# Framework and Regret.

Ref: Chap 4, Bandit Algorithms.

1. Basic setting — unstructured / structured
2. Regret — frequentist / Bayesian

3. Regret decomposition.

Notation: 1. $A$ is the set of arms (aka actions)

a. Choose any $a \in A$
b. Receive reward drawn from an (unknown) distribution $P_a$.

2. $\nu =$ joint distribution across arms.

e.g: arms are independent $\Rightarrow$ $\nu$ is a product measure
$= P_{a_1} \times P_{a_2} \times \cdots$ if arms are discrete.

i.e., $\nu = (P_a, a \in A)$ $\rightarrow$ notation for product dist.

| Name | Symbol | Definition |
|---|---|---|
| Bernoulli | $\mathcal{E}_{\mathcal{B}}^k$ | $\{(\mathcal{B}(\mu_i))_i : \mu \in [0,1]^k\}$ |
| Uniform | $\mathcal{E}_{\mathcal{U}}^k$ | $\{(\mathcal{U}(a_i, b_i))_i : a, b \in \mathbb{R}^k \text{ with } a_i \le b_i \text{ for all } i\}$ |
| Gaussian (known var.) | $\mathcal{E}_{\mathcal{N}}^k(\sigma^2)$ | $\{(\mathcal{N}(\mu_i, \sigma^2))_i : \mu \in \mathbb{R}^k\}$ |
| Gaussian (unknown var.) | $\mathcal{E}_{\mathcal{N}}^k$ | $\{(\mathcal{N}(\mu_i, \sigma_i^2))_i : \mu \in \mathbb{R}^k \text{ and } \sigma^2 \in [0, \infty)^k\}$ |
| Finite variance | $\mathcal{E}_{\mathbb{V}}^k(\sigma^2)$ | $\{(P_i)_i : \mathbb{V}_{X \sim P_i}[X] \le \sigma^2 \text{ for all } i\}$ |
| Finite kurtosis | $\mathcal{E}_{\text{Kurt}}^k(\kappa)$ | $\{(P_i)_i : \text{Kurt}_{X \sim P_i}[X] \le \kappa \text{ for all } i\}$ |
| Bounded support | $\mathcal{E}_{[a,b]}^k$ | $\{(P_i)_i : \text{Supp}(P_i) \subseteq [a,b]\}$ |
| Subgaussian | $\mathcal{E}_{\text{SG}}^k(\sigma^2)$ | $\{(P_i)_i : P_i \text{ is } \sigma\text{-subgaussian for all } i\}$ |

**Table 4.1** Typical environment classes for stochastic bandits. $\text{Supp}(P)$ is the (topological) support of distribution $P$. The kurtosis of a random variable $X$ is a measure of its tail behavior and is defined by $\mathbb{E}[(X - \mathbb{E}[X])^4]/\mathbb{V}[X]^2$. Subgaussian distributions have similar properties to the Gaussian and will be defined in Chapter 5.

Bandit Algorithms, Lattimore and Szepesvari, 2019 (pp. 58)

3.  $\mathcal{E}$ = environment = set of possible distributions

i.e., $\nu \in \mathcal{E}$

Stochastic Bandit Problem!

* nature/adversary chooses a fixed $\nu \in \mathcal{E}$

→ this is unknown to the player.

* At each (discrete) time $t$, player chooses

arm $A_t \in \mathcal{A}$. Nature then uses predetermined $\nu$ to determine arm dist $P_{A_t}$, and draws a sample $X_t \sim P_{A_t}$

\* $X_t =$ reward at time $t$

\* Player thus observes $(A_t, X_t)$ at time $t$.

\* Player uses all _causal_ information:

$$\left\{ (A_s, X_s), \; s = 1, 2, \cdots, t \right\}$$

to determine action $A_{t+1} \in \mathcal{A}$.

Two settings:

① Unstructured environment $\mathcal{E}$ : $\nu$ is a product measure, i.e; arms <u>do not</u> provide any information about each other.

$$\nu = (P_a, \; a \in \mathcal{A}) = P_{a_1} \times P_{a_2} \cdots \quad \text{(discrete arms)}$$

② Structured environment $\mathcal{E}$: arms "leak" information.

e.g: linear bandits: $\theta \in \mathbb{R}^d$ fixed but unknown

$$\mathcal{E} = \{ \nu_\theta, \theta \in \mathbb{R}^d \}$$

where $\nu_\theta = \{ P_{a,\theta}, a \in \mathcal{A} \}$.

with $P_{a,\theta} = N(\langle a, \theta \rangle, \sigma^2)$

or $\text{Bernoulli}(\langle a, \theta \rangle)$
.

Here, given $\theta$, arm distributions are known for all arms.

Since $\theta \in \mathbb{R}^d$, there are efficient ways to estimate $\theta$ from $\tilde{O}(d)$ samples under suitable assumptions.

**Regret:** The performace metric.

(i) **Frequentist viewpoint:** $\nu \in \mathcal{E}$ picked by nature.

A _genie_ is given this information. Then, the genie uses this to determine best action/arm.

Define $\{X_t^*, t = 1, 2, \cdots\}$ the associated reward for the genie.

e.g.: Finite $K$ arms, unstructured environment, with $\nu = P_{a_1} \times P_{a_2} \cdots \times P_{a_K}$

with $P_{a_i} \sim \text{Bernoulli}(\mu_i)$, $\mu_i \in [0, 1]$.

Then, WLOG, let $\mu_1 \geq \mu_2 \geq \mu_3 \cdots$.

In this case, genie is given knowledge of $\{\mu_1, \cdots \mu_K\}$ and the arm dist. (Bernoulli)

Thus, $X_t^* \sim \text{Bernoulli}(\mu_1)$.

The **PLAYER** does not have knowledge of $\mu_i$, $i = 1, 2, \cdots, K$.

Plays action $A_s$ at time $s$, and observes $X_s \sim P_{A_s}$, $\quad s = 1, 2, \cdots, t$.

$$\text{Regret} = R_t(\pi, v) = E\left[\sum_{s=1}^{t} X^*_s\right] - E\left[\sum_{s=1}^{t} X_s\right]$$

environment

time horizon

player policy

Genie expected reward until time $t$

Player expected reward until time $t$

_e.g:_

Continuing our example,

$$E\left[\sum_{s=1}^{t} X^*_s\right] = \mu_1 t \quad.$$

Notation: $\mu^*(v) = E\left[X^*_1\right]$ under genie policy that knows $v$.

Sub-optimality gap: $\Delta_a(v) = \mu^*(v) - \mu_a(v)$

where $\qquad \mu_a(v) = E_{P_a}[Z] \qquad Z \sim P_a$

$T_a(t) =$ Number of times arm $a$ has been played until time $t$ under policy $\pi$.

$$= \sum_{s=1}^{t} \chi_{\{A_s = a\}}$$

indicator function

---

Lemma  (Regret Decomposition)

Suppose that $|A|$ is countable. Then:

$$R_t(\pi, v) = \sum_{a \in A} \Delta_a(v) \cdot T_a(t).$$

---

Bayesian Regret: Average regret when we have some prior knowledge on how $v \in \mathcal{E}$ is chosen by the environment.

More specifically, suppose $\nu \in \mathcal{E}$ is chosen using some distribut $Q$ with support over $\mathcal{E}$.

Then, for some player policy $\pi$,

$$BR_t(\pi, Q) = \text{Bayesian regret using prior } Q$$

$$= E_{\nu \sim Q}\left[R_t(\pi, \nu)\right]$$

From now on (until towards end of semester), we will focus on **FREQUENTIST REGRET.**

Notation abuse: We will use $R_t$ and hide the dependence on $(\pi, \nu)$ in most cases