# Contextual Bandits.

$n$ arms

for $t = 1, 2, \ldots$

    Nature reveals context $c_t \in \mathcal{C}$

    Player chooses $a_t \in [n]$

    Receives loss $\ell(c_t, a_t)$

Objective: Minimize loss $\sum_{t=1}^{T} \ell(c_t, a_t)$

First idea: ignore context and run EXP-3.

$$\Rightarrow \quad \text{regret} \quad \max_a \sum_{t=1}^{T} \ell(c_t, a_t) - \ell(c_t, a) \leq \sqrt{nT}$$

$$\text{Total loss} = \sum_{t=1}^{T} \ell(c_t, a_t) \leq \left( \min_a \sum_{t=1}^{T} \ell(c_t, a) \right) + \sqrt{nT}$$

Second idea: Assume $|\mathcal{C}|$ is small and instantiate an independent copy of EXP-3 for each context.

$$\text{Regret} = \sum_{c \in \mathcal{C}} \max_{a \in [n]} \sum_{t: c_t = c} \left( \ell(c_t, a) - \ell(c_t, a_t) \right)$$

$$\leq \sum_{c \in \mathcal{C}} \sqrt{T_c \cdot n} = (1, \ldots, 1) \begin{pmatrix} \sqrt{T_1 n} \\ \vdots \\ \sqrt{T_k n} \end{pmatrix}$$

$$\leq \sqrt{|\mathcal{C}| \cdot \sum_{c \in \mathcal{C}} T_c n} = \sqrt{n \cdot |\mathcal{C}| \cdot T}$$

$$\text{Total loss} \quad \sum_{t=1}^{T} \ell(c_t, a_t) \leq \left( \sum_{c \in \mathcal{C}} \min_a \sum_{t: c_t = c} \ell(c, a) \right) + \sqrt{n \cdot |\mathcal{C}| \cdot T}$$

We have a collection of policies $\Pi$ s.t.

$\pi \in \Pi$ $\qquad \pi : \mathcal{C} \to \Delta_n$ (if stochastic) $\pi : \mathcal{C} \to [n]$ (if deterministic)

Policy regret $\displaystyle \max_{\pi \in \Pi} \sum_{t=1}^{T} \ell(c_t, \pi) - \ell(c_t, \pi_t)$

where at each time player chooses $\pi_t \in \Pi$

then (optionally) draws $a_t \sim \pi_t(c_t)$, $\ell(c, \pi) = \mathbb{E}_{a_t \sim \pi_t(c_t)}[\ell(c, a_t)]$

<span style="color:red">In the case of all contexts to all actions</span> $\underline{|\Pi| = n^{|c|}}$

Enumerate policies $\bar{\pi}_1, \ldots, \bar{\pi}_{|\Pi|}$ and define

$$\ell_{t,i} = \mathbb{E}_{a \sim \bar{\pi}_i(c_t)}\left[\ell(c_t, a)\right]$$

$$(\pi_t \sim a_t)$$

We observe $\ell(c_t, a_t)$, $a_t \sim \pi_{I_t}(c_t)$, $I_t \sim q_t$

$$\hat{\ell}_{t,i} = \frac{\pi_i(a_t \mid c_t)}{\displaystyle\sum_{k=1}^{|\Pi|} q_{t,k} \, \bar{\pi}_k(a_t \mid c_t)} \, \ell(c_t, a_t)$$

$$\mathbb{E}\left[\hat{\ell}_{t,i}\right] = \sum_{j} \mathbb{P}(a_t = j \mid c_t) \frac{\pi_i(j \mid c_t)}{\sum_{n} q_{s,n} \pi_n(j \mid c_t)} \ell(c_t, j)$$

$$= \sum_{j} \pi_i(j \mid c_t) \ell(c_t, j) = \mathbb{E}_{a \sim \pi_i(c_t)}\left[\ell(c_t, a)\right]$$

$$= \ell_{t,i}$$

$$\mathbb{P}(a_t = j \mid c_t) = \sum_{s=1}^{|\pi|} \mathbb{P}(a_t = j, I_t = s \mid c_t)$$

$$= \sum_{s=1}^{|\pi|} \mathbb{P}(a_t = j \mid I_t = s \mid c_t) \mathbb{P}(I_t = s \mid c_t)$$

$$= \sum_{s=1}^{|\pi|} \pi_s(j \mid c_t) \cdot q_{t,s}$$

<div style="border:1px solid">

**EXP3($\gamma$): Exponential Weights for Exploration Exploitation**

**Input:** Time horizon $T$, $n$ arms, $\eta > 0$, $\gamma \in [0,1]$

**Initialize:** Player sets $p_1 = (1/n, \ldots, 1/n) \in \triangle_n$. Adversary chooses $\{\ell_t\}_{t=1}^T \subset [-1,1]^n$.

**for:** $t = 1, \cdots, T$   ← Nature reveals context $C_t$

     Player defines $\lambda_t \in \triangle_n$ and plays $I_t \sim q_t := (1-\gamma)p_t + \gamma\lambda_t$,   $a_t \sim \pi_{I_t}(C_t)$

     Player suffers (and observes) loss $\ell_{t,I_t}$ (but does *not* observe $\ell_{t,i}$ for $i \neq I_t$)

     Player computes loss estimator $\widehat{\ell}_{t,i}$ for all $i \in [n]$

     Update iterates:

$$w_{t+1,i} = w_{t,i}\exp(-\eta\widehat{\ell}_{t,i}) \qquad p_{t+1,i} = w_{t+1,i}/\sum_{j=1}^n w_{t+1,j}.$$

</div>

**Theorem 17.** *Fix any sequence* $\ell_t \in [-1,1]$ *for all* $t$. *For any* $\widehat{\ell}_{t,i}$ *and* $\eta, \gamma \geq 0$ *that satisfy* $\mathbb{E}[\widehat{\ell}_{t,i}|\mathcal{F}_{t-1}] = \ell_{t,i}$ *and* $-\eta\widehat{\ell}_{t,i} \leq 1$ *for all* $i,t$ *we have*

$$\max_{i=1,\ldots,n} \mathbb{E}\left[\sum_{t=1}^T \ell_{t,I_t} - \ell_{t,i}\right] \leq 2\gamma T + \frac{\log(\overset{|\pi|}{n})}{\eta} + (1-\gamma)\eta\mathbb{E}\left[\sum_{t=1}^T\sum_{j=1}^{\overset{|\pi|}{n}} p_{t,j}\widehat{\ell}_{t,j}^2\right].$$

After opt.

$\leq nT$   for $\eta$ get

$$-\eta\widehat{\ell}_{t,i} \leq \eta \frac{\pi_i(a_t|C_t)}{\sum_h \pi_h(a_t|C_t)\,q_{t,h}}$$

$$\leq \eta\sum_{j=1}^n \frac{\pi_i(j|C_t)}{\sum_h \pi_h(j|C_t)\,q_{t,h}}$$

$$\leq \eta\sum_{j=1}^n \frac{\pi_i(j|C_t)}{\sum_h \pi_h(j|C_t)\,\gamma\lambda_{t,h}} \leq \frac{\eta n}{\gamma} \leq 1 \text{ if } \gamma = \eta n$$

$\sqrt{nT\log|\pi|}$

**Claim** For $P_i \in \triangle_n$ for $i = 1,\ldots,m$ we have

$$\min_{\lambda \in \triangle_m} \max_{i=1,\ldots,m} \sum_{j=1}^n \frac{P_{i,j}}{\sum_h \lambda_h P_{h,j}} = n$$