

Introduction to Reinforcement Learning (RL)

AppliedAIcourse.com

Agenda:

- Applications
- Problem Setup
- Terminology
- Bellman Equation
- Markov Decision Process
- Q-learning
- Temporal-Difference
- Deep Q-Learning
[intuition]

Assumptions :

- comfortable with math in ML & DL

Applications :

- Robotics
- Game-play
- Trading Systems
- Self-driving Cars

④ Modern RL = Deep RL
↓
DL + RL

Problem Formulation:

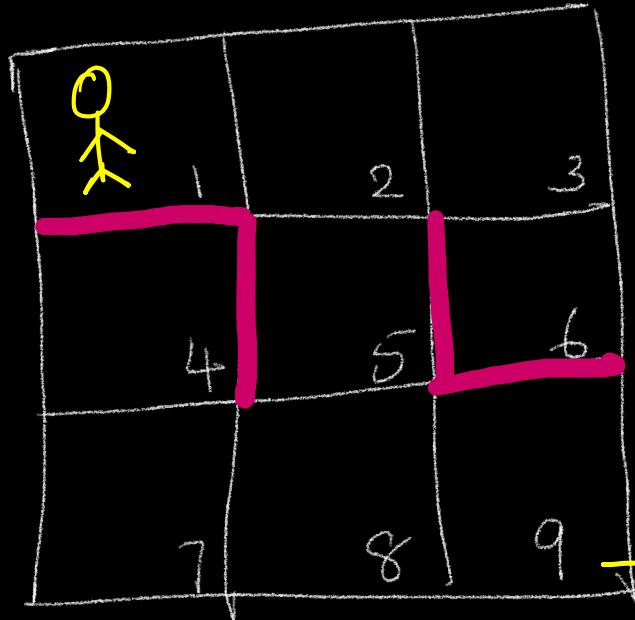
Agent : 

environment :

states : L₁-L₉

Actions : → ← ↓ ↑

Reward : R(s,a)



Locations 1-9
— : walls

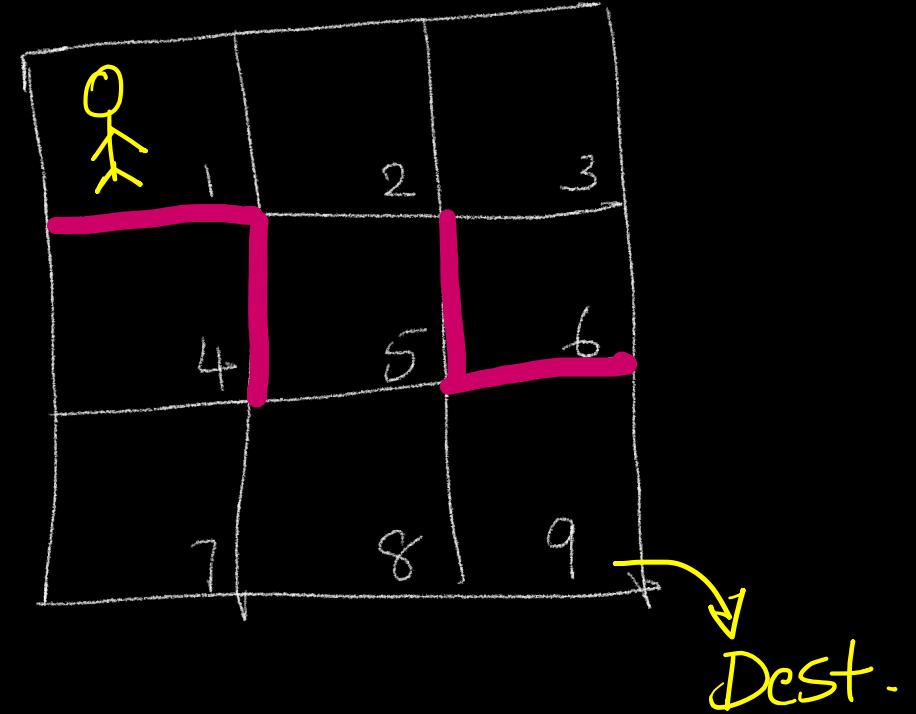
destination

Value-funktion: $V(s)$

$$s \xrightarrow{a} s'$$

$$1 \xrightarrow{} 2$$

$$2 \xrightarrow{} 5$$



$$S \xrightarrow{a} S'$$

$$V(S) = \max_a (R(S,a) + \gamma \cdot V(S'))$$



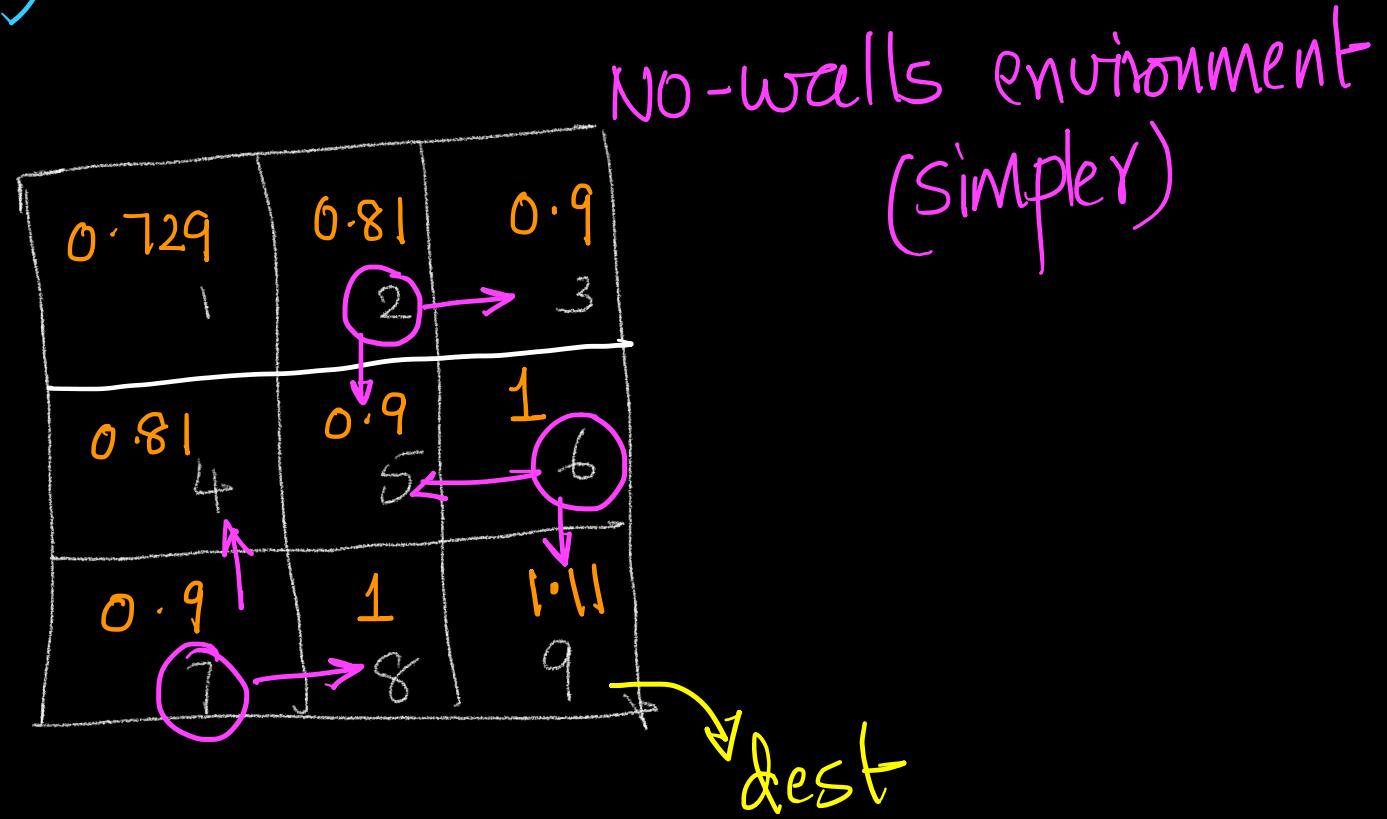
Bellman-Equation

$0 \leq \text{discounting factor} \leq 1$
Let $\gamma = 0.9$

$$V(S) = \max_a (R(S,a) + \gamma V(S'))$$

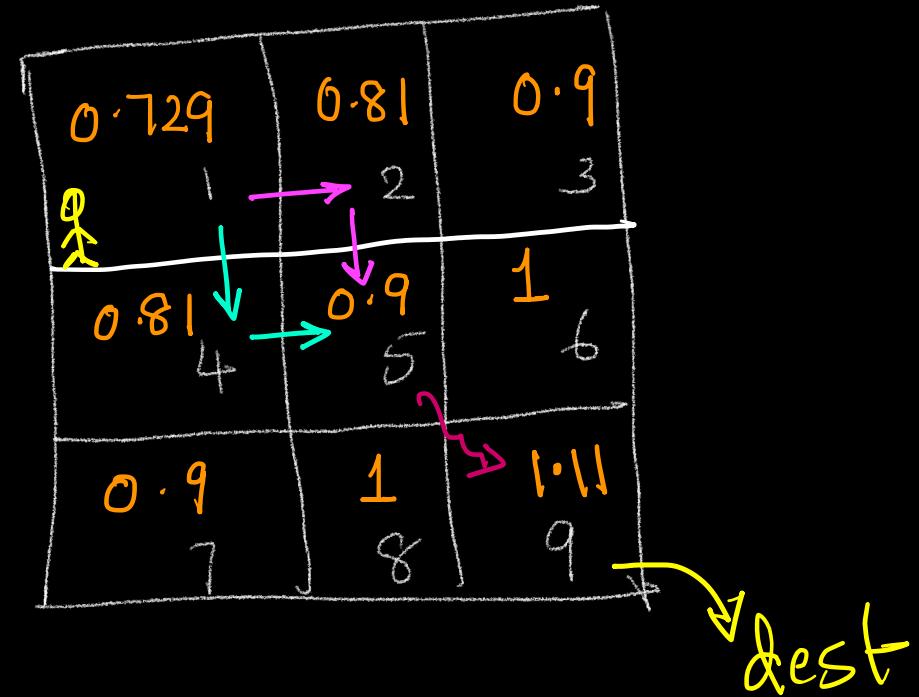
Let $\gamma = 0.9$

$$R(S,a) = 0$$



Dynamic Programming in CS:

- recursion
- overlapping Sub-problems



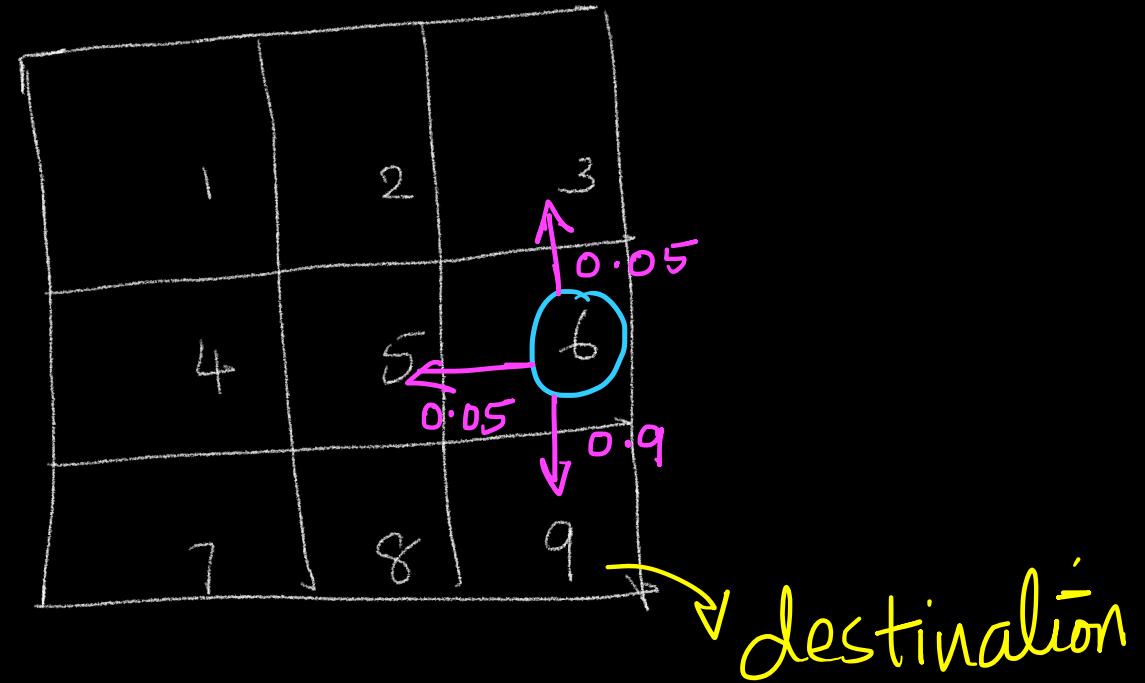
$$V(S) = \max_a R(S, a) + \gamma \cdot V(S')$$

Dynamic Programming → Bellman
by.

(Q) why not solve it trivially?

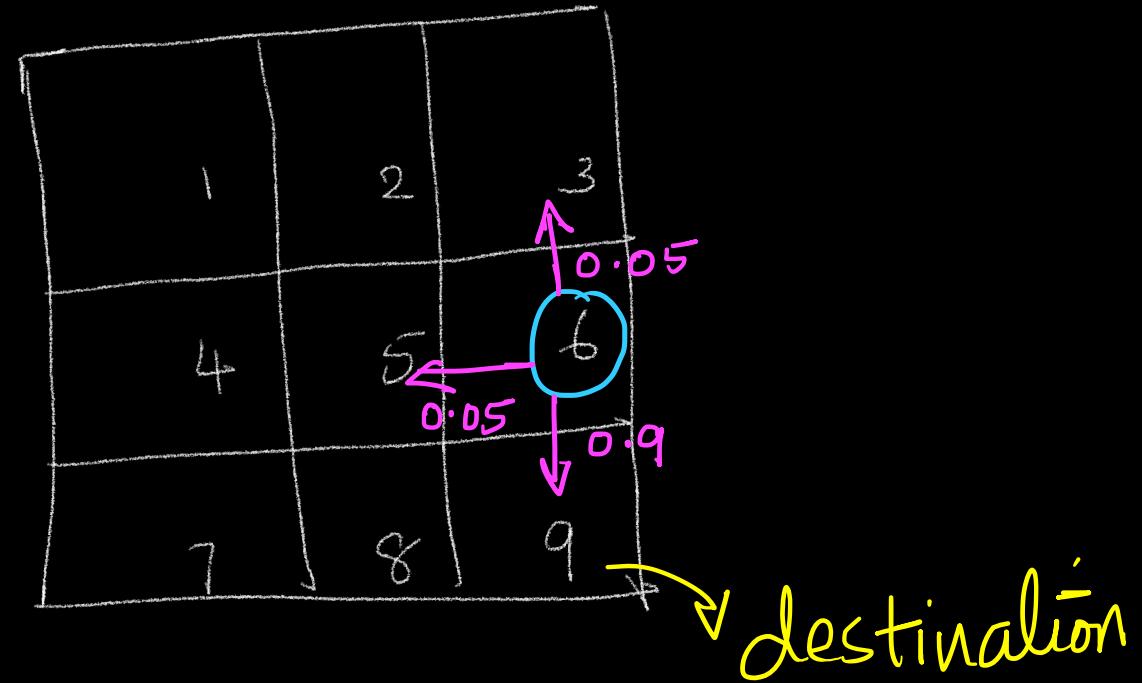
Stochastic - process: Markov Decision process [MDP]

Roomba:



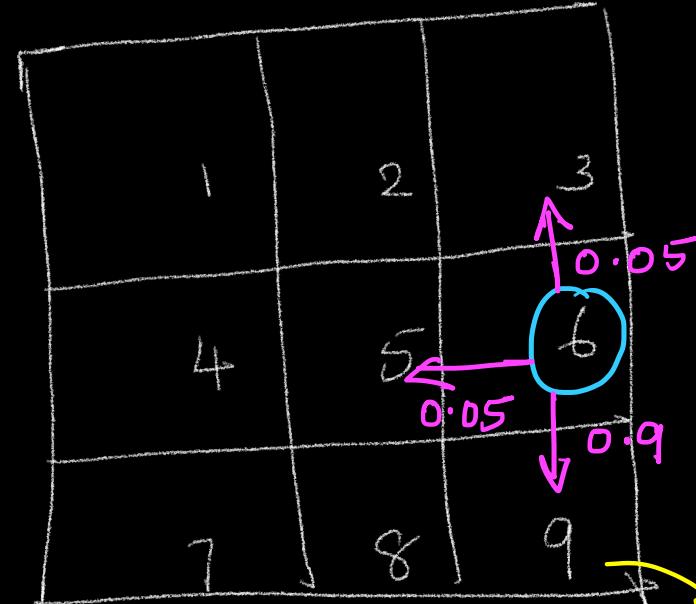
MDP:

↓
partly stochastic / random
decision made



MDP:

$$V(s) = \max_a [R(s,a) + \gamma V(s')]$$



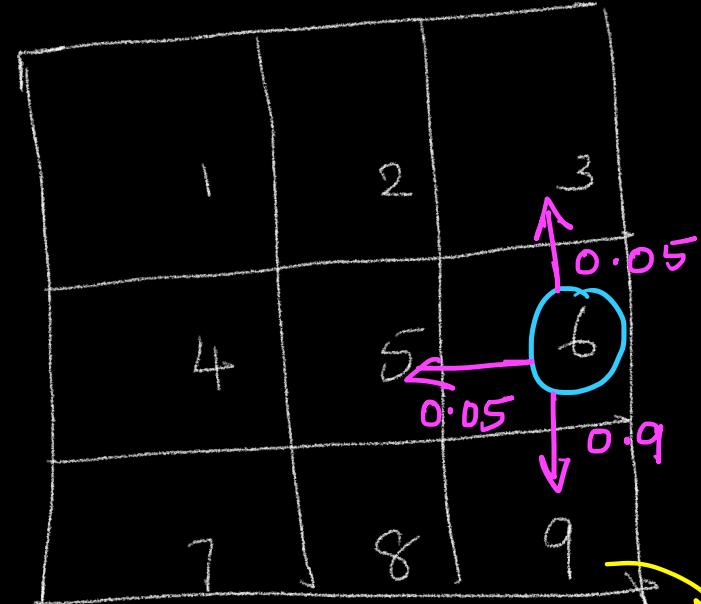
destination

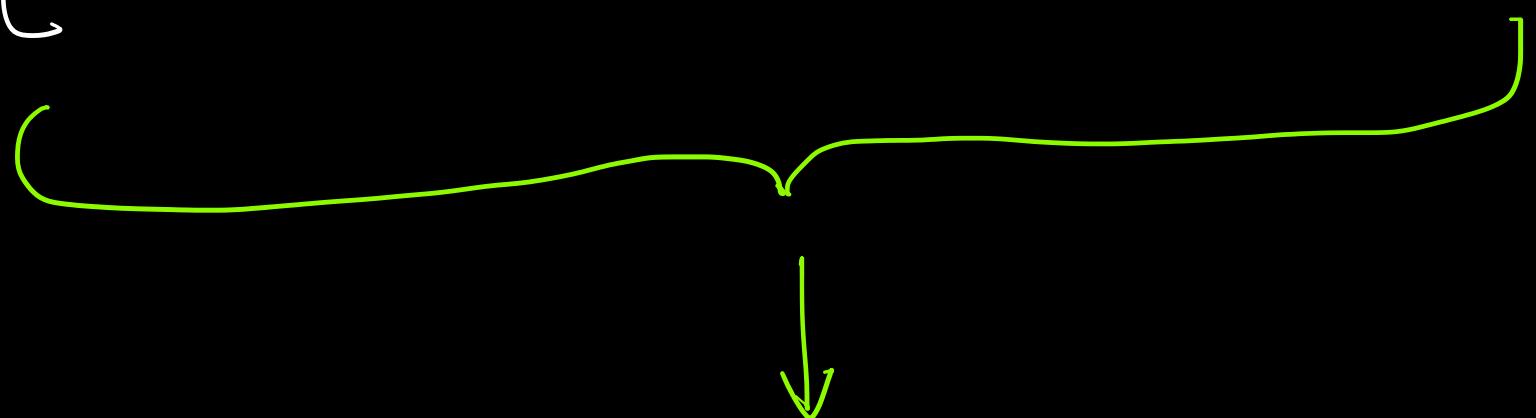
$$V(s) = \max_a R(s,a) + \gamma \sum_{s'} P(s,a,s') V(s')$$

MDP:

$$V(s) = \max_a R(s, a) + \gamma \cdot \sum_{s'} P(s, a, s') V(s')$$

Expected value of $V(s')$



$$V(s) = \max_a \left[R(s,a) + \gamma \sum_{s'} p(s,a,s') \cdot V(s') \right]$$


$Q(s,a)$: Q - function

$$\Rightarrow V(s) = \max_a Q(s,a)$$

Q - function:

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} p(s, a, s') \cdot V(s')$$



$$\max_{a'} Q(s', a')$$

$$Q(s, a) = R(s, a) + \gamma \cdot \sum_{s'} \left(p(s, a, s') \cdot \max_{a'} Q(s', a') \right)$$

Intuition:

$$Q(s, a) = R(s, a) + \gamma \cdot \sum_{s'} \left(p(s, a, s') \cdot \max_{a'} Q(s', a') \right)$$

Let $\gamma = 1$

Q :

| | a_1 | a_2 | \dots | a_K |
|----------|----------|-------|---------|-------|
| s_1 | | | \dots | |
| s_2 | | | \dots | |
| s_3 | | | \dots | |
| \vdots | \vdots | | \dots | |
| s_n | | | \dots | |

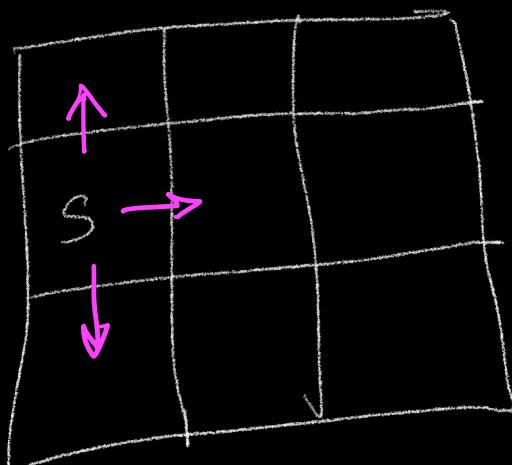
Temporal-difference : TD(s,a)

what if the environment is also changing

$$Q(s, \uparrow)$$

$$Q(s, \rightarrow)$$

$$Q(s \rightarrow \downarrow)$$



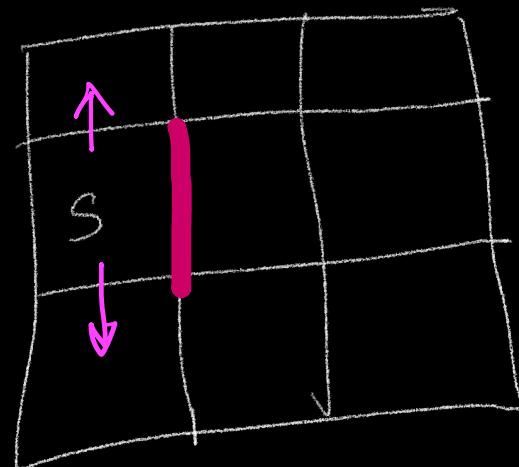
Time t-1:



$$Q(s, \uparrow)$$

$$Q(s, \rightarrow)$$

$$Q(s \rightarrow \downarrow)$$



Time t

$$Q_{t-1}(s, a) \quad \text{vs} \quad Q_t(s, a)$$

$$\begin{aligned}
 TD(s, a) &= Q_{\text{new}}(s, a) - Q_{\text{old}}(s, a) \\
 &= [R(s, a) + \gamma \sum_{s'} p(s, a, s') \max_{a'} Q(s', a')] \\
 &\quad - Q_{t-1}(s, a)
 \end{aligned}$$

$$Q_t(s, a) = Q_{t-1}(s, a) + \alpha \cdot TD_t(a, s)$$



how quickly can the robot
adapt to env. change

a.k.a. learning rate

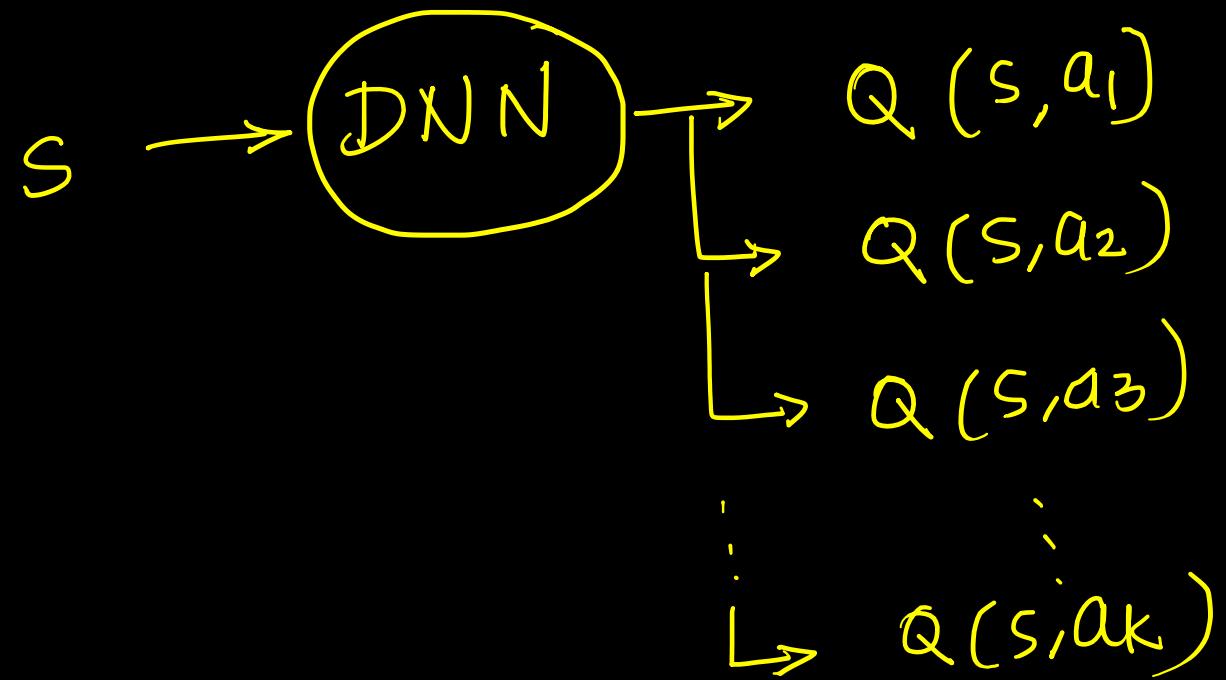
$$Q_t(s, a) = Q_{t-1}(s, a) + \alpha \cdot \left(R(s, a) + \gamma \sum_{s'} p(s, a, s') \max_{a'} Q(s', a') - Q_{t-1}(s, a) \right)$$

Simplified form:

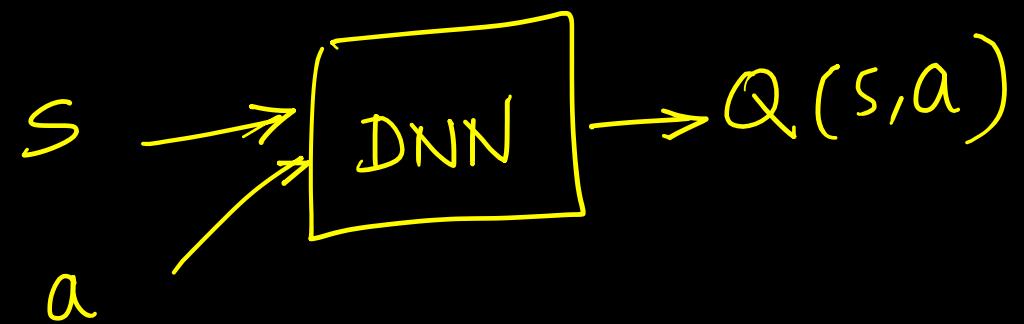
$$Q(s, a) = R(s, a) + \gamma \cdot \max_{a'} Q(s', a')$$

Deep-Q-learning : Intuition

$Q(s, a)$: Table



not



inefficient

Next Session : DQNN + code

Q & A

