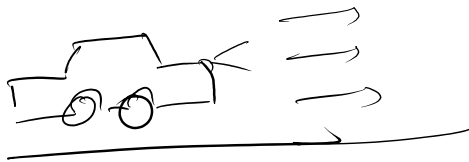


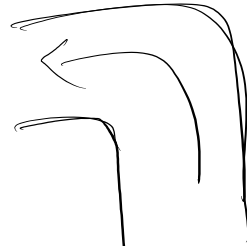
Reinforcement Learning
Introduction to Reinforcement Learning

- Subdomain of ML
- Train models (Agents in RL)
- Sequential decisions.

Driving



Acc.



left



States in RL



Sudden Brake

Waymo
↓
Alphabet

Self driving
tech

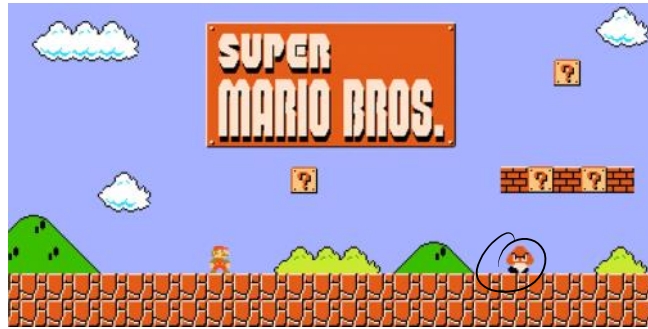


1 2 3 4 5

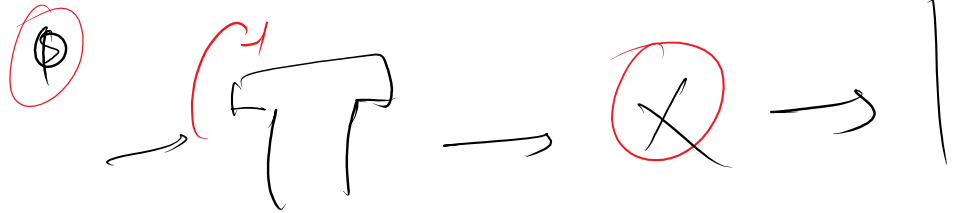


↓ techniques from RL

Super Mario



Actions in RL



Alpha Go
Alpha Star



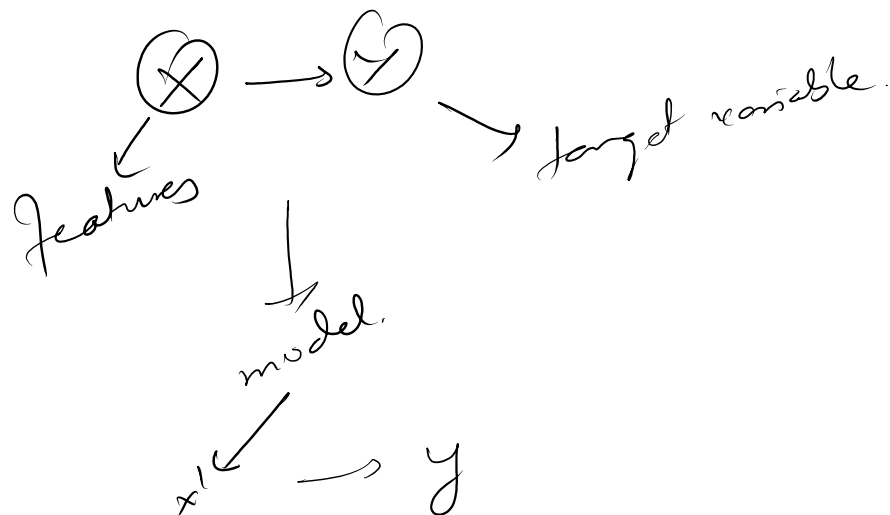
Go



Star Craft

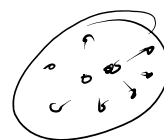
Supervised

Supervised



Unsupervised

X



Reinforcement

Situations \rightarrow Actions

States \rightarrow Actions

What to do??

Interacting

Learning from interaction.

General
discussion
on RL

\rightarrow MDP

(i) Tabular

\rightarrow DP

\rightarrow MC

\rightarrow TD

(ii) Approx

\rightarrow DRL

Papers.

1.

comp.

- quizzes
- coding.

Capstone
project

- ML
- Stats/prob
- Python

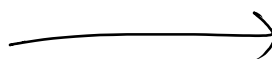
"Reinforcement Learning: An Introduction"
→ Sutton & Barto.

Reinforcement Learning Building blocks Reinforcement Learning

- Agent → "the model"
- Environment & States → Snapshot of env.
↳ Representation of the problem.
- Actions → "moves"



S₁



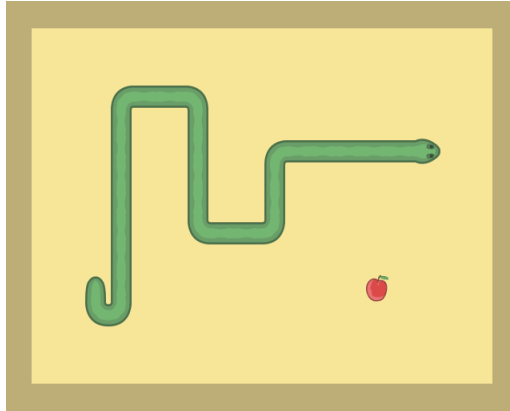
S₂

(R)



(R)

→ Reward: → "numerical feedback"



+1

+3

maximize the total reward

→ Value Function

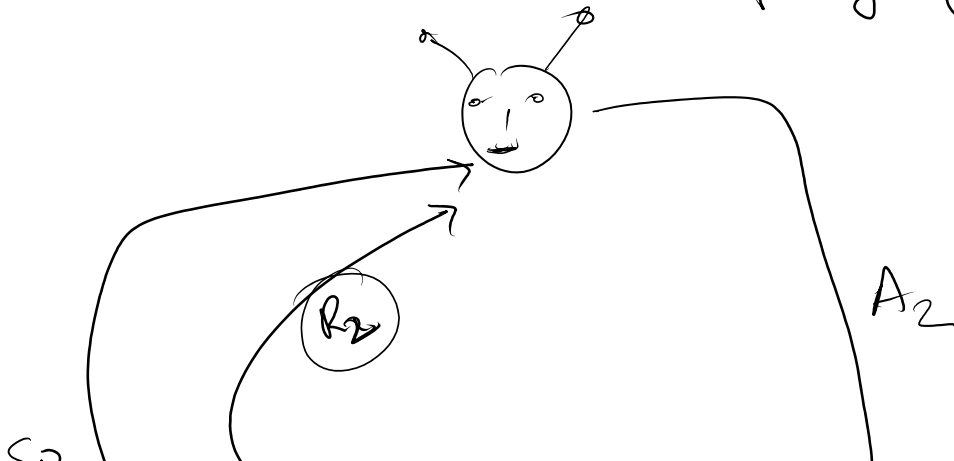
State
Action

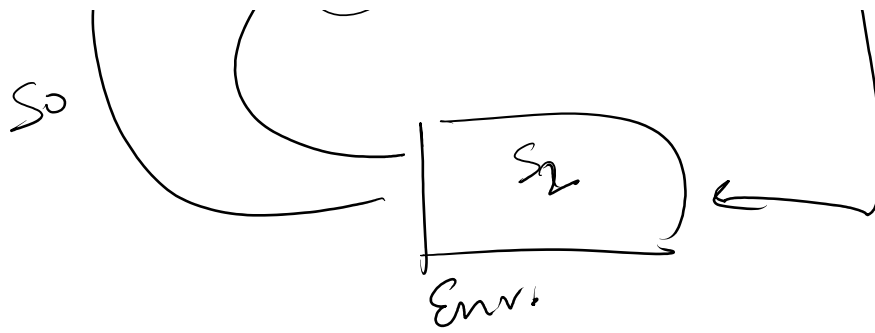
Rewards

→ Policy -

"how to behave?"

policy → (π)





Reinforcement Learning Building blocks Reinforcement Learning - Examples

Env., agent, state, action, values, policies.

→ Self driving car in a simulation

policy: well trained agent

Env :- Simulator

State :- Front, rear, side
speed, accel, road
angle.

Reward :-

- 1 each passing timestep
- 5 go off the lane
- 10 jump off the sign
- 100 hit someone

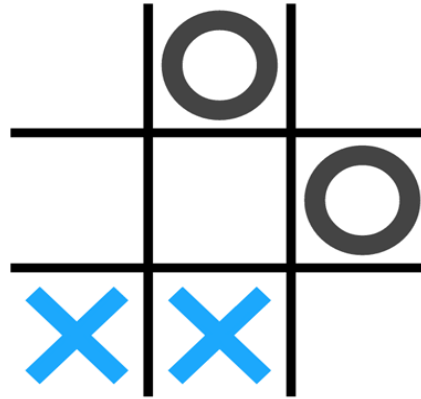


→ Tic-tac-toe

Game board

→ Tic-tac-toe

Optimal policy.



Env. :- Game board day with the opp. playing (X)

State :- X's or O's on board.

Reward :- +1 win
0 lose

Action :- placing O.

Policy - Random.

→ Mobile rag picking robot.

Actions :- Go ahead, return back to charging station.

Policy :- Go ahead if battery > 25%



Env :- The complete floor.

State :- Front, side, rear, current battery %, space in rag bag.

Reward :- +1 for rag collected.

-100 : stuck.

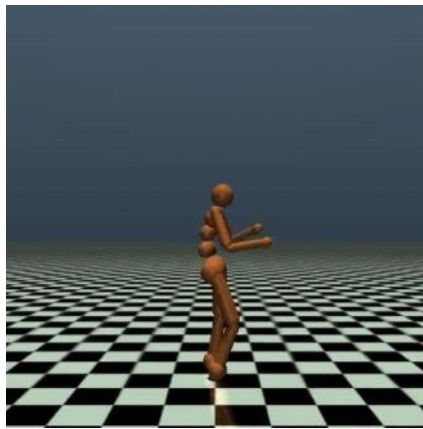
→ robot learn to walk

Env. 1 - Simulator.

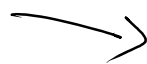
State 2 - Pos. of joint, ang. vel.
of joints, joint position
x, y vel.

Reward 2 - +1 → forward
-100 → falling.

Actions 2 - Torque of
joints.
policy 2 - optimal.



Reinforcement Learning
Introduction to the K-Armed Bandit Problem



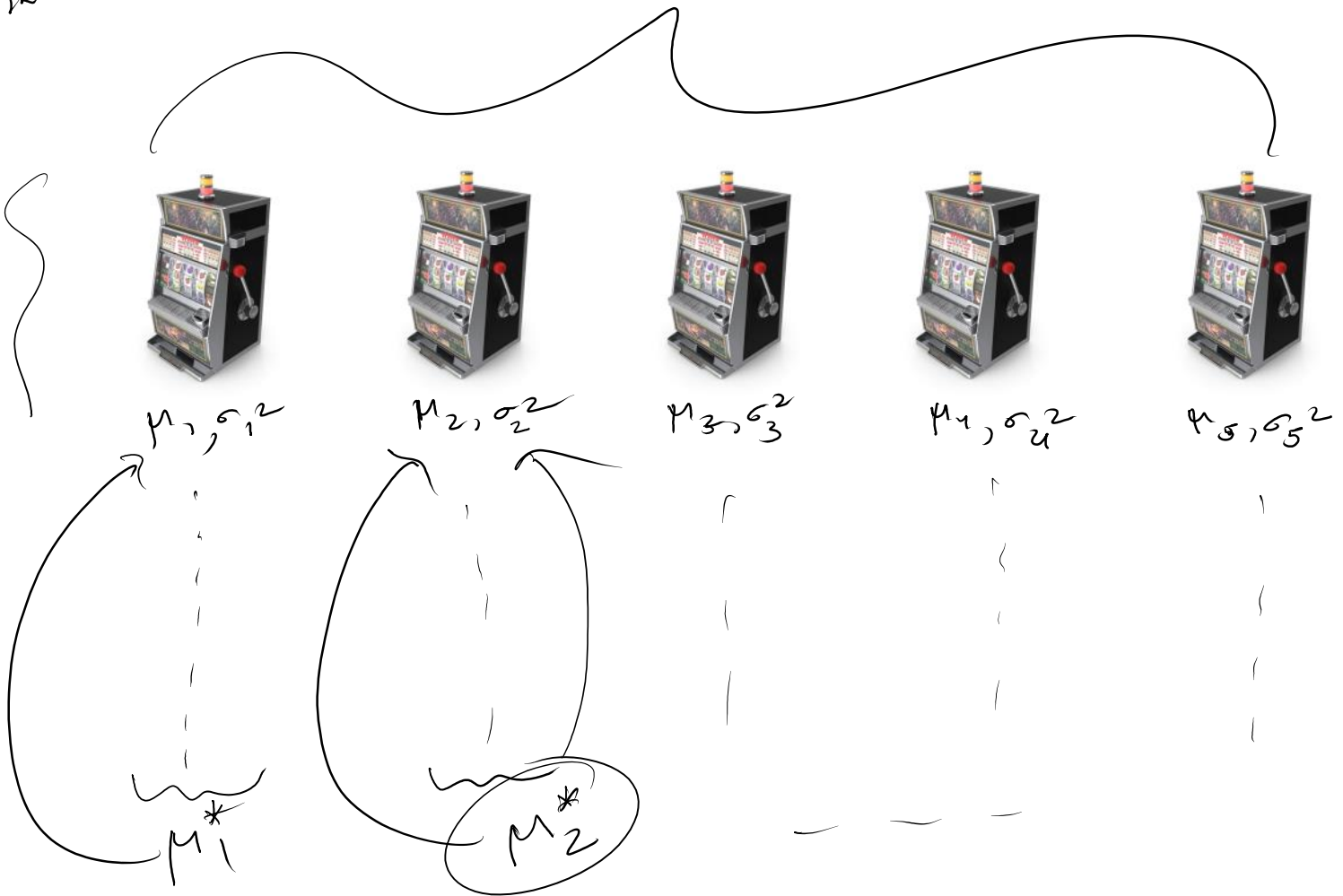
$$X \rightarrow \mathcal{N}(\mu, \sigma^2)$$

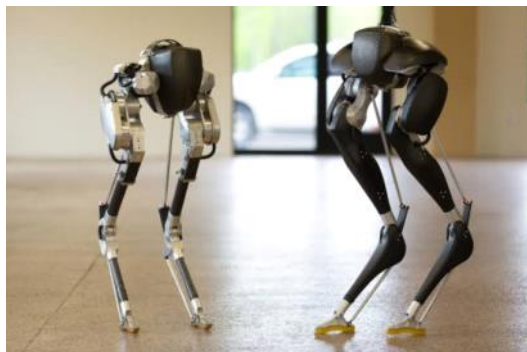
$$p(x|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

(σ² (μ))
↑ ↑



125





Reinforcement Learning

Solving K-armed bandit - Greedy method

K actions

Action values

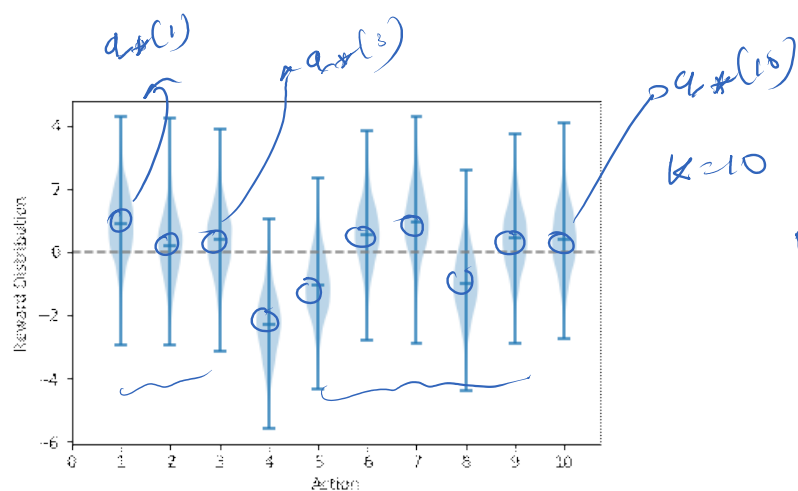
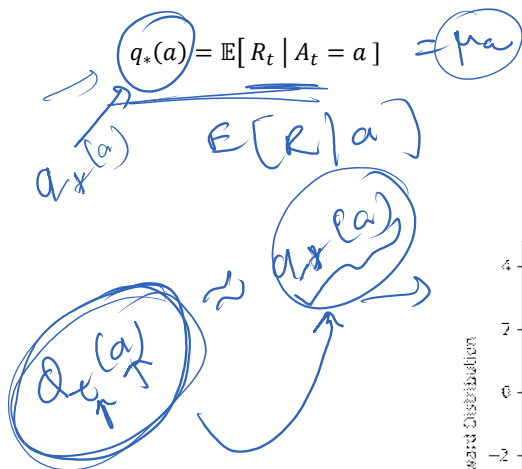
μ, σ^2

true value of action

$q_*(a)$

t, A_t, R_t

$$q_*(a) = E[R_t | A_t = a]$$



$K=10$

μ
 $\sigma^2 = 1$

Calculating action value estimates

1, 2

Calculating action value estimates

t	$Q_t(A)$	$Q_t(B)$	$Q_t(C)$	A	B	C
1	0	0	0	2	1	
2	2	0	0	1		
3	2	1	0			
4	1.5	1	0			3
5	1.5	1	3			-1
6	1.5	1	1			

$K=3$
Randomly.
A, B, C

$$\frac{2+1}{2}$$

$$Q_t(a) = \frac{\sum \text{reward, choosing } a, t-1}{\text{no. of times } a \text{ was chosen before } t}$$

A, B, C
Greedy.

t	$Q_t(A)$	$Q_t(B)$	$Q_t(C)$	A	B	C
1	0	0	0	-1		
2	-1	0	0		2	
3	-1	2	0		-5	
4	-1	-1.5	0			1
5	-1	-1.5	1			1
6	-1	-1.5	1			

$$\frac{2+(-5)}{2} = -1.5$$

Calculating action value estimates incrementally

action a
(n+1) $\rightarrow r_1, r_2, r_3, r_4, \dots, r_n$

$$Q_{n+1} = \frac{r_1 + r_2 + \dots + r_n}{n}$$

$$= \frac{1}{n} \sum_{i=1}^n r_i$$

$$\begin{aligned}
 &= \frac{1}{n} \left[R_n + \sum_{i=1}^{n-1} R_i \right] \\
 &= \frac{1}{n} \left[R_n + (n-1) \frac{1}{(n-1)} \sum_{i=1}^{n-1} R_i \right] \rightarrow \text{den} = \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \\
 &= \frac{1}{n} \left[R_n + (n-1) \text{den} \right] \\
 &= \frac{1}{n} \left[R_n + n \text{den} - \text{den} \right] \\
 \rightarrow \text{New Estimate } Q_{n+1} &= \underbrace{\text{den}}_{\text{old estimate}} + \underbrace{\left(\frac{1}{n} \right)}_{\text{step size}} \left[\underbrace{R_n - \text{den}}_{\text{Target} - \text{old estimate}} \right] \rightarrow (0.1)
 \end{aligned}$$

Diagram illustrating the update rule for the Greedy Algorithm for K-Armed bandit. A horizontal line represents the value space, with points 0, n, and T marked. A bracket labeled 0.5 [0.1] is shown below the line.

Greedy Algorithm for K-Armed bandit

Initialize, for $a = 1$ to k :

- $\rightarrow Q(a) \leftarrow 0$
- $\rightarrow N(a) \leftarrow 0$

Repeat forever:

- $A \leftarrow \text{argmax}_a Q(a)$
- $R \leftarrow \text{bandit}(A)$
- $N(A) \leftarrow N(A) + 1$
- $Q(A) \leftarrow \frac{Q(A)}{N(A)} + \frac{1}{N(A)} [R - Q(A)]$

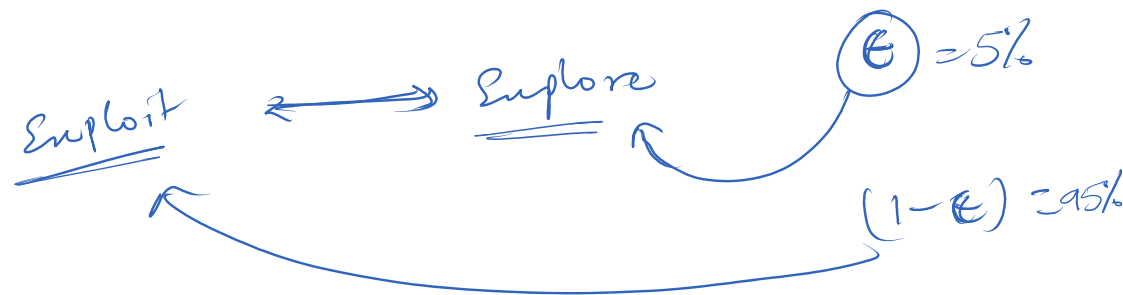
no of Bandits.

$Q(a/s)$

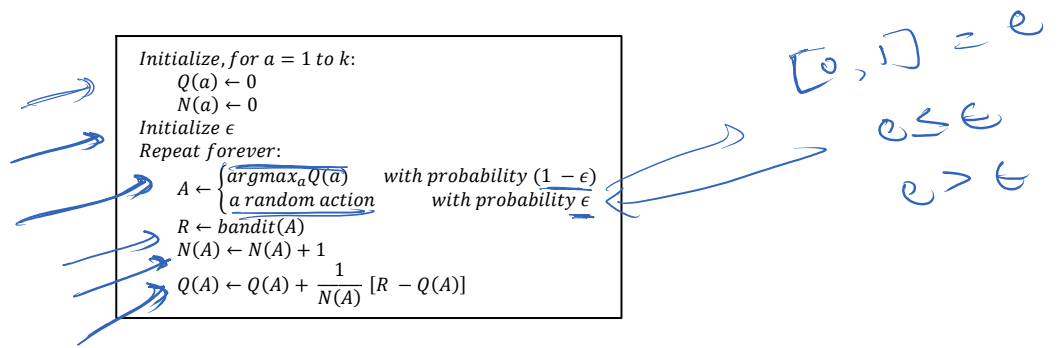
Reinforcement Learning
Solving K-armed bandit - Epsilon Greedy method

Exploration and exploitation

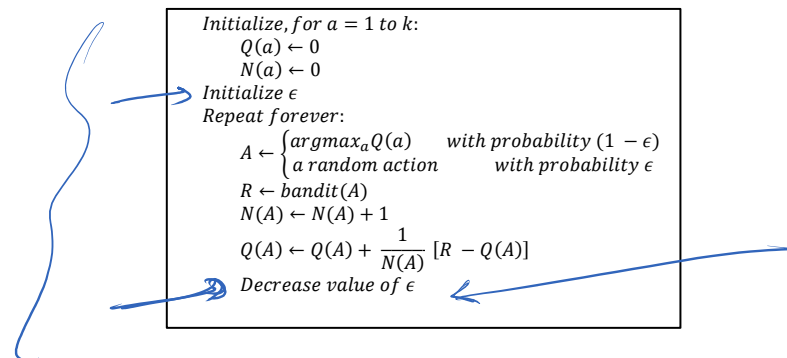
Exploiting → Best use of current knowledge
Explore → Non-optimal move



Epsilon Greedy Algorithm



Decaying - Epsilon Greedy Algorithm



Reinforcement Learning

Solving K-armed bandit - Upper Confidence Bounds

Upper Confidence Bounds:

$$\underbrace{Q_t(a)}_{\downarrow} + \underbrace{U_t(a)}_{\text{uncertainty}} \geq \underbrace{q_*(a)}_{\text{optimal value}}$$

UCB1

$$U_t(a) = \sqrt{\frac{2 \ln t}{N_t(a)}}$$

$t \rightarrow$ total time elapsed.
 $N_t(a) \rightarrow$ No. of times action 'a' is taken.

$N_t(a) \uparrow$	\rightarrow	$U_t(a) \downarrow$	
$t \uparrow$	\rightarrow	$U_t(a) \uparrow$	
$U_t(a) \uparrow$		$U_t(a) \downarrow$	$q_*(a)$

UCB1 Algorithm:

Initialize, for $a = 1$ to k :
 $Q(a) \leftarrow 0$
 $N(a) \leftarrow 0$
 Repeat forever:
 $A \leftarrow \operatorname{argmax}_a [Q(a) + \sqrt{\frac{2 \ln t}{N(a)}}]$
 $R \leftarrow \text{bandit}(A)$
 $N(A) \leftarrow N(A) + 1$
 $Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)]$

Performance comparison:

2000 \rightarrow $k=10$

$k=10$

$k=10$

1000

