# Formal Setting: MAB

Simplified RL: no states & no transitions

$$\mathcal{A}: 1, \ldots, K \quad \text{"arms"}$$

$$r : \mathcal{A} \to \Delta(\mathbb{R}) \quad \text{noisy} \quad r_t \sim r(a_t)$$
$$r(r_t | a_t) \in [0, 1]$$

$$\mathbb{E}(r(a)) = \mu_a$$

$$T : \mathbb{Z}_+ \quad \text{integer time} \quad \text{horizon}$$

__Goal__: maximize cumulative expected reward

$$\mathbb{E}\left[\sum_{t=1}^{T} r(a_t)\right] = \sum_{t=1}^{T} \mu_{a_t}$$

__Optimal action:__
$$a^* = \underset{a=1,\ldots,K}{\arg\max} \; \mu_a$$

Devise an algorithm for balancing __exploration__ and __exploitation__

## Definition (Regret)

The regret of an algorithm which choosed $a_1, \ldots, a_T$

$$R(T) = \mathbb{E}\left[\sum_{t=1}^{T} r(a^*) - r(a_t)\right]$$

$$= \sum_{t=1}^{T} \mu^* - \mu_{a_t}$$

Want to find an algorithm __sublinear__ in regret
$$R(T) \underset{\sim}{\propto} T^p \qquad p < 1$$

$$\lim_{T \to \infty} \frac{R(T)}{T} \to 0 \quad \text{in this case}$$

# Balancing exploration & exploitation

### Alg 1: Random
for $t = 1, \dots, T$:
$$a_t \sim \text{unif}(1, \dots, K)$$

### Alg 2: Greedy
Try each arm once
compute $\hat{y}_{a_t} = r_t$
for $t = K+1, \dots T$
$$a_t = \arg\max_a \hat{y}_a$$

Both suffer from <u>linear</u> regret

Why?
$$R(T) = \sum_{t=1}^{T} \mathbb{E}\left[ r(a^*) - r(a_t) \right]$$

$$= \sum_{t=1}^{T} \mathbb{E}\left[ (r(a^*) - r(a_t)) \mathbb{1}\underline{\{a_t \neq a_*\}} \right]$$

$$= \sum_{\substack{t: \\ a_t \neq a^*}} (y^* - y_{a_t}) \mathbb{P}\{ a_t \neq a_* \}$$

$$\geq \sum_{t=1}^{T} \left[ \underbrace{\min_{a \neq a^*}(y^* - y_a)}_{\substack{\text{constant} \\ \text{wrt } t}} \right] \cdot \underline{\mathbb{P}\{ a_t \neq a_* \}} \geq cT$$

## Alg 3: Explore-then-commit

Pull each arm $N$ times $(t = 1, \dots, NK)$  $\Big\}$ exploration
compute $\hat{y}_a$ as average observed reward

for $t = NK+1, \dots, T$  $\Big\}$ exploitation (commit)
$$a_t = \arg\max_a \hat{y}_a = \hat{a}^*$$

$$R(T) = \underbrace{\sum_{t=1}^{NK} \mathbb{E}(r(a^*) - r(a_t))}_{R_1} + \underbrace{\sum_{NK+1}^{T} \mathbb{E}(r(a^*) - r(a_t))}_{R_2}$$

Claim: $R_1 \leq NK$ for reward bounded $[0, 1]$