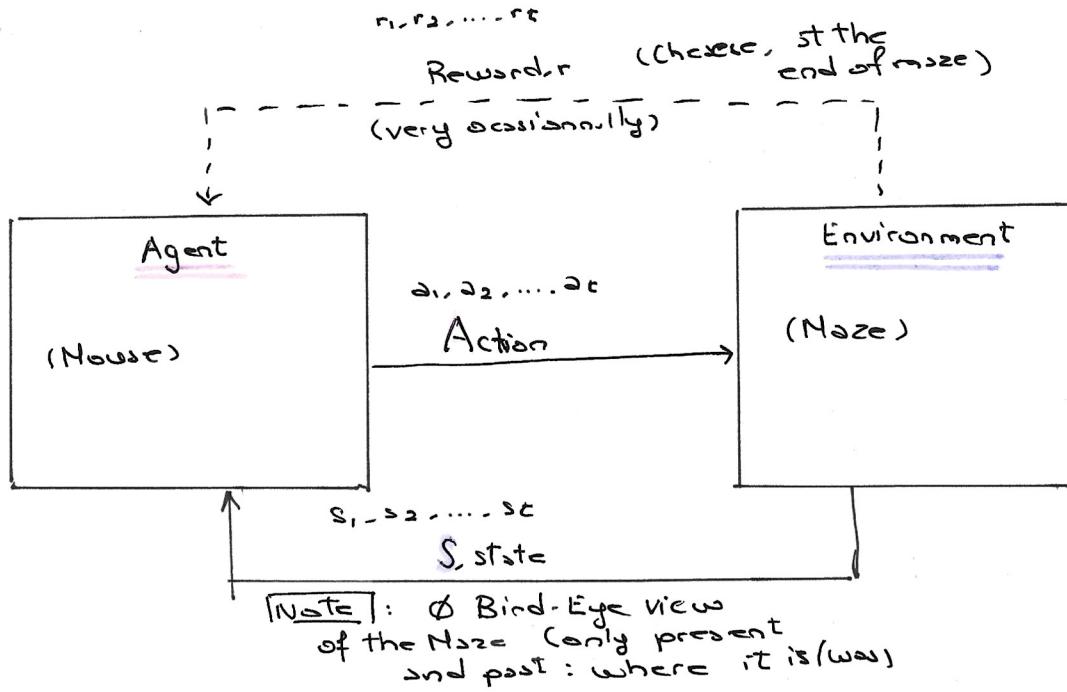


- A framework for Learning how to Interact with the Environment from Experience.
- Biologically - Inspired Idea (Trial/Error / Experiences)

## Reinforcement Learning Overview

→ NULL most of the time



- If every correct turn, it got a piece of cheese (reward), this is supervised learning. Rewards = LABELS.
- Because the reward is time-delayed (At the end of maze), not depending on each move of the mouse, Reward = TIME-DELAYED LABEL
  - ↳ Semi-supervised Learning.
- Since we have less labels, requires much more data to train these agents and more trial/error.
  - ↳ Harder optimization problem.

⇒ Another example:

Agent: Chess  
Environment: Rules of the game

## Mathematics of RL.

• Example: Chess:

→ Challenge: design a policy of what actions to take given a state  $s$  to max. my chance of getting a future reward.

↳  $\boxed{\text{policy } \pi(s, a)}$ :  $\Pr\{A=a | s=s\}$ ,  $\emptyset$  deterministic.

→ Only reward once you best chess. If I lose, do I remove the entire sequence of actions? How do you figure which actions are good/bad?

↳  $\boxed{\text{VALUE}}$ :  $V_{\pi}(s) = \mathbb{E} \left\{ \sum_t \gamma^t r_t | s_0 = s \right\}$  if I start  
↳ Expectation of how much reward I'll get in the future in that state  $s$  and I enact that policy

• Discount Rate  $0 \leq \gamma \leq 1$

→  $\boxed{\text{Goal}}$ : Optimize  $\pi(s, a)$ , the policy, to maximize future rewards.

## Markov Decision Process.

→ Environment contains a stochastic component

→ IP of going from current state/action to the next.

## Credit Assignment Problem.

→ Since the rewards are sparse  $\Rightarrow$  very difficult to find actions that are good/bad

→ 1960s Minsky

→ Dense vs. Sparse Rewards

→ Sample Efficiency Many trials/examples/games.

→ Reward Shaping.

## Optimization Techniques for RL.

- Differential Programming.
- Monte-Carlo
- Temporal Difference (Model-free)
- Bellman 1957
- Exploration vs. Exploitation.
- Policy Iteration.
- Gradient Descent, Evolutionary Opt., Simulated Annealing.

$$V(s_t) = \sum r^t \text{ future rewards}$$

$$Q(a_t, s_t) = r_{t+1} + \sum r^t$$

$$= r_{t+1} + V(s_t)$$

## Q-Learning.

- $Q(s, a)$ : quality of state / Action pair

$$Q^{\text{update}}(s_t, a_t) = Q^{\text{old}}(s_t, a_t) + \alpha \{ r_t + \gamma \max_a Q(s_{t+1}, a) - Q^{\text{old}}(s_t, a_t) \}$$

→ Given a state, we can choose the action that has the highest Quality.

$$\pi(a|s) \Rightarrow Q_\pi(\pi(s))$$

## indsight Replay.

- Data efficient
- Save «mistakes» for the future, ~~reduced sample efficient~~.

## Overview of Methods

2

Model-based RL → Based on history, only based on current/future states.

Markov Decision Process  $P(s', s, a)$

- Policy Iteration  $\pi(s, a)$
- Value Iteration  $V(s)$

Satisfies...  
Dynamic Programming & Bellman Optimality.

Non Linear Dynamics

$$\frac{d}{dt} x = f(x(t), u(t), t) dt$$

→ Optimal Control & HJB

↓  
Hamelin/Jacobi/  
Bellman

Brute force search  
Curse of Dimensionality.

Actor  
Critic

Deep  
NPC

Uses Quality +  
function

DQN

OFF Policy (\*)

• Run game more  
occasionally →  
helpful for  
learning the  
environment.

$$Q(s, a)$$

→ Q-Learning.  
(faster)

Deep  
Policy  
Network

$$\theta^{\text{new}} = \theta^{\text{old}} + \alpha P_\theta R_{s,a}$$

→ Policy Gradient Optimization

Model-free RL

Gradient-Free

On Policy

Always use the best  
policy every game

TD (Temporal Δ)

$$TD(0) \dots TD(\infty) \equiv NC$$

→ SARSA (conservative)

Gradient-based

DEEP L.

**3** MODEL-BASED RL: MARKOV DECISION PROCESS

• We assume that we have a model for how the env. works / how the system works. follows the MDP.  $\Rightarrow$  for  $P(s'|s, a)$ . We also assume a model for the reward  $(R(s', s, a))$ .

$$\cdot R(s', s, a) = \text{IP} \{ r_{k+1} | S_{k+1} = s', s_k = s, a_k = a \}$$

$$\cdot P(s', s, a) = \text{IP} \{ S_{k+1} = s' | s_k = s, a_k = a \}$$

$\Rightarrow$  Assumption: knows the rule of the game (Ex: knows how to checkmate someone in a chess game).

$\Rightarrow$  Dynamic Programming instead of Brute Force Search.

Value Function:  $V_\pi(s) = \mathbb{E} \left\{ \sum_k \gamma^k r_k | s_0 = s \right\}$  Expected value of future rewards given an initial state  $s_0 = s$ .

Optimized

$$\rightarrow V(s) = \max_\pi \mathbb{E} \left\{ \sum_{k=0}^{\infty} \gamma^k r_k | s_0 = s \right\}$$

$$= \max_\pi \mathbb{E} \left\{ r_0 + \sum_{k=1}^{\infty} \gamma^k r_k | s_1 = s' \right\} \Leftrightarrow \max_\pi \mathbb{E} [r_0 + \gamma V(s')]$$

$\downarrow$  current       $\uparrow$  future

$$V(s) = \max_\pi \mathbb{E} [r_0 + \gamma V(s')] \Rightarrow \text{Bellman's Equation. (On Policy)}$$

$\Downarrow$

$$\pi = \arg \max_\pi \mathbb{E} [r_0 + \gamma V(s')]$$

$\rightarrow$  Bellman: how to solve multi-step optimization problems by breaking into smaller recursive sub-problems. TD: list out all sub-problems, until you solve the whole thing (Table Approach)

$\rightarrow$  Bottom-Up vs Top-Down Approach.

TD: Always start with smaller sub-problems (Reward guaranteed  $\downarrow$  step / 2 steps... ahead)

Value Iteration  $\Rightarrow$  Dynamic Programming. Discount Rate.

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

$\Downarrow$   
Find the best Action

$\Downarrow$   
My current reward

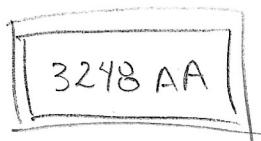
$\Downarrow$   
Next step given that Action a.

To get this optimal future reward

Assumes that  $P(s'|s, a)$ , and  $V(s')$  are given.

$R(s', s, a)$

TPB  $\Rightarrow$  continuous actions  
DQN  $\Rightarrow$  discrete actions

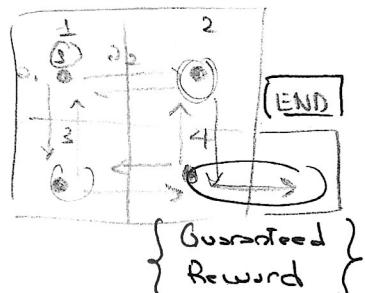


- Create a table of all possible states. (Initialize  $\pi$ ).
- Randomly take a state  $s$ .  
→ Maximize value function gives an optimized Action.
- Take a new state  $s'$  and repeat the same process.



↳ Iterate until the value function converges.

$$\pi(s, a) = \underset{a}{\operatorname{argmax}} \sum_{s'} P(s'|s, a) [R(s, s, a) + \gamma V(s')]$$



### Policy Iteration

→ directly onto the rules of the game (not to maximize the reward based on what actions to take)

- $V_\pi(s) = \mathbb{E} [R(s', s, \pi(s)) + \gamma V_\pi(s')] = \sum_{s'} P(s'|s, \pi(s)) \cdot [R(s, s, \pi(s)) + \gamma V_\pi(s')]$
- $\pi(s) = \underset{a}{\operatorname{argmax}} \mathbb{E} [R(s', s, \pi(s)) + \gamma V_\pi(s')]$

→ Typically converges in fewer iterations.



### Quality Function

$$Q(s, a) = \mathbb{E} [R(s', s, a) + \gamma V(s')] = \sum_{s'} P(s'|s, a) \cdot [R(s, s, a) + \gamma V(s')]$$

↳ we don't need a new model for the future states  $s'$  ⇒ implicit in the  $Q$  function  
⇒ becomes Model-free.

Therefore:

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

⇒ Expected value of future rewards for the maximized quality function given an action  $a$  and a state  $s$ .

$$\pi(s, a) = \underset{a}{\operatorname{argmax}} Q(s, a)$$

↳ Policy optimized: the action  $a$  that maximizes our future rewards (the quality function).

## 4) Q - LEARNING (GRADIENT-FREE)

$Q(s, a)$ : quality of state/action pair

$\Rightarrow V(s)$ : value of being in a current state  $s \Rightarrow$  assuming you take the best Action.

$\Rightarrow Q(s, a)$ : for any action  $a$  I may take (more information than the value function)

$$Q(s, a) = \mathbb{E} [R(s', s, a) + \gamma V(s')] = \sum_{s'} P(s'|s, a) [R(s', s, a) + \gamma V(s')]$$

↓                      ↓  
 current reward      future rewards.  
 (from current states  $s$  to  $s'$ )

↔

$$\left\{ \begin{array}{l} V(s) = \max_a Q(s, a) \\ \pi(s, a) = \arg \max_a Q(s, a) \end{array} \right\}$$

$\Rightarrow$  Generally, we don't know the models for  $P(s'|s, a)$  and  $R(s', s, a)$ .  
 $\hookrightarrow$  Stuck in Trial and Error Learning.

Monte-Carlo Learning (Model-Free)

(i) # steps (Episodes)  $\Rightarrow$  for one state.

$$R_\Sigma = \sum_{k=1}^n \gamma^k r_k \quad \text{TOTAL REWARD OVER EPISODE}$$

Estimated Reward ( $V^\text{old}$ )  
 Actual Reward ( $R_\Sigma$ )

$$\Rightarrow V^\text{new}(s_k) = V^\text{old}(s_k) + \frac{1}{n} \{ R_\Sigma - V^\text{old}(s_k) \} \quad \forall k \in \{1, \dots, n\}$$

$$\Rightarrow Q^\text{new}(s_k, a_k) = Q^\text{old}(s_k, a_k) + \frac{1}{n} \{ R_\Sigma - Q^\text{old}(s_k, a_k) \}$$

$\Rightarrow$  Not recommended, a lot of episodes required.  
 However, no biases.

The right action could've happened at any point in the past leading up to the reward now.

↳ TD: there's an optimal time (past) that is well correlated to the current reward.

## Temporal Difference Learning: TD(0)

•  $V(s_k) = \mathbb{E}[r_k + \gamma V(s_{k+1})]$   $\Rightarrow$  Bellman's Optimality.

$$\Rightarrow V^{\text{new}}(s_k) = V^{\text{old}}(s_k) + \alpha \left\{ r_k + \gamma V^{\text{old}}(s_{k+1}) - V^{\text{old}}(s_k) \right\}$$

$\overbrace{\qquad\qquad\qquad}^{\begin{cases} \text{TD Target} \\ \text{Estimate } R_\Sigma \end{cases}}$

$\overbrace{\qquad\qquad\qquad}^{\text{One step delay.}}$

$\hookrightarrow R_\Sigma^{(n)} = r_k + \gamma r_{k+1} + \gamma^2 r_{k+2} + \dots + \gamma^n r_{k+n} + \gamma^{n+1} V(s_{k+n+1})$

$$= \sum_{j=0}^n \gamma^j r_{k+j} + \gamma^{n+1} V(s_{k+n+1}) \quad \Rightarrow \text{for TD}(N)$$

If  $N \rightarrow \infty$ :  $\text{TD}(\infty) = \text{Monte-Carlo}$

TD-λ where  $0 \leq \lambda \leq 1$

$$R_\Sigma^\lambda = (\lambda - \lambda) \sum_{n=0}^{\infty} \lambda^{n-1} R_\Sigma^{(n)}$$

$$\Rightarrow V^{\text{new}}(s_k) = V^{\text{old}}(s_k) + \alpha \{ R_\Sigma^\lambda - V^{\text{old}}(s_k) \}$$



## Q-Learning

$\mapsto$  NOT on the Policy Functions...

• Temporal Difference on Quality functions...

$$Q^{\text{new}}(s_k, a_k) = Q^{\text{old}}(s_k, a_k) + \alpha \left\{ r_k + \gamma \max_a Q(s_{k+1}, a) - Q^{\text{old}}(s_k, a_k) \right\}$$

$\overbrace{\qquad\qquad\qquad}^{\begin{cases} \text{TD Error} \\ \text{TD Estimate } R_\Sigma \\ \text{Target} \end{cases}}$

$\Rightarrow$  Off-Policy TD(0) learning of the quality function  $Q$ .

$\Rightarrow$  You can still learn even if you take sub-optimal actions  $a$ .

SARSA: State - Action - Reward - State - Action

$$Q^{\text{new}}(s_k, a_k) = Q^{\text{old}}(s_k, a_k) + \alpha \left\{ r_k + \gamma Q^{\text{old}}(s_{k+1}, a_{k+1}) - Q^{\text{old}}(s_k, a_k) \right\}$$

$\overbrace{\qquad\qquad\qquad}^{\begin{cases} \text{Need to } \uparrow, \\ \text{we want to converge} \end{cases}}$

ON-POLICY

$\Rightarrow$  On-Policy TD(0) learning  
on the  $Q(s_k, a_k)$ .

- Q-Learning :  $\max_a Q(s_{k+1}, a) \Rightarrow$  off-policy

- { ① Better/Faster Learning  
Higher Variance
- ② Can Learn from Imitation  
and Exp. Replay.

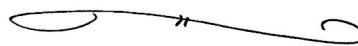
$\Rightarrow$  Epsilon-Greedy:  $\epsilon \uparrow \rightarrow$  cools down to 0

$\epsilon$ : % of exploration (random actions) into off-policy strategies.

- SARSA:  $Q^{\text{old}}(s_{k+1}, a_{k+1}) \Rightarrow$  On-Policy.

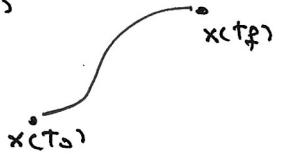
- { ① Often Safer
- ② Better Total Reward While Learning.

$\Rightarrow$  Q-Learning and SARSA: Approximate Dynamic Programming.  
 ↳ for Deep RL.



## 5 | OPTIMAL NON-LINEAR CONTROL

$\frac{d}{dt} x = f(x(t), u(t)) dt$ , where: info that characterizes the state of my system  
 (Ex: pendulum  $\Rightarrow$  position of the cart, b velocity, etc.)  
 .  $f$ : dynamics of that system  
 • Goal: design a control function  $u(t)$  to follow a state  $x(t)$   
 to minimize the cost  $J$ .



$$\Rightarrow J(x(t), u(t), t_0, t_f) = Q(x(t_f), t_f) + \int_{t_0}^{t_f} L(x(\tau), u(\tau)) d\tau$$

↳ Ex:  $x(\tau)^2 + u(\tau)^2$  for example.

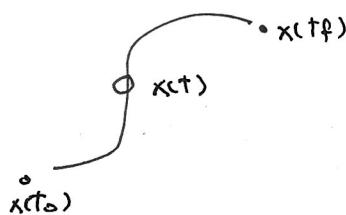
$$\Rightarrow V(x(t_0), t_0, t_f) = \min_{u(t)} J(x(t), u(t), t_0, t_f)$$

$$-\frac{\partial V}{\partial t} = \min_{u(t)} \left\{ \left( \frac{\partial V}{\partial x} \right)^T f(x(t), u(t)) + L(x(t), u(t)) \right\}$$

HAMILTON JACOBIAN  
BELLMAN (HJB)  
EQUATION.

$$V(x(t_0), t_0, t_f) = V(x(t_0), t_0, t) + V(x(t), t, t_f)$$

Bellman's  
Optimality.



## Deriving HJB Equation

$$\begin{aligned}
 \frac{\partial}{\partial t} V(x(t), t, t_f) &= \frac{\partial V}{\partial t} + \left[ \frac{\partial V}{\partial x} \right]^T \frac{dx}{dt} \\
 &= \min_{u(t)} \frac{d}{dt} \left\{ \int_0^{t_f} \mathcal{L}(x(\tau), u(\tau)) d\tau + Q(x(t_f), t_f) \right\} \\
 &= \min_{u(t)} \left\{ \underbrace{\frac{d}{dt} \int_0^{t_f} \mathcal{L}(x(\tau), u(\tau)) d\tau}_{-L(x(t), u(t))} \right\} \\
 \Rightarrow -\frac{\partial V}{\partial t} &= \min_{u(t)} \left\{ \left( \frac{\partial V}{\partial x} \right)^T f(x, u) + L(x, u) \right\}
 \end{aligned}$$

## Discrete-Time HJB

$$x_{k+1} = F(x_k, u_k)$$

$$J(x_0, \{u_k\}_{k=0,n}) = \sum_{k=0}^n \mathcal{L}(u_k, x_k) + Q(x_n, t_n)$$

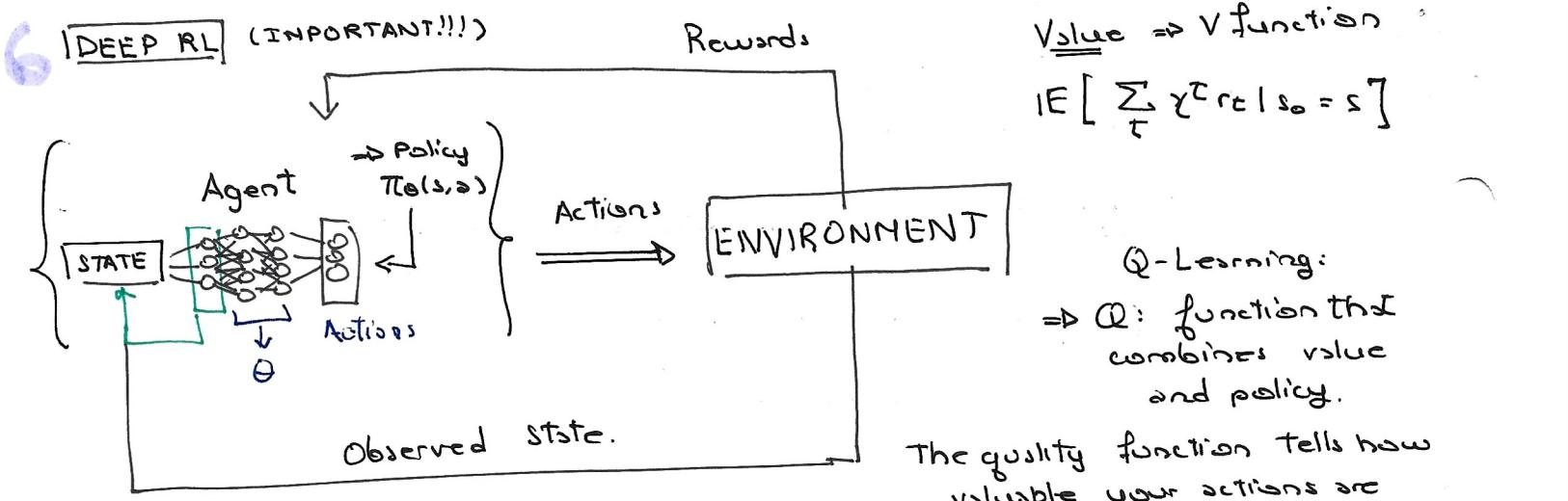
$$V(x_0, n) = \min_{\{u_k\}_{k=0}^n} J(x_0, \{u_k\}_{k=0,n})$$

$$V(x_0, n) = V(x_0, k) + V(x_k, n-1) \quad \forall k \in (0, n) \quad \Rightarrow \text{Optimality Condition}$$

$$\Rightarrow V(x_k, n) = \min_{u_k} \left\{ \underbrace{\mathcal{L}(x_k, u_k)}_{\substack{\text{current cost}}} + \underbrace{V(x_{k+1}, n-1)}_{\substack{\text{future steps} \\ \text{(its value function)}}} \right\}, \text{ where } x_{k+1} = F(x_k, u_k)$$

Therefore ...

$$\left\{
 \begin{array}{l}
 V(x) = \min_u (\mathcal{L}(x, u) + V(F(x, u))) \\
 \pi(x) = \underset{u}{\operatorname{argmin}} (\mathcal{L}(x, u) + V(F(x, u)))
 \end{array}
 \right.$$

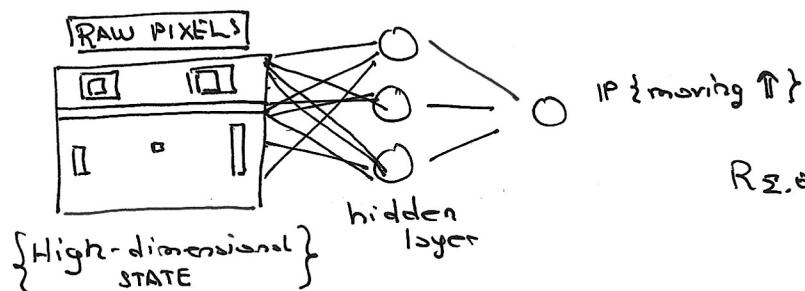


- However, the inconveniences related to RL are still influenced to Q-Learning.

## 7 DEEP RL METHODS (IMPORTANT!!!)

### 1. Deep Policy Network

$\Rightarrow \pi_\theta(s, a)$ : parametrized as a neural network (weights + bias)



{through back propagation}

$$\text{Policy } \pi_\theta(s, a) \\ \Rightarrow \theta^{\text{new}} = \theta^{\text{old}} + \alpha \nabla_{\theta} R_{\Sigma, \theta}$$

$R_{\Sigma, \theta}$ : cumulated expected reward  
 (we can use a value network for example)...

Policy Gradient Opt: (more general than NN)

$$R_{\Sigma, \theta} = \left[ \sum_{s \in S} \mu_\theta(s) \right] \left[ \sum_{a \in A} \pi_\theta(s, a) Q(s, a) \right]$$

↳ quality      ↳ probability

Randomized      Expected prob. of having  $s$       Randomized

{Asymptotic distribution}

$\Rightarrow$  Consistent with  $R_{\Sigma, \theta}$

$$\nabla_{\theta} R_{\Sigma, \theta} = \sum_{s \in S} \mu_\theta(s) \sum_{a \in A} Q(s, a) \nabla_{\theta} \pi_\theta(s, a)$$

↓  
 Move towards MAX Rewards!

why not differentiate this variable?

Assumed stationary  
 in an ergodic environment

$$\begin{aligned}
 \nabla_{\theta} R_{Z,\theta} &= \sum_{s \in S} \mu_{\theta}(s) \sum_{a \in A} \pi_{\theta}(s, a) Q(s, a) \cdot \frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)} \\
 &= \sum_{s \in S} \underbrace{\mu_{\theta}(s)}_{\text{Policy}} \sum_{a \in A} \underbrace{\pi_{\theta}(s, a) Q(s, a)}_{\text{Value}} \nabla_{\theta} \log \{ \pi_{\theta}(s, a) \} \\
 &= \mathbb{E} \left[ Q(s, a) \underbrace{\nabla_{\theta} \log \{ \pi_{\theta}(s, a) \}}_{\text{back propagation...}} \right] \\
 \theta^{\text{new}} &= \theta^{\text{old}} + \alpha \nabla_{\theta} R_{Z,\theta}
 \end{aligned}$$

$$\left\{ \begin{array}{l} \frac{(x^2)'}{x^2} = \frac{2x}{x^2} = \frac{2}{x} \\ \log \{ x^2 \} = \frac{1}{x^2} \cdot (x^2) \end{array} \right.$$

Example with  $x^2$

## 2. Deep Q-Learning (DQN)

• Off-Policy TD(·)

$$Q^{\text{new}}(s_k, a_k) = Q^{\text{old}}(s_k, a_k) + \alpha \{ r_k + \gamma \max_a Q(s_{k+1}, a) - Q^{\text{old}}(s_k, a_k) \}$$

$\Rightarrow Q(s, a) \approx Q(s, a, \theta) \Rightarrow$  parameterize Q Function with NN.

$\Downarrow$   
↓ dimension  
(Extract low dimensional features)

$$\begin{aligned}
 \text{• NN Cost Function: } \mathcal{L} &= \mathbb{E} \left[ \{ r_k + \gamma \max_a Q(s_{k+1}, a_{k+1}, \theta) - Q(s_k, a_k, \theta) \}^2 \right] \\
 &\text{minimize temporal diff-error by optimizing } \theta \text{ in NN.}
 \end{aligned}$$

## 3. Advantage Network (DDQN: Deep-Dueling)

$$Q(s, a, \theta) = V(s, \theta_1) + A(s, a, \theta_2)$$

$\downarrow$   
Value Network.  
(Explanation from state)       $\downarrow$  Advantage Network (how effective the action is in that state).

• Splits into 2 networks

## 4. Actor-Critic Network

$$\begin{cases} \text{ACTOR: Policy based} \rightarrow \pi(s, a) = \pi(s, a, \theta) \\ \text{CRITIC: VALUE based} \rightarrow V(s_k) = \mathbb{E} [r_k + \gamma V(s_{k+1})] \end{cases}$$

$$\theta_{k+1} = \theta_k + \alpha \{ r_k + \gamma V(s_{k+1}) - V(s_k) \}$$

$\underbrace{\theta_{k+1} - \theta_k}_{\text{Policy}}$        $\underbrace{r_k + \gamma V(s_{k+1}) - V(s_k)}_{\text{Value (Error)}}$

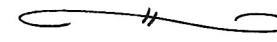
$\Rightarrow$  Best of Both worlds  
 $\Downarrow$   
Using TD signals  
from critic to update policy parameters

## 5. Advantage Actor-Critic Network

- { ACTOR: Deep Policy Network  $\rightarrow \pi(s, a) \approx \pi(s, a, \theta)$   
CRITIC: Deep Dueling Q Network  $\rightarrow Q(s_k, a_k, \theta_2)$

$$\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} \left\{ (\log \pi(s_k, a_k, \theta)) \cdot Q(s_k, a_k, \theta_2) \right\} \rightarrow \text{UPDATE-}$$

Policy   DDQN:  
  quality.



## 6

### DEEP RL IN FLUID DYNAMICS

- Biological RL: The rewards are internal (release dopamine)
- Usually, rewards are in the environment.