

Reinforcement Learning

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
لعلوم والتكنولوجيا
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

ML

supervised ML

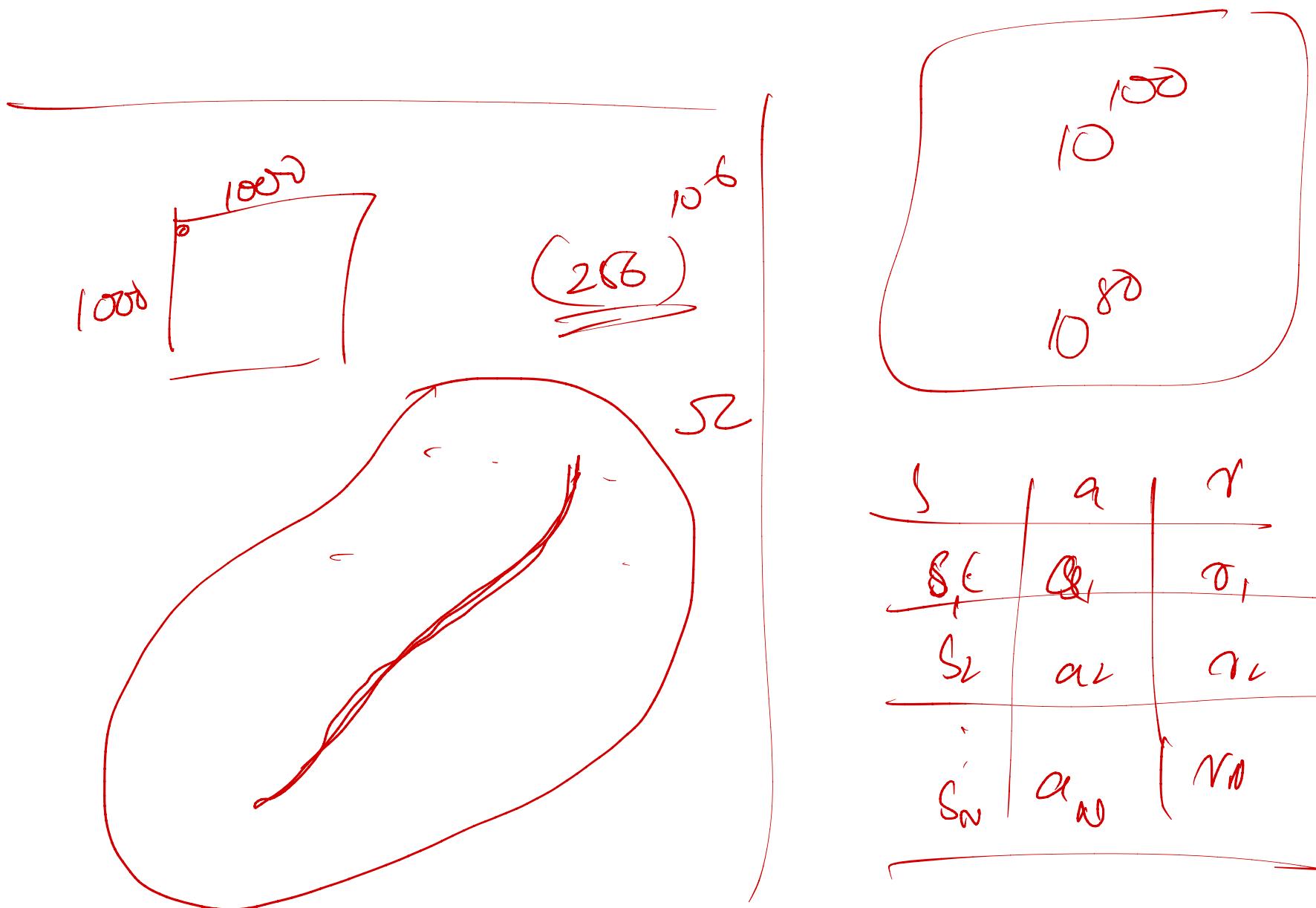
$x, y?$

$\hat{y} \in f_\theta(x)$

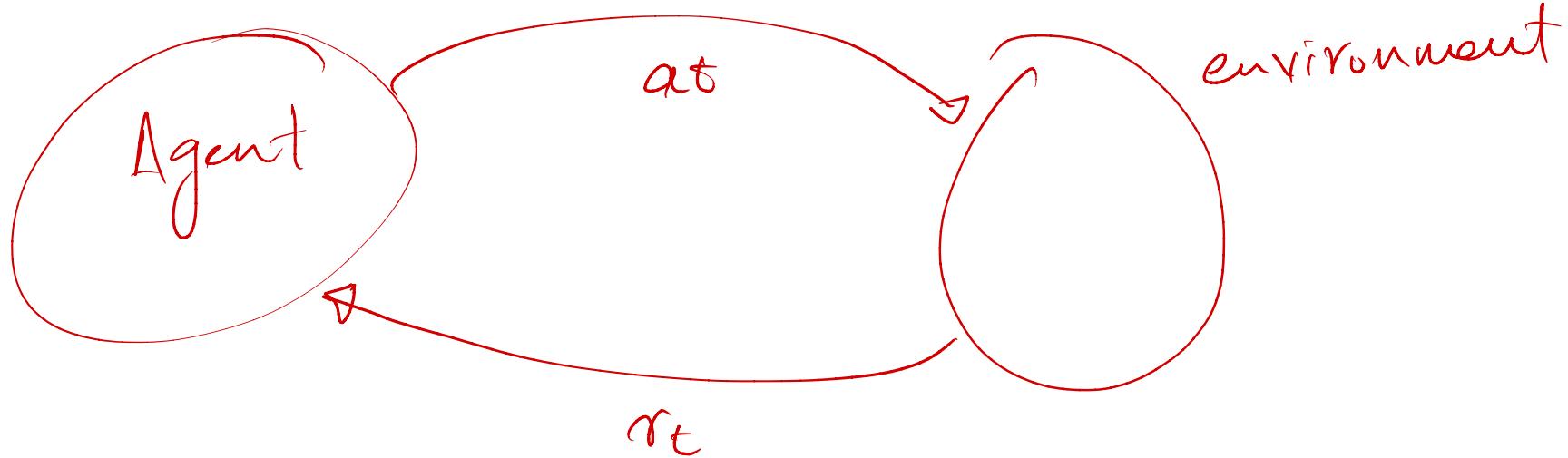
unsupervised ML
 X

$x \sim P_x(\lambda)$

playing chess

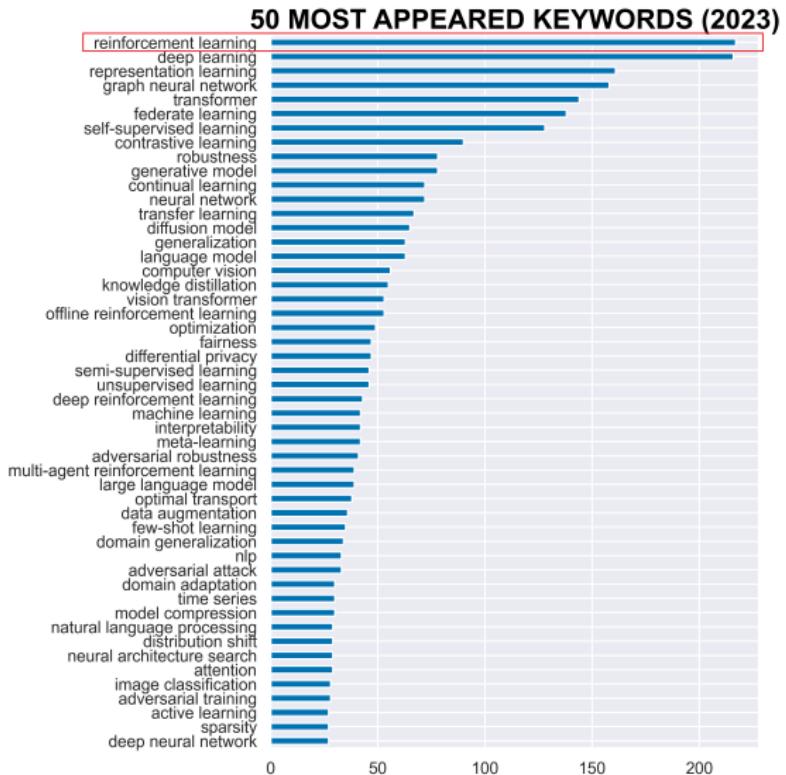


s	a	r
s_L	a_L	r_L
S_2	a_2	r_2
s_n	a_n	r_n



state s_t
 action a_t
 transition dynamics $\mathbb{P}[s_{t+1} | (s_t, a_t)]$
 reward r_{t+1} $p(r_t | (s_t, a_t))$

The probability problems involved are formidable. . . [but] the theory of sequential design will be of the greatest importance to mathematical statistics and to science... – Robins 1952

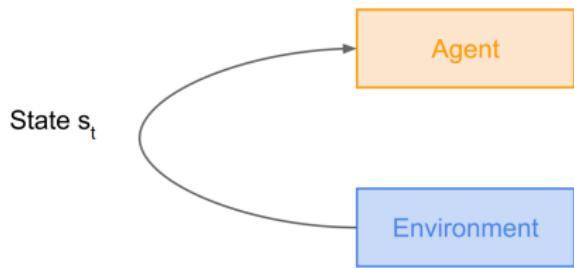


What is Reinforcement Learning?

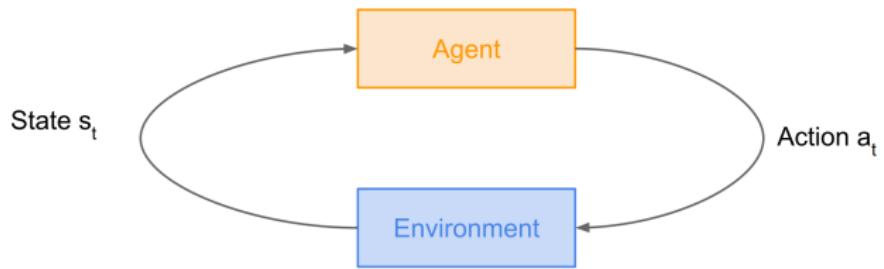
Agent

Environment

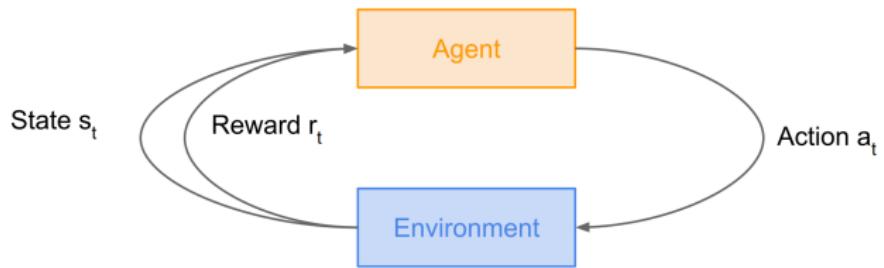
What is Reinforcement Learning? (cont.)



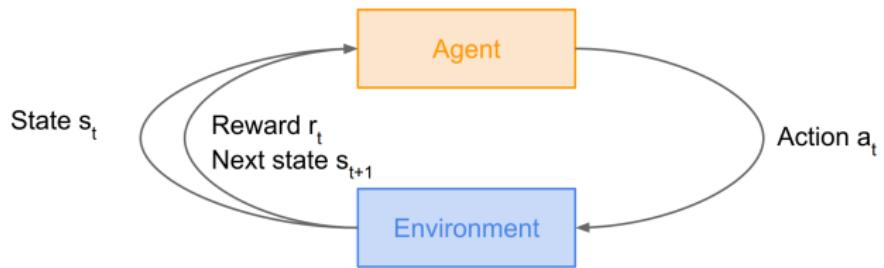
What is Reinforcement Learning? (cont.)

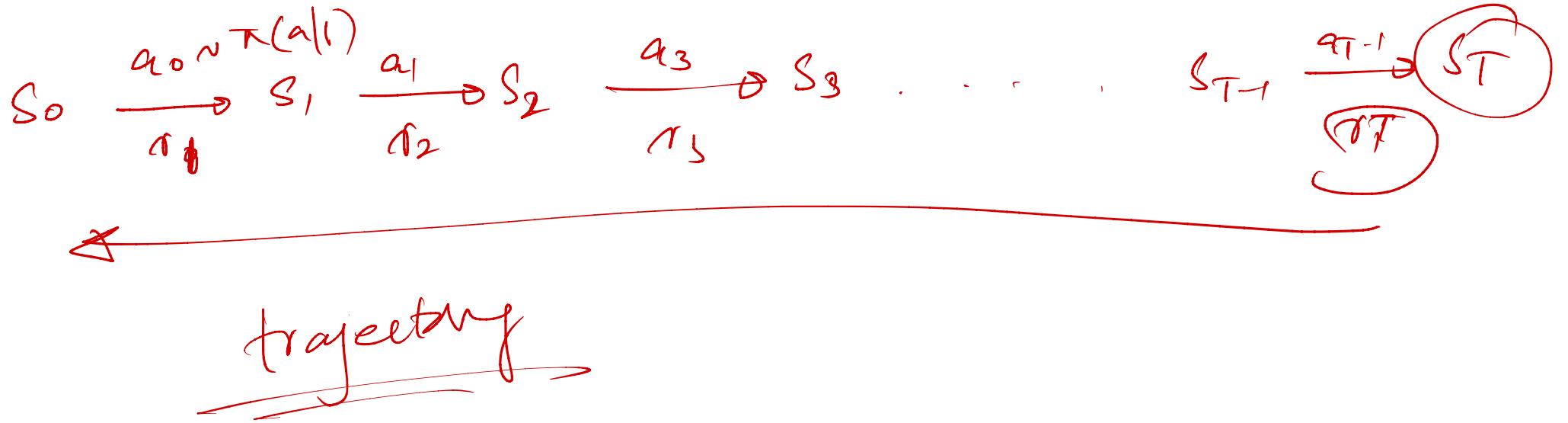


What is Reinforcement Learning? (cont.)



What is Reinforcement Learning? (cont.)





$$G = \sum_{t=0}^T r_t$$

maximize
 gain
 Expect Return

Reinforcement Learn

Value function
based methods

Policy Gradient
Based methods

Supervised Learning vs. Reinforcement Learning

Supervised Learning

- Given labeled data: $\{(x_i, y_i)\}$
learn $f(x) \approx y$

- Directly told what to output

- Inputs x are independently, identically distributed (i.i.d.)

$\sum_{in} \log P(a|x_i, y_i)$

$\pi \rightarrow policy$
~~no labels~~

Reinforcement Learning

- Learn behavior $\underline{\pi(a|s)}$

- from experience, indirect feedback

- data not i.i.d.: actions affect the future observations

Some Examples - ChatGPT

Step 1

Collect demonstration data and train a supervised policy.

A prompt is sampled from our prompt dataset.

Explain reinforcement learning to a 6 year old.

A labeler demonstrates the desired output behavior.

We give treats and punishments to teach...

This data is used to fine-tune GPT-3.5 with supervised learning.

SFT

Step 2

Collect comparison data and train a reward model.

A prompt and several model outputs are sampled.

Explain reinforcement learning to a 6 year old.

A labeler ranks the outputs from best to worst.

A: Explain reinforcement learning to a 6 year old.
B: Explain rewards...
C: We give treats and punishments to teach...
D: In machine learning...

This data is used to train our reward model.

RM

D > C > A > B

Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.

Write a story about otters.

PPO

The PPO model is initialized from the supervised policy.

Once upon a time...

RM

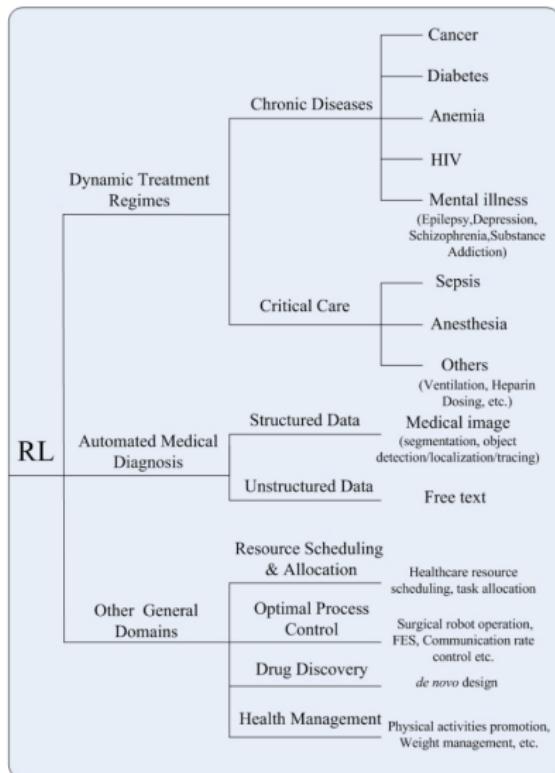
The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

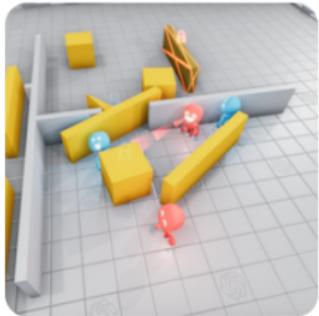
r_k

Some Examples - Healthcare

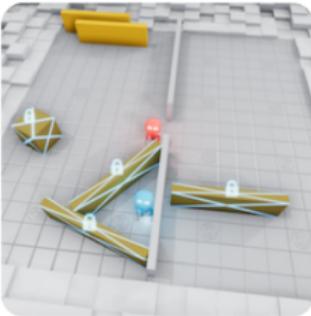


Some Examples - Multiagent Hide & Seek

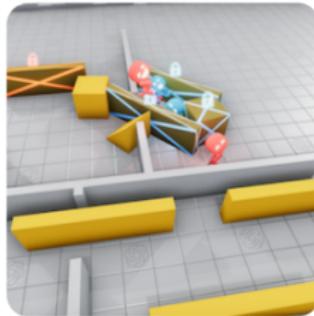
(a) Running and Chasing



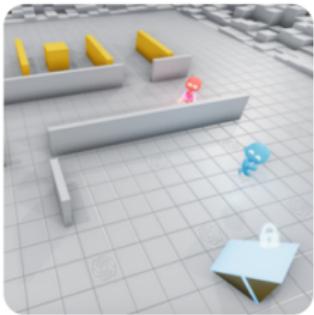
(b) Fort Building



(c) Ramp Use



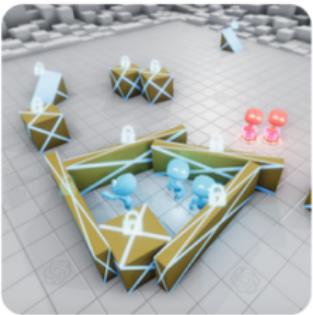
Some Examples - Multiagent Hide & Seek (cont.)



(d) Ramp Defense



(e) Box Surfing



(f) Surf Defense

Some Examples - Multiagent Hide & Seek (cont.)



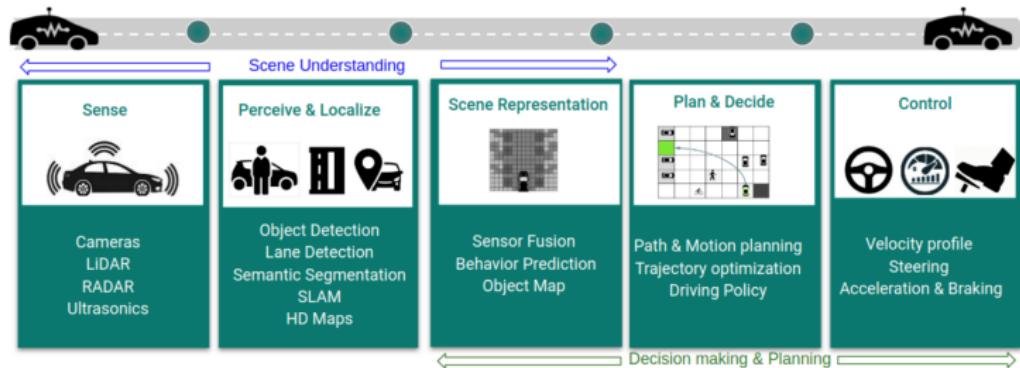
<https://www.youtube.com/watch?v=kopoLzvh5jY>

Some Examples - Grand Turismo GTSophie



[0https://www.gran-turismo.com/us/gran-turismo-sophy/technology/](https://www.gran-turismo.com/us/gran-turismo-sophy/technology/)

Some Examples - Autonomous Driving



'AlphaGo' AI Scores Narrow Win Against Ke Jie, World's Top 'Go' Player

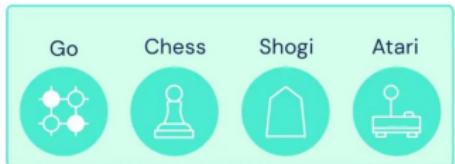
By [Lucian Armasu](#) published May 23, 2017

       [Comments \(4\)](#)



⁰<https://www.tomshardware.com/news/alphago-narrow-win-kejie,34486.html>

Some Examples - MuZero



MuZero learns the rules of the game, allowing it to also master environments with unknown dynamics.
(Dec 2020, Nature)

⁰<https://www.deepmind.com/blog/muzero-mastering-go-chess-shogi-and-atari-without-rules>

Some Examples - Minecraft

AI learns how to play Minecraft by watching videos



Open AI has trained a neural network to play Minecraft by Video PreTraining (VPT) on a massive unlabeled video dataset of human Minecraft play, while using just a small amount of labeled contractor data.

With a bit of fine-tuning, the AI research and deployment company is confident that its model can learn to craft diamond tools, a task that usually takes proficient humans over 20 minutes (24,000 actions). Its model uses the native human interface of keypresses and mouse...



29 June 2022 | Artificial Intelligence

⁰<https://www.artificialintelligence-news.com/2022/06/29/ai-learns-how-to-play-minecraft-by-watching-videos/>

Some Examples

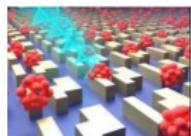


Reinforcement learning: From board games to protein design

Scientists have successfully applied reinforcement learning to a challenge in molecular biology. The team of researchers developed powerful new protein design software adapted from a strategy proven adept at board games like ...

BIOTECHNOLOGY

⌚ APR 20, 2023 🗂 0 📁 37



Chiral detection of biomolecules based on reinforcement learning

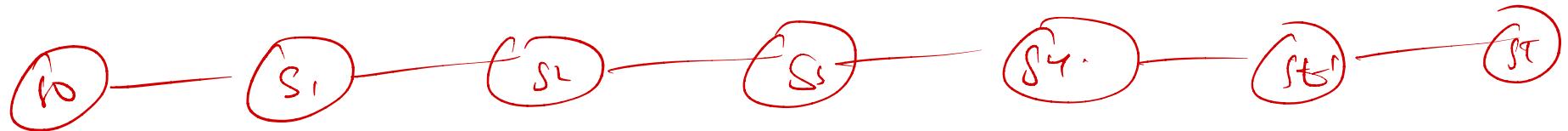
As one of the basic physical properties, chirality plays an important role in many fields. Especially in biomedical chemistry, the discrimination of enantiomers is a very important research subject. Most biomolecules exhibit ...

OPTICS & PHOTONICS

⌚ FEB 22, 2023 🗂 0 📁 4

states
actions

- ▶ Mathematical formulation of the RL problem
- ▶ **Markov property:** Current state completely characterises the state of the world i.e., future is independent of the past given the present
- ▶ Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$
 - \mathcal{S} : set of possible states
 - \mathcal{A} : set of possible actions
 - \mathcal{R} : distribution of reward given (state, action) pair
 - \mathbb{P} : transition probability i.e. distribution over next state given (state, action) pair
 - γ : discount factor



$$p(s_t | s_{t-1}) = p(s_t | s_{t-1}, s_{t-2}, \dots, s_0)$$



\downarrow
Markov process

process

Markov Decision

$$p(s_{t+1} | s_t, a_t) = p(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0)$$

Markov Decision Process (cont.)

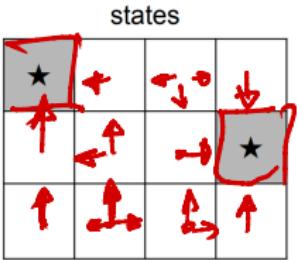
- ▶ At time step $t = 0$, environment samples initial state $s_0 \sim p(s_0)$
- ▶ Then, for $t=0$ until done:
 - Agent selects action a_t
 - Environment samples reward $r_t \sim R(\cdot | s_t, a_t)$
 - Environment samples next state $s_{t+1} \sim P(\cdot | s_t, a_t)$
 - Agent receives reward r_t and next state s_{t+1}
- ▶ A policy π is a function from S to A that specifies what action to take in each state
- ▶ **Objective:** Find policy π^* that maximizes cumulative discounted reward: $\sum_{t \geq 0} \gamma^t r_t$

Reinforcement Learning

- ▶ The Agent and the Environment interact in a time-sequenced loop
- ▶ Agent responds to [State, Reward] by taking an Action
- ▶ Environment responds by producing next step's State
- ▶ Environment also produces a scalar denoted as Reward
- ▶ Each State is assumed to have the Markov Property, meaning:
 - Next State/Reward depends only on Current State (for a given Action)
 - Current State captures all relevant information from History
 - Current State is a sufficient statistic of the future (for a given Action)
- ▶ Goal of Agent is to maximize Expected Sum of all future Rewards
- ▶ By controlling the (Policy : State → Action) function
- ▶ This is a dynamic (time-sequenced control) system under uncertainty

A simple MDP: Grid World

actions = {
1. right →
2. left ←
3. up ↑
4. down ↓
}

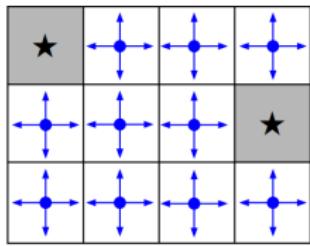


Set a negative “reward”
for each transition
(e.g. $r = -1$)

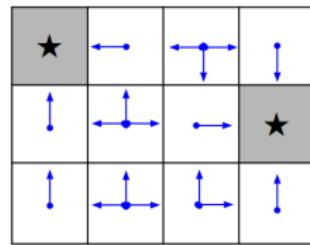
Objective: reach one of terminal states (greyed out) in least number of actions

$$P(w_t | w_1, \dots, w_{t-1})$$

A simple MDP: Grid World (cont.)



Random Policy



Optimal Policy



Real World Problems that fit MDP Framework

- ▶ Self-driving vehicle (speed/steering to optimize safety/time)
- ▶ Game of Chess (Boolean Reward at end of game)
- ▶ Complex Logistical Operations (eg: movements in a Warehouse)
- ▶ Make a humanoid robot walk/run on difficult terrains
- ▶ Manage an investment portfolio
- ▶ Control a power station
- ▶ Optimal decisions during a football game
- ▶ Strategy to win an election (high-complexity MDP)

Why are these problems hard?

- ▶ State space can be large or complex (involving many variables)
- ▶ Sometimes, Action space is also large or complex
- ▶ No direct feedback on “correct” Actions (only feedback is Reward)
- ▶ Time-sequenced complexity (Actions influence future States/Actions)
- ▶ Actions can have delayed consequences (late Rewards)
- ▶ Agent often doesn't know the Model of the Environment
- ▶ “Model” refers to probabilities of state-transitions and rewards
- ▶ So, Agent has to learn the Model AND solve for the Optimal Policy
- ▶ Agent Actions need to tradeoff between “explore” and “exploit”

- ▶ Policy π determines how the agent chooses actions
- ▶ $\pi : S \rightarrow A$, mapping from states to actions
- ▶ Deterministic Policy:

$$\pi(s) = a$$

- ▶ Stochastic Policy:

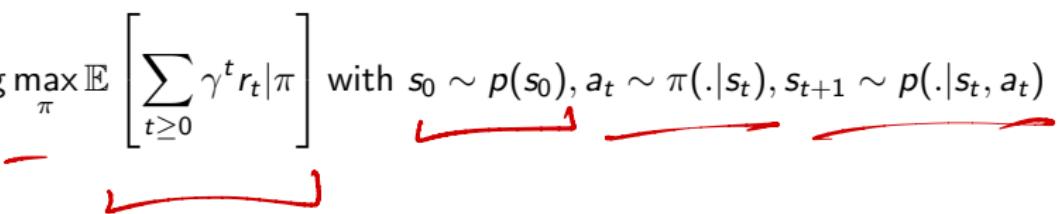
$$\pi(a|s) = Pr(a_t = a | s_t = s)$$

The Optimal Policy

- ▶ We want to find optimal policy π^* that maximizes the sum of reward
- ▶ But, how do we handle the randomness (initial state, transition probability...)?

The Optimal Policy

- ▶ We want to find optimal policy π^* that maximizes the sum of reward
- ▶ But, how do we handle the randomness (initial state, transition probability...)?
- ▶ **Solution:** Maximize the expected sum of rewards!
- ▶ Formally,

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi \right] \text{ with } s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$$


- ▶ How good is a state?

$$V^\pi(s) = E[\xi] = E\left[\beta^t r_t \mid a_t \sim \pi(a_t|s_t), s_0 = s\right]$$

Bellman Equations

$$\hat{V}(s) = \mathbb{E}_{\pi}[r(s, \underline{\pi(a|s)}) + \gamma V(s')]$$

$$\begin{aligned}\hat{r} &= E[r] \\ r &\sim p(s'|s,a)\end{aligned}$$

✓ factor

$$\begin{aligned}\hat{V}(s) &= \mathbb{E}_{\pi}[r(s, \underline{\pi(a|s)}) + \gamma V(s')] \\ &\quad + \gamma \mathbb{E}_{s' \sim P(s'|s,a)}[V(s')]\end{aligned}$$

~~θ function~~

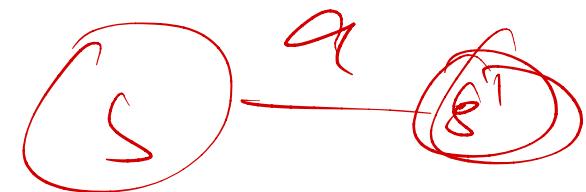
$$V^\pi(s) =$$

$$E \left[\sum_{t \geq 0} \beta^t r_t \mid s_0 = s, a_0 \sim \pi(a_0 | s) \right]$$

$$\tilde{\theta}^{\pi}(s, a) = \pi(s, a) + \beta \mathbb{E}_{\substack{s' \sim p(s'|s, a)}} [\tilde{V}^{\pi}(s')]$$

$$\tilde{V}^{\pi}(s) = \sum_a \pi(a|s) \tilde{\theta}^{\pi}(s, a)$$

$$\theta^{\pi}(s, a) \geq \hat{J}(s, a) + \beta \left(\sum_{s' \in \text{Sup}(\mathcal{C}(s, a)) \cap A(s)} \pi(a|s) \cdot \theta(s', a') \right).$$

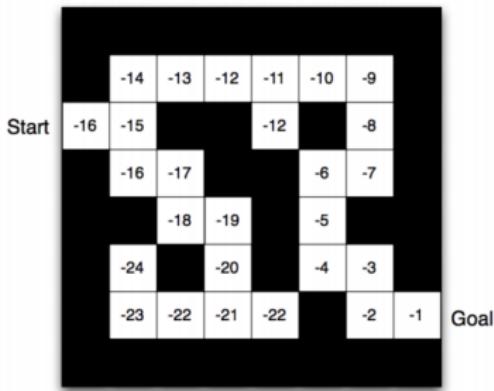
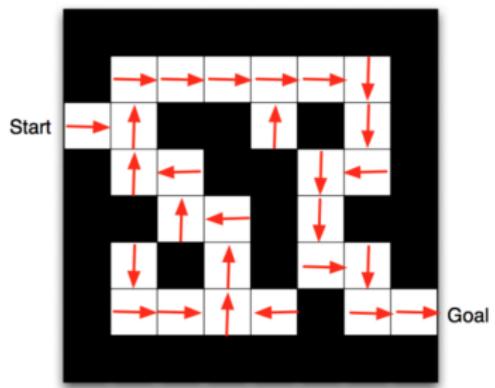


- ▶ How good is a state?
- ▶ The **value function** at state s , is the expected cumulative reward from following the policy from state s :

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi \right]$$

- ▶ Computing the value function is generally impractical, but we can try to approximate (learn) it
- ▶ The benefit is credit assignment: see directly how an action affects future returns rather than wait for rollouts

Value Function



Q-Value Function

- ▶ Can we use a value function to choose actions?

$$\arg \max_a r(s_t, a) + \gamma \mathbb{E}_{p(s_{t+1}|s_t, a_t)} [v^\pi(s_{t+1})]$$

- ▶ Can we use a value function to choose actions?

$$\arg \max_a r(s_t, a) + \gamma \mathbb{E}_{p(s_{t+1}|s_t, a_t)} [v^\pi(s_{t+1})]$$



- ▶ **Problem:** this requires taking the expectation with respect to the environment's dynamics, which we don't have direct access to!
- ▶ Instead, learn a Q-value function.

- ▶ The Q-value function at state s and action a , is the expected cumulative reward from taking action a in state s and then following the policy:

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

- ▶ Relationship:

$$V^\pi(s) = \sum_a \pi(a|s) Q^\pi(s, a)$$

- ▶ Optimal Action:

$$\arg \max_a Q^\pi(s, a)$$

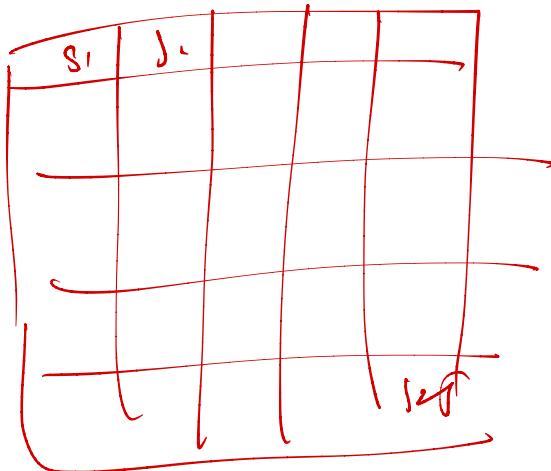
$$\hat{Q}^{\pi}(s, a) = r(s, a) + \beta \mathbb{E}_{\substack{s' \sim p(s'|s, a)}} [\hat{V}^{\pi}(s')]$$

$$\hat{Q}^{\pi}(s, a) = r(s, a) + \beta \mathbb{E}_{\substack{s' \sim p(s'|s, a)}} \left[\sum_{a' \sim \pi(a'|s')} \hat{Q}^{\pi}(s', a') \right]$$

$$\underline{Q}^{\pi}(s, a) \geq r(s, a) + \beta \mathbb{E}_{\substack{s' \sim p(s'|s, a)}} \left[\sum_{\substack{a' \in A(s) \\ a' \neq a}} \pi(a'|s) \hat{Q}^{\pi}(s', a') \right]$$

$\theta(s, a)$

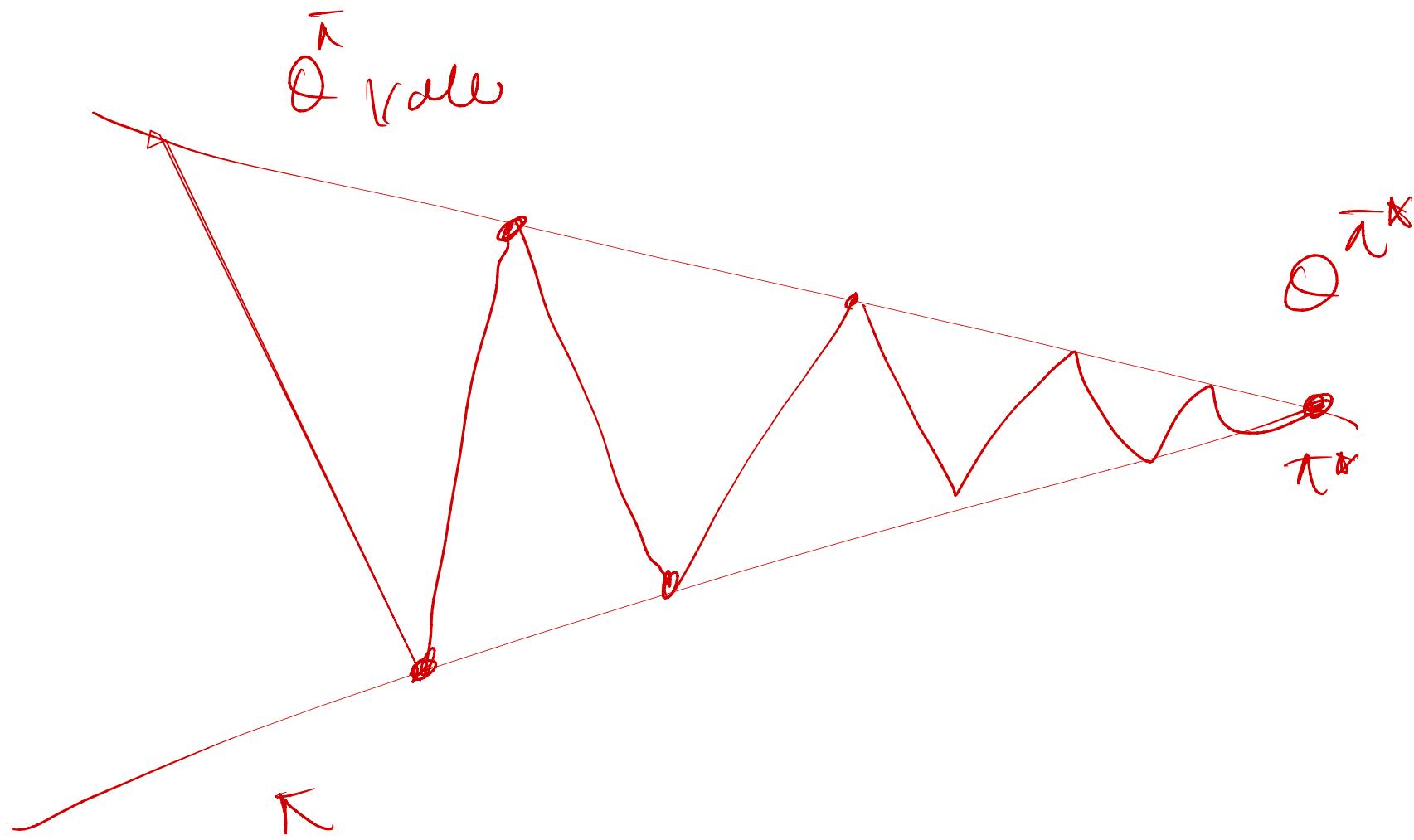
	a_1	a_2	a_3	a_4
s_1	0	0	0	0
s_2	0	0	0	0
s_3	0	0	0	0
s_4	0	0	0	0
s_5	0	0	0	0
s_6	0	0	0	0
s_7	0	0	0	0
s_8	0	0	0	0
s_9	0	0	0	0
s_{10}	0	0	0	0
s_{11}	0	0	0	0
s_{12}	0	0	0	0
s_{13}	0	0	0	0
s_{14}	0	0	0	0
s_{15}	0	0	0	0
s_{16}	0	0	0	0
s_{17}	0	0	0	0
s_{18}	0	0	0	0
s_{19}	0	0	0	0
s_{20}	0	0	0	0



$$s_0 \sim p(s_0)$$

$$\pi(a_t | s_t) = \max_a \theta(s_t, a)$$

$$\pi(a_t | s_t) = \begin{cases} 1 & \arg \max_a \theta(s_t, a) \\ 0 & \text{otherwise} \end{cases}$$



Q Learning Algorithm

$\pi \rightarrow$ uniform Random

Repeat until convergence

calculate $\hat{Q}(s, a)$

optimize $\pi(a_t | s_t) =$

$$\begin{cases} 1-\epsilon & \text{if } a_t = \arg\max \\ \epsilon & \text{otherwise} \\ \frac{1}{|A(s)|} & \end{cases}$$

$$\theta^\pi(s, a) = r(s, a) + \beta \mathbb{E}_{\substack{s' \sim p(s'|s, a)}} [\sqrt{\pi}(s')]$$

$$\theta^\pi(s, a) = r(s, a) + \beta \mathbb{E}_{\substack{s' \sim p(s'|s, a)}} \left[\sum_{a' \in A(s')} \pi(a'|s') \theta^\pi(s', a') \right]$$

$$\theta^t(s, a) \leftarrow \theta^{t-1}(s, a)(1 - \alpha) + \alpha \hat{\theta}^*(s, a)$$

- ▶ The Bellman Equation is a recursive formula for the Q-value function:

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)\pi(a'|s')} [Q^\pi(s', a')]$$

- ▶ There are various Bellman equations, and most RL algorithms are based on repeatedly applying one of them.

Optimal Bellman Equation

- ▶ The optimal policy π^* is the one that maximizes the expected discounted reward, and the optimal Q-value function Q^* is the Q-value function for π^*
- ▶ The Optimal Bellman Equation gives a recursive formula for Q^* :

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} \left[\max_{a'} Q^\pi(s_{t+1}, a') | s_t = s, a_t = a \right]$$

- ▶ This system of equations characterizes the optimal Q-value function. So maybe we can approximate Q^* by trying to solve the optimal Bellman equation!
- ▶ **Intuition:** If the optimal state-action values for the next time-step $Q^*(s', a')$ are known, then the optimal strategy is to take the action that maximizes the expected value of $r(s, a) + \gamma Q^\pi(s', a')$

- ▶ Let Q be an action-value function which hopefully approximates Q^*
- ▶ Q-learning is an algorithm that repeatedly adjusts Q to minimize the Bellman error
- ▶ Each time we sample consecutive states and actions (s_t, a_t, s_{t+1})

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \underbrace{\left[r(s, a) + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]}_{\text{Bellman Error}}$$

Exploration-Exploitation Tradeoff

- ▶ Notice: Q-learning only learns about the states and actions it visits.
- ▶ Exploration-exploitation tradeoff: the agent should sometimes pick suboptimal actions in order to visit new states and actions.
- ▶ Simple solution: ϵ -greedy policy
 - With probability $1 - \epsilon$, choose the optimal action according to Q
 - With probability ϵ , choose a random action
- ▶ ϵ -greedy policy still widely used today regardless of simplicity

Q-Learning Algorithm

Initialize $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 Take action A , observe R, S'

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

$S \leftarrow S'$;

 until S is terminal

$$\Theta(s, a) \rightarrow \Theta(s, a) (1 - \alpha) + \alpha \{ R + \gamma \max_a Q(s', a) \}$$

Deep Q-Learning (DQN)

- ▶ So far, we've been assuming a tabular representation of Q : one entry for every state/action pair
- ▶ What's the problem with this?

Deep Q-Learning (DQN)

- ▶ So far, we've been assuming a tabular representation of Q : one entry for every state/action pair
- ▶ What's the problem with this?
- ▶ Not scalable. Must compute $Q(s, a)$ for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

Deep Q-Learning (DQN)

- ▶ So far, we've been assuming a tabular representation of Q : one entry for every state/action pair
- ▶ What's the problem with this?
- ▶ Not scalable. Must compute $Q(s, a)$ for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!
- ▶ **Solution:** use a function approximator to estimate $Q(s, a)$. E.g. a neural network!

$$Q(s, a, ; \theta) \approx Q^*(s, a)$$

Deep Q-Learning (DQN)

- ▶ So far, we've been assuming a tabular representation of Q : one entry for every state/action pair
- ▶ What's the problem with this?
- ▶ Not scalable. Must compute $Q(s, a)$ for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!
- ▶ **Solution:** use a function approximator to estimate $Q(s, a)$. E.g. a neural network!

$$Q(s, a, ; \theta) \approx Q^*(s, a)$$

- ▶ If the function approximator is a deep neural network \Rightarrow Deep Q-Learning (DQN)

Deep Q-Learning (DQN)

- ▶ Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{p(s'|s, a)} \left[r(s, a) + \gamma \max_{a'} Q^*(s', a') | s, a \right]$$

- ▶ Forward Pass:

$$\mathcal{L}(\theta) = \mathbb{E}_{s, a \sim p(\cdot)} [(y - Q(s, a; \theta))^2]$$

where $y = \mathbb{E}_{p(s'|s, a)} \left[r(s, a) + \gamma \max_{a'} Q(s', a'; \theta) | s, a \right]$

- ▶ Backward Pass:

- Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta} \mathcal{L}(\theta) = \mathbb{E}_{s, a \sim p(\cdot), p(s'|s, a)} \left[(r(s, a) + \gamma \max_{a'} Q(s', a'; \theta) - Q(s, a; \theta)) \nabla_{\theta} Q(s, a; \theta) \right]$$

Training the Q-network: Experience Replay

- ▶ Learning from batches of consecutive samples is problematic:
 - Samples are correlated \Rightarrow inefficient learning
 - Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand side) \Rightarrow can lead to bad feedback loops

Training the Q-network: Experience Replay

- ▶ Learning from batches of consecutive samples is problematic:
 - Samples are correlated \Rightarrow inefficient learning
 - Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand side) \Rightarrow can lead to bad feedback loops
- ▶ Address these problems using experience replay
 - Continually update a replay memory table of transitions (s_t, a_t, r_t, s_{t+1}) as game (experience) episodes are played
 - Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N
Initialize action-value function Q with random weights

for episode = 1, M **do**
 Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
 for $t = 1, T$ **do**
 With probability ϵ select a random action a_t
 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
 Execute action a_t in emulator and observe reward r_t and image x_{t+1}
 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}
 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
 end for
end for

DQN with Target Network

- ▶ DQN, as stated above, will have learning problems because the target and the prediction are not independent as they both rely on the same network.
- ▶ It's like a dog chasing its own tail.

DQN with Target Network

- ▶ DQN, as stated above, will have learning problems because the target and the prediction are not independent as they both rely on the same network.
- ▶ It's like a dog chasing its own tail.
- ▶ **Solution:** Use two separate Q -value estimators, each of which is used to update the other.
- ▶ The target values are calculated using a target Q-network. The target Q-network's parameters are updated to the current networks every C time steps.
- ▶ Target network prevents the network from spiraling around.

DQN with Target Network Algorithm

Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory D to capacity N

Initialize action-value function Q with random weights θ

Initialize target action-value function \hat{Q} with weights $\theta^- = \theta$

For episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

For $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in D

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from D

 Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ with respect to the network parameters θ

 Every C steps reset $\hat{Q} = Q$

End For

End For

⁰Minh et al. <https://www.nature.com/articles/nature14236>

- ▶ Q-learning suffers from a maximization bias.
- ▶ This is because the update target is $r + \max_a Q^*(s, a)$. If Q -value is slightly overestimated then this error gets compounded.

- ▶ Q-learning suffers from a maximization bias.
- ▶ This is because the update target is $r + \max_a Q^*(s, a)$. If Q -value is slightly overestimated then this error gets compounded.
- ▶ **Solution:** Decompose the max operation in the target into action selection and action evaluation.
- ▶ Use the current network to select the max action for the next state and then use the target network to get the target Q -value for that action.
- ▶ Using these independent estimators, we can have unbiased Q -value estimates of the actions selected using the opposite estimator.
- ▶ We can thus avoid maximization bias by disentangling our updates from biased estimates.

DDQN Algorithm

Algorithm 1: Double DQN Algorithm.

```
input :  $\mathcal{D}$  – empty replay buffer;  $\theta$  – initial network parameters,  $\theta^-$  – copy of  $\theta$ 
input :  $N_r$  – replay buffer maximum size;  $N_b$  – training batch size;  $N^-$  – target network replacement freq.
for episode  $e \in \{1, 2, \dots, M\}$  do
    Initialize frame sequence  $\mathbf{x} \leftarrow ()$ 
    for  $t \in \{0, 1, \dots\}$  do
        Set state  $s \leftarrow \mathbf{x}$ , sample action  $a \sim \pi_B$ 
        Sample next frame  $x^t$  from environment  $\mathcal{E}$  given  $(s, a)$  and receive reward  $r$ , and append  $x^t$  to  $\mathbf{x}$ 
        if  $|\mathbf{x}| > N_f$  then delete oldest frame  $x_{t_{min}}$  from  $\mathbf{x}$  end
        Set  $\mathbf{x}' \leftarrow \mathbf{x}$ , and add transition tuple  $(s, a, r, s')$  to  $\mathcal{D}$ ,
            replacing the oldest tuple if  $|\mathcal{D}| \geq N_r$ 
        Sample a minibatch of  $N_b$  tuples  $(s, a, r, s') \sim \text{Unif}(\mathcal{D})$ 
        Construct target values, one for each of the  $N_b$  tuples:
        Define  $a'^{\max}(s'; \theta) = \arg \max_{a'} Q(s', a'; \theta)$ 
         $y_j = \begin{cases} r & \text{if } s' \text{ is terminal} \\ r + \gamma Q(s', a'^{\max}(s'; \theta); \theta^-), & \text{otherwise.} \end{cases}$ 
        Do a gradient descent step with loss  $\|y_j - Q(s, a; \theta)\|^2$ 
        Replace target parameters  $\theta^- \leftarrow \theta$  every  $N^-$  steps
    end
end
```

On-Policy vs. Off-Policy Algorithms



- ▶ **On Policy:** In this, the learning agent learns the value function according to the current action derived from the policy currently being used.
- ▶ **Off Policy:** In this, the learning agent learns the value function according to the action derived from another policy.

- ▶ SARSA algorithm is a slight variation of the Q-Learning algorithm.
- ▶ Q-Learning technique is an Off-Policy technique and uses the greedy approach to learn the Q-value. SARSA technique, on the other hand, is an On-Policy and uses the action performed by the current policy to learn the Q-value.

- ▶ SARSA algorithm is a slight variation of the Q-Learning algorithm.
- ▶ Q-Learning technique is an Off-Policy technique and uses the greedy approach to learn the Q-value. SARSA technique, on the other hand, is an On-Policy and uses the action performed by the current policy to learn the Q-value.

Q-Learning:
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r(s, a) + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]$$

- ▶ SARSA algorithm is a slight variation of the Q-Learning algorithm.
- ▶ Q-Learning technique is an Off-Policy technique and uses the greedy approach to learn the Q-value. SARSA technique, on the other hand, is an On-Policy and uses the action performed by the current policy to learn the Q-value.

Q-Learning:
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r(s, a) + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]$$

SARSA:
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r(s, a) + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

- ▶ The update equation for SARSA depends on the current state, current action, reward obtained, next state and next action.
- ▶ This observation lead to the naming of the learning technique as **SARSA** stands for **State Action Reward State Action** which symbolizes the tuple (s, a, r, s', a') .

- ▶ The update equation for SARSA depends on the current state, current action, reward obtained, next state and next action.
- ▶ This observation lead to the naming of the learning technique as **SARSA** stands for **State Action Reward State Action** which symbolizes the tuple (s, a, r, s', a') .
- ▶ Just like DQN, we also have a Deep SARSA.

SARSA Algorithm

Input : States, $s \in S$, Actions $a \in A(s)$, Initialize $Q(s, a)$, α , γ , π to an arbitrary policy (non-greedy)

Output: Optimal action value $Q(s, a)$ for each state-action pair
while *True* do

```
    for ( $i = 0$ ;  $i \leq \# \text{ of episodes}$ ;  $i++$ ) do
        Initialize  $s$ 
        Choose  $a$  from  $s$ , using policy derived from  $Q$ 
        Repeat(for each step of episodes):
            Take action  $a$ ; observe reward,  $r$ , and next state,  $s'$ 
            Choose action  $a'$  from state  $s'$  using policy derived from  $Q$ 
             $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$ 
             $s \leftarrow s'; a \leftarrow a'$ ;
            until  $s$  is terminal
```

end
end

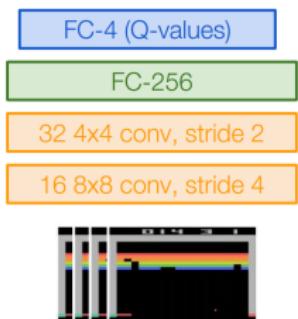
Case Study: Playing Atari Games



- ▶ **Objective:** Complete the game with the highest score
- ▶ **State:** Raw pixel inputs of the game state
- ▶ **Action:** Game controls e.g. Left, Right, Up, Down
- ▶ **Reward:** Score increase/decrease at each time step

Case Study: Playing Atari Games

$Q(s, a; \theta)$:
neural network
with weights θ



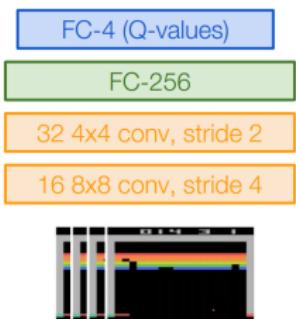
Last FC layer has 4-d output (if 4 actions), corresponding to $Q(s_t, a_1), Q(s_t, a_2), Q(s_t, a_3), Q(s_t, a_4)$

Number of actions between 4-18 depending on Atari game

Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

Case Study: Playing Atari Games

$Q(s, a; \theta)$:
neural network
with weights θ



Last FC layer has 4-d output (if 4 actions), corresponding to $Q(s_t, a_1), Q(s_t, a_2), Q(s_t, a_3), Q(s_t, a_4)$

Number of actions between 4-18 depending on Atari game

Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

<https://www.youtube.com/watch?v=V1eYniJ0Rnk>

- ▶ Markov Decision Process is the mathematical formulation of the Reinforcement Learning Problem defined by $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$
- ▶ Each state satisfies the Markov Property i.e., future is independent of the past given the present.
- ▶ The Agent and the Environment act in a time sequenced loop. Policy π determines how the agent chooses actions
- ▶ Value function approximates how good a state is while Q-value function approximates the goodness of a state action pair.
- ▶ Bellman Equation is a recursive formula for the Q-value function
- ▶ Q-learning is an algorithm that repeatedly adjusts Q to minimize the Bellman error
- ▶ If the Q-value function approximator is Deep Neural Network, we get Deep Q-Learning
- ▶ SARSA is an on-policy variation of Q-learning

These slides have been adapted from

- ▶ Emma Brunskill, Stanford CS234: [Reinforcement Learning](#)
- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: [Deep Learning for Computer Vision](#)
- ▶ Ashwin Rao, Stanford CME241: [Foundations of Reinforcement Learning with Applications in Finance](#)
- ▶ Jimmy Ba & Bo Wang, UofT CSC413/2516: [Neural Networks and Deep Learning](#)