

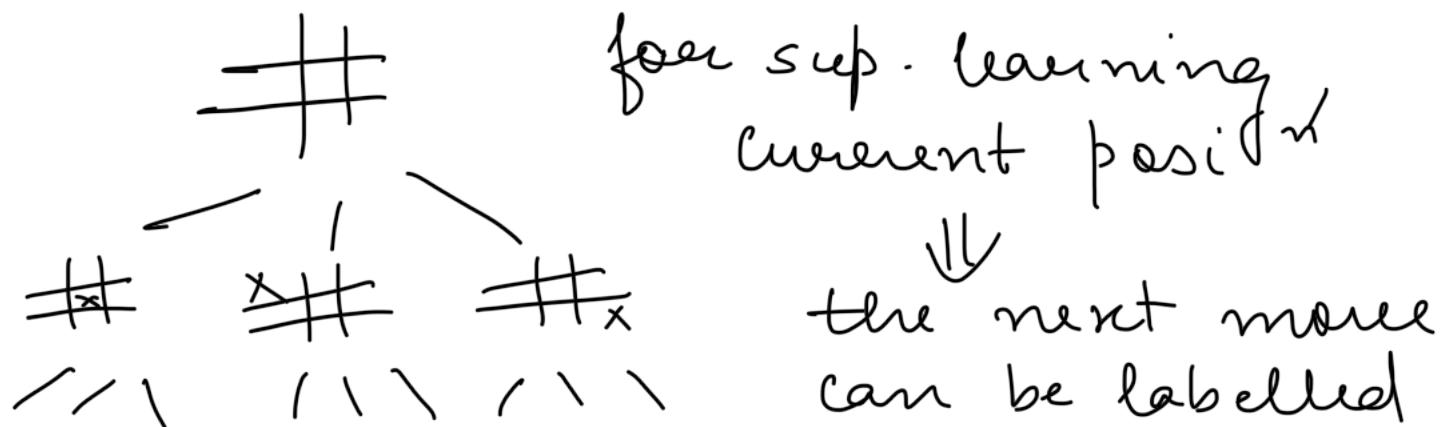
# RL-Week-1

## 1 Intro

- ML  $\rightarrow$  learn func<sup>n</sup> from input to output from data
- familiar model  $\rightarrow$  learning from data
- cycling, Talk  $\rightarrow$  trial & error, evalua<sup>n</sup>  
(not instanc<sup>n</sup>), RL (feedback mechanism)  
 $\hookrightarrow$  walking also (people cheer if you do good, see if you fall you cry)
- RL  $\rightarrow$  trial & error learning paradigm  
(rewards & punishments)
- learn abt. a sys. thru interac<sup>n</sup>, by behavioural psychology.
- Pavlov's dog  $\rightarrow$  bell, food, salivate, dog
- from 80s, chatgpt (human feedback), complex dynamics (helicopter, humanoid), complex workspace, stochastic sensing & actua<sup>n</sup> (visual atten<sup>n</sup>), human in the loop learning, cognitively motivated learning, news customization, ad selec<sup>n</sup>, recommeda<sup>n</sup> engines  $\rightarrow$  finally learn thru self-play
- TD Gammmon (learn completely by self play), Arcade games (using complex neural network, learn to play from video input), AlphaGo, AlphaZero, 3 years

- AlphaFold (predicted protein structure), Protein structure database

## 12 Tic-Tac-Toe & TD Learning



- by RL  $\rightarrow$  learn from real (+1 / -1 | 0), you have to learn from repeated play
- MENACE (matchbox engine) - if you make a valid / winning move, leads increase that means probab. of that move has increased.
- assume opponent is an imperfect opponent (it makes mistakes)
- there are certain positions in the game, from where you will certainly win.  
 $\therefore$ , increase the probab. of such moves.
- TD (Temporal difference learning) -  $\Delta$  the probabilities by the timeline.
- "pudic" of outcome at time  $t+1$  is better-than the "pudic" at time  $t$ . So, now use the latter "pudic" to adjust the earlier "pudic".

- RL  $\rightarrow$  Deep RL . Great  $\rightarrow$  omniscient learning (consume any info. to learn)

### L3 Immediate RL & Bandits

- limits of explore & exploit
- IRL  $\rightarrow$  payoff accrues immediately after an act<sup>n</sup> is chosen . Bandit problem
- explore to find profitable act<sup>n</sup>. exploits to act according to the best observation so far. Always doing str. is not optimal.
- MAB (n-arm) selects a particular ac<sup>n</sup> from a set of n act<sup>n</sup> (1, 2, 3 ... n). Each selec<sup>n</sup> gives Rewards from respective probab. distribu<sup>n</sup>.

Arm i has a distibn. with mean  $\mu_i$

$$\mu^* = \max \{\mu_i\}$$

#### - Objectives

1. Identify the correct arm eventually.
- 2.
3.  $a_{i,k}$  reward when i<sup>th</sup> arm for k<sup>th</sup> time

$$Q(a_i) = \frac{\sum_{k=1}^{n_i} r_{i,k}}{\sum_{k: r_{i,k}} 1} \quad \begin{array}{l} \text{true} \\ \text{reward} = \mu_i \\ (\text{this is } r_{i,k}) \end{array}$$

↓  
 expected  
 payoff for  
 ac<sup>n</sup>  $a_i$

$$\text{expected reward} \stackrel{\text{max.}}{\overbrace{Q(a^*)}} = \max_i \{Q(a_i)\}$$

$$Q_{k+1}(a_i) = Q_k(a_i) + \alpha(u_{i,k} - Q_k(a_i))$$

$$\alpha = \frac{1}{k_i + 1} \rightarrow \text{yields the avg.}$$

$$\bar{x}_k = \frac{x_1 + \dots + x_k}{k}$$

$$\bar{x}_{k+1} = \frac{x_1 + \dots + x_{k+1}}{k+1} = \frac{k\bar{x}_k + x_{k+1}}{k+1}$$

$$\bar{x}_{k+1} = \bar{x}_k + \frac{1}{k+1} [x_{k+1} - \bar{x}_k]$$

Exploratory methods →

- Epsilon Greedy → arm  $a^* = \operatorname{argmax}_i \{Q_k(a_i)\}$   
with probab.  $1 - \epsilon$  and other arm  
with prob.  $\epsilon$

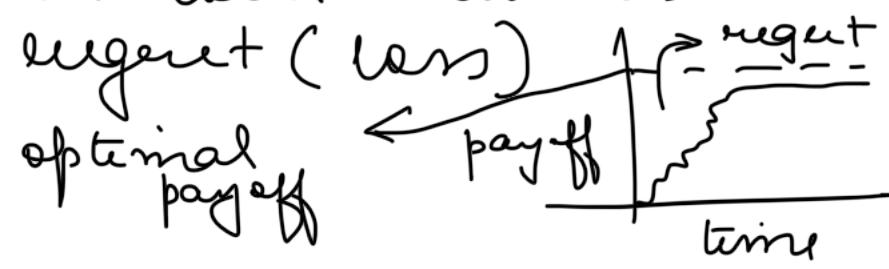
- Softmax → select max with probab.  
 $\propto$  to current value estimates.

$$\pi_k(a_i) = \frac{\exp\left(\frac{Q_k(a_i)}{\tau}\right)}{\sum_j \exp\left(\frac{Q_k(a_j)}{\tau}\right)} \rightarrow \text{temp.}$$

both give asymptotic convergence  
(always).

## L4 Regrets & PAC Frameworks

- Bandit use case → Yahoo News

- Obj. 2  $\rightarrow$  max. the total rewards obtained or min. regret (less) 

how to learn while keeping regret as less as possible.
- PAC Framework (Probably Approximately Correct)
  1. identification of an  $\epsilon$ -optimal arm with probab. ( $1 - \delta$ )
  2.  $\epsilon$  optimal: mean of the selected arm satisfies  $\Rightarrow \mu > \mu^* - \epsilon$
  3. min sample complexity - order of samples req. for such an arm identification  $\rightarrow$  Round based algo.
- PAC Guarantees  $\rightarrow$  Median Elimination, Upper Confidence Bounds (UCB), Thompson Sampling  $\hookrightarrow$  Regret Optimality (Bayesian Approach)  $\rightarrow$

L5 UCB  $\rightarrow$  Upper confidence bound

- The arm with the best estimate  $\hat{\mu}^*$  so far serves as a benchmark & other arms are played only if the upper bound of a suitable confidence interval is atleast  $\hat{\mu}^*$ .

$\hookrightarrow$  So be greedy w.r.t. upper confidence bounds!

- Play machine  $j$
- $$\max \left\{ Q(a_j) + \sqrt{\frac{2 \ln n}{n_j}} \right\} \rightarrow \begin{array}{l} \text{overall no.} \\ \text{of plays} \end{array}$$
- ↓
- sub optimal arm  
 $j$  is played less  
 - than  $(8/\Delta_j) \ln n$  times.

## L6 Contextual Bandits

- challenge for "customizing" & "ad selecting" for each user, not a generic 1-timer.
- diff. ads for diff. user. 1 bandit  $\rightarrow$  1 user.  
 ↳ hard to train for 1 person
- assume that the params of the reward distribution themselves are determined by set of params. (groups are made for users)  $\rightarrow$  linear parameterization of the expecta".
- each user is represented by a set of features.  
 (joint features of user & arm)
- the "static"  $\rightarrow$  of choosing now depends on these features.  
 ↳ check for the presence or absence of diff. signals.

- Lin UCB - most popular continual bandit
  - ↳ predicted expected reward assumed to be a linear func<sup>n</sup> of features.
- 1. use ridge deg. to fit
- 2. get the upper confidence
- 3. UCB like ac<sup>n</sup> selec<sup>n</sup>.
- 4. better performance with lesser training data.

Practice Assignment