n alternatives / treatments / arms

Input: $n$ arms described by distributions $\nu_i$, $i = 1, \ldots, n$

for $t = 1, 2, \ldots, T$     user $t$ arrives at nytimes.com

Player chooses $I_t \in [n] := \{1, \ldots, n\}$

Nature reveals $X_{I_t, t} \overset{iid}{\sim} \nu_{I_t}$

Example: $\nu_i = \mathcal{N}(\theta_i^*, 1)$     $X_{i,t} \in \mathbb{R}$     $\forall t$

Example: $\nu_i = \text{Bernoulli}(\theta_i^*)$     $X_{I_t, t} \in \{0, 1\}$     $\forall t$

We do not know $\{\nu_i\}_{i=1}^{n}$, but we typically assume some property. For example $\text{support}(\nu_i) \in [0,1]$, $\nu_i$ is 1-sub-Gaussian.

Regret Minimization / Reward Accumulation Maximization

Goal: maximize $\mathbb{E}\left[ \sum_{t=1}^{T} X_{I_t, t} \right]$

Def) Regret $= \max_{i=1,\ldots,n} \mathbb{E}\left[ \sum_{t=1}^{T} X_{i,t} - \sum_{t=1}^{T} X_{I_t, t} \right]$

Best action in hindsight        What player recieves

Goal: minimize regret.

Obtain sub-linear regret:

$\text{Regret}(T) = o(T)$     $\left( e.g. \; R_T = O(\sqrt{T}) \right)$

**Proposition** $\Big]$ $R_T := \max\limits_{i=1,\ldots,n} \mathbb{E}\left[ \sum\limits_{t=1}^{t} X_{i,t} - \sum\limits_{t=1}^{T} X_{I_t,t} \right]$

$$= \sum\limits_{i=1}^{n} \Delta_i \, \mathbb{E}\left[ T_i \right]$$

where $\Delta_i := \max\limits_{s=1,\ldots,n} \theta_s^* - \theta_i^*$, $T_i$ is # times arm

$i$ is played up to time $T$. $\theta_i^* = \mathbb{E}\limits_{X \sim \nu_i}\left[ X \right]$.

**Proof** $\Big]$

$\max\limits_{i=1,\ldots,n} \mathbb{E}\left[ \sum\limits_{t=1}^{t} X_{i,t} - \sum\limits_{t=1}^{T} X_{I_t,t} \right]$

$= \max\limits_{i} \sum\limits_{t=1}^{T} \theta_i^* - \sum\limits_{t=1}^{T} \mathbb{E}\left[ X_{I_t,t} \right]$

$= T \cdot \max\limits_{i} \theta_i^* - \sum\limits_{t=1}^{t} \mathbb{E}\left[ \underbrace{\mathbb{E}\left[ X_{I_t,t} \mid I_t = i \right]}_{\theta_{I_t}^*} \right]$

$= T \cdot \max\limits_{i} \theta_i^* - \sum\limits_{t=1}^{t} \sum\limits_{i=1}^{n} \mathbb{E}\left[ \underline{\mathbb{1}\{I_t = i\}} \, \theta_i^* \right]$

$= T \cdot \max\limits_{i} \theta_i^* - \sum\limits_{i=1}^{n} \theta_i^* \underbrace{\sum\limits_{t=1}^{T} \mathbb{E}\left[ \mathbb{1}\{I_t = i\} \right]}$

$= \mathbb{E}\left[ \sum\limits_{t=1}^{T} \mathbb{1}\{I_t = i\} \right]$

$= \mathbb{E}\left[ T_i \right]$

$$= T \cdot \max_i \theta_i^* - \sum_{i=1}^n \theta_i^* \, \mathbb{E}[T_i]$$

$$= \sum_{i=1}^n \mathbb{E}[T_i]$$

$$= \sum_{i=1}^n \left( \max_j \theta_j^\emptyset - \theta_i^* \right) \mathbb{E}[T_i]$$

$$= \Delta_i$$

# Best arm identification

## (adaptive A/B/n testing)

Input confidence $\delta \in (0,1)$.

Objective is for an algorithm to pull arms and stop ASAP and output $\hat{i} \in [n]$

where $\mathbb{P}\left( \hat{i} = \operatorname*{argmax}_{i \in \{1,\dots,n\}} \theta_i^\emptyset \right) \geq 1 - \delta$.

$(\varepsilon, \delta)$- PAC identification

Identify an arm $\hat{i}$ : $\max_j \theta_j^* - \theta_{\hat{i}}^\emptyset \leq \varepsilon$

Top-$k$ : identify top $k$ means

Multiple testing / threshold bandits

Combinatorial bandits

$c \in C$ $\qquad C \subset [n]$

$\underset{c \in C}{argmax} \sum_{i \in c} \theta_i^*$

Top-$k$
Threshold bandits
Matching
Routing

Max-bandits $\qquad$ Pull $k$ arms and observe $\underset{i \in S}{max} X_{i,t}$ $\qquad |S| = k$