# K-Armed bandit problem

We face repeatly with a choice among K options and after every choice we receive a Reward from a probability distribution.

GOAL → MAX total reward over some time period

ACTION VALUE: Value of selecting on action a is:

$$q_*(a) = E[R_t \mid A_t = a]$$

## How to estimate action value?

## 1. SAMPLE-AVG METHOD

We don't know the reward distribution:

$$Q_t(a) = \frac{\sum_{i=1}^{t-1} R_i \, \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}} = \frac{\Sigma \, R \text{ when take action } a}{\Sigma \text{ we take action } a}$$

## How to select action?

- RANDOM
- GREEDY → Take the max estimated value → SUBOPTIMAL

$$A_t = \arg\max_a Q_t(a)$$

- $\varepsilon$- GREEDY $\rightarrow$ Behave greedy most of the time but sometimes select randomly from all action wit equal probability

$$A_t = \begin{cases} \underset{a}{\arg\max} \; Q_t(a) & \text{with prob} = 1-\varepsilon \\ a \sim U(\{a_1,...,a_t\}) & \text{oth.} \end{cases}$$

# 2.1 INCREMENTAL METHOD (Stationary)

$$Q_n = \frac{R_1 + ... + R_{n-1}}{n-1} = \begin{array}{l} \text{estimate of the action value} \\ \text{of the action that has } R_i \end{array}$$

$$Q_{n+1} = \frac{1}{n} \sum_{i=1}^{n} R_i = ... = \boxed{Q_n + \frac{1}{n} [R_n - Q_n]}$$

New Estimate = Old Estimate + StepSize $[\underbrace{\text{Target} - \text{OldEst}}_{\text{error estimate}}]$

Bandit Algo: $Q(a) = 0 = N(a)$

While True ()

$$A = \begin{cases} \arg\max_a Q(a) & \text{if } P = 1-\varepsilon \\ a \text{ random} & \text{else} \end{cases}$$

$R = $ bandit $(a) \rightarrow$ take on action and return a reward

$N(A) += 1$

$Q(A) = Q(A) + 1/N(A) \; [R - Q(A)]$

# 2.2 INCREMENTAL METHOD (Non Stationary)

Non stationary → the reward probability **changes** over time

- Mantain a costant step size:

$$Q_{n+1} = Q_n + \alpha [R_n - Q_n]$$
$$\hookrightarrow (0,1]$$

- OPTIMISTIC INITIAL VALUE: Initial action values can be used to _IMPROVE EXPLORATION_ so that the system will do a good amount of exploration.