April 22, 2020 :

Adversarial bandits: High probability bound
(from Bubeck & Bianchi)

- Issue - the variance of loss $\tilde{Y}_{i,t}$ is
$$O\left(1/p_{i,t}\right).$$

- One idea: Ensure that $p_{i,t}$ is not
too small, by playing exp-3 policy with
prob $(1-\varepsilon)$ & uniform with prob $\varepsilon$, so
that each arm is played with $p_{i,t}$
$> \varepsilon/k$.

= Want a regret of $\sqrt{n}$ ; Since in
each round you may collect
constant regret, $\varepsilon$ can't be too large.

In fact $\varepsilon \leq \frac{1}{\sqrt{n}}$, so that over

$n$ rounds the uniform dist will

Contribute at most $O(\sqrt{n})$ regret!

.But then the variance of cumulative

regret can be $\sqrt{n}$ per round & so

$\sqrt{n} \cdot n = n^{3/2}$; So this doesn't work;

= We work with gain $-g_{i,t} = 1 - l_{i,t}$

= Introduce a bias in the gain
estimate.

Lemma: Let $\beta \leq 1$ & Set

$$\tilde{g}_{i,t} = \frac{g_{i,t} \, \mathbb{1}_{E_t = i} + \beta}{p_{i,t}}$$

Then with prob $1-\delta$,

$$\sum_{t=1}^{n} g_{i,t} \leq \sum_{t=1}^{n} \widetilde{g}_{i,t} + \frac{\ln(1/\delta)}{\beta}$$

Proof: $\mathbb{E}_t$ - be conditional exp, conditional

on $Z_1, Z_2, \ldots, Z_{t-1}$

$$\underbrace{\mathbb{E}_t \exp\left(\overset{\leq 1}{\overbrace{\beta g_{it}}} - \beta\left(\boxed{\frac{g_{it} \mathbb{1}_{I_t=i} + \beta}{p_{it}}}\right)\right)}$$

$$= \mathbb{E}_t \exp\left(\beta g_{it} - \frac{\beta g_{it} \mathbb{1}_{I_t=i}}{p_{it}} - \frac{\beta^2}{p_{it}}\right)$$

$$= \mathbb{E}_t\left[\exp\left(\beta g_{it} - \frac{\beta g_{it} \mathbb{1}_{I_t=i}}{p_{it}}\right)\exp\left(\frac{-\beta^2}{p_{it}}\right)\right]$$

$$= \exp\left(-\beta^2/p_{it}\right) \mathbb{E}_t\left[\exp\left(\beta g_{it} - \underbrace{\frac{\beta g_{it} \mathbb{1}_{I_t=i}}{p_{it}}}_{X}\right)\right]$$

$$\exp(x) \le 1 + x + x^2 \quad \forall x \le 1.$$

$$\le \exp\left(\frac{-\beta^2}{p_{it}}\right)\left(1 + \mathbb{E}_t(x) + \mathbb{E}_t(x^2)\right)$$

$$\le \exp\left(\frac{-\beta^2}{p_{it}}\right)\left(1 + 0 + \frac{\beta^2 \, g_{it}^2}{p_{it}} - \underset{\underset{+ve}{\uparrow}}{(\ )}\right)$$

$$\le \exp\left(\frac{-\beta^2}{p_{it}}\right)\underbrace{\left(1 + \frac{\beta^2 \, g_{it}^2}{p_{it}}\right)}_{\le \; e\left(\frac{\beta^2 \, g_{it}^2}{p_{it}}\right)}$$

$$\le 1.$$

$$\therefore \; \mathbb{E}\left[\exp\left(\beta \overbrace{\sum_{t=1}^{n} g_{it} - \beta \sum_{t=1}^{m}\left(\frac{g_{it}\mathbb{1}_{I_t = i} + \beta}{p_{it}}\right)}^{X}\right)\right]$$

$$\le 1.$$

$$\Pr\left[x > \ln(\delta^{-1})\right] = \Pr\left[e^x > 1/\delta\right]$$

$$\le \delta \, \mathbb{E}[e^x] \le \delta.$$

$$\rightarrow \quad \mathbb{P}\left[ \beta \sum_{t=1}^{n} \tilde{g}_{it} - \beta \sum_{t=1}^{n} \frac{g_{it} \, \mathbb{1}_{I_t=i} + \beta}{p_{i,t}} \right.$$

$$\left. \leq \ln(1/\delta) \right]$$

$$\geq 1 - \delta;$$

---

<u>Exp 3.P:</u>

Input: $\eta, \delta, \beta \in (0,1)$

$p_1 \stackrel{a}{=}$ uniform dist on $1 \cdots, K$

for $t = 1, \cdots n$

  ① Draw arm $I_t$ from prob dist $p_t$

  ② Compute estimated gain,

$$\tilde{g}_{i,t} = \frac{g_{it} \, \mathbb{1}_{I_t=i} + \beta}{p_{i,t}}$$

- Update estimated cum gain

$$\widetilde{G}_{i,t} = \sum_{s=1}^{t} \tilde{g}_{i,s}$$

(3) $\quad p_{t+1} = \left( p_{t+1,1}, -- \quad , p_{t+1,k} \right)$

$$p_{i,t+1} = \frac{(1-\gamma)\, exp\left( \eta\, \widetilde{G}_{i,t} \right)}{\sum\limits_{k=1}^{K} exp\left( \eta\, \widetilde{G}_{k,t} \right)} + \frac{\gamma}{k}$$

and for:

$$\underline{\qquad\qquad} \times \overline{\qquad\qquad}$$

THM: Set $\beta = \sqrt{\dfrac{ln(K\delta^{-1})}{nK}}$ ; $\gamma = 1.05 \sqrt{\dfrac{K\, ln\, K}{n}}$

$$\eta = 0.95 \sqrt{\dfrac{ln\, K}{K}} ;$$

Get:

With prob $1-\delta$,

$$R_n \leq 5.15 \sqrt{n K \ln\left(K/\delta\right)}$$

## Contextual bandits:

- learner has access to extra info.

ex: movie recommendation:
- we should look at contextual info, past history of movies, and also the content/type of movie when making a recommendation!

- Need to devise algo's which use this contextual info.

- **Basic example:**

Bandits with side info;

A fixed set of contexts $\mathcal{C}$;
rounds are marked by contexts $c_1, \cdots, \in \mathcal{C}$;

learner must learn a mapping

$$g: \mathcal{C} \longrightarrow \{1, \cdots, K\}.$$

**Idea:** Run a different EXP3 on each context!

=

—

$$c = |\mathcal{E}|; \checkmark$$

Run $\overset{one}{\wedge}$ exp3 on each context. $\checkmark$

$n_c = \#$ times context $c$ is played;

$$R_n = \mathbb{E}\left[ \sum_{c \in \mathcal{C}} \max_{k \in K} \left( \sum_{\substack{t: \\ c_t = c}} l_{I_t,t} - l_{k,t} \right) \right]$$

$$= \sum_{\mathcal{E} \in \mathcal{E}} \max_{k=1\cdots K} \sum_{t: c_t = c} \left( l_{I_t,t} - l_{k,t} \right)$$

$$\leq \sum_{c \in \mathcal{E}} \sqrt{2 n_c \, K \ln K}$$

concave $\sqrt{\phantom{x}}$

$$= |\mathcal{E}| \sum_{c \in \mathcal{E}} \frac{1}{|\mathcal{C}|} \sqrt{2 n_c K \ln K}$$

$$\boxed{\begin{array}{c} \theta(g(x)) \\ = \sqrt{x}. \\ \leq g\,\theta(x). \end{array}}$$

$$\leq |\mathcal{E}| \sqrt{\sum_c \frac{2 n_c K \ln K}{|\mathcal{E}|}}$$

$$= |\mathcal{G}| \sqrt{\frac{2n K \ln K}{|\mathcal{G}|}} = \sqrt{2n K |\mathcal{E}| \ln K}.$$

- Playing with experts:

when $|\mathcal{E}|$ is large - bad idea!

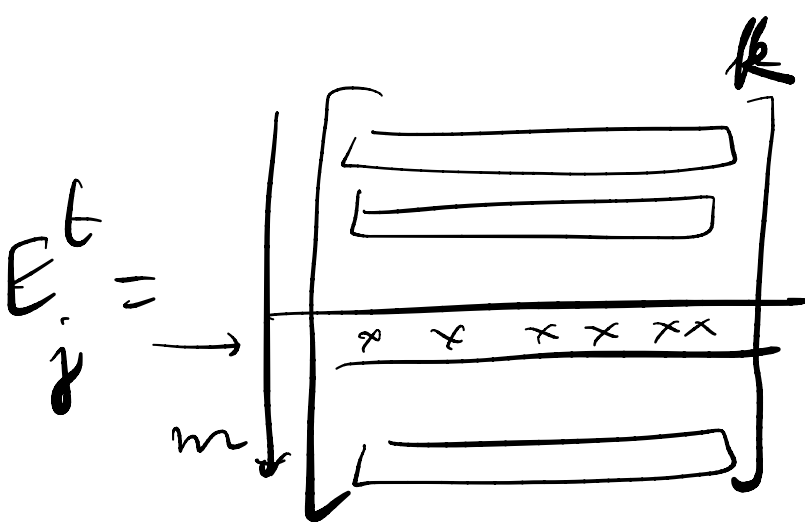- Users with similar demographics like similar movies!
∴ Contexts are structured!

set $R_n = \mathbb{E}\left[\max_{\phi \in \Phi} \sum_{t=1}^{n}\left(x_{t,\phi(t)} - X_t\right)\right]$ $\Phi = \{f: \mathcal{E} \to k\}$

- At the beginning of each round experts announce their predictions!

- In fact experts give a prob dist over actions, (experts are randomized)

- The expert advice is sound ⊢

$$E^t_j = \quad\xrightarrow{\hspace{1cm}}$$



$R_n$ — measured w.r.t best export in hindsight;

$$R_n = \mathbb{E}\left[\max_{m \in M} \sum_{t=1}^{n} E^t_m x_t - \sum x_t\right]$$

$$x_t = (x_{t1}, \ldots, x_{tk})$$

Exp 4:

Expert; Exp with, Exp beg Explicit

Input: $n, k, M, \eta, \tau$

2) $Q_1 = \left( 1/M, \quad \cdots, \quad 1/M \right)$

3) for $t = 1, \cdots n$

4) Receive advice $E^t$

5) Choose $A_t \sim P_t,$ $\quad P_t = Q_t E^t$
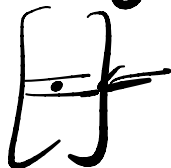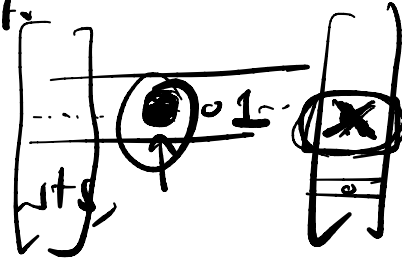
6) Receive reward $X_t = a_{t A_t}$

7) Estimate $a_t!$

$$\hat{X}_{ti} = 1 - \frac{\mathbb{1}\{A_t = i\}(1 - X_t)}{P_{ti} + \tau}$$

8) Propagate $\cdots$ to experts

$$\tilde{X}_t = \boxed{E^t \hat{X}_t}$$

9) Update $Q_t$ using exp wts

$$\rho_{t+1,i} = \frac{\exp\left(\eta \, \tilde{x}_{t,i}\right) \rho_{t,i}}{\sum_{j=1}^{m} \exp\left(\eta \, \breve{x}_{t,j}\right) \rho_{t,j}}$$

to find:

—— $\times$ —— $\to$ $\lambda$ —— $\hat{} $ —— $\hat{}$ ——

$$\sqrt{nK \ln K \, |\mathcal{E}|}$$

$\uparrow$

$\log |\mathcal{E}|$.

—— $\times$ ——