

<sup>th</sup>  
Feb 14, 2017

All ANNs we have learned so far are useless if we have small datasets or even worse, we have none.

Can we learn from scratch?

↳ Mechanism? Learning by doing.  
online learning

Online learning is necessary for several reasons:

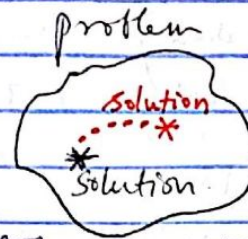
\* data not available.

\* non-stationary.

Learning by doing • No teacher

• But there is indirect feedback. (either punishment or reward.)

• the possibility to get permanent / constant / continuous feedback.

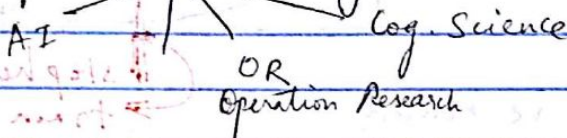


Stationary: Solution is fixed.

Location of solution in feature space is fixed. "Stationary"

Non stationary: solution moves.

## Reinforcement Learning



RL has nothing to do with LK.

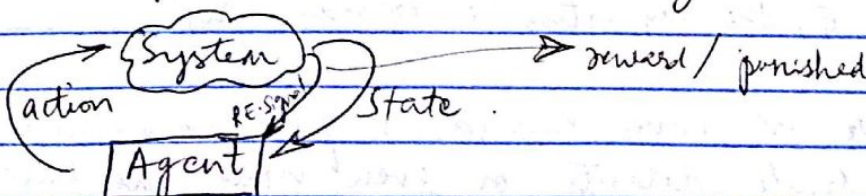
## Typical Applications

① Robotics - ~~elevator~~ elevator control  
- Robo-soccer

② Board games - Backgammon  
- Checkers  
- Chess

③ Scheduling - Dynamic channel allocation  
- inventory problems.

## Idea of RL (Reinforcement Learning)



Smart = find a trade off between exploitation & exploration.

Taking actions and observing reward/punishment.

↳ take the right (exploitation)

↳ or experiment (exploration).



# RL math Model

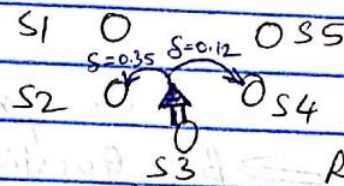
- Stochastic/non-stationary : uncertainty
- decision making under uncertain conditions.
- weakly supervised : maximize reward.

Markov  
Decision  
Process (MDP)

## MDP

- Set of states  $S$
- Set of actions  $A$
- Transition probability  $\delta$
- Reward function  $R$

Environment : Set of states



taking actions:  $\delta(S, a, S')$

$$R(S, a) = -1 \text{ (with } \delta = 0.5)$$

$$= +10 \text{ (with } \delta = 0.35)$$

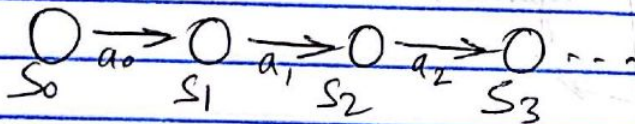
$$= +20 \text{ (with } \delta = 0.15)$$

Reward:  $R(S, a)$

Cardinality of  $A \ll$  Cardinality of  $S$

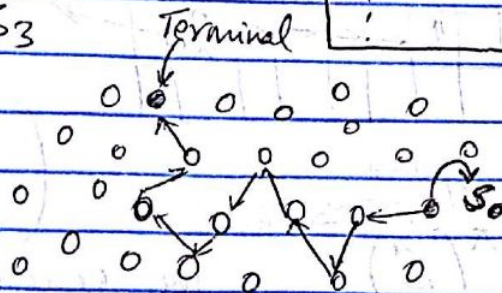
States	Actions			
	$a_1$	$a_2$	$a_3$	$a_4$
$S_1$	0.1	0.1	0.08	0.01
$S_2$	0.4	0.69	0.17	0.15
$S_3$	0.2	0.4	0.14	
$\vdots$				

## MDP Model : Trajectory



(it will be a long matrix)

Finding trajectory = learning



Any scheme working

based on reward has

to answer this question : immediate or long-term ?

"greedy" approach.

Another question : Are early rewards more important than later ones ?

Early ones

→ discount the rewards?



Return the function (sum of rewards)

① Finite horizon return =  $\sum_i R(s_i, a_i)$

② Infinite horizon  
(2a) discounted, return =  $\sum_{i=0}^{\infty} \gamma^i R(s_i, a_i)$

(2b) undiscounted, return =  $\frac{1}{N} \sum_{i=0}^{N-1} R(s_i, a_i)$   $0 < \gamma < 1$

→ one episode with  $N$  iterations/observations.

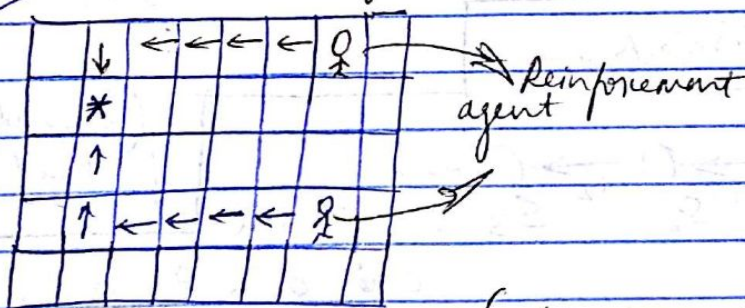
⇒ Big Question: Can we observe the system?

→ Yes, fully observable

→ No, partially observable (maybe we have some hidden variables)

The Goal of R.L. is to find a policy  $\pi$  that can guide the agent along an optimal trajectory.

Grid world problem



policy → state →  $V^{\pi}(S)$  = the expected return starting at state  $S$  following  $\pi$ .  
value

(find the value of which state under the conditions of policy  $\pi$ )

policy → state → action →  $Q^{\pi}(S, a)$  = the expected return starting at state  $S$  with action  $a$  and then following  $\pi$ .  
substitution (the value)