# Lecture 21

**Instructor: Haipeng Luo**

## 1  Softening Policy Elimination

In this lecture we are finally ready to discuss the state-of-the-art algorithm for the i.i.d. contextual bandit problem, which is both optimal and oracle-efficient [Agarwal et al., 2014]. Recall that the idea of Policy Elimination is to find $P_t \in \Delta(\Pi_t)$ such that $V(P_t, \pi) \le 2K$ for all $\pi \in \Pi_t$ where

$$V(P, \pi) = \mathbb{E}_x \left[ \frac{1}{P^\mu(\pi(x)|x)} \right]$$

is essentially the variance of the loss estimators that we want to control. To obtain an efficient algorithm, we need to forget about the idea of removing policies from $\Pi$. So is it possible to ensure $V(P_t, \pi) \le 2K$ for all $\pi \in \Pi$ while at the same time $P_t$ puts most of the weights on good policies?

Unfortunately this is too strong of a requirement. For example, if there is a bad policy $\pi_{\text{bad}}$ which always picks a bad action $a_{\text{bad}}$ with loss 1 and no other policy ever picks $a_{\text{bad}}$, then

$$2K \ge V(P_t, \pi_{\text{bad}}) = \mathbb{E}_x \left[ \frac{1}{P_t^\mu(a_{\text{bad}}|x)} \right] = \frac{1}{(1 - K\mu)P_t(\pi_{\text{bad}}) + \mu}$$

which implies that $P_t(\pi_{\text{bad}})$ will be pretty large assuming $\mu$ is small. This is clearly not a good algorithm.

From this example, however, we can see that the condition $V(P_t, \pi) \le 2K$ should be somehow relaxed for bad policies. Just as in Policy Elimination, whether a policy is good or bad can be roughly determined by its empirical performance compared to the empirically best policy. Specifically, recall the notation $\bar{\ell}_t(\pi) = \frac{1}{t} \sum_{\tau=1}^{t} \widehat{\ell}_\tau(\pi(x_\tau))$ for the empirical average loss and $\pi_t^\star = \operatorname{argmin}_{\pi \in \Pi} \bar{\ell}_t(\pi)$ for the empirically best policy up to time $t$. Define empirical average regret for a policy $\pi$ to be

$$\text{Reg}_t(\pi) = \bar{\ell}_t(\pi) - \bar{\ell}_t(\pi_t^\star).$$

We now relax the low-variance condition as: find $P_t$ such that

$$V(P_t, \pi) \le 2K + \beta \text{Reg}_{t-1}(\pi) \quad \forall \, \pi \in \Pi$$

for some parameter $\beta > 0$ to be specified later. Now there is hope to impose exploitation simultaneously. Specifically, we want $\sum_{\pi \in \Pi} P_t(\pi) \text{Reg}_{t-1}(\pi)$ to be as small as possible. How small can it be? The following lemma answers this question.

**Lemma 1.** *For any $\beta > 0$, there always exists a distribution $P \in \Delta(\Pi)$ such that*

$$\sum_{\pi \in \Pi} P(\pi) \text{Reg}_{t-1}(\pi) \le \frac{2K}{\beta}$$

$$V(P, \pi) \le 2K + \beta \text{Reg}_{t-1}(\pi) \quad \forall \, \pi \in \Pi.$$

*Proof.* Define function $F_t : \Delta(\Pi) \to \mathbb{R}_+$ as

$$F_t(P) = \sum_{\pi \in \Pi} P(\pi) \text{Reg}_{t-1}(\pi) + \frac{2}{\beta} \mathbb{E}_x \left[ \sum_{a=1}^{K} \ln \frac{1}{P^\mu(a|x)} \right].$$

The claim is that the minimizer of $F_t(P)$, which always exists due to compactness of $\Delta(\Pi)$ and continuousness of $F_t$, satisfies both conditions. To see this, first notice that we can extend the function to a set of "sub-distributions" $\Delta(\Pi)' = \{P \in \mathbb{R}_+^N : \sum_{\pi \in \Pi} P(\pi) \leq 1\}$ and still have

$$\min_{P \in \Delta(\Pi)} F_t(P) = \min_{P \in \Delta(\Pi)'} F_t(P).$$

This is because for any sub-distribution $P \in \Delta(\Pi)'$, one can make it a distribution by increasing the weight for policy $\pi_{t-1}^\star$ until the weights sum up to 1. This will only decrease the function value since $\text{Reg}_{t-1}(\pi_{t-1}^\star) = 0$ and the second term of $F_t$ is decreasing in any coordinate of $P$.

Next note that the derivate of $F_t$ with respect to a policy $\pi$ is

$$\nabla F_t(P)(\pi) = \text{Reg}_{t-1}(\pi) - \frac{2(1 - K\mu)}{\beta} V(P, \pi).$$

Let $P^*$ be a minimizer of $F_t$ over $\Delta(\Pi)'$. By KKT conditions, we have

$$\text{Reg}_{t-1}(\pi) - \frac{2(1 - K\mu)}{\beta} V(P^*, \pi) - \lambda_\pi + \lambda = 0 \tag{1}$$

for some Lagrangian multipliers $\lambda_\pi \geq 0$ and $\lambda \geq 0$. Multiply both sides by $P^*(\pi)$ and sum over $\pi \in \Pi$ gives

$$\sum_{\pi \in \Pi} P^\star(\pi) \text{Reg}_{t-1}(\pi) = \frac{2(1 - K\mu)}{\beta} \sum_{\pi \in \Pi} P^*(\pi) V(P^*, \pi) + \sum_{\pi \in \Pi} P^*(\pi)\lambda_\pi - \lambda \quad (P^* \in \Delta(\Pi))$$

$$= \frac{2(1 - K\mu)}{\beta} \sum_{\pi \in \Pi} P^*(\pi) V(P^*, \pi) - \lambda \qquad \text{(complementary slackness)}$$

$$\leq \frac{2}{\beta} \mathbb{E}_x \left[ \sum_{\pi \in \Pi} \frac{P^*(\pi)}{P^*(\pi(x)|x)} \right] - \lambda = \frac{2K}{\beta} - \lambda \leq \frac{2K}{\beta},$$

showing that $P^*$ satisfies the first condition. Moreover, the last equality above also implies $\lambda \leq \frac{2K}{\beta}$ since $\text{Reg}_{t-1}(\pi) \geq 0$. Rearranging Eq. (1) thus gives

$$V(P^*, \pi) \leq \frac{\beta}{2(1 - K\mu)} \left( \text{Reg}_{t-1}(\pi) + \lambda \right) \leq 2K + \beta \text{Reg}_{t-1}(\pi),$$

where we assume $\mu \leq \frac{1}{2K}$ so that $2(1 - K\mu) \geq 1$ (since otherwise we trivially have $V(P^*, \pi) \leq 1/\mu \leq 2K$). This shows that $P^\star$ satisfies the second condition too. $\qquad \square$

The question is now what $\beta$ we should use. Assuming $\text{Reg}_{t-1}(\pi)$ concentrates well around the actual expected regret of $\pi$ compared to $\pi^\star$,

$$\text{Reg}(\pi) \stackrel{\text{def}}{=} \bar{\ell}(\pi) - \bar{\ell}(\pi^\star),$$

which is exactly what we hope for, $\text{Reg}_{t-1}(\pi)$ should be at most a constant. A reasonable choice of $\beta$ would then be of order $1/\mu$, since $V(P, \pi)$ is trivially bounded by $1/\mu$. In other words, when a policy $\pi$ is good, which means $\text{Reg}_{t-1}(\pi)$ is close to zero, we still require $V(P, \pi)$ to be close to $2K$, while when the policy is bad, which means $\text{Reg}_{t-1}(\pi)$ is a large constant, then there is almost no requirement on $V(P, \pi)$ with this choice of $\beta$.

On the other hand, this means that the exploitation constraint is $\sum_{\pi \in \Pi} P(\pi)\text{Reg}_{t-1}(\pi) = \mathcal{O}(K\mu)$, which also makes sense because $\mu$ should be of order $1/\sqrt{T}$, and if the per round regret is of order $1/\sqrt{T}$, then the overall regret over $T$ rounds is of order $\sqrt{T}$. With some specific constant (chosen based on the analysis), this leads to the final algorithm called ILOVETOCONBANDITS (see Algorithm 1).

## 2 Oracle-Efficiency

To discuss oracle-efficiency, keep in mind that as in Policy Elimination, the true context distribution in the definition of $V$ can be replaced by the empirical distribution of observed contexts, that is, a

---

**Algorithm 1:** ILOVETOCONBANDITS (colloquially referred as Mini-monster)

---

**Input**: failure probability $\delta \in (0,1)$

**Initialization**: let $\mu = \min \left\{ \frac{1}{K}, \sqrt{\frac{\ln(TN/\delta)}{TK}} \ln T \right\}$

**for** $t = 1, \ldots, T$ **do**

$\quad$ find $P_t$ such that

$$\sum_{\pi \in \Pi} P_t(\pi) \text{Reg}_{t-1}(\pi) \leq 20K\mu$$

$$V(P_t, \pi) \leq 2K + \frac{\text{Reg}_{t-1}(\pi)}{10\mu} \quad \forall \, \pi \in \Pi.$$

$\quad$ play $a_t \sim P_t^\mu(\cdot|x_t)$

---

uniform distribution over $x_1, \ldots, x_{t-1}$ at time $t$. (For simplicity, the analysis of next section will assume that the true context distribution is known instead.)

According to the proof of Lemma 1, to find distribution $P_t$ it suffices to solve the optimization problem $\text{argmin}_{P \in \Delta(\Pi)} F_t(P)$ (in fact an approximate solution is enough). This is in fact very similar to FTRL with a special regularizer. To see how to solve it efficiently with the oracle, notice that the derivative of $F_t(P)$ with respect to a policy $\pi$ can be written as (with $\beta = 1/(10\mu)$)

$$\nabla F_t(P)(\pi) = \frac{1}{t-1} \sum_{\tau=1}^{t-1} \left( \widehat{\ell}_\tau(\pi(x_\tau)) - \widehat{\ell}_\tau(\pi_{t-1}^\star(x_\tau)) \right) - \frac{20\mu(1-K\mu)}{(t-1)} \sum_{\tau=1}^{t-1} \frac{1}{P^\mu(\pi(x_\tau)|x_\tau)}.$$

Since the part involving $\pi_{t-1}^\star$ is independent of $\pi$, if we feed the oracle with a training set

$$\mathcal{S} = \left\{ \left( x_1, \widehat{\ell}_1 - \frac{20\mu(1-K\mu)}{P^\mu(\cdot|x_1)} \right), \cdots, \left( x_{t-1}, \widehat{\ell}_{t-1} - \frac{20\mu(1-K\mu)}{P^\mu(\cdot|x_{t-1})} \right) \right\},$$

we have $\text{ERM}(\mathcal{S}) = \text{argmin}_{\pi \in \Pi} \nabla F_t(P)(\pi)$. In other words, the oracle can tell us the minimum coordinate of the gradient of $F_t(P)$ for any $P$, which opens up many possibilities to utilize the theory of optimization to find $P_t$. For example, one can directly apply the Frank-Wolfe algorithm (also known as conditional gradient method). Specifically, for a constraint convex optimization problem $\min_{w \in \Omega} f(w)$, the Frank-Wolfe algorithm performs the following iterative updates (staring with an arbitrary $w_1 \in \Omega$):

$$v_k = \underset{v \in \Omega}{\text{argmin}} \, \langle v, \nabla f(w_k) \rangle$$

$$w_{k+1} = (1 - \gamma_k)w_k + \gamma_k v_k$$

for some step-size $\gamma_k$ (default choice is $2/(k+1)$). When $\Omega$ is the simplex, the first step is exactly to find the minimum coordinate of the gradient. Therefore, with the oracle we can implement the Frank-Wolfe algorithm to solve $\text{argmin}_{P \in \Delta(\Pi)} F_t(P)$ efficiently. We omit the details on how many iterations are needed but it will be polynomial in $T$, $K$, and $\ln N$.

Importantly, notice that unlike gradient descent, Frank-Wolfe leads to a sparse solution: when $\Omega$ is the simplex, after $k$ rounds $w_k$ has only $k$ non-zero coordinates (assuming $w_1$ concentrates on one element to start with). This means that $P_t$'s are all sparse distributions and operations involving $P_t$, such as constructing the training set $\mathcal{S}$ and sampling $a_t$, are all efficient.

Instead of using Frank-Wolfe, another possibility is to do some kind of coordinate descent: iteratively use the oracle to fine the coordinate with minimum derivative and adjust the weight for this coordinate appropriately. This is exactly the method taken in [Agarwal et al., 2014]. In fact, with additional tricks that are specialized for this task, it was shown that over $T$ rounds only $\mathcal{O}(\sqrt{T})$ oracle calls are needed, which also implies that all $P_t$'s are $\mathcal{O}(\sqrt{T})$-sparse.

## 3 Regret Analysis

Finally in this section we prove that ILOVETOCONBANDITS enjoys optimal regret. The key is to show the following concentration results on regret.

**Lemma 2.** *With probability $1 - \delta/2$, Algorithm 1 ensures that for all $t \in [T]$ and all $\pi \in \Pi$,*

$$\text{Reg}(\pi) \leq 2\text{Reg}_t(\pi) + \epsilon_t \quad and \quad \text{Reg}_t(\pi) \leq 2\text{Reg}(\pi) + \epsilon_t$$

*where $\epsilon_t = \frac{20C}{\mu t} + 15K\mu$ and $C = \ln\left(\frac{4NT}{\delta}\right)\ln T$.*

*Proof.* By Freedman's inequality and a union bound, we have with probability $1 - \delta/2$, for all $t \in [T]$, all $\pi \in \Pi$, and any $\lambda \in [0, \mu]$,

$$|\bar{\ell}_t(\pi) - \bar{\ell}(\pi)| \leq \frac{\lambda}{t}\sum_{\tau=1}^{t} V(P_\tau, \pi) + \frac{\ln\left(\frac{4NT}{\delta}\right)}{\lambda t}.$$

Specifically picking $\lambda = \frac{\mu}{\ln T}$ gives

$$|\bar{\ell}_t(\pi) - \bar{\ell}(\pi)| \leq \frac{\mu}{t\ln T}\sum_{\tau=1}^{t} V(P_\tau, \pi) + \frac{C}{\mu t}. \tag{2}$$

Now we use induction to prove the lemma. The base case $t = 0$ is trivial. Assuming the statement holds for all rounds before time $t$, we have by the algorithm

$$V(P_\tau, \pi) \leq 2K + \frac{\text{Reg}_{\tau-1}(\pi)}{10\mu} \leq 2K + \frac{\text{Reg}(\pi)}{5\mu} + \frac{\epsilon_{\tau-1}}{10\mu} \tag{3}$$

for all $\tau = 2, \ldots, t$ and $V(P_1, \pi) \leq 2K$. Therefore, we have

$$
\begin{aligned}
\text{Reg}(\pi) - \text{Reg}_t(\pi) &= \bar{\ell}(\pi) - \bar{\ell}(\pi^\star) - \bar{\ell}_t(\pi) + \bar{\ell}_t(\pi_t^\star) \\
&\leq \bar{\ell}(\pi) - \bar{\ell}(\pi^\star) - \bar{\ell}_t(\pi) + \bar{\ell}_t(\pi^\star) && \text{(by optimality of } \pi_t^\star\text{)} \\
&\leq \frac{2C}{\mu t} + \frac{\mu}{t\ln T}\sum_{\tau=1}^{t}(V(P_\tau, \pi) + V(P_\tau, \pi^\star)) && \text{(by Eq. (2))} \\
&\leq \frac{2C}{\mu t} + \frac{4K\mu}{\ln T} + \frac{\text{Reg}(\pi)}{5\ln T} + \frac{1}{5t\ln T}\sum_{\tau=2}^{t}\epsilon_{\tau-1} && \text{(by Eq. (3) and } \text{Reg}(\pi^\star) = 0\text{)} \\
&\leq \frac{2C}{\mu t} + \frac{4K\mu}{\ln T} + \frac{\text{Reg}(\pi)}{5\ln T} + \frac{8C}{\mu t} + \frac{3K\mu}{\ln T} && \text{(by plugging in } \epsilon_\tau\text{)} \\
&\leq \frac{10C}{\mu t} + 7K\mu + \frac{\text{Reg}(\pi)}{2} \leq \frac{\epsilon_t}{2} + \frac{\text{Reg}(\pi)}{2}.
\end{aligned}
$$

Rearranging proves $\text{Reg}(\pi) \leq 2\text{Reg}_t(\pi) + \epsilon_t$. Similarly, we also have

$$
\begin{aligned}
\text{Reg}_t(\pi) - \text{Reg}(\pi) &= \bar{\ell}_t(\pi) - \bar{\ell}_t(\pi_t^\star) - \bar{\ell}(\pi) + \bar{\ell}(\pi^\star) \\
&\leq \bar{\ell}_t(\pi) - \bar{\ell}_t(\pi_t^\star) - \bar{\ell}(\pi) + \bar{\ell}(\pi_t^\star) && \text{(by optimality of } \pi^\star\text{)} \\
&\leq \frac{2C}{\mu t} + \frac{\mu}{t\ln T}\sum_{\tau=1}^{t}(V(P_\tau, \pi) + V(P_\tau, \pi_t^\star)) && \text{(by Eq. (2))} \\
&\leq \frac{2C}{\mu t} + \frac{4K\mu}{\ln T} + \frac{\text{Reg}(\pi)}{5\ln T} + \frac{\text{Reg}(\pi_t^\star)}{5\ln T} + \frac{1}{5t\ln T}\sum_{\tau=2}^{t}\epsilon_{\tau-1} && \text{(by Eq. (3))} \\
&\leq \frac{2C}{\mu t} + \frac{4K\mu}{\ln T} + \frac{\text{Reg}(\pi)}{5\ln T} + \frac{\epsilon_t}{5\ln T} + \frac{8C}{\mu t} + \frac{3K\mu}{\ln T} && (4) \\
&\leq \frac{14C}{\mu t} + 10K\mu + \text{Reg}(\pi) \leq \epsilon_t + \text{Reg}(\pi),
\end{aligned}
$$

where Step (4) uses the fact $\text{Reg}(\pi) \leq 2\text{Reg}_t(\pi) + \epsilon_t$ just proven above with $\pi$ set to $\pi_t^\star$, and also the fact $\text{Reg}_t(\pi_t^\star) = 0$. Rearranging then proves $\text{Reg}_t(\pi) \leq 2\text{Reg}(\pi) + \epsilon_t$ as well. $\qquad\square$

The final regret bound is now a simple application of this lemma and the exploitation constraint of the algorithm.

**Theorem 1.** *Algorithm 1 ensures that with probability $1 - \delta$, we have $\mathcal{R}_T = \widetilde{\mathcal{O}}\left(\sqrt{TK\ln(N/\delta)}\right)$.*

*Proof.* The first step is exactly the same as analyzing Policy Elimination: by Azuma's inequality we have with probability $1 - \delta/2$,

$$\sum_{t=1}^{T} \ell_t(a_t) \leq \sum_{t=1}^{T} \sum_{\pi \in \Pi} P_t(\pi)\bar{\ell}(\pi) + TK\mu + \mathcal{O}\left(\sqrt{T\ln(1/\delta)}\right).$$

Conditioning on this event and the event stated in Lemma 2, which happen simultaneously with probability $1 - \delta$, we have

$$\mathcal{R}_T \leq \sum_{t=1}^{T} \sum_{\pi \in \Pi} P_t(\pi)\text{Reg}(\pi) + TK\mu + \mathcal{O}\left(\sqrt{T\ln(1/\delta)}\right)$$

$$\leq 2 \sum_{t=1}^{T} \sum_{\pi \in \Pi} P_t(\pi)\text{Reg}_{t-1}(\pi) + \sum_{t=2}^{T} \epsilon_{t-1} + TK\mu + \mathcal{O}\left(\sqrt{T\ln(1/\delta)}\right)$$

$$\leq 56TK\mu + \frac{40C\ln T}{\mu} + \mathcal{O}\left(\sqrt{T\ln(1/\delta)}\right), \qquad \text{(by the exploitation constraint)}$$

which is of order $\widetilde{\mathcal{O}}\left(\sqrt{TK\ln(N/\delta)}\right)$ with the optimal tuning of $\mu$. $\qquad\qquad\square$

## References

Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.