\Theta$ in lemma 1 should have this particular form.

## 1.2   Bounding $S_2$

To bound $S_2$ and $S_3$, we first prove the following lemmas.

**Lemma 2.**

$$\sum_{t=1}^{T} \mathbb{P}\left(A_t = a, E_{a,t}^{\mu}, \overline{E_{a,t}^{\theta}}\right) \leq L_a(T) + 1$$

*where $L_a(T) = \frac{\log T}{d(x_a, y_a)}$.*

From [AG12]: "This follows from the observation that $\theta_{a,t}$ is well-concentrated around its mean when $N_t(a)$ is large, that is, larger than $L_a(T)$". We will formally illustrate this in the proof.

*Proof.* First, we make sure that we have taken enough samples of $a$:

$$\sum_{t=1}^{T} \mathbb{P}\left(A_t = a, E_{a,t}^{\mu}, \overline{E_{a,t}^{\theta}}\right) \le L_a(T) + \sum_{t=1}^{T} \mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}}, E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T)\right). \qquad (1)$$

For each term in the sum in the RHS above, we condition on the history

$$\mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}}, E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T)\right) = \mathbb{E}\left[\mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}}, E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T) \mid \mathcal{H}_t\right)\right] \qquad (2)$$

Using the Bayes rule

$$\mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}}, E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T) \mid \mathcal{H}_t\right)$$
$$= \mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}} \mid E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T), \mathcal{H}_t\right) \mathbb{P}\left(E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T) \mid \mathcal{H}_t\right),$$

the fact that $\mathbb{P}\left(E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T) \mid \mathcal{H}_t\right) \le 1$ and (2), we get

$$\mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}}, E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T)\right) \le \mathbb{E}\left[\mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}} \mid E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T), \mathcal{H}_t\right)\right].$$

We now claim that

$$\mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}} \mid E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T), \mathcal{H}_t\right) \le \frac{1}{T}. \qquad (3)$$

Given the claim, the RHS of (1) $\le L_a(T) + \sum_{t+1}^{T} 1/T = L_a(T) + 1$, which proves the lemma. So it remains to show that (3) holds. Now,

$$\mathbb{P}\left(A_t = a, \overline{E_{a,t}^{\theta}} \mid E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T), \mathcal{H}_t\right) \qquad (4)$$

$$\le \mathbb{P}\left(\overline{E_{a,t}^{\theta}} \mid E_{a,t}^{\mu}, N_{t-1}(a) > L_a(T), \mathcal{H}_t\right) \qquad (5)$$

$$= \mathbb{P}\left(\theta_{a,t} > y_a \mid \hat{\mu}_{a,t-1} \le x_a, N_{t-1}(a) > L_a(T), \mathcal{H}_t\right) \quad \text{by definition of } E_{a,t}^{\theta} \text{ and } E_{a,t}^{\mu} \qquad (6)$$

$$= \mathbb{P}\left(\theta_{a,t} > y_a \mid \underbrace{S_{a,t-1} \le x_a(N_{t-1}(a) + 1), N_{t-1}(a) > L_a(T), \mathcal{H}_t}_{=:\star}\right) \quad \text{by definition of } \hat{\mu}_{a,t-1} \qquad (7)$$

$$= \mathbb{P}\left(Beta(S_{a,t-1} + 1, N_{t-1}(a) - S_{a,t-1} + 1) > y_a \mid \star\right) \qquad (8)$$

where the last equality follows from the definition of $\theta_{a,t}$, and $Beta(\alpha, \beta) > y_a$ denotes the event that a $Beta(\alpha, \beta)$ distributed random variable is greater than $y_a$. We have a general fact about the Beta distribution:

**Lemma 3.** *For $\alpha' > \alpha$ and $y \in [0, 1]$, we have*

$$\mathbb{P}\left(Beta(\alpha, C - \alpha) > y\right) \le \mathbb{P}\left(Beta(\alpha', C - \alpha') > y\right)$$

*Proof.* Recall the identity

$$1 - \text{CDF}_{Beta(\alpha,\beta)}(y) = \text{CDF}_{Binom(\alpha+\beta-1,y)}(\alpha - 1). \tag{9}$$

Using the identity, we have

$$
\begin{aligned}
\mathbb{P}\left(Beta(\alpha, C - \alpha) > y\right) &= 1 - \text{CDF}_{Beta(\alpha,C-\alpha)}(y) \\
&= \text{CDF}_{Binom(C-1,y)}(\alpha - 1) \\
&\leq \text{CDF}_{Binom(C-1,y)}(\alpha' - 1) \\
&= 1 - \text{CDF}_{Beta(\alpha',C-\alpha')}(y) = \mathbb{P}\left(Beta(\alpha', C - \alpha') > y\right)
\end{aligned}
$$

This proves lemma 3. □

Going back to (8), we have

$$\mathbb{P}\left(Beta(S_{a,t-1} + 1, N_{t-1}(a) - S_{a,t-1} + 1) > y_a \mid \bigstar\right) \tag{10}$$

$$\leq \mathbb{P}\left(Beta(x_a(N_{t-1}(a) + 1) + 1, (1 - x_a)(N_{t-1}(a) + 1)) > y_a \mid \bigstar\right) \quad \because \text{lemma 3,} \tag{11}$$

$$= \text{CDF}_{Binom(N_{t-1}(a)+1,y_a)}(x_a(N_{t-1}(a) + 1)) \quad \because \text{identity (9).} \tag{12}$$

To bound the quantity above, we need the following lemma.

**Lemma 4** (KL-divergence version of Chernoff-Hoeffding). *Let $X_1, \ldots, X_n$ be independent Bernoulli random variables. Let $p_i = \mathbb{E}[X_i]$, $X = \frac{1}{n}\sum_{i=1}^n X_i$, $\mu = \mathbb{E}[X]$. Then*

$$\mathbb{P}\left(X \geq \mu + \lambda\right) \leq \exp\left(-nd(\mu + \lambda, \mu)\right), \quad \forall 0 < \lambda \leq 1 - \mu.$$
$$\mathbb{P}\left(X \leq \mu - \lambda\right) \leq \exp\left(-nd(\mu - \lambda, \mu)\right), \quad \forall 0 < \lambda < \mu$$

*where $d(a, b) = a\log(\frac{a}{b}) + (1 - a)\log(\frac{1-a}{1-b})$.*

See the supplementary material of [AG13]. Now continuing from (12)

$$
\begin{aligned}
\text{CDF}_{Binom(N_{t-1}(a)+1,y_a)}(x_a(N_{t-1}(a) + 1)) &\leq e^{-(N_{t-1}(a)+1)d(x_a,y_a)} \quad \because \text{lemma 4} \\
&\leq e^{-L_a(T)d(x_a,y_a)} \quad \because N_{t-1}(a) > L_a(T) \text{ by } \bigstar \\
&= \frac{1}{T} \quad \because \text{definition of } L_a(T).
\end{aligned}
$$

This proves (3) which concludes the proof of lemma 2. □

## 1.3 Bounding $S_3$

**Lemma 5.**

$$\sum_{t=1}^T \mathbb{P}\left(A_t = a, \overline{E_{a,t}^\mu}\right) \leq 1 + \frac{1}{d(x_a, y_a)}.$$

*Proof.* Define $\tau_k$ as the time index at which action $a$ is taken for the $k$-th time and $\tau_0 = 0$.

$$\sum_{t=1}^{T} \mathbb{P}\left(A_t = a, \overline{E_{a,t}^\mu}\right) \leq \mathbb{E}\left[\sum_{k=0}^{T-1} \sum_{t=\tau_k+1}^{\tau_{k+1}} \mathbb{1}(A_t = a)\mathbb{1}(\overline{E_{a,t}^\mu})\right] \quad \because \tau_T \geq T$$

$$= \mathbb{E}\left[\sum_{k=0}^{T-1} \mathbb{1}(\overline{E_{a,\tau_{k+1}}^\mu})\right]$$

$$\leq 1 + \mathbb{E}\left[\sum_{k=1}^{T-1} \mathbb{1}(\overline{E_{a,\tau_{k+1}}^\mu})\right] \quad \text{(Shifting the start of the summation)}$$

$$\leq 1 + \sum_{k=1}^{T-1} e^{-kd(x_a,\mu_a)} \quad \because \text{lemma 4 and } \overline{E_{a,\tau_{k+1}}^\mu} = \{\hat{\mu}_{a,\tau_{k+1}} > x_a\}$$

$$\leq 1 + \frac{e^{-d(x_a,\mu_a)}}{1 - e^{-d(x_a,\mu_a)}} \quad \because \text{geometric series}$$

$$\leq 1 + \frac{1}{d(x_a,\mu_a)} \quad \because \frac{e^{-x}}{1-e^{-x}} \leq \frac{1}{x} \iff e^x \geq 1 + x$$

$\square$

## 1.4 Putting it all together

$$\mathbb{E}\left[N_T(a)\right] \leq \frac{24}{\Delta_a'^2} + \sum_{j=0}^{T-1} \Theta\left(e^{-j\Delta_a'^2/2} + \frac{1}{(j+1)\Delta_a'^2}e^{-jD_a} + \frac{1}{e^{j\Delta_a'^2/4}-1}\right) \quad \because \text{lemma 1}$$

$$+ \frac{\log T}{d(x_a,y_a)} + 1 \quad \because \text{lemma 2}$$

$$+ \frac{1}{d(x_a,y_a)} + 1 \quad \because \text{lemma 5}$$

$$\leq (1+\epsilon)^2 \frac{\log T}{d(\mu_a,\mu_0)} + O\left(\frac{1}{\epsilon^2}\right) \quad \because d(x_a,y_a) = \frac{d(\mu_a,\mu_0)}{(1+\epsilon)^2} \text{ by construction.}$$

where the $O\left(\frac{1}{\epsilon^2}\right)$ term hides the "constants" that depend on the distribution.

# References

[AG12] Shipra Agrawal and Navin Goyal. Analysis of Thompson sampling for the multi-armed bandit problem. In *COLT*, pages 39–1, 2012.

[AG13] Shipra Agrawal and Navin Goyal. Further optimal regret bounds for Thompson sampling. In *AISTATS*, pages 99–107, 2013.