

INTRODUCTION TO REINFORCEMENT LEARNING

Nathanaël FISTALKOW

$$RL \subseteq ML$$

ML:



Supervised learning:

Training data is labelled

Objective: predict unlabelled data

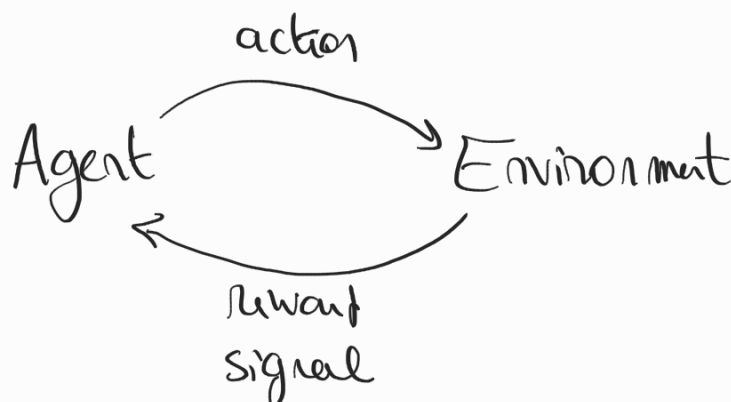
Unsupervised learning:

Training data is not labelled

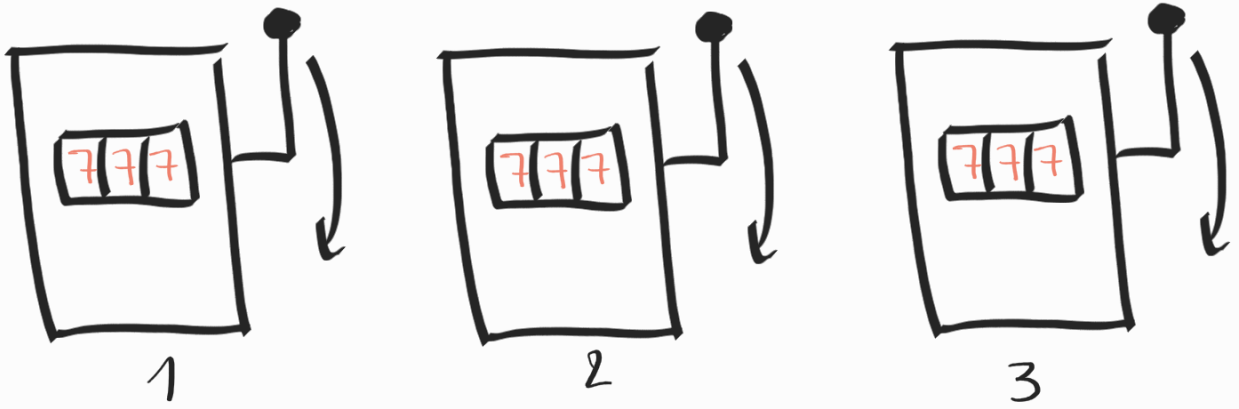
Objective: finding structure / pattern

RL is neither supervised nor unsupervised!

RL: | active learning
| reward-based



MULTI ARMED BANDITS



arm = machine = action

K machines : $1, \dots, K$

we assume that each machine i has a
reward distribution δ_i

Bernoulli distribution

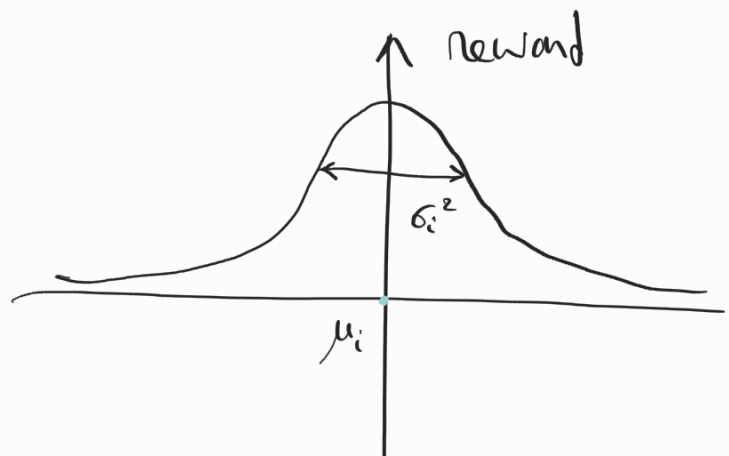
$B(p)$

$\begin{cases} 1 \\ 0 \end{cases}$

with probability p
with probability $1-p$

Normal distribution

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



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

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 δ_i
- Two (similar) goals:
 - (1) identify the best arm: $\operatorname{argmax}_i \mu_i$
 - (2) maximise the total reward:

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

$$K=2$$

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

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

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

Greedy strategy:

play

$\operatorname{argmax}_{i \in [1, K]}$

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

EXPLOITATION

ϵ . Greedy strategy:

ϵ fixed

$\epsilon = 0.1$

exploration

uniformly at random over all actions

w. probability ϵ

exploitation

greedy:

$\operatorname{argmax}_{i \in [1, K]}$

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

w. probability $1 - \epsilon$

Notations

$\arg \max_{i \in [1, k]} \mu_i \stackrel{\text{def}}{=} i^* \in [1, k]$ such that

$$\mu_{i^*} = \max_{i \in [1, k]} \mu_i$$

Regret analysis

$\text{Regret}(T) = R(T) =$ difference between
"best a posteriori"

and

"what we achieved"

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

⚠ μ_i is the actual expectation

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

$$\max \sum_{t=1}^T r(t)$$

(\Leftrightarrow)

$$\min R(T)$$

We use the regret for comparing different strategies

So far:

Greedy

$$\operatorname{argmax}_i \hat{\mu}_i(T)$$

ϵ -Greedy

$$\begin{cases} \operatorname{argmax}_i \hat{\mu}_i(T) & \text{with probability } 1-\epsilon \\ \text{uniform at random} & \text{with probability } \epsilon \end{cases}$$

- 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

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

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

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

Chernoff-Hoeffding bounds

let y_1, \dots, y_n iid samples of Y

Law of large numbers: $\frac{1}{n} \sum_{i=1}^n y_i \rightarrow \mathbb{E}[Y]$

Better: this happens fast!

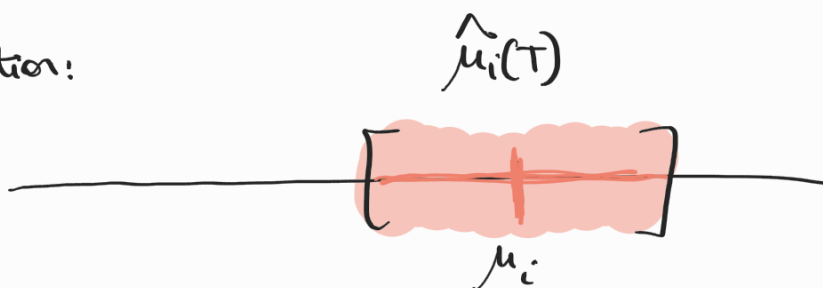
$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n y_i - \mathbb{E}[Y]\right| \geq c\right) \leq 2e^{-2c^2 n}$$

We apply it to our case:

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

$$c(i, T) = \sqrt{\frac{\log(T)}{n(i, T)}} \quad \text{gives} \quad \frac{2}{T^2} \xrightarrow{T \rightarrow +\infty} 0$$

Intuition:



with high probability

UPDATES

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

After choosing i :

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

Small calculations

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

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

we call this α step size

Two possible updates :

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

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

THOMPSON SAMPLING \equiv POSTERIOR SAMPLING

Key idea: instead of choosing the best action from the model, look at the model as a distribution of which action is best

Bernoulli bandit

θ_i : probability of success for arm i

$\left. \begin{array}{l} \alpha(i, t) : \text{number of successes for } i \\ \beta(i, t) : \text{number of failures for } i \end{array} \right\} \text{up to time } t$

Greedy

$$\left[\hat{\theta}(i, T) = \frac{\alpha(i, T)}{\alpha(i, T) + \beta(i, T)} \right]$$

Thompson

$$\hat{\theta}(i, T) = \text{Sample}(\alpha(i, T), \beta(i, T))$$

Choose $\arg\max_{i \in [1, k]} \hat{\theta}(i, T)$

Update $\alpha(i, t), \beta(i, t)$

Update:

$$\alpha(i, t+1), \beta(i, t+1) =$$

$$\begin{cases} \alpha(i, t), \beta(i, t) & \text{if } i \text{ not chosen} \\ (\alpha(i, t), \beta(i, t)) + (r(t), 1-r(t)) & \text{o/w} \end{cases}$$

$$r(t) = \begin{cases} 1 & \text{if success} \\ 0 & \text{if failure} \end{cases}$$

β . distribution:

$$P(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

Here Bayesian inference is easy
(computationally). In general it's
hard \rightarrow approximation efforts