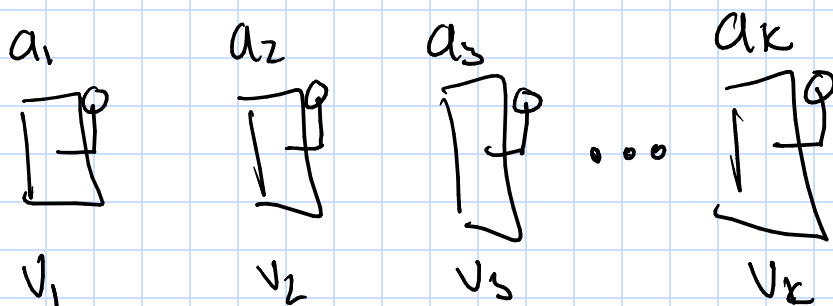## Admin

- Homework #5 due 04/30/2021   (Semi Supervised learning)

- Final Project Offically due  05/05/2021

    - Both partners must submit he project
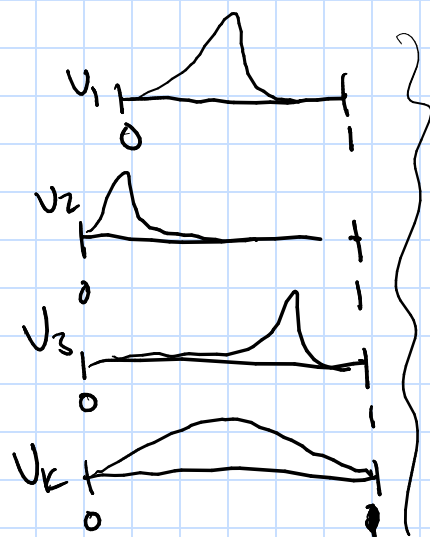
    - IEEE format

- Exam 04/09 - 04/12

## Multi arm bandits

The multi arm bandit (MAB) is a sequential allocation problem defined by a set of actions. Our goal is to maximize our reward.

$a_1$      $a_2$      $a_3$          $a_k$

$V_1$      $V_2$      $V_3$          $V_k$

$$M_1 = 0.9, \quad \sigma^2 = 0.01$$

$$M_2 = 0.92, \quad \sigma^2 = 0.1$$

## Approaches

- Randomized

- Round Robin

# The Stochastic Bandit Problem

Known Parameters: number of arms $k$, number of rounds $n > k$

Unknown Parameters: $k$ probability distributions $v_1, \ldots, v_k$ on $[0,1]$

for $t = 1, \ldots, n$

   ① The forecaster chooses $I_t \in \{1, \ldots, k\}$

   ② Given $I_t$, the environment draws reward $X_{I_t, t} \sim V_{I_t}$ independently from the past and reveals it to the forecaster

## Defs:

$I_t$: arm sampled at time $t$ from $\{1, \ldots, k\}$

$t$: round of play

$v$: distributions

$X_{I_t, t}$: the reward sampled from $V_{I_t}$ at time $t$

$V_i$ has a mean $\mu_i$

$\mu^* = \max\limits_{i \in [k]} \mu_i$, $\quad i^* = \arg\max\limits_{i \in [k]} \mu_i$

## Regret

$$\max_{i \in [k]} \sum_{t=1}^{n} X_{i,t} - \sum_{t=1}^{n} X_{I_t,t} = R_n$$

## Pseudo Regret

$$\max_{i \in [k]} \mathbb{E}\left[ \sum_{t=1}^{n} X_{i,t} - \sum_{t=1}^{n} X_{I_t,t} \right] = \tilde{R}_n$$

In a stochastic setting, it is easy to show that

$$\bar{R}_n = n\mu^* - \sum_{t=1}^{n} \mathbb{E}[\mu_{I_t}] \qquad \text{\#}$$

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

## UCB1

$$\bar{X}_i + \sqrt{\frac{2\log(t)}{n_i}}$$

$\bar{X}_i \rightarrow$ average samples from machine $i$

$n_i \rightarrow$ # of times we sampled machine $i$