# INTRODUCTION TO REINFORCEMENT LEARNING

Nathanaël Fijalkow

$$RL \subseteq ML$$

ML:

$$\left.\begin{array}{l} \text{Data} \\ \text{Background} \\ \text{knowledge} \end{array}\right\} \longrightarrow \text{Model} \longrightarrow \left\{\begin{array}{l} \text{Action} \\ \text{Prediction} \end{array}\right.$$

## Supervised learning:

Training data is **labelled**
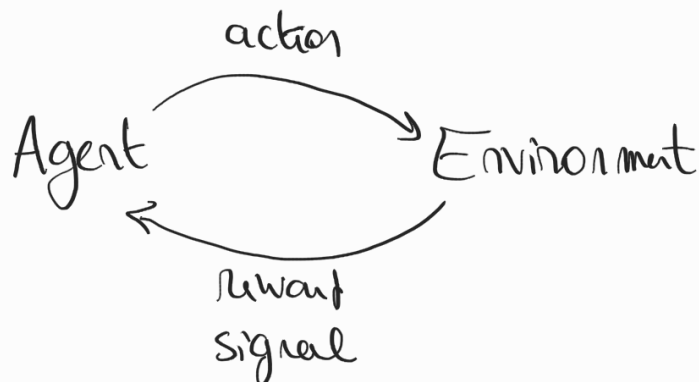
Objective: predict unlabelled data
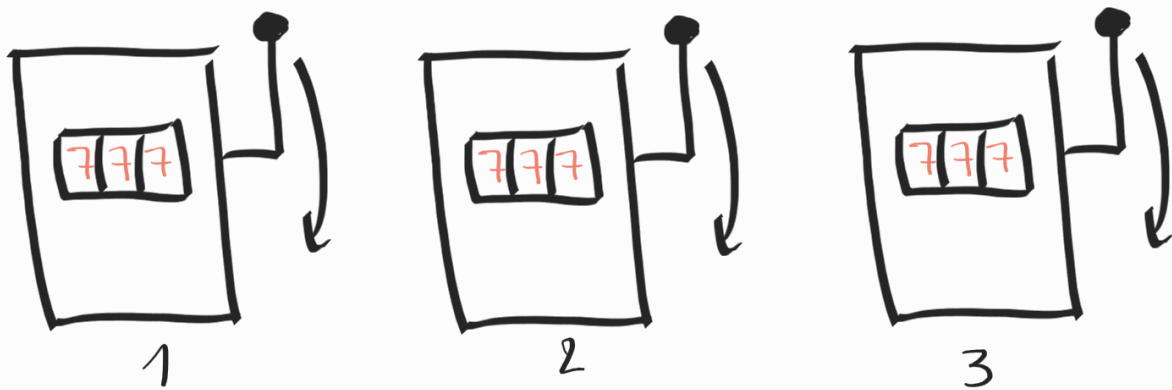
## Unsupervised learning:

Training data is **not labelled**

Objective: finding structure / pattern

**RL** is neither supervised nor unsupervised!

RL: | active learning
    | reward-based

action
Agent ⟶ Environment
reward
signal

# MULTI ARMED BANDITS



1          2          3

arm = machine = action

K machines : $1, \dots, K$

We assume that each machine $i$ has a
reward distribution $\delta_i$

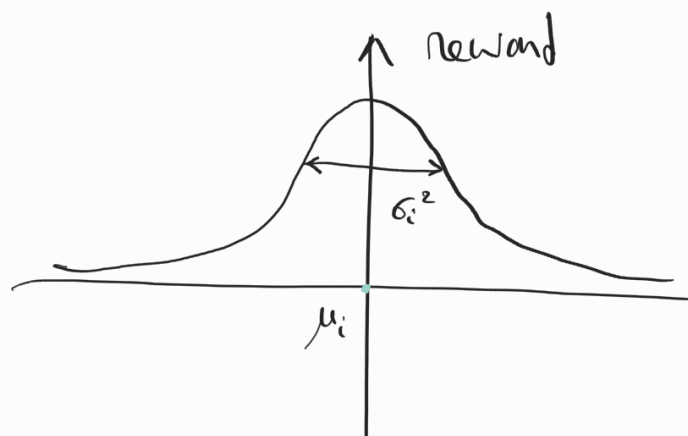Bernoulli distribution $\qquad B(p)$

$$\begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1-p \end{cases}$$

Normal distribution $\qquad \mathcal{N}(\mu, \sigma^2)$



Notation: $\quad \mu = \mathbb{E}[\delta] \qquad \sigma^2 = \text{Var}(\delta)$

**Scenario:**

- We know the number of machines but
  NOT THE REWARD DISTRIBUTIONS

- at each time step $t$ we pick an arm $i$
  and draw a reward $r(t) \in \mathbb{R}$ from $\delta_i$

- Two (similar) goals:

  $\rightarrow$ (1) identify the best arm: $\underset{i}{\text{argmax}} \; \mu_i$

  $\rightarrow$ (2) maximise the total reward:

$$\sum_{t \geq 1} r(t) \longleftarrow \text{reward at time } t$$

$K = 2$

$$\delta_1 = \begin{cases} 2 & \text{prob.} \quad 1/2 \\ -1 & \text{prob} \quad 1/2 \end{cases}$$

$$\delta_2 = \begin{cases} 100 & \text{prob.} \quad 1/10 \\ 0 & \text{prob} \quad 9/10 \end{cases}$$

| time | machine | |
|------|---------|---|
| 1 | 1 $\longrightarrow$ | $r(1) = 2$ |
| 2 | 2 $\longrightarrow$ | $r(2) = 0$ |
| 3 | 2 $\longrightarrow$ | $r(3) = 0$ |
| 4 | 2 $\longrightarrow$ | $r(4) = 0$ |

**Difficulty:** we have to make choices based on incomplete statistics !

**Trade off** between:

- getting **good information** on all machines:
  EXPLORATION

- getting good rewards
  EXPLOITATION

after T steps

for each machine $i \in [1, K]$

we have a set of samples

$\hookrightarrow \hat{\mu}_i(T)$ empirical expectation

play $\quad \underset{i \in [1, K]}{\text{argmax}} \quad \hat{\mu}_i(T)$ $\qquad$ EXPLOITATION

E. Greedy strategy: $\qquad \varepsilon$ fixed $\quad \varepsilon = 0.1$

exploration $\begin{cases} \text{Uniformly at} \\ \text{random over} \\ \text{all actions} \end{cases}$ $\qquad$ w. probability $\varepsilon$

exploitation $\begin{cases} \text{greedy:} \\ \underset{i \in [1, K]}{\text{argmax}} \quad \hat{\mu}_i(T) \end{cases}$ $\qquad$ w. probability $1 - \varepsilon$

$$\underset{i \in [1,k]}{\arg\max} \ \mu_i \ \overset{def}{=} \ * \in [1,K] \ \text{such that}$$

$$\mu_* = \underset{i \in [1,K]}{\max} \ \mu_i$$

## Regret analysis

$$\text{Regret}(T) = R(T) = \text{difference between}$$

"best a posteriori"

and

"what we achieved"

$$R(T) = \boxed{T \cdot \mu_k} - \boxed{\sum_{t=1}^{T} r(t)}$$

⚠ $\mu_i$ is the actual expectation

$\hat{\mu}_i(T)$ is the empirical expectation at time $T$

$$\max_{t \geq 1} \sum^{T} r(t)$$

$$(\Longleftrightarrow)$$

$$\min \quad R(T)$$

We use the regret for comparing different strategies

So far:

Greedy $\qquad \qquad \operatorname*{argmax}_{i} \hat{\mu}_i(T)$

ε.Greedy $\qquad \begin{cases} \operatorname*{argmax}_{i} \hat{\mu}_i(T) & \text{with probability } 1-\varepsilon \\ \text{uniform at} & \text{with probability } \varepsilon \\ \text{random} \end{cases}$

Issues: • exploration never stops : at least ε is lost

• exploration does not take existing info into account

• may take a long time to converge

# UPPER CONFIDENCE BOUNDS (UCB)

$r(t)$: reward at time $t$

$r(i,t) : \begin{cases} r(t) & \text{if } i \text{ was chosen at time } t \\ 0 & \text{o/w} \end{cases}$

$$R(T) = T \cdot \mu_* - \sum_{t=1}^{T} r(t)$$

$$\hat{\mu}_i(T) = \frac{1}{n(i,T)} \sum_{t=1}^{T} r(i,t)$$

**UCB**

$$\underset{i}{\text{argmax}} \ \hat{\mu}_i(T) + c(i,T)$$

$$c(i,T) = \sqrt{\frac{\log(T)}{n(i,T)}}$$

**intuitions:**

- if $n(i,T)$ is small (little information) then $c(i,T)$ is large : we need to explore

- if $n(i,T)$ is large then $c(i,T)$ is small, except when $T$ grows, but exponentially less often

Why $c(i,T) = \sqrt{\dfrac{\log(T)}{m(i,T)}}$ ?

## Chernoff - Hoeffding bounds

let $y_1, \ldots, y_n$ iid samples of $Y$

Law of large numbers: $\dfrac{1}{n} \sum_{i=1}^{n} y_i \longrightarrow \mathbb{E}[Y]$
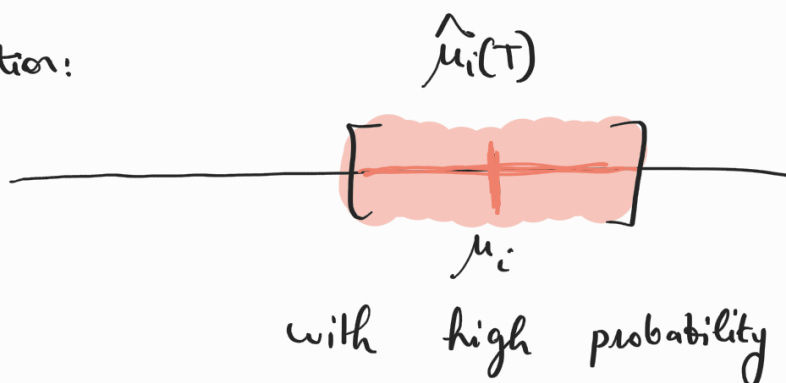
Better: this happens fast!

$$\mathbb{P}\left( \left| \frac{1}{n} \sum_{i=1}^{n} y_i - \mathbb{E}[Y] \right| \geq c \right) \leq 2e^{-2c^2 n}$$

We apply it to our case:

$$\mathbb{P}\left( \left| \hat{\mu}_i(T) - \mu_i \right| \geq c(i,T) \right) \leq 2e^{-2c(i,T)^2 m(i,T)}$$

$c(i,T) = \sqrt{\dfrac{\log(T)}{m(i,T)}}$    gives    $\dfrac{2}{T^2} \xrightarrow[T \to +\infty]{} 0$

Intuition:      $\hat{\mu}_i(T)$



$\mu_i$

with high probability

# UPDATES

$$\hat{\mu}_i(T) = \frac{1}{m(i,T)} \sum_{t=1}^{T} r(i,t)$$

After choosing $i$ :

$$\hat{\mu}_i(T+1) = \frac{1}{\underbrace{m(i,T+1)}_{m(i,T)+1}} \sum_{t=1}^{T+1} r(i,t)$$

Small calculations

$$\hat{\mu}_i(T+1) = \hat{\mu}_i(T) + \frac{1}{m(i,T)+1} \left[ X(i,T+1) - \hat{\mu}_i(T) \right]$$

$$\text{NEW} = \text{OLD} + \alpha \left[ \text{CURRENT} - \text{OLD} \right]$$

We call this step size

Two possible updates :

- empirical mean

- constant step size     ($\alpha = 0.1$ for instance)

$$\hat{\mu}_i(T+1) = \hat{\mu}_i(T) + \alpha \left( \hat{\mu}_i(T) - X(i,T+1) \right)$$