



Work Integrated Learning Programmes

Reinforcement Learning

Reinforcement learning (RL) is based on rewarding desired behaviors or punishing undesired ones. Instead of one input producing one output, the algorithm produces a variety of outputs and is trained to select the right one based on certain variables – Gartner

When to use RL?

RL can be used in large environments in the following situations:

- 1.A model of the environment is known, but an analytic solution is not available;
- 2.Only a simulation model of the environment is given (the subject of simulation-based optimization)
- 3.The only way to collect information about the environment is to interact with it.

3

BITSPilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

(Deep) Reinforcement Learning

Paradigm	Supervised Learning	Unsupervised Learning	Reinforcement Learning
Objective	$p_{\theta}(y x)$	$p_{\theta}(x)$	$\pi_{\theta}(a s)$
Applications	→ Classification → Regression	→ Inference → Generation	→ Prediction → Control

4

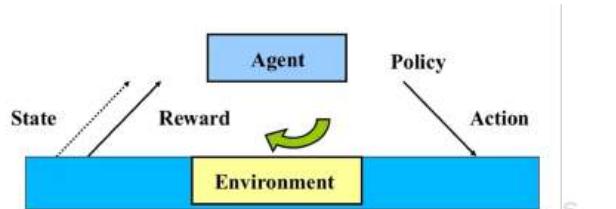
BITSPilani, Deemed to be University under Section 3 of UGC Act, 1956

Criteria	Definition
Type of data	Learning in 3D
Type of problem	Classification
Supervision	Supervised Learning
Algorithms	Decision Tree
Aim	W
Application	3D Reconstruction

Criteria	Supervised ML	Unsupervised ML	Reinforcement ML
Definition	Learns by using labelled data	Trained using unlabelled data without any guidance.	Works on interacting with the environment
Type of data	Labelled data	Unlabelled data	No – predefined data
Type of problems	Regression and classification	Association and Clustering	Exploitation or Exploration
Supervision	Extra supervision	No supervision	No supervision
Algorithms	Linear Regression, Logistic Regression, SVM, KNN etc.	K – Means, C – Means, Apriori	Q – Learning, SARSA
Aim	Calculate outcomes	Discover underlying patterns	Learn a series of action
Application	Risk Evaluation, Forecast Sales	Recommendation System, Anomaly Detection	Self Driving Cars, Gaming, Healthcare

5

Elements of Reinforcement Learning



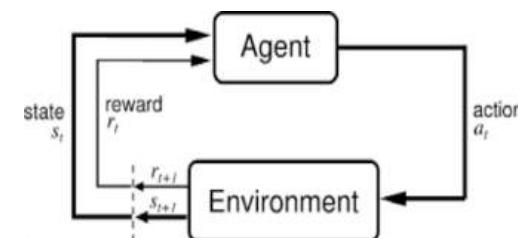
1

Characteristics of RL

- No supervision, only a real value or reward signal
 - Decision making is sequential
 - Time plays a major role in reinforcement problems
 - Feedback isn't prompt but delayed



Elements of Reinforcement Learning



Beyond the agent and the environment, one can identify four main sub-elements of a reinforcement learning system: *a policy*, *a reward*, *a value function*, and, optionally, *a model* of the environment.

Elements of Reinforcement Learning

•Agent

- An **entity** that tries to learn the best way to perform a specific task.
- In our example, the child is the agent who learns to ride a bicycle.

•Action (A) -

- **What the agent does** at each time step.
- In the example of a child learning to walk, the action would be “walking”.
- A is the set of all possible moves.
- In video games, the list might include running right or left, jumping high or low, crouching or standing still.

9

Elements of Reinforcement Learning

•Environment

- Outside world of an agent or physical world in which the agent operates.

Formal Definition - Reinforcement learning (RL) is an area of machine learning concerned with how intelligent agents ought to take actions in an environment in order to maximize the notion of cumulative reward.

11

Elements of Reinforcement Learning

•State (S)

- **Current situation** of the agent.
- After doing performing an action, the agent can move to different states.
- In the example of a child learning to walk, the child can take the action of taking a step and move to the next state (position).

•Rewards (R)

- Feedback that is given to the agent based on the action of the agent.
- If the action of the agent is good and can lead to winning or a positive side then a positive reward is given and vice versa.

10

Tic-Tac-Toe



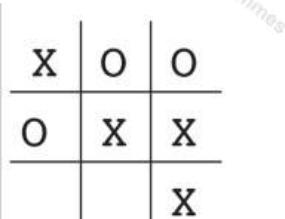
X	O	O
O	X	X
		X



Tic-Tac-Toe

States	Initial Values
$\begin{bmatrix} X \\ \end{bmatrix}$	0.5
$\begin{bmatrix} X & O & O \\ & X & \end{bmatrix}$	0.5
$\begin{bmatrix} X & O & O \\ & X & X \\ & & X \end{bmatrix}$	1.0
$\begin{bmatrix} X & O \\ X & O \\ & X O \end{bmatrix}$	0
...	...

Learning Task: Play as many times against the opponent, and learn the values



Set up a table of states initial values

Tic-Tac-Toe

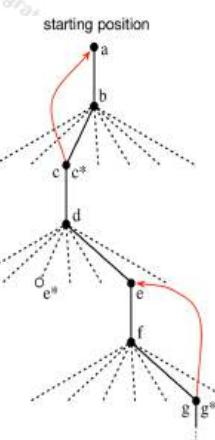
States	Initial Values
$\begin{bmatrix} X \\ \end{bmatrix}$	0.5
$\begin{bmatrix} X & O & O \\ & X & \end{bmatrix}$	0.5
$\begin{bmatrix} X & O & O \\ & X & X \\ & & X \end{bmatrix}$	1.0
$\begin{bmatrix} X & O \\ X & O \\ & X O \end{bmatrix}$	0
...	...

S_t - state before greedy move
 S_{t+1} - state after greedy move

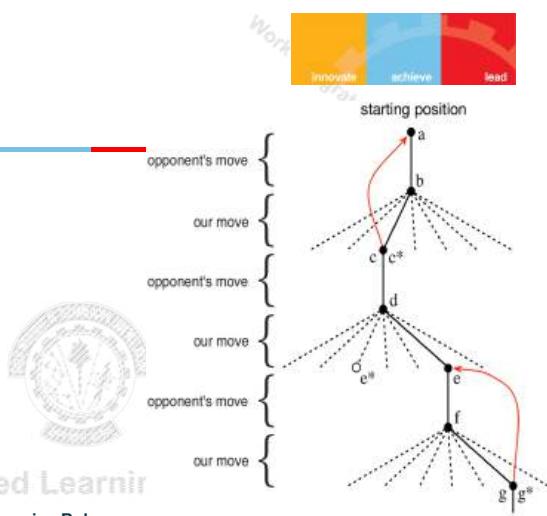
$$V(S_t) \leftarrow V(S_t) + \alpha [V(S_{t+1}) - V(S_t)]$$



starting position



Tic-Tac-Toe



Temporal Difference Learning Rule

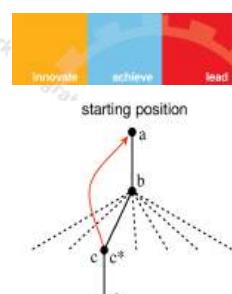
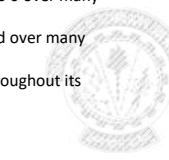
$$V(S_t) \leftarrow V(S_t) + \alpha [V(S_{t+1}) - V(S_t)]$$

α - Step Size Parameter

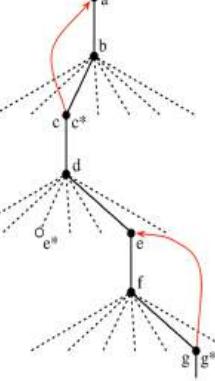
Tic-Tac-Toe

Questions:

- (1) What happens if α is gradually made to 0 over many games with the opponent?
- (2) What happens if α is gradually reduced over many games, but never made 0?
- (3) What happens if α is kept constant throughout its life time?



starting position



Temporal Difference Learning Rule

$$V(S_t) \leftarrow V(S_t) + \alpha [V(S_{t+1}) - V(S_t)]$$

α - Step Size Parameter

Work Integrated Learning Programmes

Tic-Tac-Toe

Key Takeaways:

- (1) Learning while interacting with the environment (opponent).
- (2) We have a clear goal
- (3) Our policy is to make moves that maximizes our chances of reaching goal
 - Use the values of states most of the time (exploration) and explore rest of the time.

Work Integrated Learni

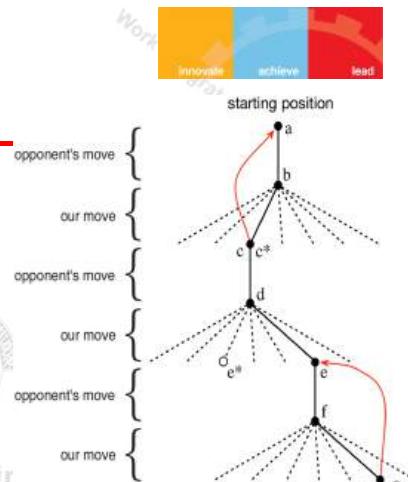
Temporal Difference Learning Rule

$$V(S_t) \leftarrow V(S_t) + \alpha [V(S_{t+1}) - V(S_t)]$$

α - Step Size
Parameter

17

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

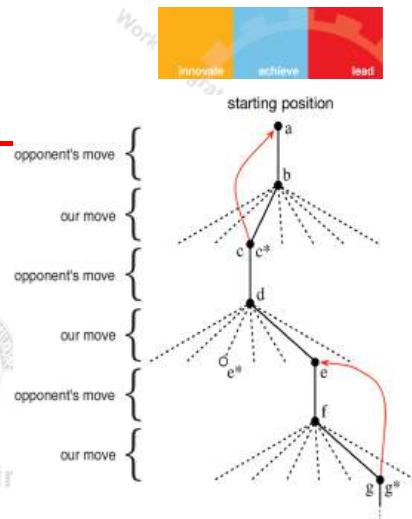


Work Integrated Learning Programmes

Tic-Tac-Toe

Reading Assigned:

Identify how this reinforcement learning solution is different from solutions using minimax algorithm and genetic algorithms.
Post your answers in the discussion forum;



Work Integrated Learni

Temporal Difference Learning Rule

$$V(S_t) \leftarrow V(S_t) + \alpha [V(S_{t+1}) - V(S_t)]$$

α - Step Size
Parameter

18

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

K-armed Bandit Problem

Problem

- You are faced repeatedly with a choice among k different options, or actions
- After each choice of actions you receive a numerical reward
 - Reward is chosen from a stationary probability distribution that depends on the selected action
- **Objective** : to maximize the expected total reward over some time period

Work Integrated Learning Programmes

19

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

K-armed Bandit Problem

Problem

- You are faced repeatedly with a choice among k different options, or actions
- After each choice of actions you receive a numerical reward
 - Reward is chosen from a stationary probability distribution that depends on the selected action
- **Objective** : to maximize the expected total reward over some time period



20

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

K-armed Bandit Problem

Problem

- You are faced repeatedly with a choice among k different options, or actions
- After each choice of actions you receive a numerical reward
 - Reward is chosen from a stationary probability distribution that depends on the selected action
- **Objective** : to *maximize the expected total reward over some time period*



Strategy:

- Identify the best lever(s)
- Keep pulling the identified ones

Questions:

- How do we define the **best ones**?
- What are the best levers?

21

K-armed Bandit Problem

Problem

- You are faced repeatedly with a choice among k different options, or actions
- After each choice of actions you receive a numerical reward
 - Reward is chosen from a stationary probability distribution that depends on the selected action
- **Objective** : to *maximize the expected total reward over some time period*



Strategy:

- Identify the best lever(s)
- Keep pulling the identified ones

Questions:

- How do we define the **best ones**?
- What are the best levers?

22

K-armed Bandit Problem

Problem

- You are faced repeatedly with a choice among k different options, or actions
- After each choice of actions you receive a numerical reward
 - Reward is chosen from a stationary probability distribution that depends on the selected action
- **Objective** : to *maximize the expected total reward over some time period*



- **Expected Mean Reward** for each action selected
→ call it **Value** of the action

$$q_*(a) \doteq \mathbb{E}[R_t | A_t = a]$$

23

$$q_*(a) \doteq \mathbb{E}[R_t | A_t = a]$$

- A_t
- $Q_t(a)$
- $q_*(a)$

- action selected on time step t
- estimated value of action a at time step t
- value of an arbitrary action a

Work Integrated Learning Programmes

Note: If you knew the value of each action, then it would be trivial to solve the k -armed bandit problem: you would always select the action with highest value :-)

24

K-armed Bandit Problem



$$\mathbb{E}[a] = -\$0.5$$



$$\mathbb{E}[b] = -\$0.2$$



$$\mathbb{E}[c] = \$0.1$$



$$\mathbb{E}[d] = \$0.11$$

K-armed Bandit Problem



$$-1, -1, 5
\\ \mathbb{E}[a] = 1$$



$$-0.2, -0.2
\\ \mathbb{E}[b] = -0.2$$



$$-0.5, -0.5, -0.5
\\ \mathbb{E}[c] = -0.5$$



$$-2, -2
\\ \mathbb{E}[d] = -2$$

K-armed Bandit Problem



$$-1, -1, 5
\\ \mathbb{E}[a] = 1$$



$$-0.2, -0.2
\\ \mathbb{E}[b] = -0.2$$



$$-0.5, -0.5, -0.5
\\ \mathbb{E}[c] = -0.5$$



$$-2, -2
\\ \mathbb{E}[d] = -2$$

Keep pulling the levers; update the estimate of action values;

K-armed Bandit Problem

1. How to maintain the estimate of expected rewards for each action?

- Average the rewards actually received !!!

$$Q_t(a) \doteq \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t}$$

$$= \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

1. How to use the estimate in selecting the right action?

Greedy Action Selection

$$A_t \doteq \arg \max_a Q_t(a)$$

K-armed Bandit Problem

2. How to use the estimate in selecting the right action?

Greedy Action Selection

$$A_t \doteq \arg \max_a Q_t(a)$$

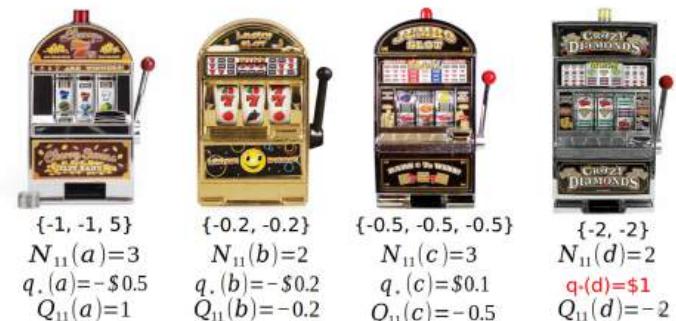
Actions which are inferior by the value estimate upto time t, could be indeed better than the greedy action at t !!!

3. Exploration vs. Exploitation?

ϵ -Greedy Action Selection / near-greedy action selection

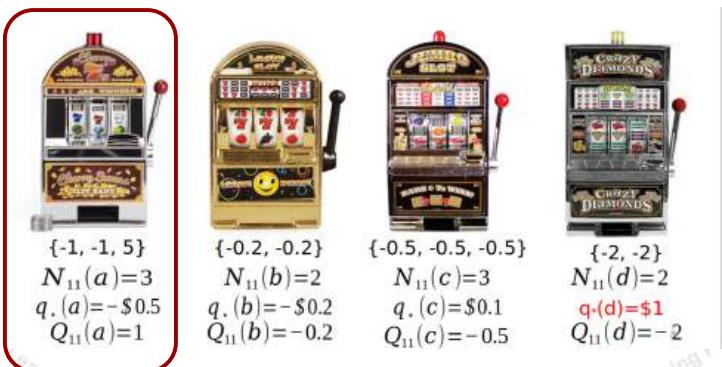
Behave greedily most of the time; Once in a while, with small probability ϵ select randomly from among all the actions with equal probability, independently of the action-value estimates.

K-armed Bandit Problem



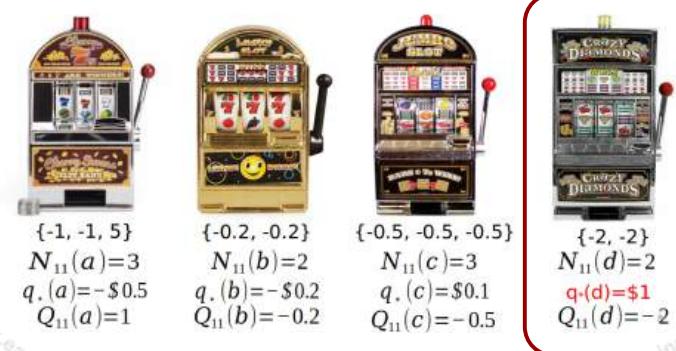
K-armed Bandit Problem

Greedy Action



K-armed Bandit Problem

Action to Explore



K-armed Bandit Problem

ϵ -Greedy Action Selection / near-greedy action selection

```
epsilon = 0.05 // small value to control exploration
def get_action():
    if random.random() > epsilon:
        return argmax(Q(a))
    else:
        return random.choice(A)
```

- In the limit as the number of steps increases, every action will be sampled by ϵ -greedy action selection an infinite number of times. This ensures that all the $Q_t(a)$ converge to $q^*(a)$.
- Easy to implement / optimize for epsilon / yields good results

33



Work Integrated Learning Programmes

34

Ex-1: In ϵ -greedy action selection, for the case of two actions and $\epsilon = 0.5$, what is the probability that **the greedy action** is selected?

$$\begin{aligned}
 p(\text{greedy action}) &= p(\text{greedy action AND greedy selection}) + p(\text{greedy action AND random selection}) \\
 &= p(\text{greedy action} | \text{greedy selection}) p(\text{greedy selection}) \\
 &\quad + p(\text{greedy action} | \text{random selection}) p(\text{random selection}) \\
 &= p(\text{greedy action} | \text{greedy selection})(1-\epsilon) + p(\text{greedy action} | \text{random selection})(\epsilon) \\
 &= p(\text{greedy action} | \text{greedy selection})(0.5) + p(\text{greedy action} | \text{random selection})(0.5) \\
 &= (1)(0.5) + (0.5)(0.5) \\
 &= 0.5 + 0.25 \\
 &= 0.75
 \end{aligned}$$

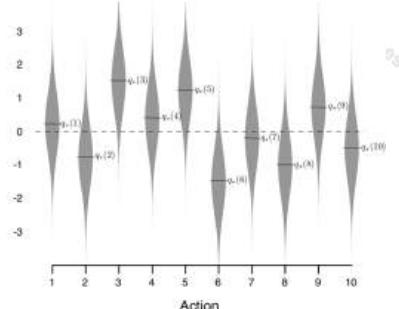
35

Ex-1: In ϵ -greedy action selection, for the case of two actions and $\epsilon = 0.5$, what is the probability that **the greedy action** is selected?

10-armed Testbed

Example:

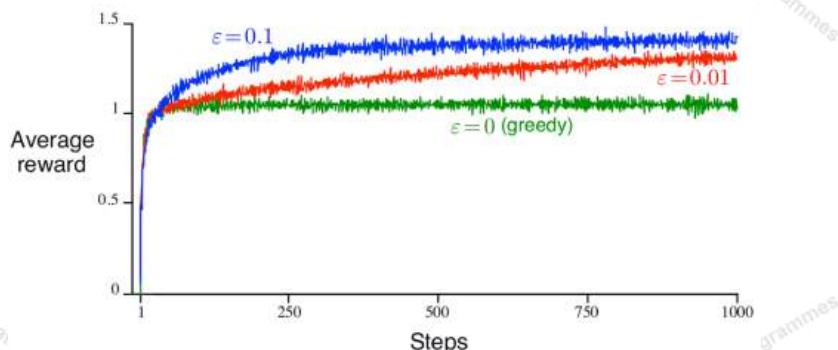
- A set of 2000 randomly generated k - armed bandit problems with $k = 10$
- Action values were selected according to a normal (Gaussian) distribution with mean 0 and variance 1.
- While selecting action A_t at time step t, the actual reward, R_t , was selected from a normal distribution with mean $q_t(A_t)$ and variance 1
- **One Run :** Apply a method for 1000 time steps to one of the bandit problems
- Perform 2000 runs, each run with a different bandit problem, to get an algorithms average behavior



An example bandit problem from the 10-armed testbed

36

Average performance of ϵ -greedy action-value methods on the 10-armed testbed



37

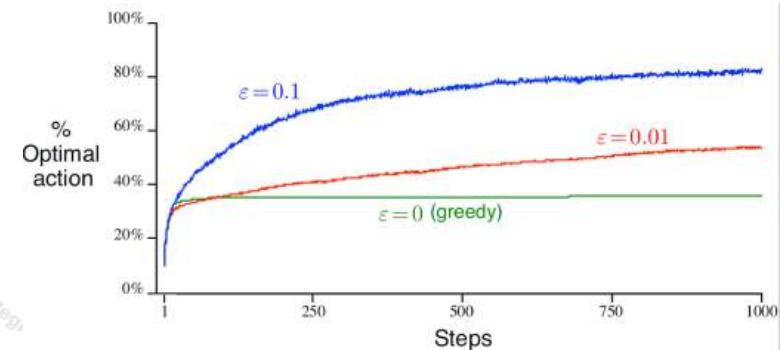
Discussion on Exploration vs. Exploitation

- 1) What if the reward variance is
 - a. larger, say 10 instead of 1?
 - b. zero? [deterministic]
- 2) What if the bandit task is non-stationary? [that is, the true values of the actions changed over time]



39

Average performance of ϵ -greedy action-value methods on the 10-armed testbed



38

Ex-2:

Consider a k-armed bandit problem with $k = 4$ actions, denoted 1, 2, 3, and 4.

Consider applying to this problem a bandit algorithm using ϵ -greedy action selection, sample-average action-value estimates, and initial estimates of $Q_1(a) = 0$, for all a .

Suppose the initial sequence of actions and rewards is $A_1 = 1, R_1 = 1, A_2 = 2, R_2 = 1, A_3 = 2, R_3 = 2, A_4 = 2, R_4 = 2, A_5 = 3, R_5 = 0$.

On some of these time steps the ϵ case may have occurred, causing an action to be selected at random.

On which time steps did this definitely occur? On which time steps could this possibly have occurred?

40

Incremental Implementation

- Efficient approach to compute the estimate of action-value;

$$Q_n = \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$$

- Given Q_n and the nth reward, R_n , the new average of all n rewards can be computed as follows

$$\begin{aligned} Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\ &= \frac{1}{n} \left(R_n + \sum_{i=1}^{n-1} R_i \right) \\ &= \frac{1}{n} \left(R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\ &= \frac{1}{n} (R_n + (n-1)Q_n) \\ &= \frac{1}{n} (R_n + nQ_n - Q_n) \\ &= Q_n + \frac{1}{n} [R_n - Q_n], \end{aligned}$$



41

Bandit Algorithm with Incremental Update/ ϵ -greedy selection

Initialize, for $a = 1$ to k :

$$\begin{aligned} Q(a) &\leftarrow 0 \\ N(a) &\leftarrow 0 \end{aligned}$$

Loop forever:

$$\begin{aligned} A &\leftarrow \begin{cases} \arg \max_a Q(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases} \quad (\text{breaking ties randomly}) \\ R &\leftarrow \text{bandit}(A) \\ N(A) &\leftarrow N(A) + 1 \\ Q(A) &\leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)] \end{aligned}$$

43

Incremental Implementation

Note:

- StepSize decreases with each update
- We use α or $\alpha_t(a)$ to denote step size (constant / varies with each step)

Discussion:

Const vs. Variable step size?



Work Integrated Learning Programmes

$$\begin{aligned} Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\ &= \frac{1}{n} \left(R_n + \sum_{i=1}^{n-1} R_i \right) \\ &= \frac{1}{n} \left(R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\ &= \frac{1}{n} (R_n + (n-1)Q_n) \\ &= \frac{1}{n} (R_n + nQ_n - Q_n) \\ &= Q_n + \frac{1}{n} [R_n - Q_n], \end{aligned}$$

$$\text{NewEstimate} \leftarrow \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}]$$

42

What are problems whose solutions are modelled as MAB?



Work Integrated Learning Programmes

44

Non-stationary Problem

- Most RL problems are non-stationary !
- Give more weight to recent rewards than to long-

$$Q_{n+1} \doteq Q_n + \alpha [R_n - Q_n]$$

Work Integrated Learning Programmes

Non-stationary Problem

- Most RL problems are non-stationary !
- Give more weight to recent rewards than to long-

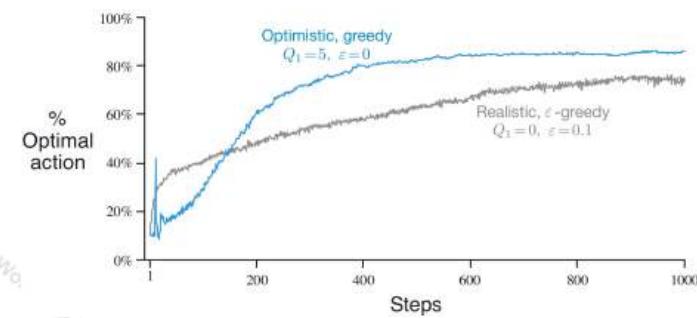
$$\begin{aligned} Q_{n+1} &= Q_n + \alpha [R_n - Q_n] \\ &= \alpha R_n + (1 - \alpha) Q_n \\ &= \alpha R_n + (1 - \alpha) [\alpha R_{n-1} + (1 - \alpha) Q_{n-1}] \\ &= \alpha R_n + (1 - \alpha) \alpha R_{n-1} + (1 - \alpha)^2 Q_{n-1} \\ &= \alpha R_n + (1 - \alpha) \alpha R_{n-1} + (1 - \alpha)^2 \alpha R_{n-2} + \\ &\quad \dots + (1 - \alpha)^{n-1} \alpha R_1 + (1 - \alpha)^n Q_1 \\ &= (1 - \alpha)^n Q_1 + \sum_{i=1}^n \alpha(1 - \alpha)^{n-i} R_i. \end{aligned}$$

Exponential recency-weighted average

Optimistic Initial Values

- All the above discussed methods are **biased** by their initial estimates
- For sample average method the bias disappears once all actions have been selected at least once
- For methods with constant α , the bias is permanent, though decreasing over time
- Initial action values can also be used as a simple way of encouraging exploration.
- In 10 armed testbed, set initial estimate to +5 rather than 0.
- This can encourage action-value methods to explore.
 - Whichever actions are initially selected, the reward is less than the starting estimates;
 - the learner switches to other actions, being disappointed with the rewards it is receiving.
 - The result is that all actions are tried several times before the value estimates converge.

Optimistic Initial Values



The effect of optimistic initial action-value estimates on the 10-armed testbed.
Both methods used a constant step-size parameter, $\alpha = 0.1$.

Caution:
Optimistic Initial Values can only be considered as a simple trick that can be quite effective on stationary problems, but it is far from being a generally useful approach to encouraging exploration.

Question:
Explain how in the non-stationary scenario the optimistic initial values will fail (to explore adequately).

Upper-Confidence-Bound Action Selection

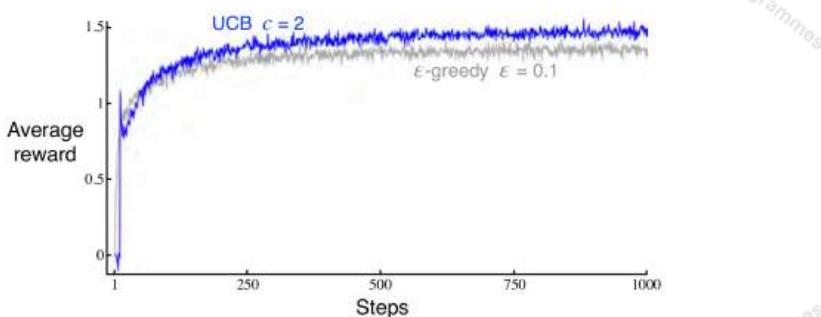
- ϵ -greedy action selection forces the non-greedy actions to be tried,
 - Indiscriminately, with no preference for those that are nearly greedy or particularly uncertain
- It would be better to select among the non-greedy actions according to their potential for actually being optimal
 - Take it being $A_t \doteq \arg \max_a \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$

Upper-Confidence-Bound Action Selection

- Each time a is selected the uncertainty is presumably reduced
- Each time an action other than a is selected, t increases but $N_t(a)$ does not; because t appears in the numerator, the uncertainty estimate increases.
- Actions with lower value estimates, or that have already been selected frequently, will be selected with decreasing frequency over time

$$A_t \doteq \arg \max_a \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

Upper-Confidence-Bound Action Selection



Policy-based algorithms

- Forget about action-value (Q) estimates, we don't really care about them
- We care about what actions to chose
 - Let's assign a preference to each action and tweak its value
- Define $H_t(a)$ as a numerical preference value associated with action a
- Which action is selected?
 - $A_t = \arg \max_a [H_t(a)]$
 - Hardmax results in no exploration -- deterministic action selection

Simple!



Softmax function

- Input: vector of preferences
- Output: vector of probabilities forming a valid distribution

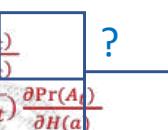
$$\Pr(a_t = a) = \frac{e^{H_t(a)}}{\sum_{a' \in A} e^{H_t(a')}} = \Pr(a)$$

$$H_t \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} = \begin{bmatrix} 6 \\ 9 \\ 2 \end{bmatrix}$$

$$\text{softmax} \begin{pmatrix} 6 \\ 9 \\ 2 \end{pmatrix} = \begin{bmatrix} 0.047 \\ 0.952 \\ 0.00087 \end{bmatrix}$$

That is, with $\Pr(0.95)$ choose a_2 , $\Pr(0.05)$ choose a_1 , and $< \Pr(0.01)$ choose a_3

Gradient ascend

- $H_{t+1}(A_t) = H_t(A_t) + \alpha(R_t - \bar{R}_t) \frac{\partial \Pr(A_t)}{\partial H(A)}$ 
- $\forall a \neq A_t, H_{t+1}(a) = H_t(a) + \alpha(R_t - \bar{R}_t) \frac{\partial \Pr(A_t)}{\partial H(a)}$ 
- Update the preferences based on observed reward and a baseline reward (\bar{R}_t)
- If the observed reward is larger than the baseline:
 - Increase the preference of the chosen action, A_t
 - Decrease the preference of all other actions, $\forall a \neq A_t$
- Else do the opposite

Softmax function

- Exploration – checked!
- Softmax provides another important attribute – a **differentiable policy**
- Say that we learn that action a_1 results in good relative performance
- Hardmax: $A_t = \underset{a}{\operatorname{argmax}}[H_t(a_1), H_t(a_2), H_t(a_3)]$
 - Change $H(a_1)$ such that $\Pr(a_1)$ is increased, $\frac{\partial \Pr(a_1)}{\partial H(a_1)} = NA$
- Softmax: $\Pr(a_1) = \frac{e^{H_t(a_1)}}{\sum_{a' \in A} e^{H_t(a')}}$
 - Change $H(a_1)$ such that $\Pr(a_1)$ is increased -> update towards $\frac{\partial \Pr(a_1)}{\partial H(a_1)}$

Update Rule

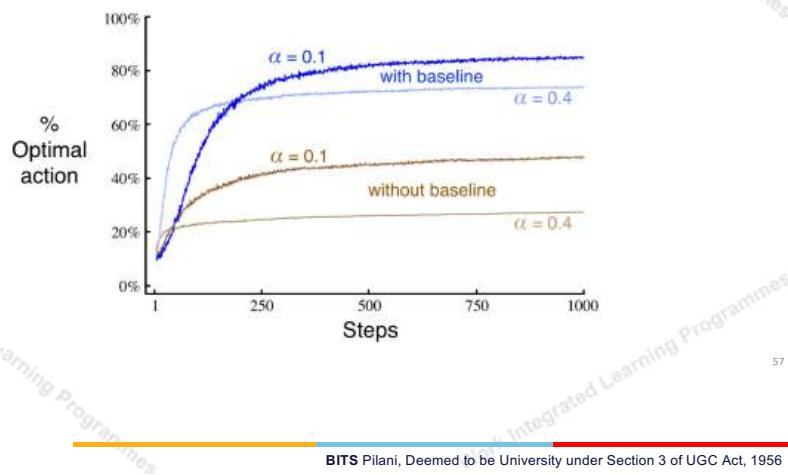
On each step, after selecting action A_t and receiving the reward R_t ,
Update the action preferences :

$$H_{t+1}(A_t) = H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \Pr(A_t))$$

$$\forall a \neq A_t, H_{t+1}(a) = H_t(a) - \alpha(R_t - \bar{R}_t)\Pr(a)$$

$$H_{t+1}(A_t) \doteq H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(A_t)), \quad \text{and}$$

$$H_{t+1}(a) \doteq H_t(a) - \alpha(R_t - \bar{R}_t)\pi_t(a), \quad \text{for all } a \neq A_t$$

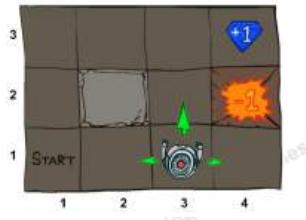


57

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

A different scenario

- Associative vs. Non-associative tasks ?
- **Policy:** A mapping from situations to the actions that are best in those situations
- (discuss) How do we extend the solution for non-associative task to an associative task?
 - **Approach:** Extend the solutions to non-stationary task to non-associative tasks
 - Works, if the true action values changes slowly
 - What if the context switching between the situations are made explicit?
 - How?
 - Need Special approaches !!!



59

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

What did we learn?

- **Problem:** choose the action that results in highest expected reward
- **Assumptions:** 1. actions' expected reward is unknown, 2. we are confronted with the same problem over and over, 3. we are able to observe an action's outcome once chosen
- **Approach:** learn the actions' expected reward through exploration (value based) or learn a policy directly (policy based), exploit learnt knowledge to choose best action
- **Methods:** 1. greedy + initializing estimates optimistically, 2. epsilon-greedy, 3. Upper-Confidence-Bounds, 4. gradient ascend + soft-max

58

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



BITS Pilani

Deep Reinforcement Learning

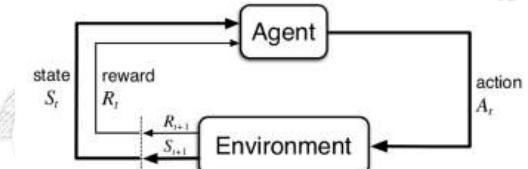
Prof. Vimal SP
CSIS



BITS Pilani

<AIMLCZG512, Deep Reinforcement Learning>

Lecture No. 3,4,5



- **Agent** - Learner & the decision maker
- **Environment** - Everything outside the agent
- **Interaction:**
 - Agent performs an action
 - Environment responds by
 - presenting a new situation (change in state)
 - presents numerical reward
- **Objective (of the interaction):**
 - Maximize the return (cumulative rewards) over time

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Grid World Example

- A maze-like problem
 - The agent lives in a grid
 - Walls block the agent's path
- Noisy movement: actions do not always go as planned
 - 80% of the time, the action North takes the agent North (if there is no wall there)
 - 10% of the time, North takes the agent West; 10% East
 - If there is a wall in the direction the agent would have been taken, the agent stays put
- The agent receives rewards each time step
 - -0.1 per step (battery loss)
 - +1 if arriving at (4,3); -1 for arriving at (4,2); 1 for arriving at (2,2)
- Goal: maximize accumulated rewards



Markov Decision Processes

- An MDP is defined by
 - A set of **states**
 - A set of **actions**
 - **State-transition probabilities**
 - Probability of arriving to s' after performing a
 - Also called the **model dynamics**
 - A **reward function**
 - The utility gained from arriving to s' after performing a
 - Sometimes just r or even $r(s, a)$
 - A **start state**
 - **Maybe a terminal state**



Work Integrated Learning Programmes

Markov Decision Processes

Model Dynamics

$$p(s', r | s, a) \doteq \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$$

State-transition probabilities

$$p(s' | s, a) \doteq \Pr\{S_t = s' \mid S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in \mathcal{R}} p(s', r | s, a)$$

Expected rewards for state-action-next-state triples

$$r(s, a, s') \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in \mathcal{R}} r \frac{p(s', r | s, a)}{p(s' | s, a)}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Markov Decision Processes – Discussion

- *MDP framework is a considerable abstraction of the problem of goal-directed learning from interaction.*
- *It proposes that whatever the details of the sensory, memory, and control apparatus, and whatever objective one is trying to achieve, any problem of learning goal-directed behavior can be reduced to three signals passing back and forth between an agent and its environment:*
 - *one signal to represent the choices made by the agent (the actions)*
 - *one signal to represent the basis on which the choices are made (the states),*
 - *and one signal to define the agent's goal (the rewards).*

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

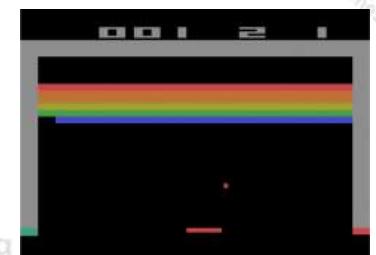
Markov Decision Processes – Discussion

- *MDP framework is abstract and flexible*
 - *Time steps* need not refer to fixed intervals of real time
 - *The actions* can be
 - at low-level controls or high-level decisions
 - totally mental or computational
 - *States* can take a wide variety of forms
 - Determined by *low-level sensations* or *high-level and abstract* (ex. symbolic descriptions of objects in a room)
- *The agent-environment boundary represents the limit of the agent's absolute control, not of its knowledge.*
 - *The boundary can be located at different places for different purposes*

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MDP Formalization : Video Games

- ***State:***
 - raw pixels
- ***Actions:***
 - game controls
- ***Reward:***
 - change in score
- ***State-transition probabilities:***
 - defined by stochasticity in game evolution



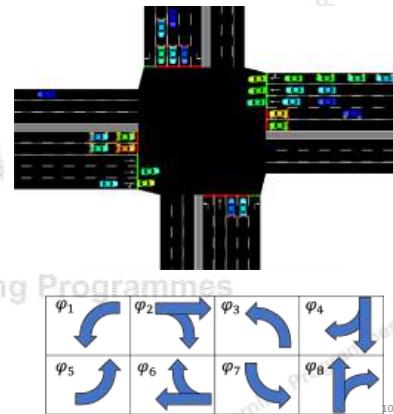
Ref: [Playing Atari with deep reinforcement learning](#), Mnih et al., 2013

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MDP Formalization : Traffic Signal Control

- State:**
 - Current signal assignment (green, yellow, and red assignment for each phase)
 - For each lane: number of approaching vehicles, accumulated waiting time, number of stopped vehicles, and average speed of approaching vehicles
- Actions:**
 - signal assignment
- Reward:**
 - Reduction in traffic delay
- State-transition probabilities:**
 - defined by stochasticity in approaching demand

Ref: "Learning an Interpretable Traffic Signal Control Policy", Ault et al., 2020



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MDP Formalization : Recycling Robot (Detailed Ex.)

- State:**
 - Assume that only two charge levels can be distinguished
 - $S = \{\text{high, low}\}$
- Actions:**
 - $A(\text{high}) = \{\text{search, wait}\}$
 - $A(\text{low}) = \{\text{search, wait, recharge}\}$
- Reward:**
 - Zero most of the time, except when securing a can
 - Cans are secured by searching and waiting, but $r_{\text{search}} > r_{\text{wait}}$
- State-transition probabilities:**
 - [Next Slide]



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MDP Formalization : Recycling Robot (Detailed Ex.)

- Robot has**
 - sensors for detecting cans
 - arm and gripper that can pick the cans and place in an onboard bin;
- Runs on a rechargeable battery**
- Its control system has components for interpreting sensory information, for navigating, and for controlling the arm and gripper**
- Task for the RL Agent: Make high-level decisions about how to search for cans based on the current charge level of the battery**



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MDP Formalization : Recycling Robot (Detailed Ex.)

- State-transition probabilities (contd...):**

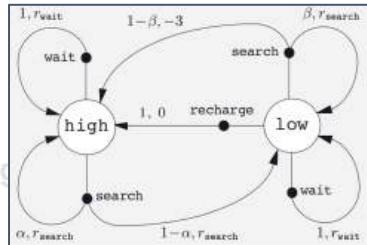
s	a	s'	$p(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
high	search	low	$1 - \alpha$	r_{search}
low	search	high	$1 - \beta$	-3
low	search	low	β	r_{search}
high	wait	high	1	r_{wait}
high	wait	low	0	-
low	wait	high	0	-
low	wait	low	1	r_{wait}
low	recharge	high	1	0
low	recharge	low	0	-

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MDP Formalization : Recycling Robot (Detailed Ex.)

- State-transition probabilities (contd...):

s	a	s'	$p(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
high	search	low	$1 - \alpha$	$r_{\text{search}} - 3$
low	search	high	$1 - \beta$	r_{search}
low	search	low	β	r_{search}
high	wait	high	1	r_{wait}
high	wait	low	0	-
low	wait	high	0	-
low	wait	low	1	r_{wait}
low	recharge	high	1	0
low	recharge	low	0	-



14

Note on Goals & Rewards

- Reward Hypothesis:

All of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward).

- The rewards we set up truly indicate what we want accomplished,
 - not the place to impart prior knowledge on how we want it to do
- Ex: Chess Playing Agent
 - If the agent is rewarded for taking opponents pieces, the agent might fall for the opponent's trap.
- Ex: Vacuum Cleaner Agent
 - If the agent is rewarded for each unit of dirt it sucks, it can repeatedly deposit and suck the dirt for larger reward

15

Returns & Episodes

- Goal is to maximize the expected return
- Return (G_t) is defined as some specific function of the reward sequence
- Episodic tasks vs. Continuing tasks
- When there is a notion of final time step, say T , return can be

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

- Applicable when agent-environment interaction breaks into episodes
- Ex: Playing Game, Trips through maze etc. [called episodic tasks]

16

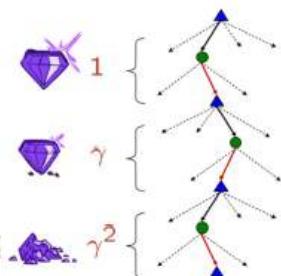
Returns & Episodes

- Generally $T = \infty$

- What if the agent receive a reward of +1 for each timestep?
- Discounted Return:

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Note: γ is a parameter, $0 \leq \gamma \leq 1$, called the discount rate



- Discount rate determines the present value of future rewards

17

Returns & Episodes

- What if γ is 0?
- What if γ is 1?
- Computing discounted rewards incrementally

$$\begin{aligned} G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

- Sum of an infinite number of terms, it is still finite if the reward is nonzero and constant and if $\gamma < 1$.
- Ex: reward is +1 constant

$$G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}.$$

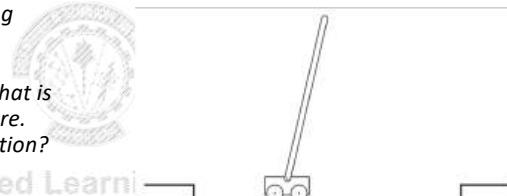
18

Returns & Episodes

→ **Objective:** To apply forces to a cart moving along a track so as to keep a pole hinged to the cart from falling over

→ **Discuss:**

- Consider the task as episodic, that is try/maintain balance until failure. What could be the reward function?
- Repeat prev. assuming task is continuous.



19

Policy

- A mapping from states to probabilities of selecting each possible action.
 - $\pi(a|s)$ is the probability that $A_t = a$ if $S_t = s$
- The purpose of learning is to improve the agent's policy with its experience



20

Defining Value Functions

State-value function for policy π

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right], \text{ for all } s \in \mathcal{S}.$$

Action-value function for policy π

$$q_\pi(s, a) \doteq \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right].$$

21

Defining Value Functions

State Value function in terms of Action-value function for policy π

$$v_\pi(s) = \sum_{a \in \mathcal{A}(s)} \pi(a|s) q_\pi(s, a)$$

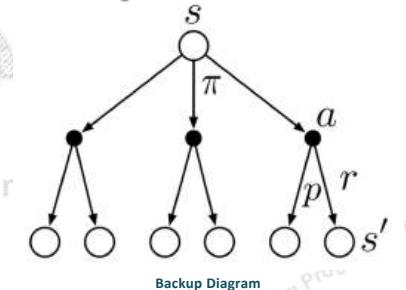
Action Value function in terms of State value function for policy π

$$q_\pi(s, a) = \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r | s, a) [r + \gamma v_\pi(s')]$$

22

Bellman Equation for V_π

$$v_\pi(s) \doteq \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]$$



Backup Diagram

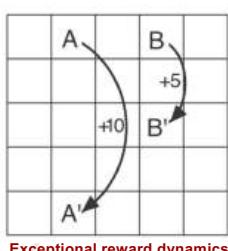
23

Value of the start state must equal
 (1) the (discounted) value of the expected next state,
 plus
 (1) the reward expected along the way

Understanding $V_\pi(s)$ with Gridworld

Reward:

- -1 if an action takes agent off the grid
- Exceptional reward from A and B for all actions taking agent to A' and B' resp.
- 0, everywhere else



Exceptional reward dynamics

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

State-value function for the equiprobable random policy with $\gamma = 0.9$

24

Understanding $V_\pi(s)$ with Gridworld

$$v_\pi(s) \doteq \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]$$

Verify $V_\pi(s)$ using Bellman equation for this state with $\gamma = 0.9$, and equiprobable random policy

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

25

Understanding $V_\pi(s)$ with Gridworld

$$v_\pi(s) \doteq \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_\pi(s')]$$

$$\begin{aligned} v_\pi(s) &= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_\pi(s')] \\ &= \sum_a 0.25 \cdot [0 + 0.9 \cdot (2.3 + 0.4 - 0.4 + 0.7)] \\ &= 0.25 \cdot [0.9 \cdot 3.0] = 0.675 \approx 0.7 \end{aligned}$$

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

Optimal Policies and Optimal Value Functions

- $\pi \geq \pi'$ if and only if $v_\pi(s) \geq v_{\pi'}(s)$ for all $s \in S$
- There is always at least one policy that is better than or equal to all other policies \rightarrow optimal policy (denoted as π_*)
 - There could be more than one optimal policy !!!

Optimal state-value function $v_*(s) \doteq \max_\pi v_\pi(s)$.

Optimal action-value function $q_*(s, a) \doteq \max q_\pi(s, a)$

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

Ex-1

Recollect the reward function used for Gridworld as below:

- -1 if an action takes agent off the grid
- Exceptional reward from A and B for all actions taking agent to A' and B' resp.
- 0, everywhere else

Let us add a constant c (say 10) to the rewards of all the actions. Will it change anything?

Optimal Policies and Optimal Value Functions

Bellman optimality equation - expresses that the value of a state under an optimal policy must equal the expected value for the best action from that state

$$\begin{aligned} v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\ &= \max_a \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a] \\ &= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\ &= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v_*(s')]. \end{aligned}$$

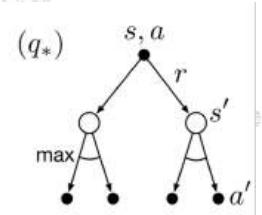
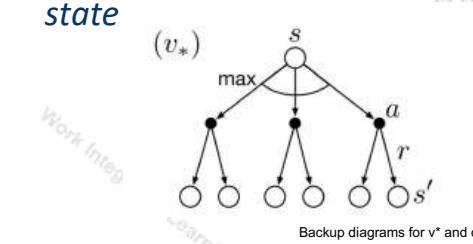
Bellman optimality equation - expresses that *the value of a state under an optimal policy must equal the expected return for the best action from that state*

Bellman optimality equation for q_*

$$\begin{aligned} q_*(s, a) &= \mathbb{E}\left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a\right] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]. \end{aligned}$$

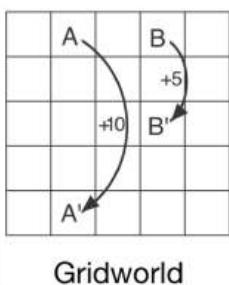
30

Bellman optimality equation - expresses that *the value of a state under an optimal policy must equal the expected return for the best action from that state*



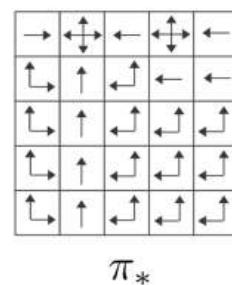
31

Optimal solutions to the gridworld example



22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

v_*

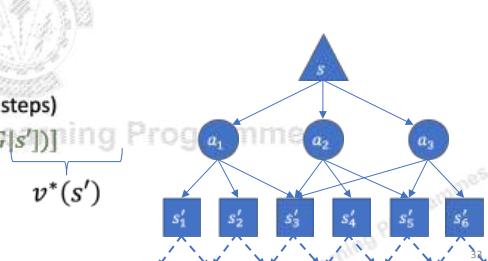


MDP - Objective

A set of states $s \in \mathcal{S}$
A set of actions $a \in \mathcal{A}$
State-transition probabilities $P(s'|s, a)$
A reward function $R(s, a, s')$

- Compute a policy: what action to take at each state
 - $\pi: S \rightarrow A$
- Compute the **optimal policy**: maximum expected reward, π^*
- $\pi^*(s) = ?$
- $= \underset{a}{\operatorname{argmax}} [\sum_{s'} P(s'|s, a) R(s, a, s')]$
- Must also optimize over the future (next steps)

$$= \underset{a}{\operatorname{argmax}} [\sum_{s'} P(s'|s, a) (R(s, a, s') + \mathbb{E}_{\pi^*}[G(s')])]$$





Notation

- π^* - a policy that yields the maximal expected sum of rewards
- G - observed sum of rewards, i.e., $\sum r_t$
- $v^*(s)$ - the expected sum of rewards from being at s then following π^*
• $= \mathbb{E}_{\pi^*}[G|s]$



Work Integrated Learning Programmes

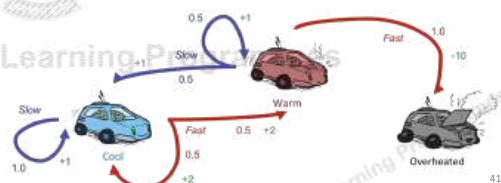
34

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Race car example

- Consider a discount factor, $\gamma = 0.9$
- What is $v^*(Cool)$
- $= \max_a [r(s, a) + \sum_{s'} p(s'|s, a) \gamma v^*(s')]$
- $= \max[1 + 0.9 \cdot 1v^*(Cool), 2 + 0.9 \cdot 0.5v^*(Cool) + 0.9 \cdot 0.5v^*(Warm)]$
 - Computing...
 - ...Stack overflow
- Work in iterations



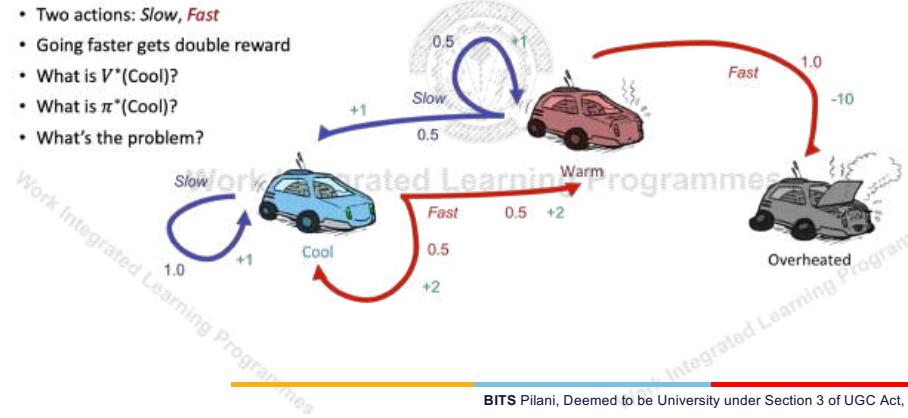
41

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Race car example

- A robot car wants to travel far, quickly
- Three states: Cool, Warm, Overheated
- Two actions: Slow, Fast
- Going faster gets double reward
- What is $V^*(Cool)$?
- What is $\pi^*(Cool)$?
- What's the problem?



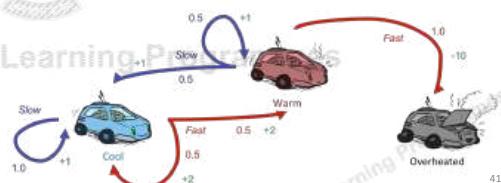
35

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Race car example

- Consider a discount factor, $\gamma = 0.9$
- What is $v^*(Cool)$
- $= \max_a [r(s, a) + \sum_{s'} p(s'|s, a) \gamma v^*(s')]$
- $= \max[1 + 0.9 \cdot 1v^*(Cool), 2 + 0.9 \cdot 0.5v^*(Cool) + 0.9 \cdot 0.5v^*(Warm)]$
 - Computing...
 - ...Stack overflow
- Work in iterations



41

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Value iteration

Value iteration

Initialize array V arbitrarily (e.g., $V(s) = 0$ for all $s \in S^+$)

Repeat

$$\Delta \leftarrow 0$$

For each $s \in S$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, $\pi \approx \pi_*$, such that

$$\pi(s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

42

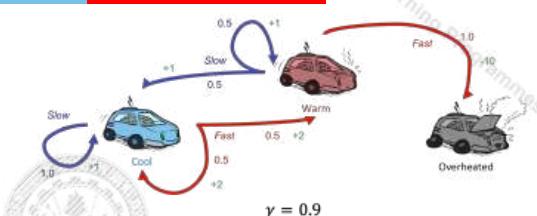
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Value Iteration

V₀

0	0	0
---	---	---



2	1	0
---	---	---

$$v_{k+1}(s) \doteq \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_k(s')]$$

3.35	2.35	0
------	------	---

Check this computation on paper.

43

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

k=0

Gridworld Display			
▲	▲	▲	0.00
0.00	0.00	0.00	0.00
▲		▲	0.00
0.00		0.00	0.00
▲	▲	▲	0.00
0.00	0.00	0.00	0.00

VALUES AFTER 0 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

45

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Example: Grid World

A maze-like problem

The agent lives in a grid

Walls block the agent's path

Noisy movement: actions do not always go as planned

80% of the time, the action North takes the agent North

10% of the time, North takes the agent West; 10% East

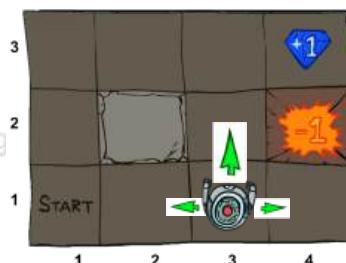
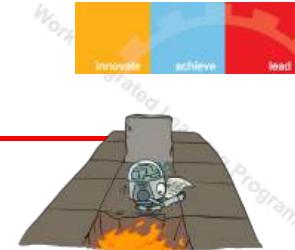
If there is a wall in the direction the agent would have been taken, the agent stays put

The agent receives rewards each time step

Small negative reward each step (battery drain)

Big rewards come at the end (good or bad)

Goal: maximize sum of (discounted) rewards



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

k=1

Gridworld Display			
▲	▲	▲	1.00
0.00	0.00	0.00	1.00
▲		▲	-1.00
0.00		0.00	0.00
▲	▲	▲	0.00
0.00	0.00	0.00	0.00

VALUES AFTER 1 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

46

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



k=2



Noise = 0.2
Discount = 0.9
Living reward = 0

47



k=3



Noise = 0.2
Discount = 0.9
Living reward = 0

48



k=4

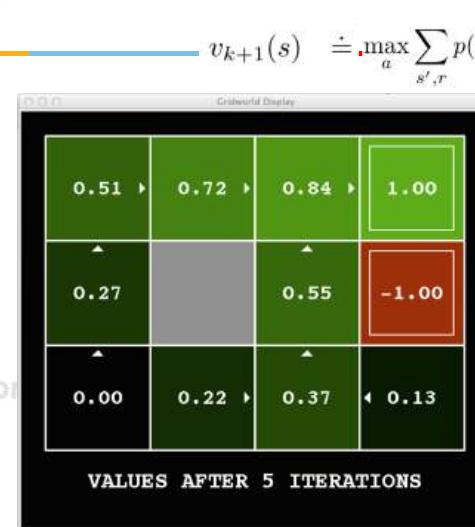


Noise = 0.2
Discount = 0.9
Living reward = 0

49



k=5



Noise = 0.2
Discount = 0.9
Living reward = 0

50



k=6

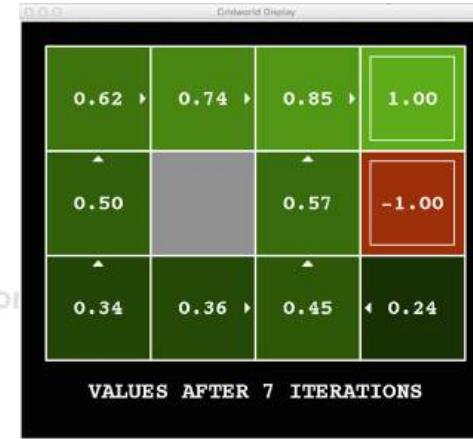


Noise = 0.2
Discount = 0.9
Living reward = 0

51



k=7



Noise = 0.2
Discount = 0.9
Living reward = 0

52



k=8

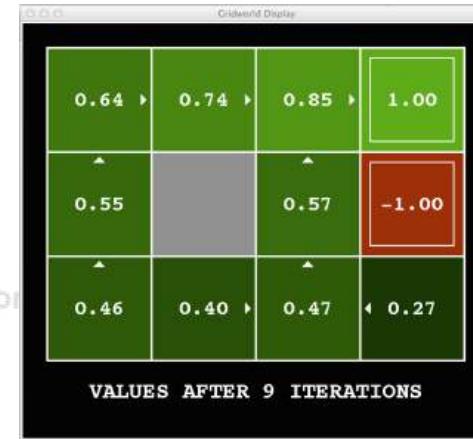


Noise = 0.2
Discount = 0.9
Living reward = 0

53



k=9

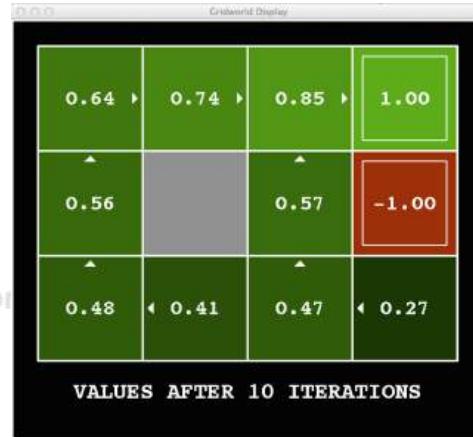


Noise = 0.2
Discount = 0.9
Living reward = 0

54



k=10



Noise = 0.2
Discount = 0.9
Living reward = 0



k=11



Noise = 0.2
Discount = 0.9
Living reward = 0



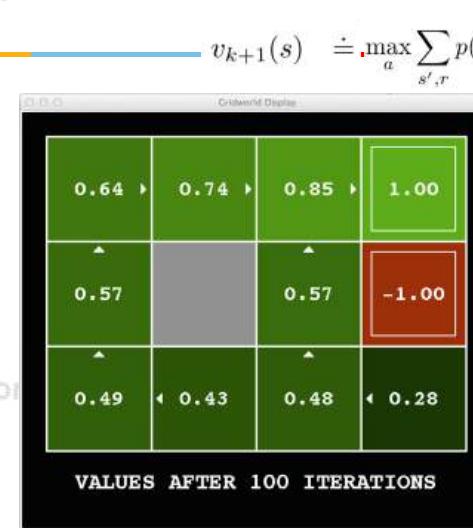
k=12



Noise = 0.2
Discount = 0.9
Living reward = 0



k=100



Noise = 0.2
Discount = 0.9
Living reward = 0



Problems with Value Iteration

- Value iteration repeats the Bellman updates:

$$V_{k+1}(s) \leftarrow \max_a [R(s, a) + \sum_{s'} P(s'|s, a) \gamma V_k(s')]$$

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

- Issue 1: It's slow – $O(S^2A)$ per iteration

- Do we really need to update every state at every iteration?

- Issue 2: A policy cannot be easily extracted
 - Policy extraction requires another $O(S^2A)$

- Issue 3: The policy often converges long before the values
 - Can we identify when the policy converged?

- Issue 4: requires knowing the model, $P(s'|s, a)$, and the reward function, $R(s, a)$

- Issue 5: requires discrete (finite) set of actions

- Issue 6: infeasible in large state spaces

59

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Solutions (briefly, more later...)

- Issue 1: It's slow – $O(S^2A)$ per iteration
 - Asynchronous value iteration

- Issue 2: A policy cannot be easily extracted
 - Learn q (action) values

- Issue 3: The policy often converges long before the values
 - Policy-based methods

- Issue 4: requires knowing the model and the reward function
 - Reinforcement learning

- Issue 5: requires discrete (finite) set of actions
 - Policy gradient methods

- Issue 6: infeasible for large (or continuous) state spaces
 - Function approximators

60

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Issue 1: It's slow – $O(S^2A)$ per iteration

- Asynchronous value iteration

- In value iteration, we update every state in each iteration

- Actually, *any sequences of Bellman updates will converge if every state is visited infinitely often regardless of the visitation order*

- Idea: prioritize states whose value we expect to change significantly

Work Integrated Learning Programmes

61

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Asynchronous Value Iteration

0.00	0.00	0.00	1.00
0.00		0.00	-1.00
0.00	0.00	0.00	0.00

VALUES AFTER 1 ITERATIONS

0.00	0.00	0.72	1.00
0.00		0.00	-1.00
0.00	0.00	0.00	0.00

VALUES AFTER 2 ITERATIONS

A single update per iteration

Algorithm 3 Prioritized Value Iteration

```

1: repeat
2:    $\rightarrow s \leftarrow \arg \max_{s \in S} H(s)$ 
3:    $V(s) \leftarrow \max_{a \in A} \{R(s, a) + \gamma \sum_{s' \in S} Pr(s'|s, a)V(s')\}$ 
4:   for all  $s' \in SDS(s)$  do
5:     // recompute  $H(s')$ 
6:   end for
7: until convergence

```

$$SDS(s) = \{s': \exists a, p(s'|s, a) > 0\}$$

For the home assignment set:

$$H(s') = \left| V(s') - \max_a \left\{ R(s', a) + \gamma \sum_{s''} Pr(s''|s', a)V(s'') \right\} \right|$$

62

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



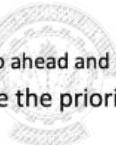
Double the work?

- Computing priority is similar to updating the state value (W.R.T computational effort)
- Why do double work?
 - If we computed the priority, we can go ahead and update the value for free
- Notice that we don't need to update the priorities for the entire state space
- For many of the states the priority doesn't change following an updated value for a single state s
- Only states s' with $\sum_a p(s'|s, a) > 0$ require update



For the home assignment set:

$$H(s') = \left| V(s') - \max_a \left\{ R(s', a) + \gamma \sum_{s''} \Pr(s''|s', a)V(s'') \right\} \right|$$

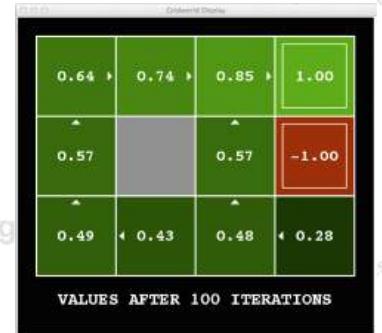


63



Issue 2: A policy cannot be easily extracted

- Given state values, what is the appropriate policy?
 - $\pi(s) \leftarrow \operatorname{argmax}_a [R(s, a) + \sum_{s'} P(s'|s, a) \gamma V_k(s')]$
 - Requires another full value sweep: $O(S^2 A)$
- Learn q (action) values instead
- $Q^*(s, a)$ - the expected sum of rewards from being at s , taking action a and then following π^*

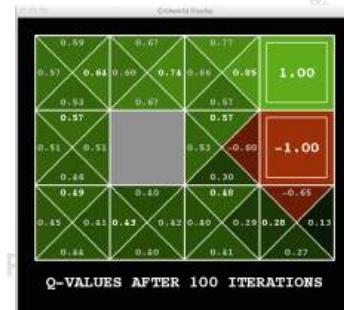


64



Q-learning

- $Q^*(s, a)$ - the expected sum of rewards from being at s , taking action a and then following π^*
- $\pi^*(s) \leftarrow \operatorname{argmax}_a [Q^*(s, a)]$
- Can we learn Q values with dynamic programming?
 - Yes, similar to value iteration



65



Q-learning as value iteration

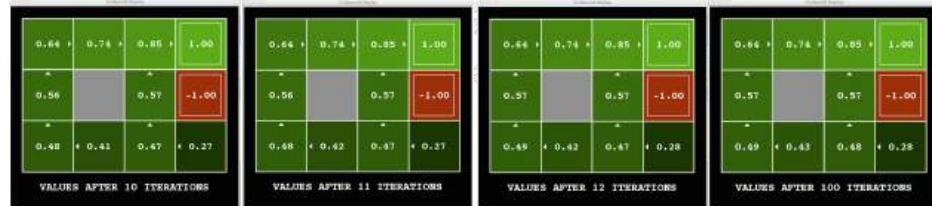
- $V^*(s) := \max_a [\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^*(s'))]$
- $V^*(s) := \max_a [Q^*(s, a)]$
- $Q^*(s, a) := \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^*(s'))$
- $Q^*(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_a [Q^*(s', a)])$
- Solve iteratively
 - $Q_{k+1}(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_a [Q_k(s', a)])$
 - Can also use Asynchronous learning

66



Issue 3: The policy often converges long before the values

- Value iteration converges to the true utility value: $V_{k \rightarrow \infty} \rightarrow V^*$
- V^* implies the optimal policy: π^*
- Can we converge directly on π^* ?
 - Improve the policy in iteration until reaching the optimal one



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

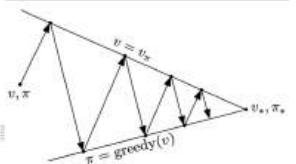


Policy Iteration

1. **Compute V_π :** calculate state value for some fixed policy (not necessarily the optimal values, $V_\pi \neq V^*$)
2. **Update π :** update policy using one-step look-ahead with the resulting (non optimal) values
3. Repeat until policy converges (optimal values and policy)

- Guaranteed converges to π^*

- $\forall s, V_{k>0}(s) \leq V_{k+1}(s)$ i.e., π_i improves monotonically with i
- A fixed point, $\forall s, V_k(s) = V_{k+1}(s)$, implies π^*



68

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Policy Evaluation

- Why is calculating V_π easier than calculating V^* ?
 - Turns non-linear Bellman equations into linear equations
- $v^*(s) = \max_a [\sum_s P(s'|s, a) (R(s, a, s') + \gamma v^*(s'))]$
- $v_\pi(s) = \sum_s P(s'|s, \pi(s)) (R(s, \pi(s), s') + \gamma v^*(s'))$
- Solve a set of linear equations in $O(S^2)$
 - Solve with Numpy (`numpy.linalg.solve`)
 - Required for your home assignment
 - See: <https://numpy.org/doc/stable/reference/generated/numpy.linalg.solve.html#numpy.linalg.solve>



69

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

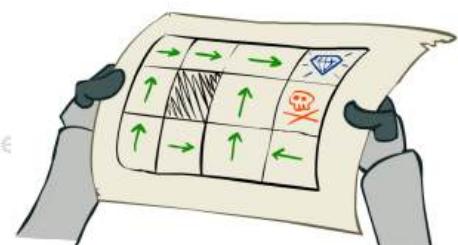


Policy value as a Linear program

- $v_{11} = 0.8 \cdot (-0.1 + 0.95 \cdot v_{12}) + 0.1 \cdot (-0.1 + 0.95 \cdot v_{21}) + 0.1 \cdot (-0.1 + 0.95 \cdot v_{11})$
- $v_{12} = 0.8 \cdot (-0.1 + 0.95 \cdot v_{13}) + 0.2 \cdot (-0.1 + 0.95 \cdot v_{12})$
- ...
- $v_{42} = -1$
- $v_{43} = 1$



Work Integrated Learning Programmes

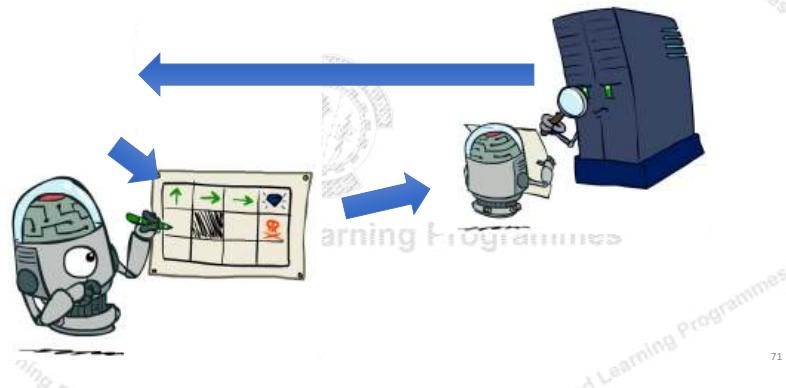


70

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Policy iteration



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

71



Issue 4: requires knowing the model and the reward function

- We will explore online learning (reinforcement learning) approaches
- How can we learn the model and reward function from interactions?
- Do we even need to learn them? Can we learn V^* , Q^* without a model?
- Can we do without V^* , Q^* ? Can we run policy iteration without a model?



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

73



Comparison

- Both value iteration and policy iteration compute the same thing (optimal state values)
- In value iteration:
 - Every iteration updates both the values and (implicitly) the policy
 - We don't track the policy, but taking the max over actions implicitly define it
- In policy iteration:
 - We do several passes that update utilities with fixed policies (each pass is fast because we consider only one action, not all of them)
 - After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
 - The new policy will be better (or we're done)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

72



Issue 5: requires discrete (finite) set of actions

- We will explore policy gradient approaches that are suitable for continuous actions, e.g., throttle and steering for a vehicle
- Can such approaches be relevant for discrete action spaces?
 - Yes! We can always define a continuous action space as a distribution over the discrete actions (e.g., using the softmax function)
- Can we combine value-based approaches and policy gradient approaches and get the best of both?
 - Yes! Actor-critic methods

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

74



Issue 6: infeasible in large (or continues) state spaces

- Most real-life problems contain very large state spaces (practically infinite)
- It is infeasible to learn and store a value for every state
- Moreover, doing so is not useful as the chance of encountering a state more than once is very small
- We will learn to generalize our learning to apply to unseen states
- We will use value function approximators that can generalize the acquired knowledge and provide a value to any state (even if it was not previously seen)



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

75



BITS Pilani

Deep Reinforcement Learning

Prof. Vimal SP
CSIS



Notation

- π^* - a policy that yields the maximal expected sum of rewards
- $V^*(s)$ - the expected sum of rewards from being at s then following π^*
- $V_\pi(s)$ - the expected sum of rewards from being at s then following π
- $Q^*(s, a)$ - the expected sum of rewards from being at s , taking action a and then following π^*
- $Q_\pi(s, a)$ - the expected sum of rewards from being at s , taking action a and then following π
- G_t - observed sum of rewards following time t , i.e., $\sum_{k=t}^T r_k$



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

76



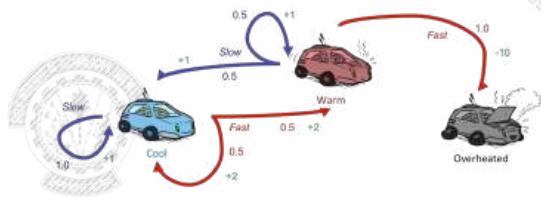
BITS Pilani

<AIMLCZG512, Deep Reinforcement Learning>
Lecture No. 6,7,8



Introduction

- Recollect the problem
 - We need to learn a policy that takes us as far and as faster possible;



Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

3



Introduction

- Still assume an underlying Markov decision process (MDP):

- A set of states $s \in S$
- A set of actions A
- A model $P(s'|s, a)$
- A reward function $R(s, a, s')$
- A discount factor γ
- Still looking for the best policy $\pi^*(S)$



Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

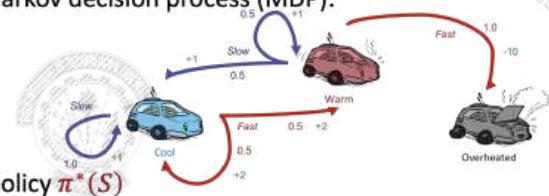
5



Introduction

- Still assume an underlying Markov decision process (MDP):

- A set of states $s \in S$
- A set of actions A
- A model $P(s'|s, a)$
- A reward function $R(s, a, s')$
- A discount factor γ
- Still looking for the best policy $\pi^*(S)$

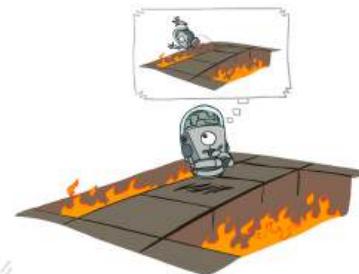


4

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



(Aside) Offline vs. Online (RL)



Offline Optimization



Online Learning

6

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Monte Carlo Methods

- Monte Carlo methods are a **broad class of computational algorithms** that *rely on repeated random sampling to obtain numerical results*
- The underlying concept is to obtain unbiased samples from a complex/unknown distribution through a random process
- They are often used in physical and mathematical problems and are most useful when it is difficult or impossible to compute a solution analytically
 - Weather prediction
 - Computational biology
 - Computer graphics
 - Finance and business
 - Sport game prediction

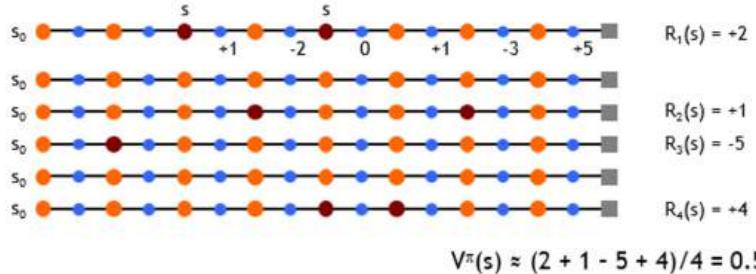


BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

7



Ex-1: First-visit Monte-Carlo Policy Evaluation [estimate $V\pi(s)$]



Acknowledgements : This example is taken from the tutorial by Peter Bodík, RAD Lab, UC Berkeley

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

9



First-visit Monte-Carlo Policy Evaluation [estimate $V\pi(s)$]

Initialize:

$\pi \leftarrow$ policy to be evaluated
 $V \leftarrow$ an arbitrary state-value function
 $Returns(s) \leftarrow$ an empty list, for all $s \in S$

Repeat forever:

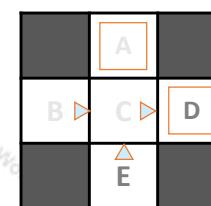
- Generate an episode using π
- For each state s appearing in the episode:
 $R \leftarrow$ return following the first occurrence of s
Append R to $Returns(s)$
 $V(s) \leftarrow$ average($Returns(s)$)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

8



Ex-2: First-visit Monte-Carlo Policy Evaluation [estimate $V\pi(s)$]

Input Policy π Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, , +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, , +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, , +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, , -10

Output Values

-10	A	
+8	B	+4
	C	+10
	D	
	E	-2

10

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956





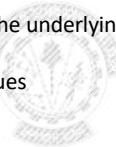
Problems with MC Evaluation

- What's good about direct evaluation?

- It's easy to understand
- It doesn't require any knowledge of the underlying model
- It converges to the true expected values

- What bad about it?

- It wastes information about transition probabilities
- Each state must be learned separately
- So, it takes a long time to learn



Output Values

	-10	
B	+8	+4
C		+10
D		-2
E		

Think: If B and E both go to C with the same probability, how can their values be different?

11

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



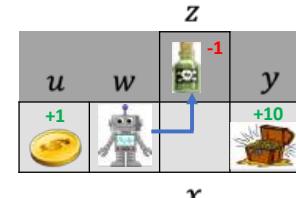
What about exploration?

- $S = \{u, w, x, y, z\}$

- $A = \{N, E, S, W, \text{exit}\}$

- Reward:

- $r(u, \text{exit}) = 1$
- $r(z, \text{exit}) = -1$
- $r(y, \text{exit}) = 10$



x

State	$\pi(s)$	$Q_\pi(N)$	$Q_\pi(E)$	$Q_\pi(S)$	$Q_\pi(W)$	$Q_\pi(\text{exit})$
u	exit	NA	NA	NA	NA	NA
w	E	0	-1	0	0	0
x	N	-1	0	0	0	NA
y	exit	NA	NA	NA	NA	0
z	exit	NA	NA	NA	NA	-1

12

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



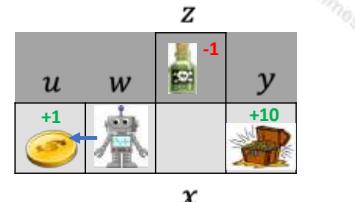
What about exploration?

- $S = \{u, w, x, y, z\}$

- $A = \{N, E, S, W, \text{exit}\}$

- Reward:

- $r(u, \text{exit}) = 1$
- $r(z, \text{exit}) = -1$
- $r(y, \text{exit}) = 10$



x

State	$\pi(s)$	$Q_\pi(N)$	$Q_\pi(E)$	$Q_\pi(S)$	$Q_\pi(W)$	$Q_\pi(\text{exit})$
u	exit	NA	NA	NA	NA	-1
w	W	0	-1	0	NA	NA
x	E	-1	0	0	0	NA
y	exit	NA	NA	NA	NA	0
z	exit	NA	NA	NA	NA	-1

13

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



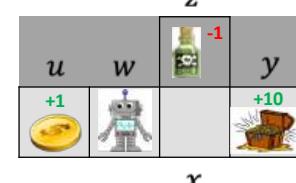
What about exploration?

- $S = \{u, w, x, y, z\}$

- $A = \{N, E, S, W, \text{exit}\}$

- Reward:

- $r(u, \text{exit}) = 1$
- $r(z, \text{exit}) = -1$
- $r(y, \text{exit}) = 10$



x

State	$\pi(s)$	$Q_\pi(N)$	$Q_\pi(E)$	$Q_\pi(S)$	$Q_\pi(W)$	$Q_\pi(\text{exit})$
u	exit	NA	NA	NA	NA	NA
w	W	0	-1	0	1	NA
x	E	-1	0	0	0	NA
y	exit	NA	NA	NA	NA	0
z	exit	NA	NA	NA	NA	-1

We converged
on a local
optimum!

14

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



ϵ -greedy MC control

On-policy first-visit MC control (for ϵ -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\epsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ϵ -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow$ average($Returns(S_t, A_t)$)

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \epsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$



MC control - example

• $Q =$

5	4,3	2,1	0
w	x	y	z

• $Returns =$

-	-,-	-,-	-
w	x	y	z

• $\pi(a|s) = (1 - \epsilon) \cdot$

• $\epsilon \cdot$ Random

exit	→	←	exit
w	x	y	z

On-policy first-visit MC control (for ϵ -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\epsilon > 0$

Initialize:

- $\pi \leftarrow$ an arbitrary ϵ -soft policy
- $Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$
- $Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

- Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$
- $G \leftarrow 0$
- Loop for each step of episode, $t = T-1, T-2, \dots, 0$:
- $G \leftarrow \gamma G + R_{t+1}$
- Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:
- Append G to $Returns(S_t, A_t)$
- $Q(S_t, A_t) \leftarrow$ average($Returns(S_t, A_t)$)
- $A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)
- For all $a \in \mathcal{A}(S_t)$:
- $\pi(a|S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \epsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$



MC control - example

• $Q =$

5	4,3	2,1	0
w	x	y	z

• $Returns =$

-	-,-	-,-	-
w	x	y	z

• $\pi(a|s) = (1 - \epsilon) \cdot$

• $\epsilon \cdot$ Random

• $\tau = x, \leftarrow, 0, w, exit, -100$

Work Integrated Learning

Programmes

On-policy first-visit MC control (for ϵ -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\epsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ϵ -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow$ average($Returns(S_t, A_t)$)

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \epsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$



MC control - example

• $Q =$

5	4,3	2,1	0
w	x	y	z

• $Returns =$

-100	-90,0	-,-	-
w	x	y	z

• $\pi(a|s) = (1 - \epsilon) \cdot$

• $\epsilon \cdot$ Random

exit	→	←	exit
w	x	y	z

On-policy first-visit MC control (for ϵ -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\epsilon > 0$

Initialize:

- $\pi \leftarrow$ an arbitrary ϵ -soft policy
- $Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$
- $Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

- Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$
- $G \leftarrow 0$
- Loop for each step of episode, $t = T-1, T-2, \dots, 0$:
- $G \leftarrow \gamma G + R_{t+1}$
- Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:
- Append G to $Returns(S_t, A_t)$
- $Q(S_t, A_t) \leftarrow$ average($Returns(S_t, A_t)$)
- $A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)
- For all $a \in \mathcal{A}(S_t)$:
- $\pi(a|S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \epsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$

Work Integrated Learning

Programmes



MC control - example

$\bullet Q =$

-100	-90,3	2,1	0
w	x	y	z

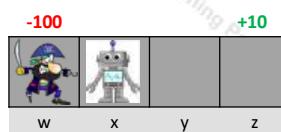
$\bullet \text{Returns} =$

-100	-90,0	-,-	-
w	x	y	z

$\bullet \pi(a|s) = (1 - \varepsilon) \cdot$

$\quad \varepsilon \cdot \text{Random}$

$\bullet \tau = x, \leftarrow, 0, w, \text{exit}, -100$



$\gamma = 0.9$

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

```

Algorithm parameter: small  $\varepsilon > 0$ 
Initialise:
     $\pi \leftarrow$  an arbitrary  $\varepsilon$ -soft policy
     $Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 
    Returns( $s, a$ )  $\leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 

Repeat forever [for each episode]:
    Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
     $G \leftarrow 0$ 
    Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
         $G \leftarrow \gamma G + R_{t+1}$ 
        Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :
            Append  $G$  to Returns( $S_t, A_t$ )
             $Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$ 
             $A^* \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken arbitrarily)
        For all  $a \in \mathcal{A}(S_t)$ :
             $\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

$\gamma = 0.9$

MC control - example

-100	-90,3	2,1	0
w	x	y	z

$\bullet \text{Returns} =$

-100	-90,0	-,-	-
w	x	y	z

$\bullet \pi(a|s) = (1 - \varepsilon) \cdot$

$\quad \varepsilon \cdot \text{Random}$



$\bullet \tau = x, \leftarrow, 0, w, \text{exit}, -100$

$\bullet A^* = [\rightarrow, \text{exit}]$

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

```

Algorithm parameter: small  $\varepsilon > 0$ 
Initialise:
     $\pi \leftarrow$  an arbitrary  $\varepsilon$ -soft policy
     $Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 
    Returns( $s, a$ )  $\leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 

Repeat forever [for each episode]:
    Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
     $G \leftarrow 0$ 
    Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
         $G \leftarrow \gamma G + R_{t+1}$ 
        Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :
            Append  $G$  to Returns( $S_t, A_t$ )
             $Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$ 
             $A^* \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken arbitrarily)
        For all  $a \in \mathcal{A}(S_t)$ :
             $\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control - example

$\bullet Q =$

-100	-90,3	2,1	0
w	x	y	z

$\bullet \text{Returns} =$

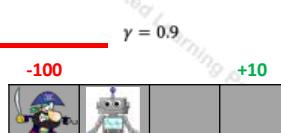
-100	-90,0	-,-	-
w	x	y	z

$\bullet \pi(a|s) = (1 - \varepsilon) \cdot$

$\quad \varepsilon \cdot \text{Random}$

$\bullet \tau = x, \leftarrow, 0, w, \text{exit}, -100$

$\bullet A^* = [\rightarrow, \text{exit}]$



$\gamma = 0.9$

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

```

Algorithm parameter: small  $\varepsilon > 0$ 
Initialise:
     $\pi \leftarrow$  an arbitrary  $\varepsilon$ -soft policy
     $Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 
    Returns( $s, a$ )  $\leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 

Repeat forever [for each episode]:
    Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
     $G \leftarrow 0$ 
    Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
         $G \leftarrow \gamma G + R_{t+1}$ 
        Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :
            Append  $G$  to Returns( $S_t, A_t$ )
             $Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$ 
             $A^* \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken arbitrarily)
        For all  $a \in \mathcal{A}(S_t)$ :
             $\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

$\gamma = 0.9$

-100	-90,3	2,1	0
w	x	y	z

$\bullet \text{Returns} =$

-100	-90,0	-,-	-
w	x	y	z

$\bullet \pi(a|s) = (1 - \varepsilon) \cdot$

$\quad \varepsilon \cdot \text{Random}$



$\bullet \tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

```

Algorithm parameter: small  $\varepsilon > 0$ 
Initialise:
     $\pi \leftarrow$  an arbitrary  $\varepsilon$ -soft policy
     $Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 
    Returns( $s, a$ )  $\leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 

Repeat forever [for each episode]:
    Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
     $G \leftarrow 0$ 
    Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
         $G \leftarrow \gamma G + R_{t+1}$ 
        Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :
            Append  $G$  to Returns( $S_t, A_t$ )
             $Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$ 
             $A^* \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken arbitrarily)
        For all  $a \in \mathcal{A}(S_t)$ :
             $\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



MC control - example

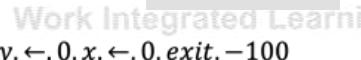
• $Q =$

-100	-90,-72.9	-81,1	0
w	x	y	z

• $\text{Returns} =$

-100	-90,-72.9	-81,-	-
w	x	y	z

• $\pi(a|s) = (1 - \varepsilon) \cdot$
 $\varepsilon \cdot \text{Random}$



• $\tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$

• $\tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$

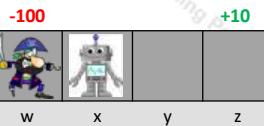
• $A^* = [\rightarrow, \rightarrow, \text{exit}]$

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



$\gamma = 0.9$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$\text{Returns}(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $\text{Returns}(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$

26

Work Integrated Learning Programmes

MC control - example

• $Q =$

-100	-90,-72.9	-81,1	0
w	x	y	z

• $\text{Returns} =$

-100	-90,-72.9	-81,-	-
w	x	y	z

• $\pi(a|s) = (1 - \varepsilon) \cdot$
 $\varepsilon \cdot \text{Random}$



• $\tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$

• $A^* = [\rightarrow, \rightarrow, \text{exit}]$

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



MC control - example

• $Q =$

-100	-90,-72.9	-81,1	0
w	x	y	z

• $\text{Returns} =$

-100	-90,-72.9	-81,-	-
w	x	y	z

• $\pi(a|s) = (1 - \varepsilon) \cdot$
 $\varepsilon \cdot \text{Random}$



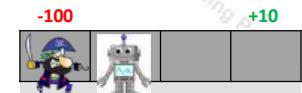
• $\tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$

• $A^* = [\rightarrow, \rightarrow, \text{exit}]$

Work Integrated Learning Programmes



$\gamma = 0.9$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$\text{Returns}(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $\text{Returns}(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

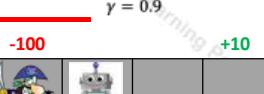
$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$

25

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



$\gamma = 0.9$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$\text{Returns}(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $\text{Returns}(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$

26

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



On-policy learning

Estimation True value

• $Q_\pi(w, \rightarrow) = q_\pi(w, \rightarrow) = -1$ (correct!)

• $q^*(w, \rightarrow) = ?$

• $q_{\pi \neq b} = ?$

• Observation drawn from π are useful for evaluating q_π

• Once the policy is changed these observations are irrelevant

• This is not sample efficient!



27



State	$\pi(s)$	$Q(s?)$	$Q(s\rightarrow)$
u	exit	NA	NA
w	\rightarrow	0	-1
x	\uparrow	-1	0
y	exit	NA	NA
z	exit	NA	NA

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Quick Recap ! On-policy vs. Off-policy Learning

Off-policy learning

- We would like to use observations drawn from some policy b to evaluate q_π where $\pi \neq b$, specifically, we want to evaluate q_{π^*}
- We strive for full utilization of previous experience
- Off-policy learning allows us to optimize a **target policy** while following another **behavior policy**
- **Pros:** sample efficient
- **Cons:** greater variance in value estimations

Importance sampling

- Given a trajectory τ drawn by running b
- We can **define** (not compute) the probability $\Pr\{\tau|b\}$
- We can also **define** $\Pr\{\tau|\pi\}$
- Define the **importance sampling ratio** as: $\rho_t = \frac{\Pr\{\tau_t|\pi\}}{\Pr\{\tau_t|b\}}$
- Can we **compute** ρ without a model, $p(S_{K+1}|S_K, A_K)$?
- $\rho_t = \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k) p(S_{K+1}|S_K, A_K)}{\prod_{k=t}^{T-1} b(A_k|S_k) p(S_{K+1}|S_K, A_K)}$ YES!



Importance sampling

- Fact: $\mathbb{E}_{\tau \sim b}[G_t | S_t = s] = v_b(s)$
- Importance sampling allows us to compute an unbiased estimate of $v_\pi(s)$ by running b
- Claim: $\mathbb{E}_{\tau \sim b}[\rho_t G_t | S_t = s] = v_\pi(s)$
- We set $v_\pi(s)$ to be a weighted average of observed returns (weighted by the importance ratio)
- Assume visiting state s over M episodes using policy b
 - s is first visited during time step t^m during each episode, $m \in M$
- $v_\pi(s) = \frac{\sum_{m \in M} \rho_t^m G_t^m}{M}$

33

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



(ordinary) Importance sampling - example

- $b(s|a) = \begin{cases} \rightarrow, p(0.5) \\ \leftarrow, p(0.5) \end{cases}$
- $\pi(s|a) = \begin{cases} \rightarrow, p(0.99) \\ \leftarrow, p(0.01) \end{cases}$
- $\tau_1 = \{w, \rightarrow, 0, x, \rightarrow, 0, y, \text{exit}, 10\}$
- $v_\pi(w) = \frac{\sum_{m \in M} \rho_t^m G_t^m}{M} = \frac{0.99 * 0.99}{0.5 * 0.5} * 10 = 3.96 * 10$
- $\tau_2 = \{w, \leftarrow, 0, u, \text{exit}, 1\}$
- $v_\pi(w) = \frac{\sum_{m \in M} \rho_t^m G_t^m}{M} = \frac{3.96 * 10 + 0.02 * 1}{2}$



39.6 ??
Ordinary Importance sampling is unbiased yet high variance

36

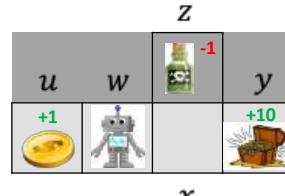
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Weighted importance sampling

- $v_\pi(s) = \frac{\sum_{m \in M} [\rho_t^m G_t^m]}{\sum_{m \in M} \rho_t^m}$
- $b(s|a) = \begin{cases} \rightarrow, p(0.5) \\ \leftarrow, p(0.5) \end{cases}$
- $\pi(s|a) = \begin{cases} \rightarrow, p(0.99) \\ \leftarrow, p(0.01) \end{cases}$
- $\tau_1 = \{w, \rightarrow, 0, x, \rightarrow, 0, y, \text{exit}, 10\}$
- $v_\pi(w) = \frac{\sum_{m \in M} [\rho_t^m G_t^m]}{\sum_{m \in M} \rho_t^m} = \frac{3.96 * 10}{3.96}$
- $\tau_2 = \{w, \leftarrow, 0, u, \text{exit}, 1\}$
- $v_\pi(w) = \frac{\sum_{m \in M} [\rho_t^m G_t^m]}{\sum_{m \in M} \rho_t^m} = \frac{3.96 * 10 + 0.02 * 1}{3.98}$

Trick: normalize by the sum of importance ratios



Ordinary Importance sampling is unbiased while the weighted version is biased (initially). Ordinary Importance sampling results in high variance while the weighted version has a bounded variance

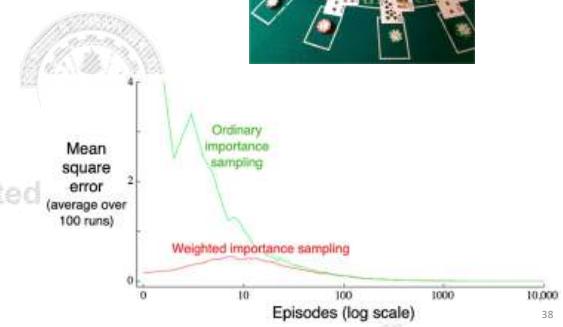
37

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Ordinary vs weighted importance sampling

- Estimating a black-jack state
- Target policy: hit on 19 or below
- Behavior policy: random (uniform)
- Both approaches converge to the true value
- weighted importance sampling is much better initially



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

38



MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

```

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :
 $Q(s, a) \leftarrow$  arbitrary
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
   $b \leftarrow$  any soft policy
  Generate an episode using  $b$ :
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T - 1, T - 2, \dots$  down to 0:
       $G \leftarrow \gamma G + R_{t+1}$ 
       $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
       $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
       $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)
      If  $A_t \neq \pi(S_t)$  then exit For loop
       $W \leftarrow W \frac{1}{h(A_t|S_t)}$ 
  
```

Discount future rewards and add immediate reward

43

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

```

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :
 $Q(s, a) \leftarrow$  arbitrary
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
   $b \leftarrow$  any soft policy
  Generate an episode using  $b$ :
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T - 1, T - 2, \dots$  down to 0:
       $G \leftarrow \gamma G + R_{t+1}$ 
       $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
       $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
       $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)
      If  $A_t \neq \pi(S_t)$  then exit For loop
       $W \leftarrow W \frac{1}{h(A_t|S_t)}$ 
  
```

Cumulative sum of IS weights affiliated with S_t, A_t (for weighted IS)

44



MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

```

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :
 $Q(s, a) \leftarrow$  arbitrary
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
   $b \leftarrow$  any soft policy
  Generate an episode using  $b$ :
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T - 1, T - 2, \dots$  down to 0:
       $G \leftarrow \gamma G + R_{t+1}$ 
       $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
       $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
       $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)
      If  $A_t \neq \pi(S_t)$  then exit For loop
       $W \leftarrow W \frac{1}{h(A_t|S_t)}$ 
  
```

Incremental update of Q values (waited moving average)

45

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

```

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :
 $Q(s, a) \leftarrow$  arbitrary
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
   $b \leftarrow$  any soft policy
  Generate an episode using  $b$ :
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T - 1, T - 2, \dots$  down to 0:
       $G \leftarrow \gamma G + R_{t+1}$ 
       $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
       $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
       $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)
      If  $A_t \neq \pi(S_t)$  then exit For loop
       $W \leftarrow W \frac{1}{h(A_t|S_t)}$ 
  
```

Update target policy (greedy)

46



MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

```
Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :
 $Q(s, a) \leftarrow$  arbitrary
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \operatorname{argmax}_a Q(S_t, a)$  (with ties broken consistently)
```

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$

$G \leftarrow 0$

$W \leftarrow 1$

For $t = T - 1, T - 2, \dots$ down to 0:

$G \leftarrow \gamma G + R_{t+1}$

$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$

$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$

$\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)

If $A_t \neq \pi(S_t)$ then exit For loop

$W \leftarrow W \frac{1}{\pi(A_t|S_t)}$

Since π is deterministic,
once we diverge from it all
IS weights of earlier
actions will be 0

47

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

```
Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :
 $Q(s, a) \leftarrow$  arbitrary
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \operatorname{argmax}_a Q(S_t, a)$  (with ties broken consistently)
```

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$

$G \leftarrow 0$

$W \leftarrow 1$

For $t = T - 1, T - 2, \dots$ down to 0:

$G \leftarrow \gamma G + R_{t+1}$

$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$

$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$

$\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)

If $A_t \neq \pi(S_t)$ then exit For loop

$W \leftarrow W \frac{1}{\pi(A_t|S_t)}$

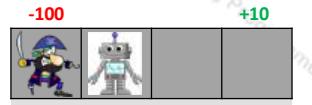
Update the joint prob by
multiplying by ρ_t . Notice
that $\pi(S_t) = 1$ in this
example (deterministic
target policy)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



MC control + IS example

-	4,3	2,1	-
w	x	y	z
-	0,0	0,0	-
w	x	y	z



Off-policy MC control, for estimating $\pi \approx \pi_*$

```
Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :
 $Q(s, a) \leftarrow$  arbitrary
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \operatorname{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
 $b \leftarrow$  any soft policy
Generate an episode using  $b$ :
 $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
 $G \leftarrow 0$ 
 $W \leftarrow 1$ 
For  $t = T - 1, T - 2, \dots$  down to 0:
 $G \leftarrow \gamma G + R_{t+1}$ 
 $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
 $\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$  (with ties broken consistently)
If  $A_t \neq \pi(S_t)$  then exit For loop
 $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
```

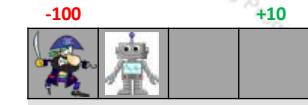
Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



MC control + IS example

-	4,3	2,1	-
w	x	y	z
-	0,0	0,0	-
w	x	y	z



Off-policy MC control, for estimating $\pi \approx \pi_*$

```
Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :
 $Q(s, a) \leftarrow$  arbitrary
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \operatorname{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
 $b \leftarrow$  any soft policy
Generate an episode using  $b$ :
 $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
 $G \leftarrow 0$ 
 $W \leftarrow 1$ 
For  $t = T - 1, T - 2, \dots$  down to 0:
 $G \leftarrow \gamma G + R_{t+1}$ 
 $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
 $\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$  (with ties broken consistently)
If  $A_t \neq \pi(S_t)$  then exit For loop
 $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



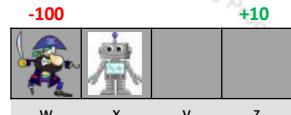


MC control + IS example

5	4,3	2,1	5
w	x	y	z
0	0,0	0,0	0
w	x	y	z



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
 $Q(s, a) \leftarrow \text{arbitrary}$ 
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

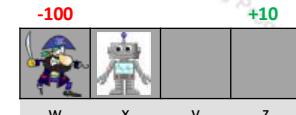
Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T-1, T-2, \dots$  down to 0:
         $G \leftarrow \gamma G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then exit For loop
         $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
    
```

MC control + IS example

5	4,3	2,1	5
w	x	y	z
0	0,0	0,0	0
w	x	y	z



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
 $Q(s, a) \leftarrow \text{arbitrary}$ 
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T-1, T-2, \dots$  down to 0:
         $G \leftarrow \gamma G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then exit For loop
         $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
    
```

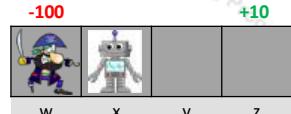


MC control + IS example

5	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
 $Q(s, a) \leftarrow \text{arbitrary}$ 
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

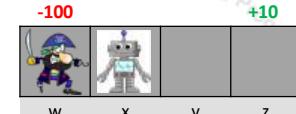
Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T-1, T-2, \dots$  down to 0:
         $G \leftarrow \gamma G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then exit For loop
         $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
    
```

MC control + IS example

-100	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Off-policy

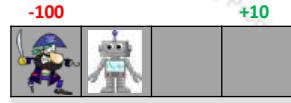
```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
 $Q(s, a) \leftarrow \text{arbitrary}$ 
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T-1, T-2, \dots$  down to 0:
         $G \leftarrow \gamma G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then exit For loop
         $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
    
```

MC control + IS example

-100	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z



Off-pol

Initialise, for all $s \in S$, $a \in A(s)$:
 $Q(s, a) \leftarrow \text{arbitrary}$
 $C(s, a) \leftarrow 0$
 $\pi(s) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)

Repeat forever:
 $b \leftarrow$ any soft policy
Generate an episode using b :
 $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$
 $G \leftarrow 0$
 $W \leftarrow 1$
For $t = T-1, T-2, \dots$ down to 0:
 $G \leftarrow \gamma G + R_{t+1}$
 $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$
 $\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)
If $A_t \neq \pi(S_t)$ then exit For loop
 $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$

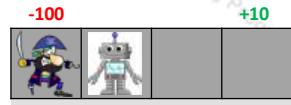
Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



MC control + IS example

-100	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z



Off-pol

Initialise, for all $s \in S$, $a \in A(s)$:
 $Q(s, a) \leftarrow \text{arbitrary}$
 $C(s, a) \leftarrow 0$
 $\pi(s) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)

Repeat forever:
 $b \leftarrow$ any soft policy
Generate an episode using b :
 $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$
 $G \leftarrow 0$
 $W \leftarrow 1$
For $t = T-1, T-2, \dots$ down to 0:
 $G \leftarrow \gamma G + R_{t+1}$
 $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$
 $\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)
If $A_t \neq \pi(S_t)$ then exit For loop
 $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$

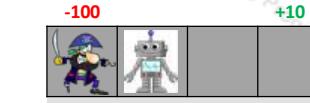
Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



MC control + IS example

-100	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z



Off-pol

Initialise, for all $s \in S$, $a \in A(s)$:
 $Q(s, a) \leftarrow \text{arbitrary}$
 $C(s, a) \leftarrow 0$
 $\pi(s) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)

Repeat forever:
 $b \leftarrow$ any soft policy
Generate an episode using b :
 $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$
 $G \leftarrow 0$
 $W \leftarrow 1$
For $t = T-1, T-2, \dots$ down to 0:
 $G \leftarrow \gamma G + R_{t+1}$
 $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$
 $\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)
If $A_t \neq \pi(S_t)$ then exit For loop
 $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$

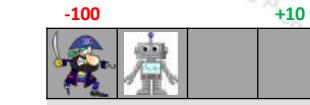
Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



MC control + IS example

-100	4,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z



Off-pol

Initialise, for all $s \in S$, $a \in A(s)$:
 $Q(s, a) \leftarrow \text{arbitrary}$
 $C(s, a) \leftarrow 0$
 $\pi(s) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)

Repeat forever:
 $b \leftarrow$ any soft policy
Generate an episode using b :
 $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$
 $G \leftarrow 0$
 $W \leftarrow 1$
For $t = T-1, T-2, \dots$ down to 0:
 $G \leftarrow \gamma G + R_{t+1}$
 $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$
 $\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken consistently)
If $A_t \neq \pi(S_t)$ then exit For loop
 $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$

Work Integrated Learning

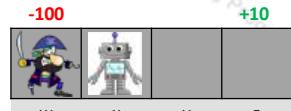
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956





MC control + IS example

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z



Off-pol

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow$  arbitrary
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
         $G \leftarrow 0$ 
         $W \leftarrow 1$ 
        For  $t = T-1, T-2, \dots$  down to 0:
             $G \leftarrow \gamma G + R_{t+1}$ 
             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
            If  $A_t \neq \pi(S_t)$  then exit For loop
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 

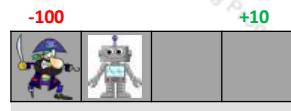
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z



Off-pol

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow$  arbitrary
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
         $G \leftarrow 0$ 
         $W \leftarrow 1$ 
        For  $t = T-1, T-2, \dots$  down to 0:
             $G \leftarrow \gamma G + R_{t+1}$ 
             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
            If  $A_t \neq \pi(S_t)$  then exit For loop
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 

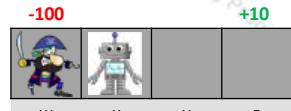
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z



Off-pol

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow$  arbitrary
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

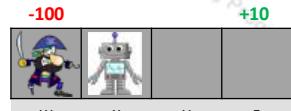
Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
         $G \leftarrow 0$ 
         $W \leftarrow 1$ 
        For  $t = T-1, T-2, \dots$  down to 0:
             $G \leftarrow \gamma G + R_{t+1}$ 
             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
            If  $A_t \neq \pi(S_t)$  then exit For loop
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z



Off-pol

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow$  arbitrary
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
         $G \leftarrow 0$ 
         $W \leftarrow 1$ 
        For  $t = T-1, T-2, \dots$  down to 0:
             $G \leftarrow \gamma G + R_{t+1}$ 
             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
            If  $A_t \neq \pi(S_t)$  then exit For loop
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z



Off-pol

```
Initialize, for all  $s \in S$ ,  $a \in A(s)$ :  

 $Q(s, a) \leftarrow$  arbitrary  

 $C(s, a) \leftarrow 0$   

 $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:  

     $b \leftarrow$  any soft policy  

    Generate an episode using  $b$ :  

         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$   

         $G \leftarrow 0$   

         $W \leftarrow 1$   

        For  $t = T-1, T-2, \dots$  down to 0:  

             $G \leftarrow \gamma G + R_{t+1}$   

             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$   

             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$   

             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)  

            If  $A_t \neq \pi(S_t)$  then exit For loop  

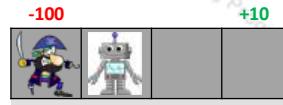
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,0	1
w	x	y	z



Off-pol

```
Initialize, for all  $s \in S$ ,  $a \in A(s)$ :  

 $Q(s, a) \leftarrow$  arbitrary  

 $C(s, a) \leftarrow 0$   

 $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:  

     $b \leftarrow$  any soft policy  

    Generate an episode using  $b$ :  

         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$   

         $G \leftarrow 0$   

         $W \leftarrow 1$   

        For  $t = T-1, T-2, \dots$  down to 0:  

             $G \leftarrow \gamma G + R_{t+1}$   

             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$   

             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$   

             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)  

            If  $A_t \neq \pi(S_t)$  then exit For loop  

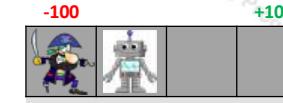
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	1
w	x	y	z



Off-pol

```
Initialize, for all  $s \in S$ ,  $a \in A(s)$ :  

 $Q(s, a) \leftarrow$  arbitrary  

 $C(s, a) \leftarrow 0$   

 $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:  

     $b \leftarrow$  any soft policy  

    Generate an episode using  $b$ :  

         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$   

         $G \leftarrow 0$   

         $W \leftarrow 1$   

        For  $t = T-1, T-2, \dots$  down to 0:  

             $G \leftarrow \gamma G + R_{t+1}$   

             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$   

             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$   

             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)  

            If  $A_t \neq \pi(S_t)$  then exit For loop  

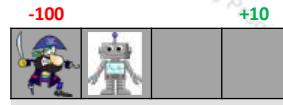
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,0	1
w	x	y	z



Off-pol

```
Initialize, for all  $s \in S$ ,  $a \in A(s)$ :  

 $Q(s, a) \leftarrow$  arbitrary  

 $C(s, a) \leftarrow 0$   

 $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:  

     $b \leftarrow$  any soft policy  

    Generate an episode using  $b$ :  

         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$   

         $G \leftarrow 0$   

         $W \leftarrow 1$   

        For  $t = T-1, T-2, \dots$  down to 0:  

             $G \leftarrow \gamma G + R_{t+1}$   

             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$   

             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$   

             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)  

            If  $A_t \neq \pi(S_t)$  then exit For loop  

             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,0	1
w	x	y	z

exit			exit
w	x	y	z

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
 $Q(s, a) \leftarrow \text{arbitrary}$ 
 $C(s, a) \leftarrow 0$ 
 $\pi(s) \leftarrow \text{argmax}_a Q(s, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ :
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T - 1, T - 2, \dots$  down to 0:
         $G \leftarrow +G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then Exit For loop
         $W \leftarrow W C(S_t, A_t)^{-1}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z

exit			exit
w	x	y	z

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
   $Q(s, a) \leftarrow$  arbitrary
   $C(s, a) \leftarrow 0$ 
 $\pi(a) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Report forever:
  b ← any soft policy
  Generate an episode using b:
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T - 1, T - 2, \dots$  down to 0:
       $G \leftarrow \gamma G + R_{t+1}$ 
       $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
       $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} (G - Q(S_t, A_t))$ 
       $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
      If  $A_t \neq \pi(S_t)$  then Exit For loop
       $W \leftarrow W \cdot \frac{1}{C(S_t, A_t)}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,0	1
w	x	y	z
exit			exit
w	x	y	z

```

Off-policy TD Control
Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow \text{arbitrary}$ 
     $C(s, a) \leftarrow 0$ 
     $\pi(a) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow \text{any soft policy}$ 
    Generate an episode using  $b$ :
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 1$ 
     $W \leftarrow 1$ 
    For  $t = T-1, T-2, \dots$  down to 0:
         $G \leftarrow r + G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then Exit For loop
     $W \leftarrow W \cdot \alpha / |S|$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
 $Q(s, a) \leftarrow \text{arbitrary}$ 
 $C(s, a) \leftarrow 0$ 
 $\pi(a) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

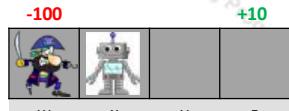
Repeat forever:
  b ← any soft policy
  Generate an episode using b:
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T - 1, T - 2, \dots, 0$  down to 0:
       $G \leftarrow \gamma G + R_{t+1}$ 
       $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
       $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
       $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
      If  $A_t \neq \pi(S_t)$  then Exit For loop
     $W \leftarrow W \frac{1}{C(S_0, A_0)}$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow \text{arbitrary}$ 
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow \text{any soft policy}$ 
    Generate an episode using  $b$ 
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
         $G \leftarrow 0$ 
         $W \leftarrow 1$ 
        For  $t = T-1, T-2, \dots$  down to 0:
             $G \leftarrow \gamma G + R_{t+1}$ 
             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
            If  $A_t \neq \pi(S_t)$  then exit For loop
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 

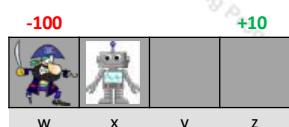
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow \text{arbitrary}$ 
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow \text{any soft policy}$ 
    Generate an episode using  $b$ 
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
         $G \leftarrow 0$ 
         $W \leftarrow 1$ 
        For  $t = T-1, T-2, \dots$  down to 0:
             $G \leftarrow \gamma G + R_{t+1}$ 
             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
            If  $A_t \neq \pi(S_t)$  then exit For loop
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 

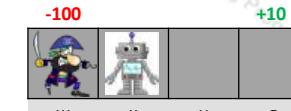
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow \text{arbitrary}$ 
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow \text{any soft policy}$ 
    Generate an episode using  $b$ 
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
         $G \leftarrow 0$ 
         $W \leftarrow 1$ 
        For  $t = T-1, T-2, \dots$  down to 0:
             $G \leftarrow \gamma G + R_{t+1}$ 
             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
            If  $A_t \neq \pi(S_t)$  then exit For loop
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 

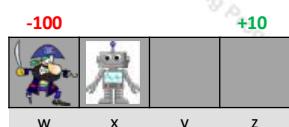
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow \text{arbitrary}$ 
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)

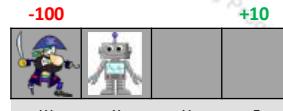
Repeat forever:
     $b \leftarrow \text{any soft policy}$ 
    Generate an episode using  $b$ 
         $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
         $G \leftarrow 0$ 
         $W \leftarrow 1$ 
        For  $t = T-1, T-2, \dots$  down to 0:
             $G \leftarrow \gamma G + R_{t+1}$ 
             $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
             $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken consistently)
            If  $A_t \neq \pi(S_t)$  then exit For loop
             $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 

```

Work Integrated Learning

MC control + IS example

-100	-90,8,1	2,9	10
w	x	y	z
1	2,4	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow \text{arbitrary}$ 
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)

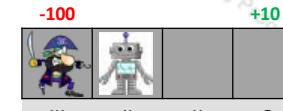
Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T-1, T-2, \dots$  down to 0:
         $G \leftarrow \gamma G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then exit For loop
         $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
    
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,8,1	2,9	10
w	x	y	z
1	2,4	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow \text{arbitrary}$ 
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T-1, T-2, \dots$  down to 0:
         $G \leftarrow \gamma G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then exit For loop
         $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
    
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MC control + IS example

-100	-90,8,1	2,9	10
w	x	y	z
1	2,4	0,2	1
w	x	y	z
exit			exit
w	x	y	z

Off-policy

```

Initialize, for all  $s \in S$ ,  $a \in A(s)$ :
     $Q(s, a) \leftarrow \text{arbitrary}$ 
     $C(s, a) \leftarrow 0$ 
     $\pi(s) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)

Repeat forever:
     $b \leftarrow$  any soft policy
    Generate an episode using  $b$ 
     $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ 
     $G \leftarrow 0$ 
     $W \leftarrow 1$ 
    For  $t = T-1, T-2, \dots$  down to 0:
         $G \leftarrow \gamma G + R_{t+1}$ 
         $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$ 
         $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$ 
         $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken consistently)
        If  $A_t \neq \pi(S_t)$  then exit For loop
         $W \leftarrow W \frac{1}{\pi(A_t|S_t)}$ 
    
```

Work Integrated Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

What did we learn?

- Online evaluation through the Monte-Carlo approach

- Update V or Q estimations based on observed returns

- On policy Monte-Carlo control

- Beware of local optimum. Must explore!
- Consider using a soft policy $\pi(a|s)$
- Assign soft policy values that mimic an ϵ -greedy strategy
- On policy learning is not sample efficient!

- Off policy Monte-Carlo control

- Define target policy (e.g., $\arg\max_a Q(S_t, a)$) that can be different than the behavior policy

- Utilize weighted importance sampling to train a target policy

$$\nu_\pi(s) = \frac{\sum_{m \in M} \rho_t m C_t^m}{\sum_{m \in M} \rho_t m}$$



Work Integrated Learning Programmes

Solving MDPs so far

Dynamic programming
 Off policy
 local learning, propagating values from neighbors
 (Bootstrapping)
 Model based

Monte-Carlo
 On-policy (though important sampling can be used)
 Requires a full episode to train on
 Model free, online learning

$y = 0.9$

-100 +10

$q^*(s, a) = \sum_{s'} p(s'|s, a) (r(s, a, s') + \gamma \max_a [q^*(s', a)])$

- $Q(z, \text{exit}) = 10$
- $Q(y, \rightarrow) = 0 + \gamma \max_a Q(z, a)$
- $Q(x, \rightarrow) = 0 + \gamma \max_a Q(y, a)$

• Episode = {x, y, z, exit}

- $Q(z, \text{exit}) = 10$
- $Q(y, \rightarrow) = 9$
- $Q(x, \rightarrow) = 8.1$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Fuse DP and MC

Dynamic programming
 Off policy
 local learning, propagating values from neighbors
 (Bootstrapping)
 Model based

Monte-Carlo
 On-policy (though important sampling can be used)
 Requires a full episode to train on
 Model free, online learning

TD Learning

Off policy
 local learning, propagating values from neighbors
 (Bootstrapping)
 Model free, online learning

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Temporal difference learning

In Monte-Carlo, G_t is an actual return from the complete episode. Now, if we replace G_t with an estimated return $R(t+1)+V(S_{t+1})$, this is what $TD(0)$ would look like:

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Where:

- $R(t+1)+V(S_{t+1})$ is called **TD target value**

- $R(t+1)+V(S_{t+1}) - V(S_t)$ is called **TD error**.

- MC uses accurate return G_t to update value, while TD uses the Bellman Optimality Equation to estimate value, and then updates the estimated value with the target value.

5

Temporal Difference learning

Tabular TD(0) for estimating v_π

Input: the policy π to be evaluated

Algorithm parameter: step size $\alpha \in (0, 1]$

Initialize $V(s)$, for all $s \in S^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

$A \leftarrow$ action given by π for S

 Take action A , observe R, S'

$V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

 until S is terminal



SARSA: On-policy TD Control

- learn an action-value function rather than a state-value function

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)].$$

- This update is done after every transition from a nonterminal state S_t . If

S_{t+1} is terminal, then $Q(S_{t+1}, A_{t+1})$ is defined as zero.

- Quintuple of events, $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$

7

SARSA: On-policy TD Control

Sarsa (on-policy TD control) for estimating $Q \approx q$.

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

Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

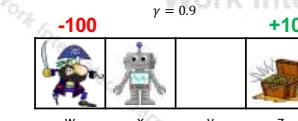
 Take action A , observe R, S'

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

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

$S \leftarrow S'; A \leftarrow A'$

 until S is terminal



8

ARSA: On-policy TD Control



Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

Initialize $Q(s, a)$, for all $s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

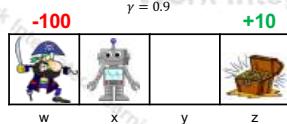
 Take action A , observe R, S'

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

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

$$S \leftarrow S'; A \leftarrow A';$$

 until S is terminal



$$Q = \begin{matrix} & 0 & 0,0 & 0,0 & 0 \\ w & & x & y & z \end{matrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

9

ARSA: On-policy TD Control



Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

Initialize $Q(s, a)$, for all $s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

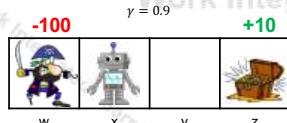
 Take action A , observe R, S'

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

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

$$S \leftarrow S'; A \leftarrow A';$$

 until S is terminal



$$Q = \begin{matrix} & 0 & 0,0 & 0,0 & 0 \\ w & & x & y & z \end{matrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

11

ARSA: On-policy TD Control



Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

Initialize $Q(s, a)$, for all $s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

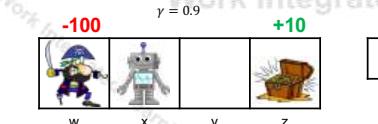
 Take action A , observe R, S'

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

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

$$S \leftarrow S'; A \leftarrow A';$$

 until S is terminal



$$Q = \begin{matrix} & 0 & 0,0 & 0,0 & 0 \\ w & & x & y & z \end{matrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

10

ARSA: On-policy TD Control



Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

Initialize $Q(s, a)$, for all $s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

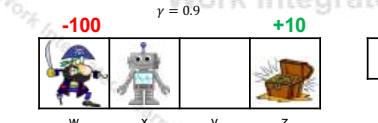
 Take action A , observe R, S'

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

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

$$S \leftarrow S'; A \leftarrow A';$$

 until S is terminal



$$Q = \begin{matrix} & 0 & 0,0 & 0,0 & 0 \\ w & & x & y & z \end{matrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

12

ARSA: On-policy TD Control



Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

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

Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

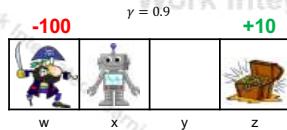
 Take action A , observe R, S'

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

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

$S \leftarrow S'; A \leftarrow A'$;

 until S is terminal



w	exit	0	w	exit
S	A	R	S'	A'

$$Q = \begin{bmatrix} 0 & 0,0 & 0,0 & 0 \\ w & x & y & z \end{bmatrix}$$

13 BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

ARSA: On-policy TD Control



Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

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

Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

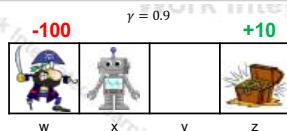
 Take action A , observe R, S'

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

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

$S \leftarrow S'; A \leftarrow A'$;

 until S is terminal



w	exit	-100	ter	.
S	A	R	S'	A'

$$Q = \begin{bmatrix} -100 & 0,0 & 0,0 & 0 \\ w & x & y & z \end{bmatrix}$$

15 BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

ARSA: On-policy TD Control



Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

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

Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

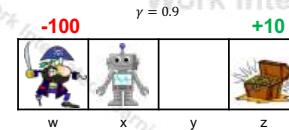
 Take action A , observe R, S'

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

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

$S \leftarrow S'; A \leftarrow A'$;

 until S is terminal



w	exit	-100	ter	exit
S	A	R	S'	A'

$$Q = \begin{bmatrix} 0 & 0,0 & 0,0 & 0 \\ w & x & y & z \end{bmatrix}$$

14 BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

ARSA: On-policy TD Control



Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

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

Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

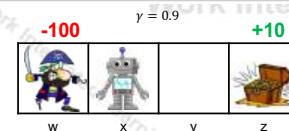
 Take action A , observe R, S'

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

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

$S \leftarrow S'; A \leftarrow A'$;

 until S is terminal



x	←	-100	ter	.
S	A	R	S'	A'

$$Q = \begin{bmatrix} -100 & 0,0 & 0,0 & 0 \\ w & x & y & z \end{bmatrix}$$

16 BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

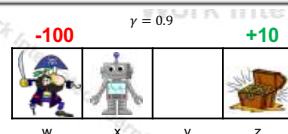
ARSA: On-policy TD Control

Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

Initialize $Q(s, a)$, for all $s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):

 Initialize S
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Repeat (for each step of episode):

 Take action A , observe R, S'
 Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)
 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$
 $S \leftarrow S'; A \leftarrow A'$;
 until S is terminal



$y = 0.9$

$+10$

x	←	0	w	exit
S	A	R	S'	A'

$$Q = \begin{matrix} -100 & 0,0 & 0,0 & 0 \\ w & x & y & z \end{matrix}$$

17 BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

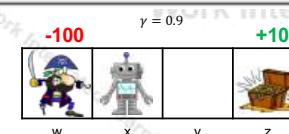
ARSA: On-policy TD Control

Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

Initialize $Q(s, a)$, for all $s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):

 Initialize S
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Repeat (for each step of episode):

 Take action A , observe R, S'
 Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)
 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$
 $S \leftarrow S'; A \leftarrow A'$;
 until S is terminal



$y = 0.9$

$+10$

x	←	0	w	exit
S	A	R	S'	A'

$$Q = \begin{matrix} -100 & -90,0 & 0,0 & 0 \\ w & x & y & z \end{matrix}$$

18 BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

ARSA: On-policy TD Control

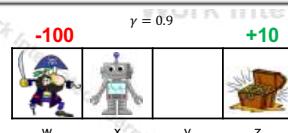
Sarsa (on-policy TD control) for estimating $Q \approx q_*$.

Initialize $Q(s, a)$, for all $s \in S, a \in A(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):

 Initialize S
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Repeat (for each step of episode):

 Take action A , observe R, S'
 Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)
 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$
 $S \leftarrow S'; A \leftarrow A'$;
 until S is terminal

And so on...



$y = 0.9$

$+10$

w	exit	0	w	exit
S	A	R	S'	A'

$$Q = \begin{matrix} -100 & -90,0 & 0,0 & 0 \\ w & x & y & z \end{matrix}$$

19 BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Q-learning: Off-policy TD Control

- Use the original TD update rule

$$\bullet \quad Q(s, a) = Q(s, a) + \alpha \left(R_{t+1} + \gamma \max_{a'} [Q(s', a')] - Q(s, a) \right)$$

- Approximates the state-action value for the optimal policy, i.e., q^*
 - Assuming that every state-action pair is visited infinitely often
- Follows from the proof of convergence for the Bellman function
 - See slides MDPs+DP

20 BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Q-learning: Off-policy TD Control



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

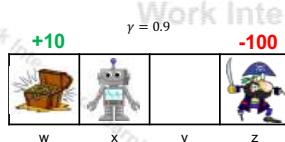
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



$$Q = \begin{matrix} & & & & \\ & 0 & 0,0 & 0,0 & 0 \\ \hline w & & x & y & z \end{matrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

21

Q-learning: Off-policy TD Control

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

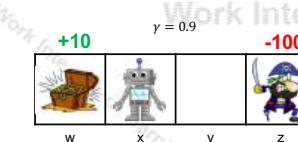
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



$$Q = \begin{matrix} & & & & \\ & 0 & 0,0 & 0,0 & 0 \\ \hline w & x & y & z \end{matrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

22

Q-learning: Off-policy TD Control



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

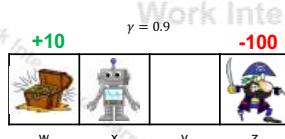
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 Take action A , observe R, S'

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

$$S \leftarrow S'$$

 until S is terminal



$$Q = \begin{matrix} & & & & \\ & 0 & 0,0 & 0,0 & 0 \\ \hline w & x & y & z \end{matrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

23

Q-learning: Off-policy TD Control

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

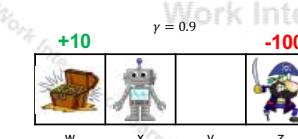
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



$$Q = \begin{matrix} & & & & \\ & 0 & 0,0 & 0,0 & 0 \\ \hline w & x & y & z \end{matrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

24

Q-learning: Off-policy TD Control



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

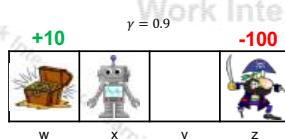
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



$$Q = \begin{bmatrix} 0 & 0,0 & 0,0 & 0 \\ w & x & y & z \end{bmatrix}$$

$\gamma = 0.9$

$\gamma = 0.9$

S	A	R	S'
w	x	y	z

y

z

w

x

y

z

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

25

Q-learning: Off-policy TD Control

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

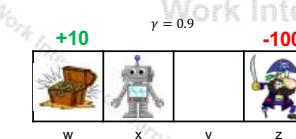
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



$\gamma = 0.9$

$\gamma = 0.9$

S	A	R	S'
w	x	y	z

z

\rightarrow

0

z

w

x

y

w	x	y	z
0	0,0	0,0	0

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

26

Q-learning: Off-policy TD Control



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

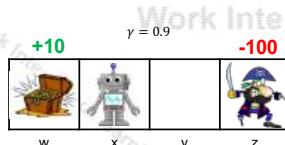
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



$\gamma = 0.9$

$\gamma = 0.9$

S	A	R	S'
w	x	y	z

z

\rightarrow

-100

$.$

w	x	y	z
0	0,0	0,0	0

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

27

Q-learning: Off-policy TD Control

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

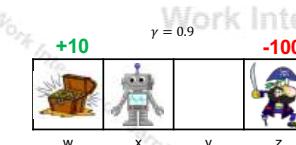
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



$\gamma = 0.9$

$\gamma = 0.9$

S	A	R	S'
w	x	y	z

z

\rightarrow

-100

$.$

w	x	y	z
0	0,0	0,0	-100

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

28

Q-learning: Off-policy TD Control



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

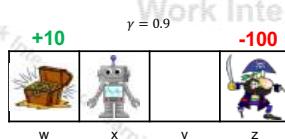
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



x	\rightarrow	0	z
S	A	R	S'

$$Q = \begin{bmatrix} 0 & 0,0 & 0,-90 & -100 \\ w & x & y & z \end{bmatrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

33

Q-learning: Off-policy TD Control

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

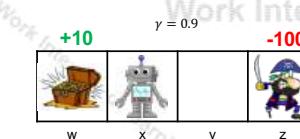
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



x	\rightarrow	0	y
S	A	R	S'

$$Q = \begin{bmatrix} 0 & 0,0 & 0,-90 & -100 \\ w & x & y & z \end{bmatrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

34

Q-learning: Off-policy TD Control



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

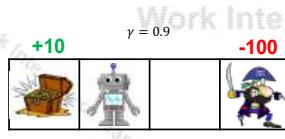
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



x	\rightarrow	0	y
S	A	R	S'

$$Q = \begin{bmatrix} 0 & 0,0 & 0,-90 & -100 \\ w & x & y & z \end{bmatrix}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

35

Q-learning: Off-policy TD Control



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

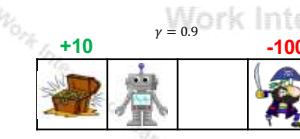
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 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



x	\rightarrow	0	y
S	A	R	S'

$$Q = \begin{bmatrix} 0 & 0,0 & 0,-90 & -100 \\ w & x & y & z \end{bmatrix}$$

And so on...

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

36

Expected SARSA

- Instead of the maximum over next state-action pairs it uses the expected value, taking into account how likely each action is under the current policy

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \mathbb{E}_\pi [Q(S_{t+1}, A_{t+1}) | S_{t+1}] - Q(S_t, A_t)] \\ &= Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \sum \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t)], \end{aligned}$$

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

n-step TD Prediction

One-step return:

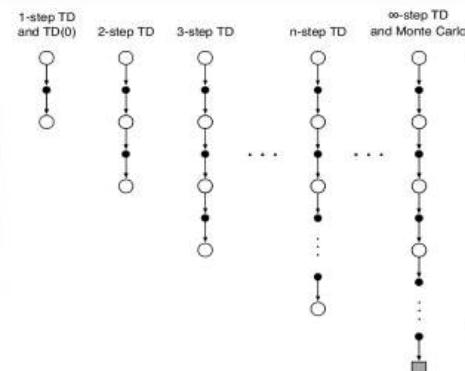
$$G_{t:t+1} = R_{t+1} + \gamma V_t(S_{t+1})$$

Two-step return:

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

Work Integrated Le

Ref Section 7.1 of TB



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Maximization Bias & Double Learning

- Maximum over estimated values is used implicitly as an estimate of the maximum value, which can lead to a significant positive bias.
- Solution – Learn two estimates – Double Learning
- Double learning doubles the memory requirements, but does not increase the amount of computation per step.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

n-step TD Prediction

One-step return:

$$G_{t:t+1} = R_{t+1} + \gamma V_t(S_{t+1})$$

Two-step return:

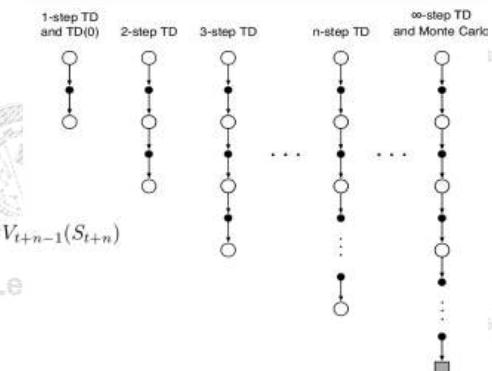
$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

n-step return:

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

Work Integrated Le

Ref Section 7.1 of TB



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

n-step TD Prediction

One-step return:

$$G_{t:t+1} \doteq R_{t+1} + \gamma V_t(S_{t+1})$$

Two-step return:

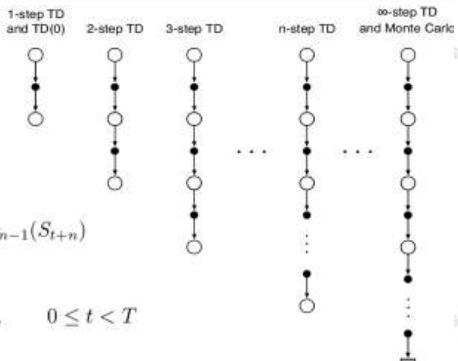
$$G_{t:t+2} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

n-step return:

$$G_{t:t+n} \doteq R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

State-value learning algorithm for using n-step returns:

$$V_{t+n}(S_t) \doteq V_{t+n-1}(S_t) + \alpha [G_{t:t+n} - V_{t+n-1}(S_t)], \quad 0 \leq t < T$$



Ref Section 7.1 of TB

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Deep Reinforcement Learning
Prof. Vimal SP
CSIS

BITS Pilani

n-step TD Prediction

n-step TD for estimating $V \approx v_\pi$

Input: a policy π
Algorithm parameters: step size $\alpha \in (0, 1]$, a positive integer n
Initialize $V(s)$ arbitrarily, for all $s \in \mathcal{S}$
All store and access operations (for S_t and R_t) can take their index mod $n + 1$

```

Loop for each episode:
  Initialize and store  $S_0 \neq$  terminal
   $T \leftarrow \infty$ 
  Loop for  $t = 0, 1, 2, \dots$ :
    If  $t < T$ , then:
      Take an action according to  $\pi(\cdot | S_t)$ 
      Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$ 
      If  $S_{t+1}$  is terminal, then  $T \leftarrow t + 1$ 
       $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose state's estimate is being updated)
      If  $\tau \geq 0$ :
         $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
        If  $\tau + n < T$ , then:  $G \leftarrow G + \gamma^n V(S_{\tau+n})$ 
         $V(S_\tau) \leftarrow V(S_\tau) + \alpha [G - V(S_\tau)]$   $(G_{\tau:n})$ 
    Until  $\tau = T - 1$ 
  
```

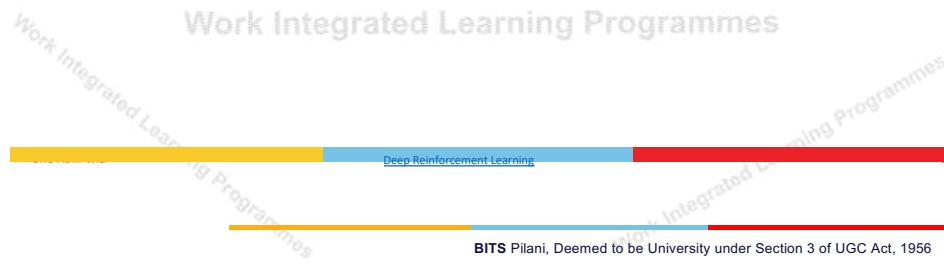
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

BITS Pilani

<AIMLCZG512, Deep Reinforcement Learning>
Lecture No. 4,5,6,7,8,9 - Incremental example

- 1 [MDP](#)
- 2 [DP](#)
- 3 [MC](#)
- 4 [TD](#)

Work Integrated Learning Programmes

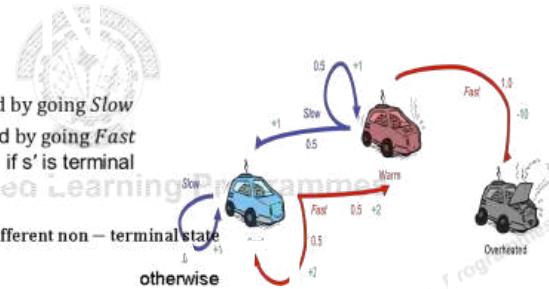


1 States (S)

2 Actions (A)

3 Rewards (R)

4 State Transition Probabilities (P)

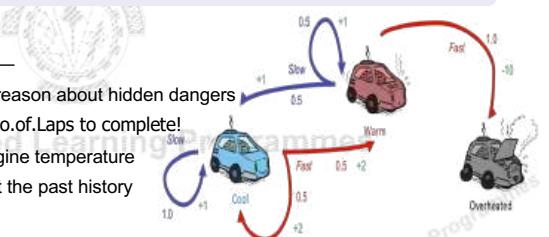


Race Car Laps

A race car completes laps while managing engine temperature and speed decisions. Here each step is a segment of a lap in the track where the car learns to trade off lap time against overheating risk. The more the lap the car completes its performance better. Car shuts down only when its engine gets overheated (optional terminal state).

Key Characteristics:

- The agent must act under **uncertainty** — percepts are local, and the agent must reason about hidden dangers
- Assume there is no fuel constraints or No. of Laps to complete!
- At every time step the next potential engine temperature depends on the current engine state not the past history



1 MDP

2 DP

3 MC

4 TD



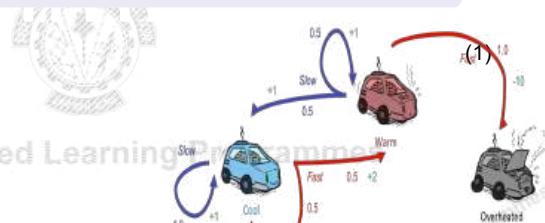
Work Integrated Learning Programmes



Scenario

In a modified Race car, there is a training track where agent is given a known map and risk factors(dynamics) in advance. Agent knows exactly where it is and its action how actions result in engine state change (KNOWN MODEL).

With full knowledge, the agent can compute optimal policies by trying through all states.



Work Integrated Learning Programmes

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Dynamic Programming Solution



Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Given

- $P(s' | s, a)$, $R(s, a, s')$, $P(s', r | s, a)$

To find

- Policy π

Solution

- Policy Iteration (PI): Prefer this in a relatively smaller state space where intermediate termination provides explicit interpretable policy
- Value Iteration (VI): Prefer this in a relatively larger state space where fast convergence to optimal policy is obtained at the end

Idea

- PI : One Step Look ahead □ Evaluate □ Improve
- VI : One Step Look ahead □ Control

VI can be thought of truncated PI where evaluation is done fully before improvement.



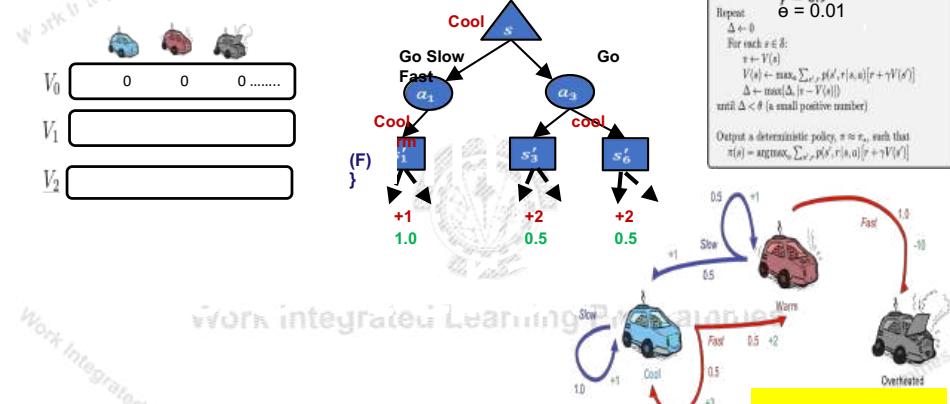
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Value Iteration

```

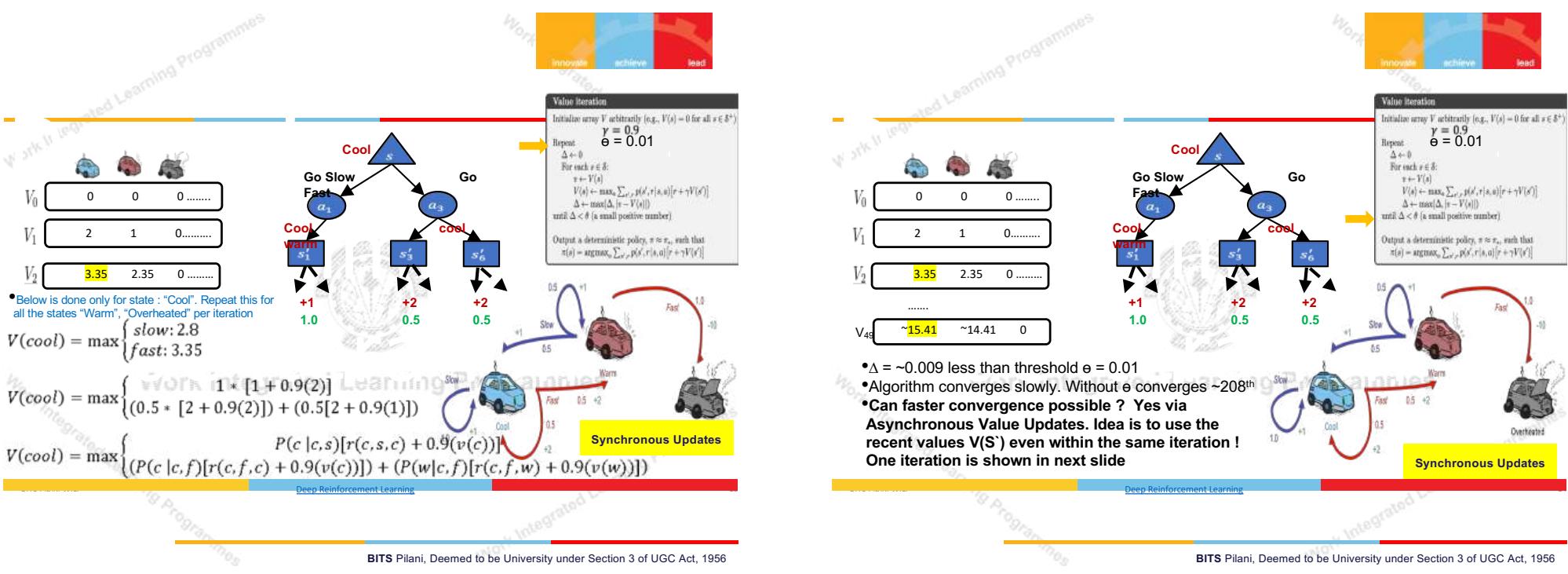
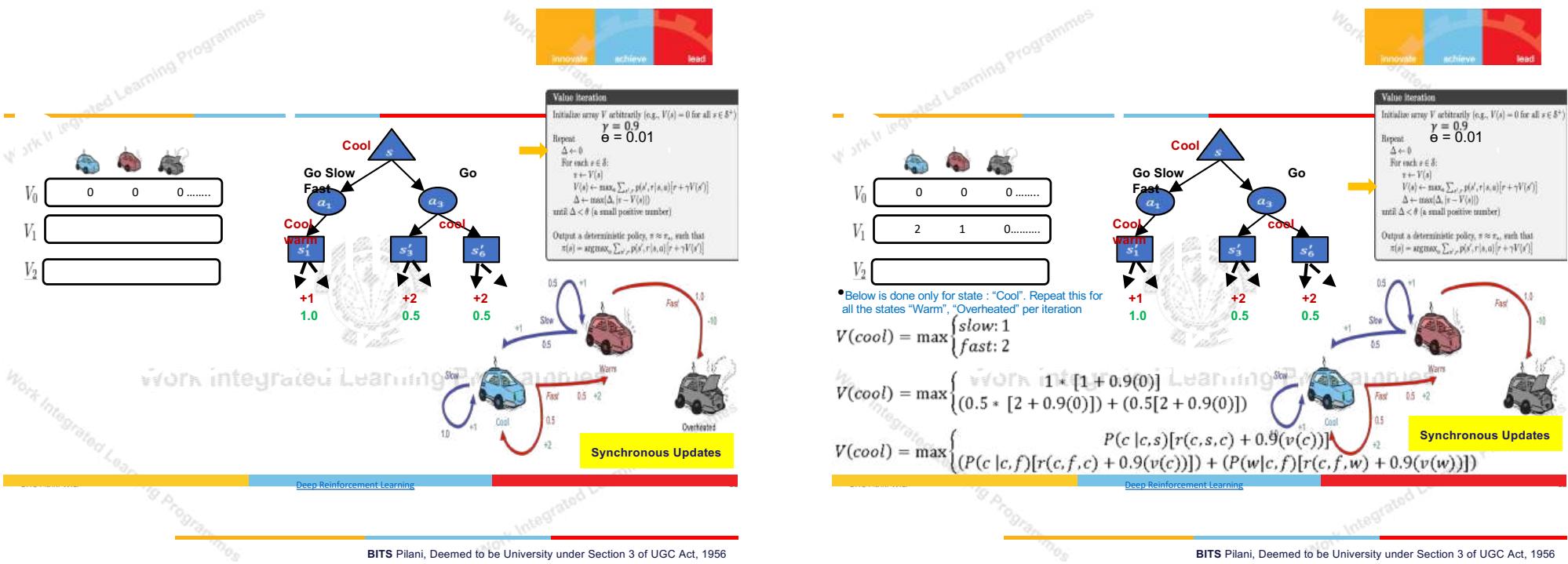
Initialize array  $V$  arbitrarily (e.g.,  $V(s) = 0$  for all  $s \in \mathcal{S}^*$ )
 $\gamma = 0.9$ 
 $\epsilon = 0.01$ 
Repeat
   $\Delta \leftarrow 0$ 
  For each  $s \in \mathcal{S}$ :
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \max_a \sum_{s'} p(s', r | s, a) [r + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
  until  $\Delta < \epsilon$  (a small positive number)
Output a deterministic policy,  $\pi \approx \pi_\star$ , such that
 $\pi(s) = \arg\max_a \sum_{s'} p(s', r | s, a) [r + \gamma V(s')]$ 
  
```



Synchronous Updates



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



	0	0	0.....
V_0	2		
V_1			

- Below is done only for state : "Cool". In this same iteration while updating the $V(\text{Warm})$ use recent value of $V(\text{Cool}) = 2$ not 0.
- $V(\text{cool}) = \max\{\text{slow: } 2, \text{fast: } 2\}$

$$V(\text{cool}) = \max \left\{ \begin{array}{l} 1 * [1 + 0.9(0)] \\ (0.5 * [2 + 0.9(0)]) + (0.5[2 + 0.9(0)]) \end{array} \right\}$$

$$V(\text{cool}) = \max \left\{ \begin{array}{l} P(c|c,s)[r(c,s,c) + 0.9(v(c))] \\ (P(c|c,f)[r(c,f,c) + 0.9(v(c))]) + (P(w|c,f)[r(c,f,w) + 0.9(v(w))]) \end{array} \right\}$$

Asynchronous Updates

Deep Reinforcement Learning

Value Iteration

```

Initialize array  $V$  arbitrarily (e.g.,  $V(s) = 0$  for all  $s \in \mathcal{S}^*$ )
 $\gamma = 0.9$ 
 $\epsilon = 0.01$ 
Repeat
   $\Delta \leftarrow 0$ 
  For each  $s \in \mathcal{S}$ :
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s', r | s, a) [r + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
  until  $\Delta < \epsilon$  (a small positive number)
Output a deterministic policy,  $\pi \approx \pi_\epsilon$ , such that
 $\pi(s) = \operatorname{argmax}_a \sum_{s' \in \mathcal{S}} p(s', r | s, a) [r + \gamma V(s')]$ 

```



Value Iteration

```

Initialize array  $V$  arbitrarily (e.g.,  $V(s) = 0$  for all  $s \in \mathcal{S}^*$ )
 $\gamma = 0.9$ 
 $\epsilon = 0.01$ 
Repeat
   $\Delta \leftarrow 0$ 
  For each  $s \in \mathcal{S}$ :
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s', r | s, a) [r + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
  until  $\Delta < \epsilon$  (a small positive number)
Output a deterministic policy,  $\pi \approx \pi_\epsilon$ , such that
 $\pi(s) = \operatorname{argmax}_a \sum_{s' \in \mathcal{S}} p(s', r | s, a) [r + \gamma V(s')]$ 

```

	0	0	0.....
V_0	2	1.9	0.....
V_1	0.....		
V_2	3.75	3.54	0.....
V_{40}	~15.44	~14.44	0

- $\Delta = \sim 0.008$ less than threshold $\epsilon = 0.01$
- Algorithm converges fast. Without ϵ converges at $\sim 157^{\text{th}}$ iteration!

Deep Reinforcement Learning

	0	0	0.....
V_0	2	1.9	0
V_1			

- Below is done only for state : "Cool". In this same iteration while updating the $V(\text{Warm})$ use recent value of $V(\text{Cool}) = 2$.
- $V(\text{warm}) = \max\{\text{fast: } -10, \text{slow: } 1.9\}$

$$V(\text{warm}) = \max \left\{ \begin{array}{l} 1 * [-10 + 0.9(0)] \\ (0.5 * [1 + 0.9(2)]) + (0.5[1 + 0.9(0)]) \end{array} \right\}$$

$$V(\text{warm}) = \max \left\{ \begin{array}{l} P(o|w,f)[r(w,f,o) + 0.9(v(o))] \\ (P(c|w,s)[r(w,s,c) + 0.9(v(c))]) + (P(w|w,s)[r(w,s,w) + 0.9(v(w))]) \end{array} \right\}$$

Asynchronous Updates

Deep Reinforcement Learning



Value Iteration

```

Initialize array  $V$  arbitrarily (e.g.,  $V(s) = 0$  for all  $s \in \mathcal{S}^*$ )
 $\gamma = 0.9$ 
 $\epsilon = 0.01$ 
Repeat
   $\Delta \leftarrow 0$ 
  For each  $s \in \mathcal{S}$ :
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s', r | s, a) [r + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
  until  $\Delta < \epsilon$  (a small positive number)
Output a deterministic policy,  $\pi \approx \pi_\epsilon$ , such that
 $\pi(s) = \operatorname{argmax}_a \sum_{s' \in \mathcal{S}} p(s', r | s, a) [r + \gamma V(s')]$ 

```

	0	0	0.....
V_0	2	1.9	0.....
V_1	0.....		
V_2	3.75	3.54	0.....
V_{40}	~15.44	~14.44	0

- After convergence extract the policy by choosing the action that lead to the converged maximum value here from V_{40} th iteration

State	Best Action
Cool	Fast
Warm	Slow
Overheated	-

Dynamic Programming Solution

In many cases, sometimes policies is observed to converge faster than value!

So why wait for entire value updates?



Work Integrated Learning Programmes

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} v_*$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

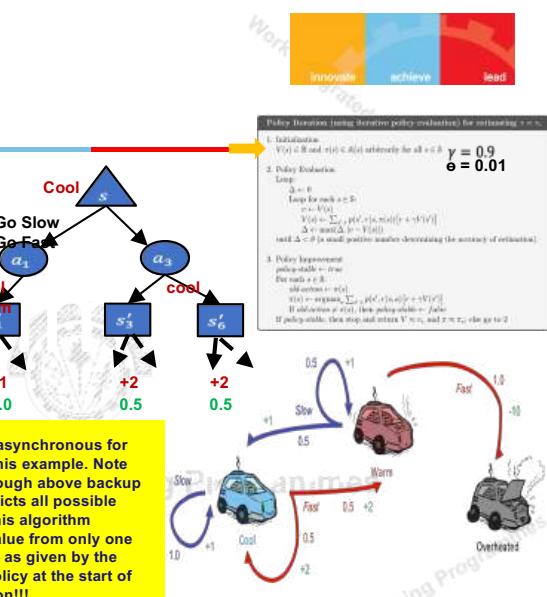


State	Best Action
Cool	Slow
Warm	Slow
Overheated	-

We will use asynchronous for updates in this example. Note that even though above backup diagram depicts all possible transition, this algorithm computes value from only one of the action as given by the initialized policy at the start of every iteration!!!

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Given

- $P(s' | s, a), R(s, a, s'), P(s', r | s, a)$

To find

- Policy π

Solution

- Policy Iteration (PI): Prefer this in a relatively smaller state space where intermediate termination provides explicit interpretable policy
- Value Iteration (VI): Prefer this in a relatively larger state space where fast convergence to optimal policy is obtained at the end

Idea

- PI : One Step Look ahead Evaluate Improve
- VI : One Step Look ahead Control

Number of state sweeps (updates) can be controlled if required (similar to mini-batch /stochastic processing). Evaluate & Improve steps are alternatively done.

VI can be thought of truncated PI where evaluation is done fully before improvement.

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Policy Evaluation Starts

V_0	0	0	0.....
V_1	1	1.45	0.....
V_2			

The value uses the initialized policy to select only one action GREEDILY. Hence below cancelled equations differentiated PI vs VI

$$V(\text{cool}) = \max \{ \text{slow}: 1, \text{fast}: 2 \}$$



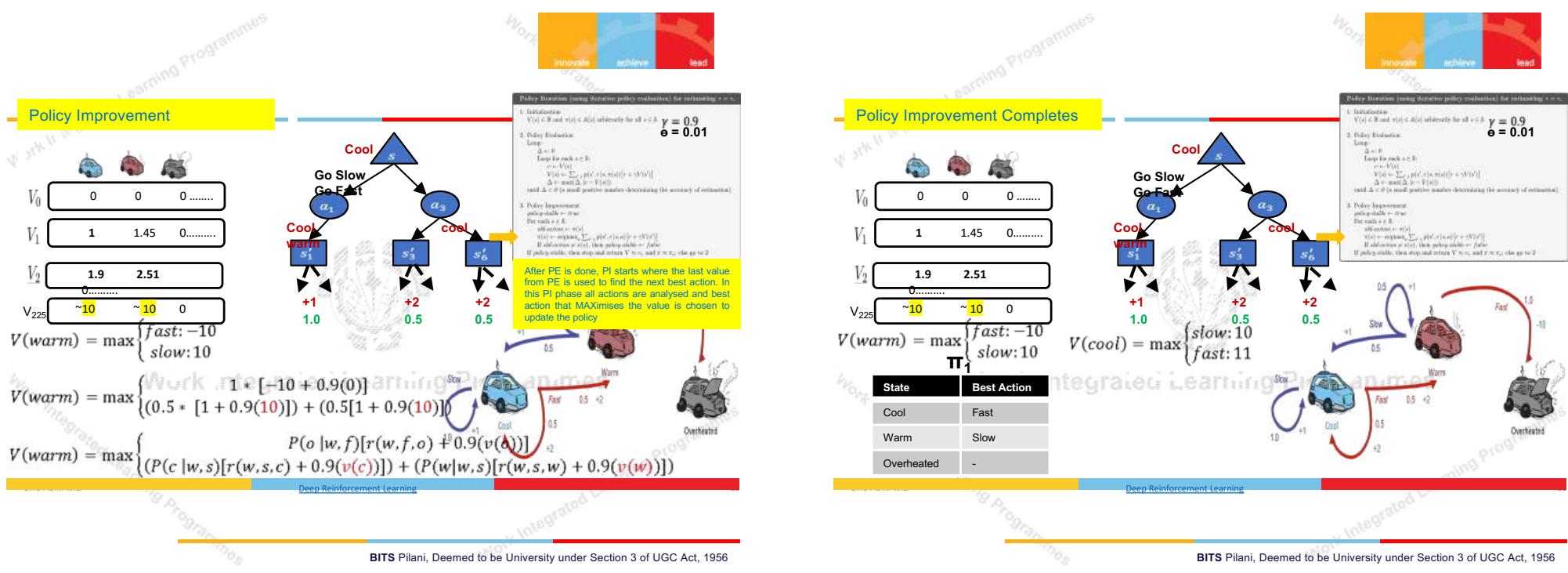
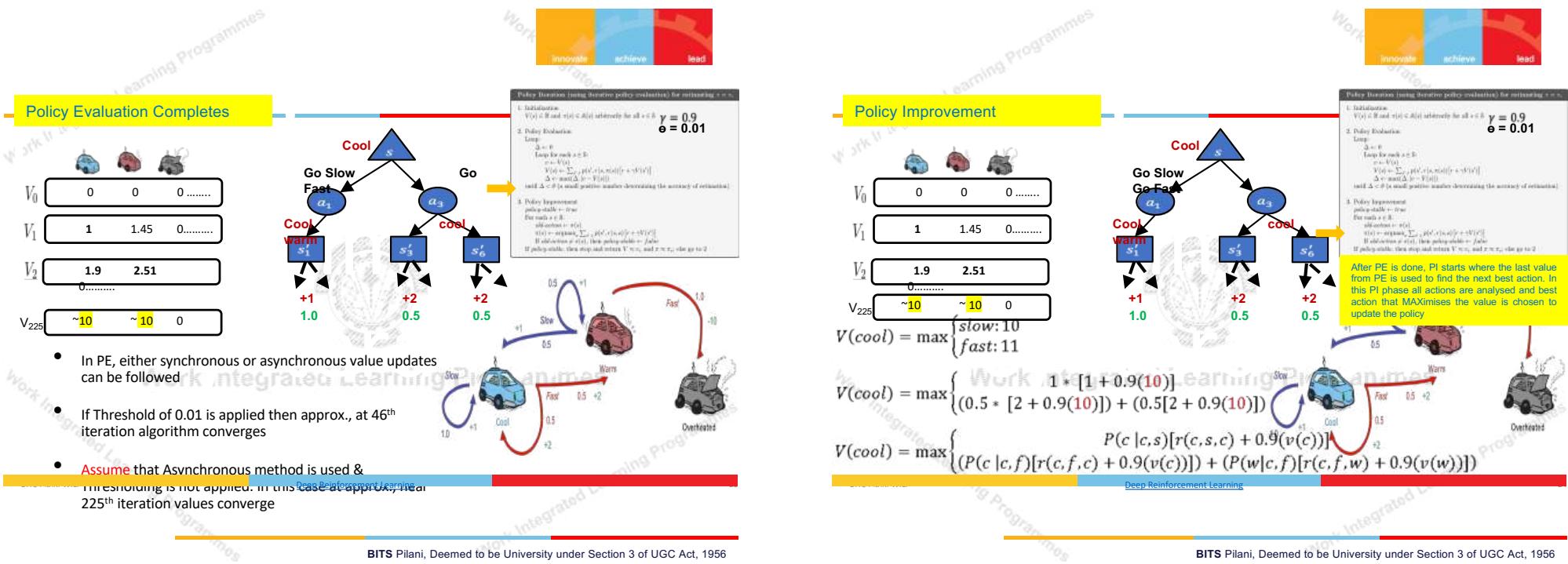
This slide show what change from Value iteration to Policy iteration algorithm in the Evaluation step. Since only action needs to be assessed, MAX of all actions outcome is no longer required in this PE cycle!

$$V(\text{cool}) = \max \{ 1 * [1 + 0.9(0)], (0.5 * [2 + 0.9(0)]) + (0.5 * [2 + 0.9(0)]) \}$$

$$V(\text{cool}) = \max \{ P(c | c, s)[r(c, s, c) + 0.9(v(c))] + (P(c | e, f)[r(e, f, c) + 0.9(v(e))] + P(w | e, f)[r(e, f, w) + 0.9(v(w))]) \}$$

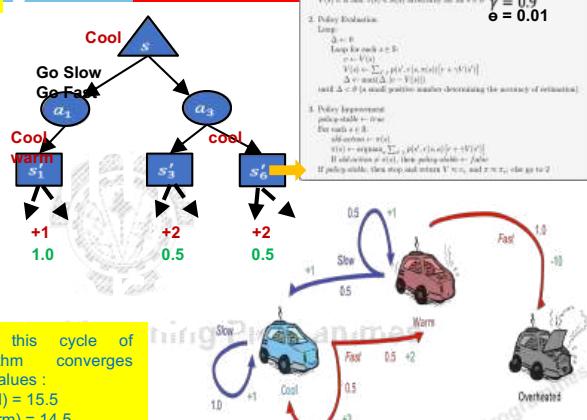
Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Continue another cycle of PE \square PI

V_0	10	10	0
V_1			
V_2			


TT1

State	Best Action
Cool	Fast
Warm	Slow
Overheated	-

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

1 Observation

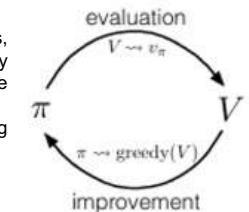
- Policy iteration consists of two simultaneous, interacting processes, one making the value function consistent with the current policy (policy evaluation), and the other making the policy greedy with respect to the current value function (policy improvement).
- In policy iteration, these two processes alternate, each completing before the other begins, **but this is not really necessary**.

Solution : GPI Generalized Policy Iteration

- Generalized policy iteration (GPI) refers to the general idea of letting policy-evaluation and policy improvement processes interact, independent of the granularity and other details of the two processes. Almost all reinforcement learning methods are well described as GPI.

Idea (Do this as an exercise for given problem)

Evaluation (ie., Value Updates) without waiting for convergence immediately follow by Policy Improvement.



Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Race Car Laps

In a modified Race car, there is a training track where agent is given a known (sometimes incomplete) map and risk factors(dynamics) may not be known in advance. Agent knows exactly where it is but its actions impact on the engine state change is difficult to predict (UNKNOWN MODEL).

With partial knowledge, the agent can compute optimal policies by trying through random experiments (Sampling). Car shuts down when its engine get overheated (optional terminal state) or Car stops once the predefined no.of.laps are completed.

Key Characteristics:

- The agent must act under **uncertainty** —
percepts are local, and the agent must reason about hidden dangers
- Assume there are fuel constraints or/and No.of.Laps to complete!
- At every time step there is **no way clear indication of Impact of next potential engine temperature depending on the current engine state. The entire past changes might lead to any possible state change. Hence immediate reward is not always available.**

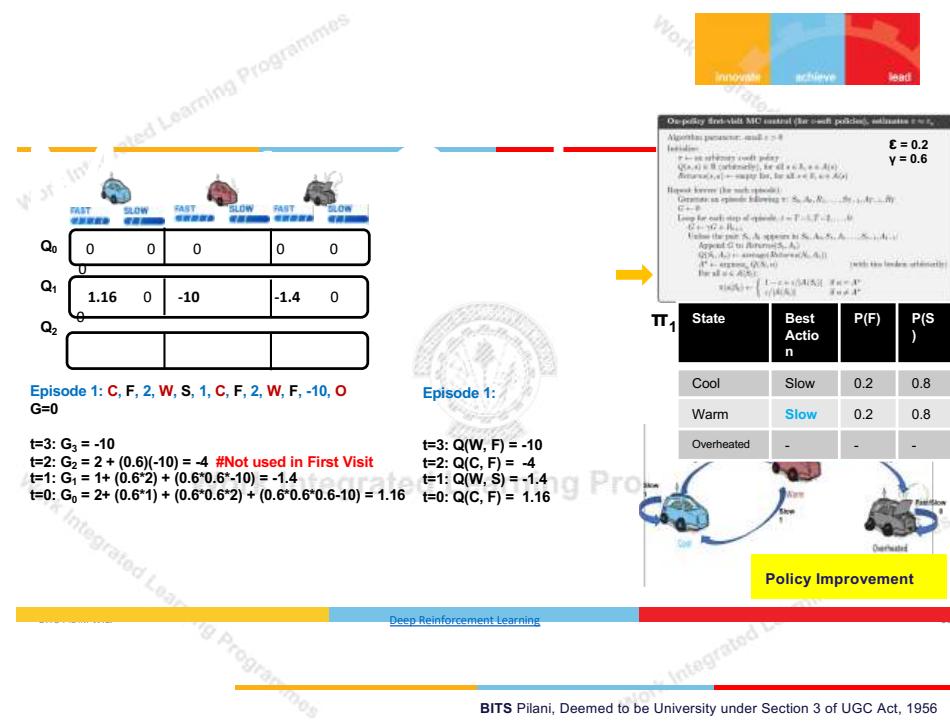
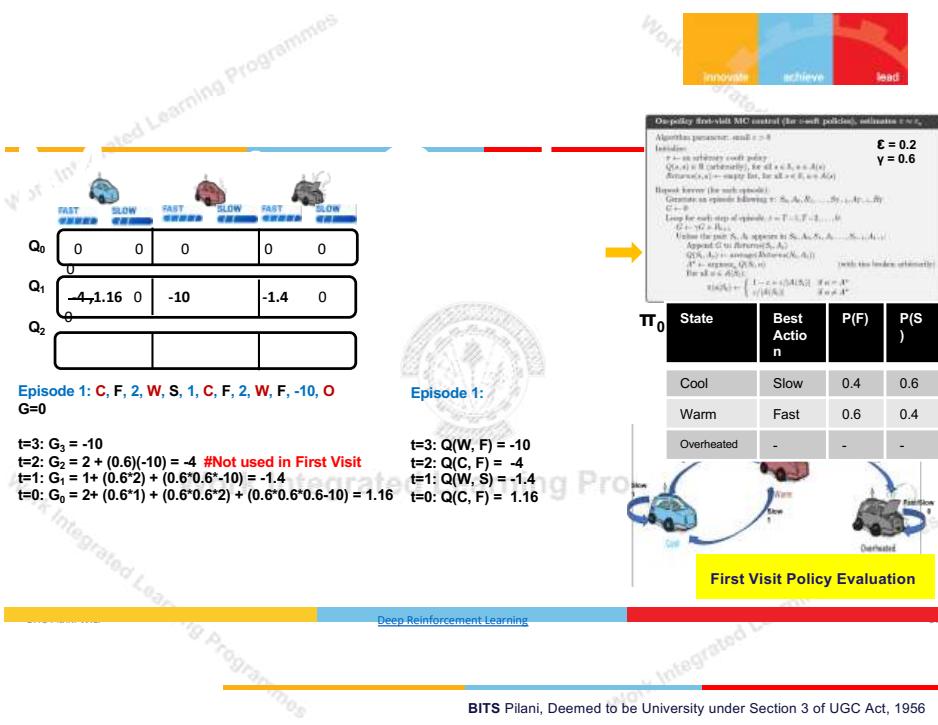
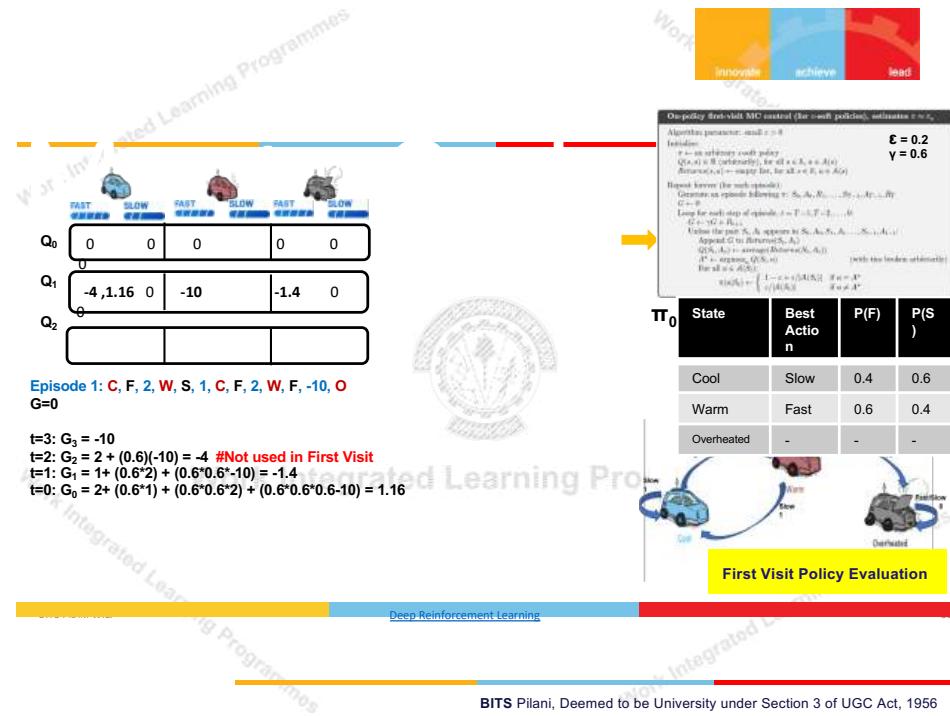
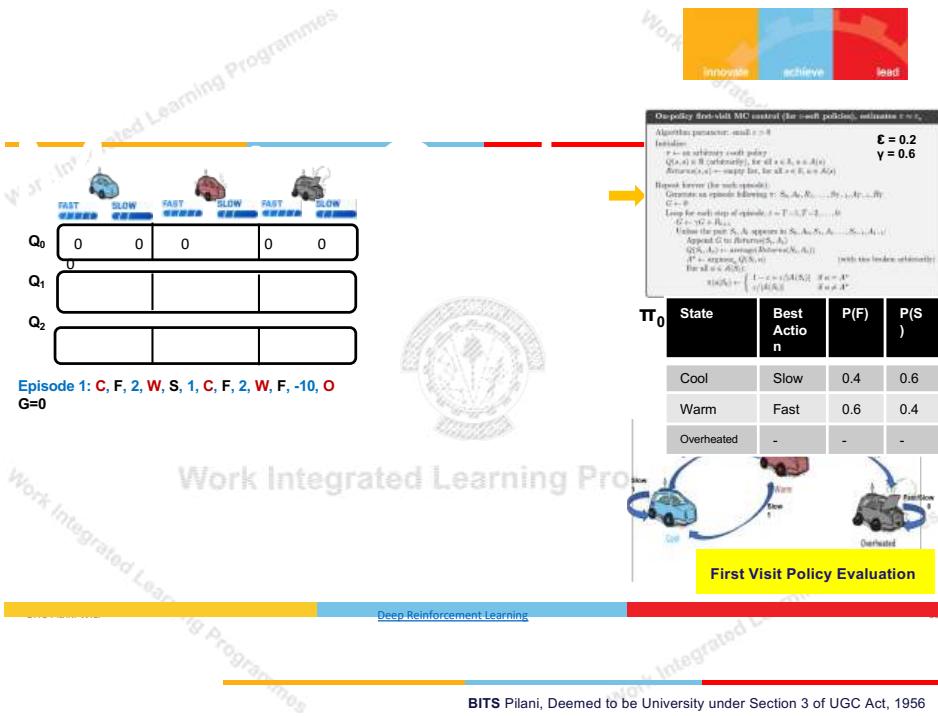


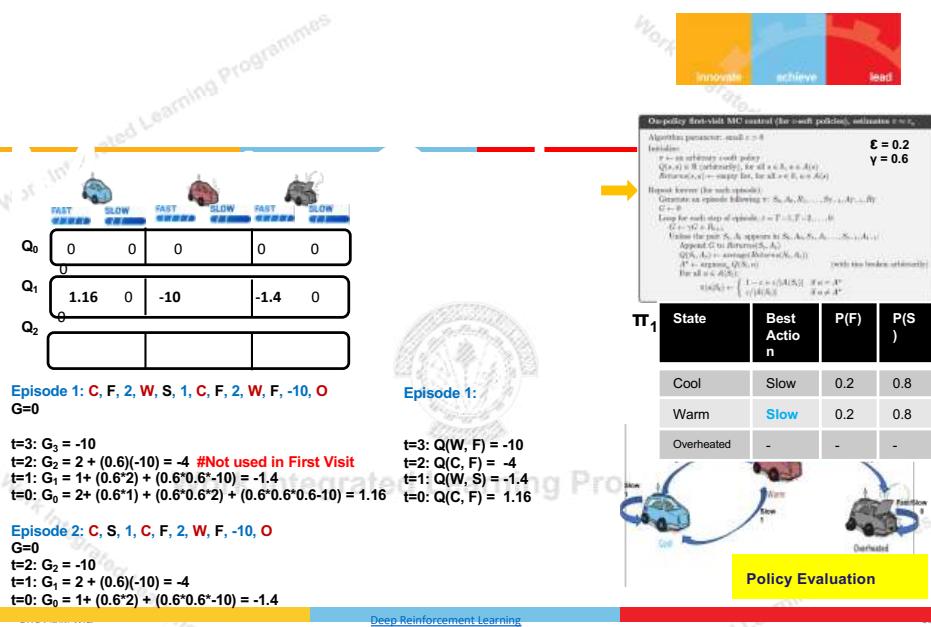
Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

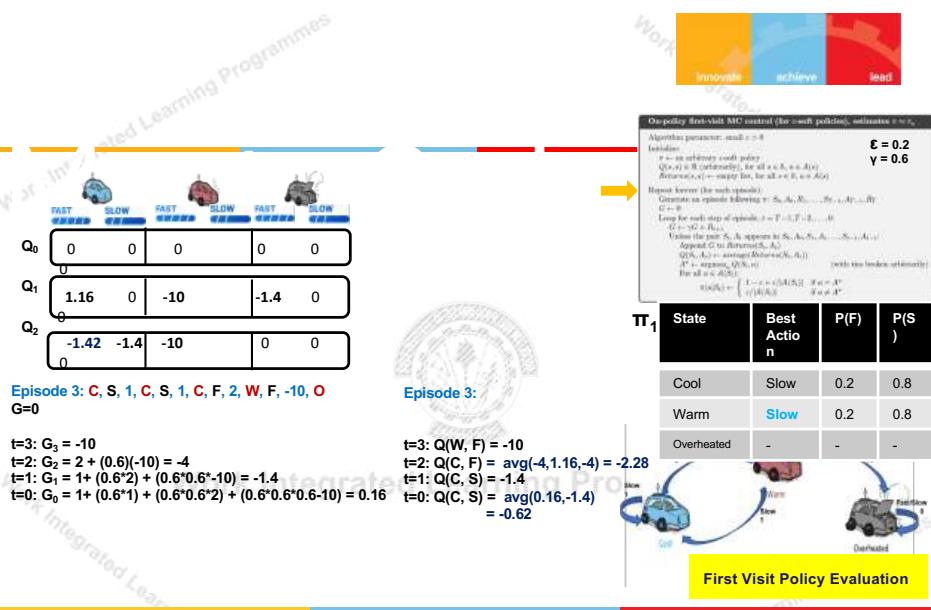
Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

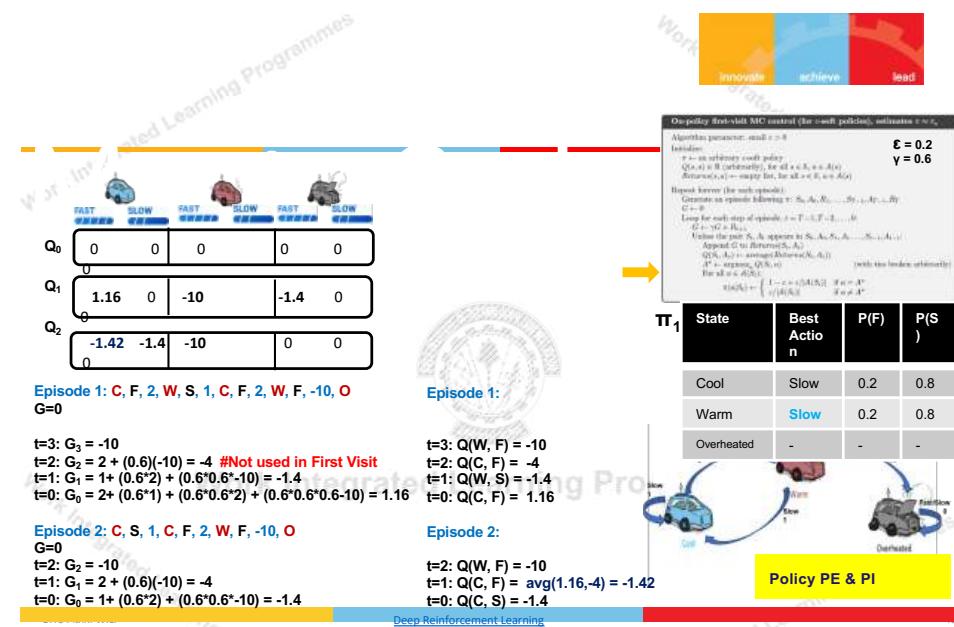




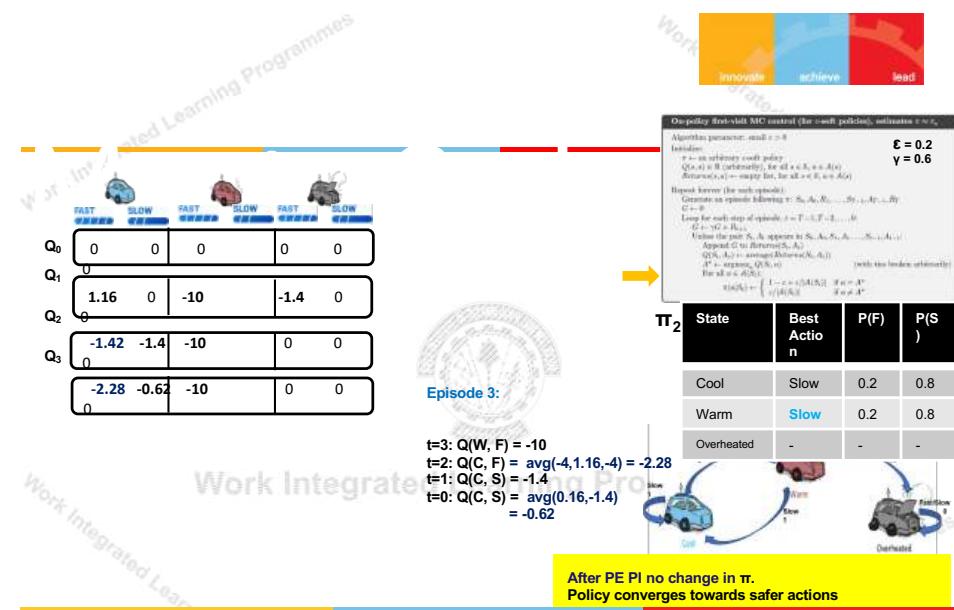
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Given

- S, A , optionally $R(s, a, s')$

To find

- Policy π . This is referred to as **TARGET** policy

Solution

- Sample Episodes
- Compute estimates of the value by discounted returns
 - Simple average
 - Weighted average
 - Normalized weighted average

Idea

- On Policy :
 - Generate Sample/Experience Evaluate using Return Improve
- Off Policy :
 - Refer to Sample/Experience generated by another reference(**BEHAVIOUR**) policy
 - Evaluate using Return Improve

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Off-policy MC control, for estimating $\pi \approx \pi_\epsilon$.

```

Initiation, for  $s \in S, a \in A(Q)$ :  

 $Q(s, a) \leftarrow \text{arbitrary}$   

 $C(s, a) \leftarrow 1$   

 $\pi(s, a) \leftarrow \text{argmax}_a Q(S, a)$  (with ties broken randomly)
  
```

Repeat forever

$\lambda = \text{any soft policy}$
 Generate an episode using λ :
 $s_0, a_0, R_1, \dots, s_{T-1}, a_{T-1}, R_T, s_T$
 $G \leftarrow 0$
 $W \leftarrow 1$
 For $t = T - 1, T - 2, \dots$ down to 0:
 $G \leftarrow \lambda G + R_t$
 $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{C(s_t, a_t)}{W} [G - Q(s_t, a_t)]$
 $Q(s_t, a_t) \leftarrow \text{argmax}_a Q(s_t, a_t) + \frac{C(s_t, a_t)}{W} [G - Q(s_t, a_t)]$
 If $a_t \neq \pi(s_t)$ then roll the loop
 $W \leftarrow W \cdot \frac{C(s_t, a_t)}{W}$

State	P(F)	P(S)
Cool	0.4	0.6
Warm	0.6	0.4
Overheated	-	-

b

Assume deterministic Target policy. Hence $\pi(a_t | s_t) = 1$
 But the behavior policy (b) is Stochastic as defined

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Monte Carlo Method

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Off-policy MC control, for estimating $\pi \approx \pi_\epsilon$.

```

Initiation, for  $s \in S, a \in A(Q)$ :  

 $Q(s, a) \leftarrow \text{arbitrary}$   

 $C(s, a) \leftarrow 1$   

 $\pi(s, a) \leftarrow \text{argmax}_a Q(S, a)$  (with ties broken randomly)
  
```

Repeat forever

$\lambda = \text{any soft policy}$
 Generate an episode using λ :
 $s_0, a_0, R_1, \dots, s_{T-1}, a_{T-1}, R_T, s_T$
 $G \leftarrow 0$
 $W \leftarrow 1$
 For $t = T - 1, T - 2, \dots$ down to 0:
 $G \leftarrow \lambda G + R_t$
 $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + W$
 $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{C(s_t, a_t)}{W} [G - Q(s_t, a_t)]$
 $Q(s_t, a_t) \leftarrow \text{argmax}_a Q(s_t, a_t) + \frac{C(s_t, a_t)}{W} [G - Q(s_t, a_t)]$
 If $a_t \neq \pi(s_t)$ then roll the loop
 $W \leftarrow W \cdot \frac{C(s_t, a_t)}{W}$

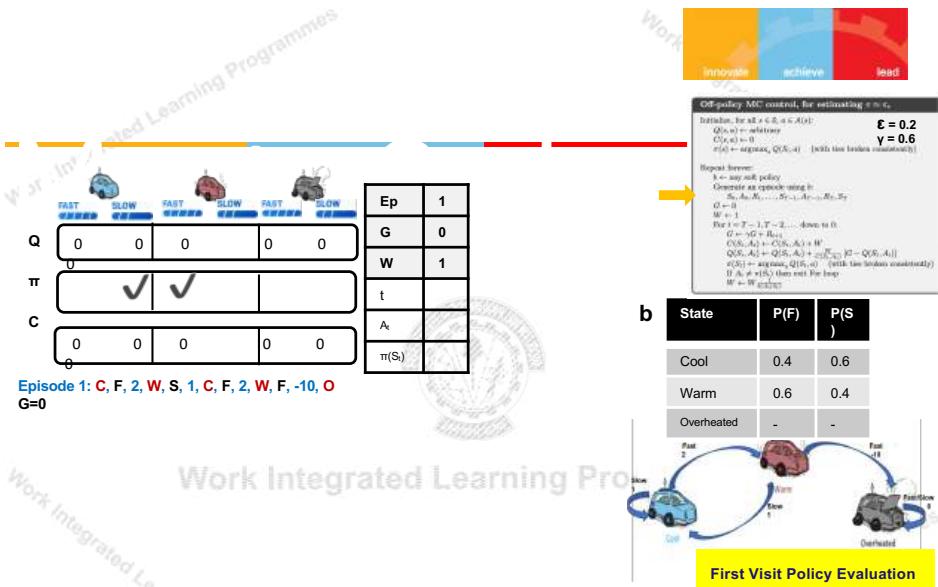
Ep	G	W
Q	0 0 0 0 0	0
π	✓ ✓	0
t		0
A_t		0
$\pi(s_t)$		0

b

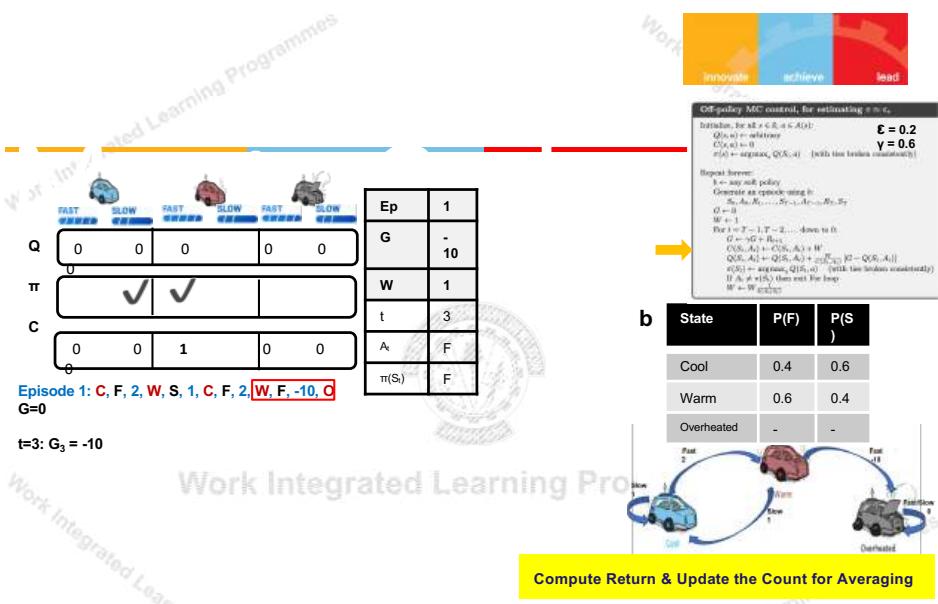
Every Visit Policy Evaluation

Deep Reinforcement Learning

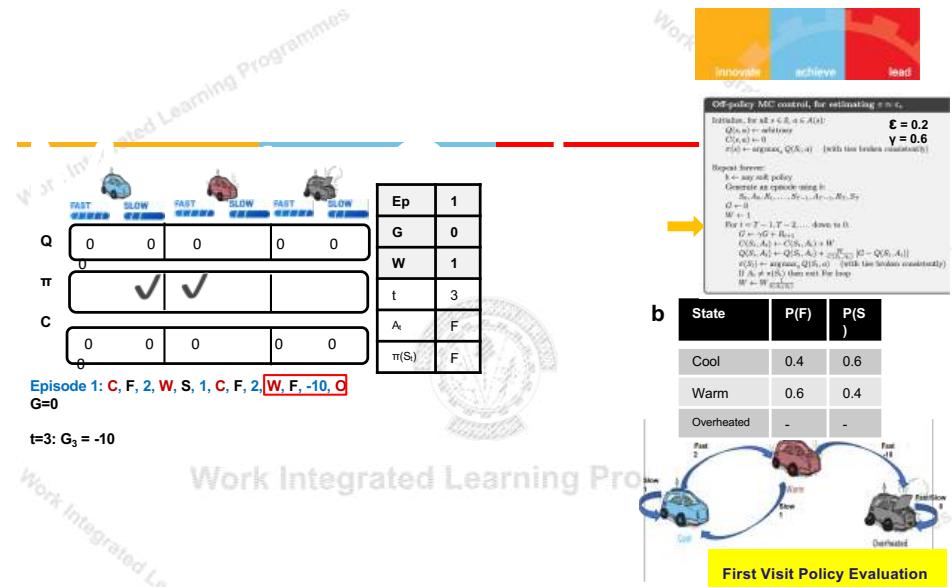
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



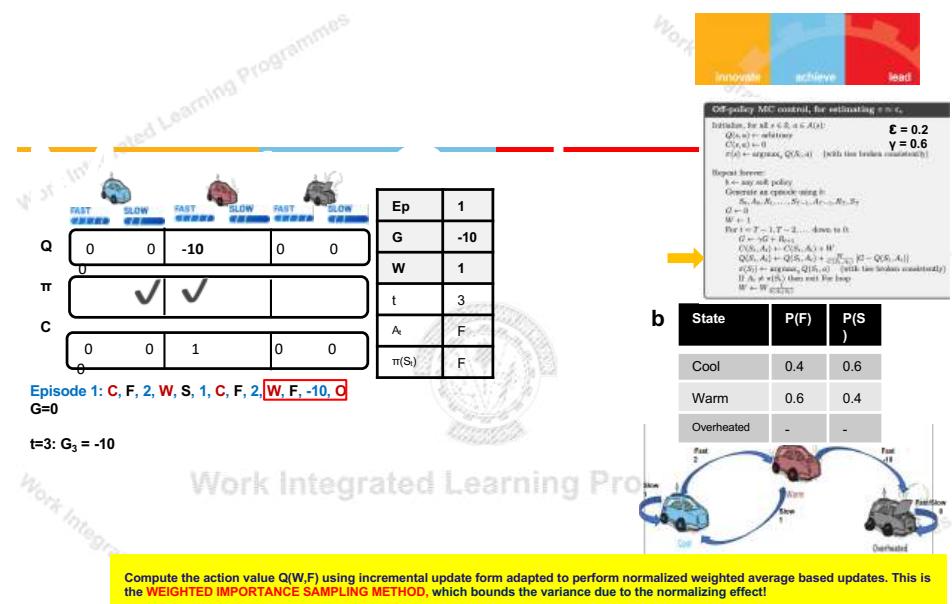
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



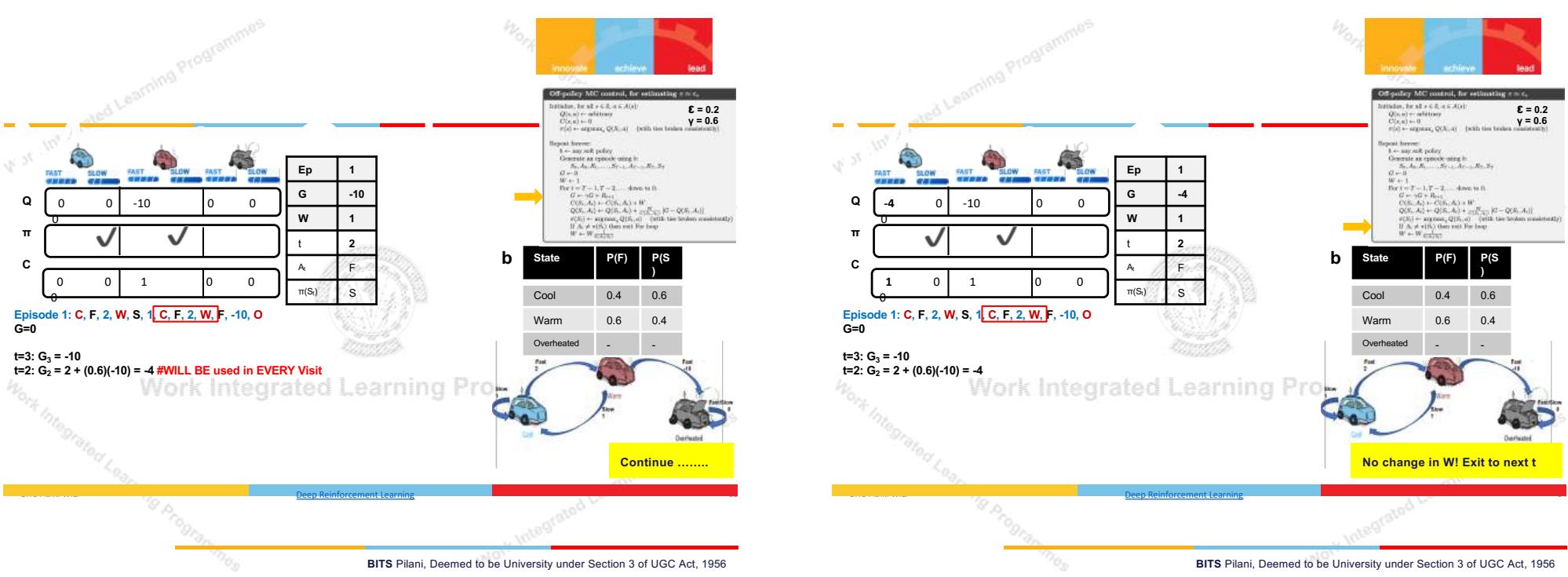
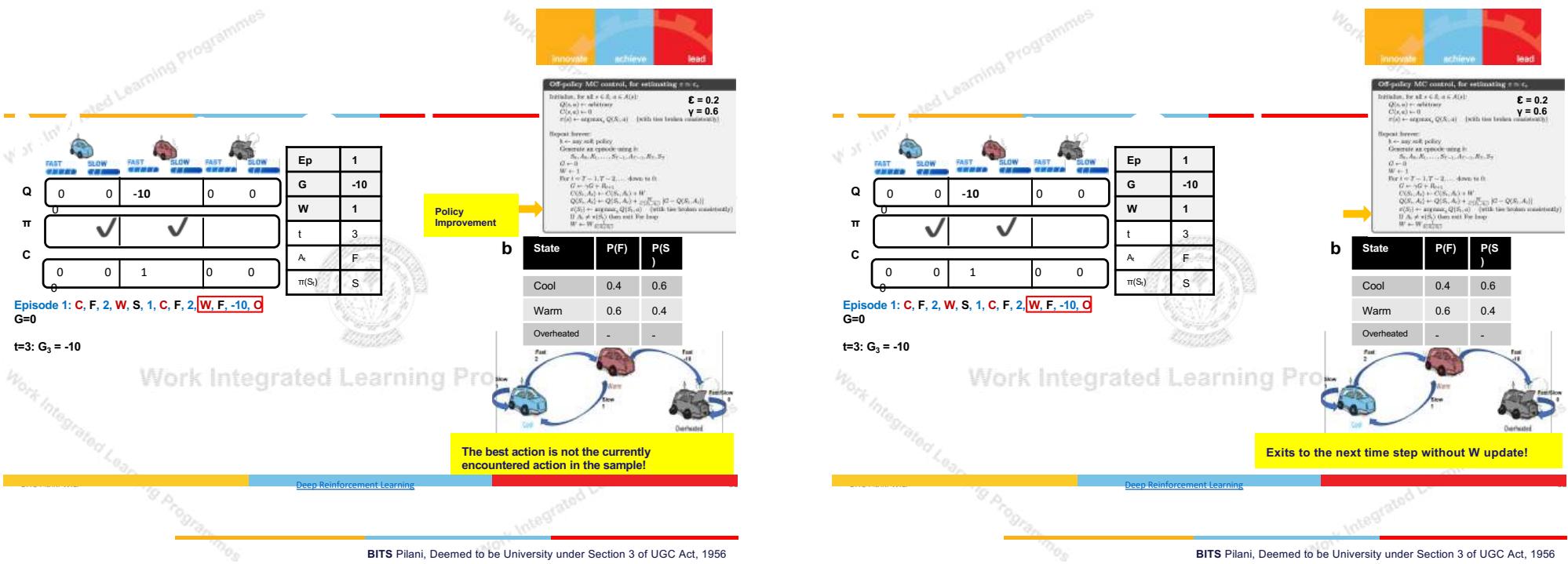
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

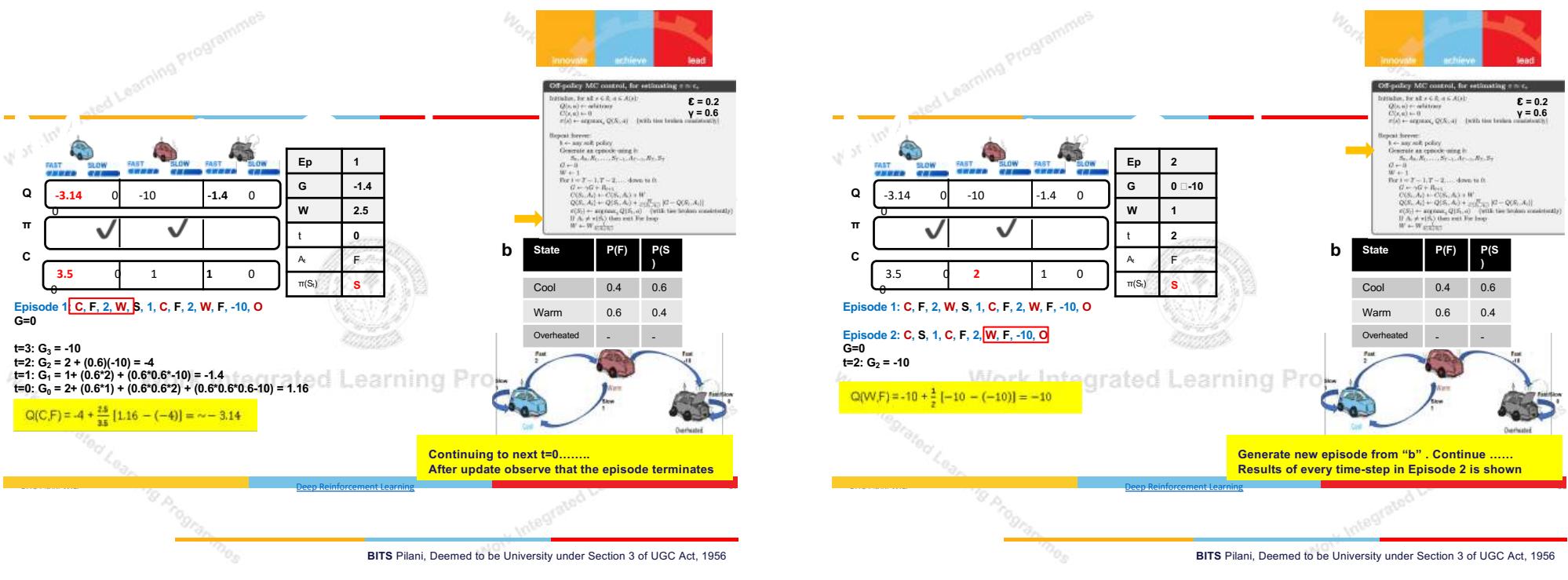
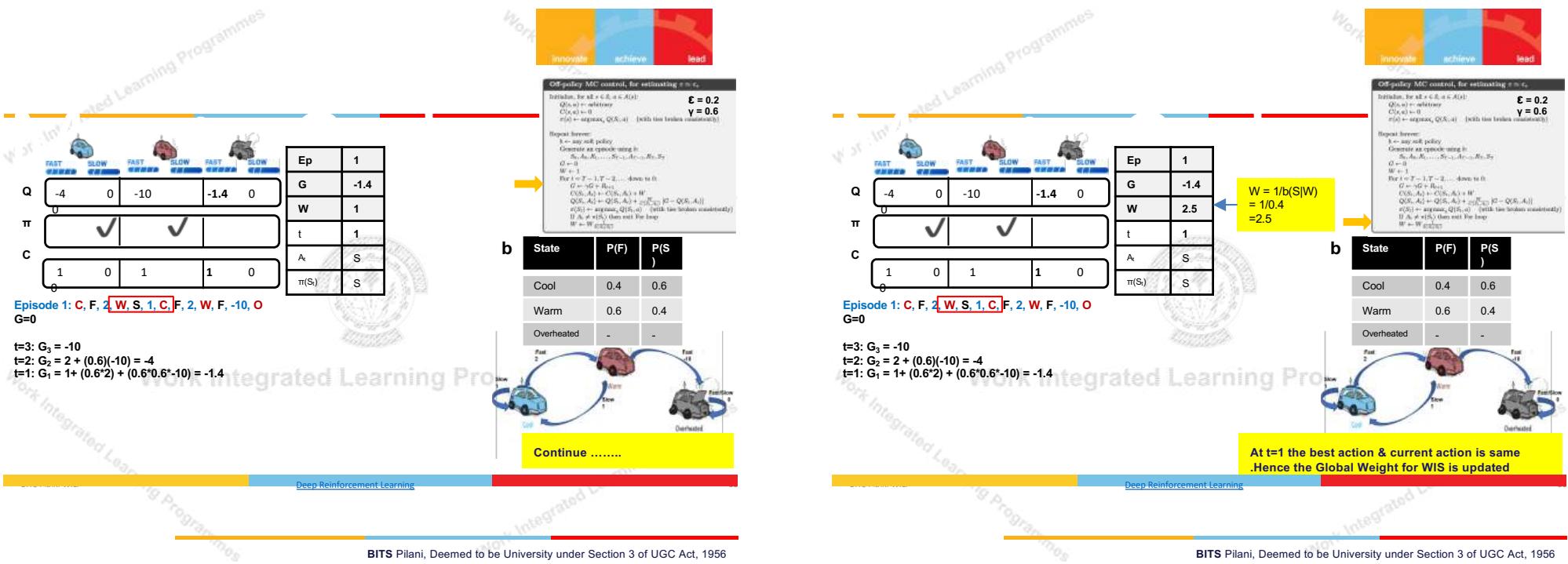


BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956





Work Integrated Learning Programmes

Off-policy MC control, for estimating $\pi \approx \pi_c$

Initiate, for $S \in \mathcal{S}, A \in \mathcal{A}(S)$:
 $Q(S, A) \leftarrow \text{arbitrary}$
 $C(S, A) \leftarrow 0$
 $\pi(S, A) \leftarrow \text{argmax}_A Q(S, A)$ (with ties broken randomly).

Episode 1: $C, F, 2, W, S, 1, C, F, 2, W, F, -10, O$

Episode 2: $C, S, 1, C, F, 2, W, F, -10, O$

$G=0$
 $t=2: G_2 = -10$
 $t=1: G_1 = 2 + (0.6)(-10) = -4$
 $Q(W, F) = -3.14 + \frac{4.5}{1} [-4 - (-3.14)] = -7.01$

b

State	P(F)	P(S)
Cool	0.4	0.6
Warm	0.6	0.4
Overheated	-	-

Continue
Results of every time-step in Episode 2 is shown

Deep Reinforcement Learning

Off-policy MC control, for estimating $\pi \approx \pi_c$

Initiate, for $S \in \mathcal{S}, A \in \mathcal{A}(S)$:
 $Q(S, A) \leftarrow \text{arbitrary}$
 $C(S, A) \leftarrow 0$
 $\pi(S, A) \leftarrow \text{argmax}_A Q(S, A)$ (with ties broken randomly).

Episode 1: $C, F, 2, W, S, 1, C, F, 2, W, F, -10, O$

Episode 2: $C, S, 1, C, F, 2, W, F, -10, O$

$G=0$
 $t=2: G_2 = -10$
 $t=1: G_1 = 2 + (0.6)(-10) = -4$
 $t=0: G_0 = 1 + (0.6)^2 + (0.6^2)(-10) = -1.4$

$Q(W, F) = 0 + \frac{1}{1} [-1.4 - (0)] = -1.4$

b

State	P(F)	P(S)
Cool	0.4	0.6
Warm	0.6	0.4
Overheated	-	-

This completes processing second episode.
Continue next episode

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

1 Given
 $\bullet S, A, \text{ optionally } R(s, a, s')$

2 To find
 $\bullet \text{Policy } \pi$. This is referred to as **TARGET** policy

3 Solution
 $\bullet \text{Sample Episodes}$
 $\bullet \text{Compute estimates of the value by discounted returns}$
Simple average
Weighted average (OIS : Ordinary Importance Sampling)
Normalized weighted average (WIS : Weighted Importance Sampling)

4 Idea
 $\bullet \text{On Policy:}$
 $\quad \bullet \text{Generate Sample/Experience } \square \text{ Evaluate using Return } \square \text{ Improve}$

Off Policy:
 $\bullet \text{Refer to Sample/Experience generated by another reference(BEHAVIOUR) policy}$
 $\quad \square \text{ Evaluate using Return } \square \text{ Improve}$

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Off-policy MC control, for estimating $\pi \approx \pi_c$

Initiate, for $S \in \mathcal{S}, A \in \mathcal{A}(S)$:
 $Q(S, A) \leftarrow \text{arbitrary}$
 $C(S, A) \leftarrow 0$
 $\pi(S, A) \leftarrow \text{argmax}_A Q(S, A)$ (with ties broken randomly).

Episode 1: $C, F, 2, W, S, 1, C, F, 2, W, F, -10, O$

Episode 2: $C, S, 1, C, F, 2, W, F, -10, O$

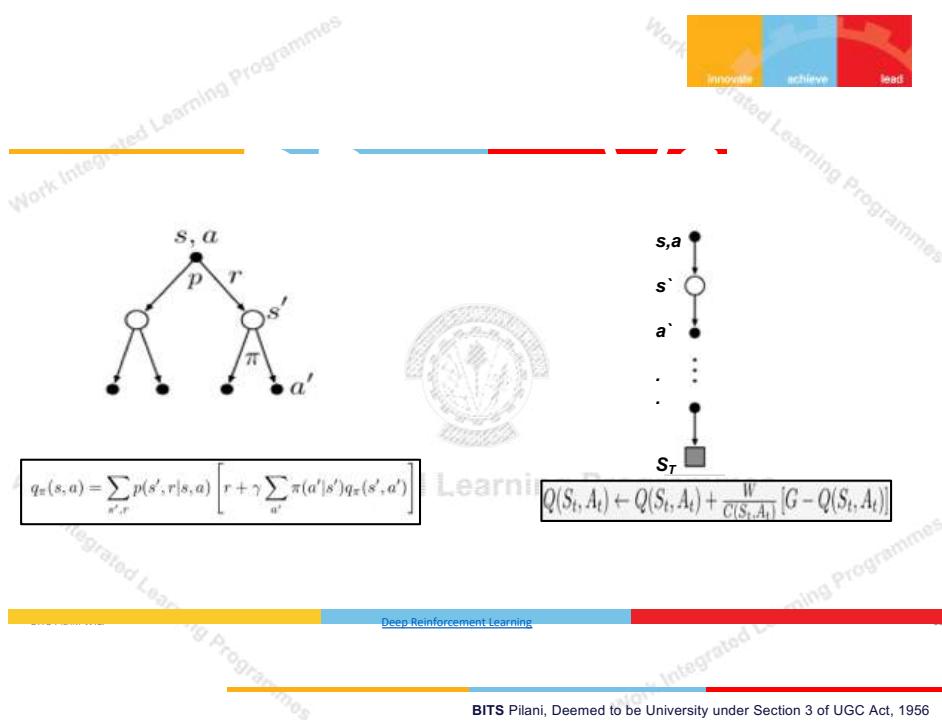
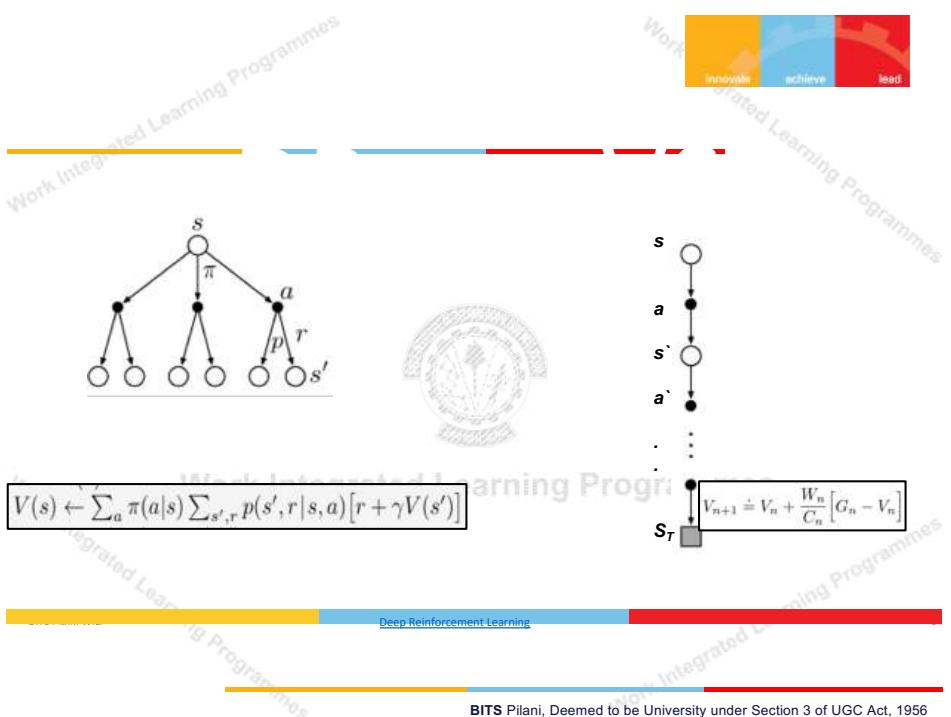
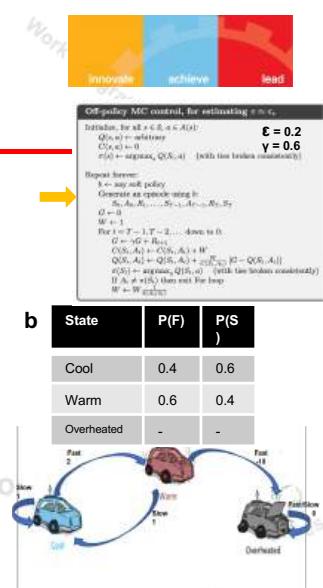
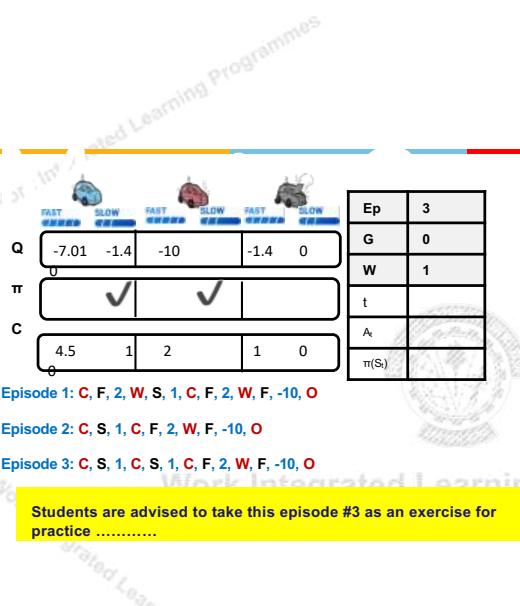
b

State	P(F)	P(S)
Cool	0.4	0.6
Warm	0.6	0.4
Overheated	-	-

Important Note : The Ordinary Importance Sampling (OIS) is unbiased but will produce high variance. The same given algorithm can be adapted to perform ordinary importance sampling by changing only the count update equation to : $C(St, At) \square C(St, At) + 1$

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956





(Modified) Race Car Example

In a modified Race car, there is a training track where agent is given a known (sometimes incomplete) map and risk factors(dynamics) may not be known in advance. Agent knows exactly where it is but its actions impact on the engine state change is difficult to predict (UNKNOWN MODEL).

With partial knowledge, the agent can compute optimal policies by trying through random experiments (Sampling). Car sometimes may shut down when its engine get overheated (optional terminal state) otherwise Race is allowed to proceed indefinitely (NO predefined no.of.laps to complete)

Key Characteristics:

- The agent must act under **uncertainty** —
percepts are local, and the agent must reason about hidden dangers
- Assume there are neither fuel constraints nor No.of.Laps to complete!
- At every time step there is an indication (direct/indirect) of
Impact of next potential engine temperature depending on the
current engine state. **Immediate reward is available.**



Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

1 Given

- $S, A, R(s, a, s')$, (optionally) $P(a | s)$

2 To find

- Policy π

3 Solution

- Sample Episodes
- Compute estimates of the value by discounted **next state estimates (bootstrapping)**
 - Simple average
 - Weighted average
 - Normalized weighted average

4 Idea

- On Policy :
 - Generate Sample/Experience Evaluate using **bootstrapping** Improve
- Off Policy :
 - Refer to Sample/Experience generated by another reference(behavior) policy
 - Evaluate using **bootstrapping** Improve

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Temporal Difference Learning Methods



Work Integrated Learning Programmes

Policy Evaluation (PE) Algorithm

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

V_0	0	0	0
V_1			
V_2			

Episode 1: C

Important Note :

- In few textbooks/ running incremental upcoming examples, dynamics are given so **only to simulate trajectories** and show updates step by step. TD itself doesn't require knowing the full transition model — it only needs sampled transitions (s, a, r, s') while training from scratch i.e., when policy is not known/learnt yet.
- The dynamics are provided for clarity and reproducibility, not because TD requires them. In practice, TD works model-free, just like Monte Carlo.

Tabular TD(0) for estimating v_π .
Input: the policy π to be evaluated.
Algorithm parameter: step size $\alpha \in [0, 1]$
Initialization: $v_\pi(s)$, for all $s \in \mathcal{S}$, arbitrarily except that $v_\pi(\text{terminal}) = 0$.
Loop for each episode:
Initialize S .
Loop for each step of episode:
 $A \leftarrow$ action given by π for S .
 Take action A ; observe R, S' .
 $V(S) \leftarrow V(S) + \alpha(R + \gamma V(S') - V(S))$.
 $S \leftarrow S'$.
until S is terminal.

$\gamma = 0.1$
 $\alpha = 0.9$



Use case : A driving policy is already learnt which needs to be evaluated for monitoring or certification/approval by safety engineers.

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes







V_0	1	0	0
V_1			
V_2			

Episode 1 [C, F, 2, W]

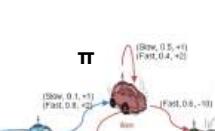
Target = $R + \gamma V(W) = 2 + (0.9 * 0) = 2$
 TD Error = Target - Old Estimate = $2 - V(C) = 2 - 0 = 2$
 Update the value : $V(C) = 0 + 0.5(2) = 1$

Tabular TD(0) for estimating v_π

Input: the policy π to be evaluated
 Algorithm parameter: step size $\alpha \in [0, 1]$
 Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:
 Initialize S
 Loop for each step of episode:
 $a \leftarrow$ action given by π for S
 Take action A , observe R, R'
 $V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$
 $S \leftarrow S'$
 until S is terminal





Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



Work Integrated Learning Programmes

V_0	1	0.95	0
V_1	1.95	2.40	1.40
V_2			

Episode 1: C, F, 2, W, S, 1, C
Episode 2: C, F, 2, C, F, 2, W, S, 1, W

$t=0 : V(C) = 1 + 0.5 [2 + (0.9 \cdot 1) - 1] = 1.95$
 $t=1 : V(C) = 1.95 + 0.5 [2 + (0.9 \cdot 0.95) - 1.95] = \sim 1.4$
 $t=2 : V(W) = 0.95 + 0.5 [1 + (0.9 \cdot 0.95) - 0.95] = \sim 1.40$

Tabular TD(0) for estimating v_π :

Input: the policy π to be evaluated
Algorithm parameter: step size $\alpha \in [0, 1]$
Initial value: $V(s)$, for all $s \in \mathcal{S}^+$; arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:
Initialize S
Loop for each step of episode:
 $A \leftarrow$ action given by π for S
 Take action A , observe R, S'
 $V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$
 $S \leftarrow S'$
until S is terminal.

Continued for 1 more episode.....

Temporal Difference Learning Methods

On Policy Control Algorithms

Work Integrated Learning Programmes

Temporal Difference Learning Methods



Work Integrated Learning Programmes

On Policy Control Algorithms

BITs Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Sarsa (on-policy TD control) for estimating $Q \approx q$

```

Initialize  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat {for each episode}
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Repeat {for each step of episode}
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal
  
```

Q_1	1.6	0	0	0
Q_2				
Q_3				

Episode 1: C, F, 2, W, S

$$Q(C,F) = 0 + 0.8 [2 + (0.9 \cdot 0) - 0] = 1.6$$

Work Integrated Learning Programmes

From State "C" & "W", Action "F" & "S" are derived respectively by applying E-greedy on current policy from Q values learnt till now.

In the shown time step : $(S, A, R, S, A) = (C, F, 2, W, S)$

Deep Reinforcement Learning

BITs Pilani, Deemed to be University under Section 3 of UGC Act, 1956

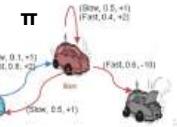
Sarsa (on-policy TD control) for estimating $Q \approx q$

```

Initialize  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat {for each episode}
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Repeat {for each step of episode}
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal
  
```

Q_0	0	0	0	0
Q_1				
Q_2				

Episode 1: C



Work Integrated Learning Programmes

Use case : A driving policy is already learnt but there are few unforeseen changes. Policy needs to be retrained to reflect the actual risk observed via exploration

Deep Reinforcement Learning

BITs Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Sarsa (on-policy TD control) for estimating $Q \approx q$

```

Initialize  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat {for each episode}
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Repeat {for each step of episode}
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal
  
```

Q_1	1.6	0	0	0
Q_2				
Q_3				

Episode 1: C, F, 2, W, S

$$Q(C,F) = 0 + 0.8 [2 + (0.9 \cdot 0) - 0] = 1.6$$

Work Integrated Learning Programmes

From State "C" & "W", Action "F" & "S" are derived respectively by applying E-greedy on current policy from Q values learnt till now.

In the shown time step : $(S, A, R, S, A) = (C, F, 2, W, S)$

Deep Reinforcement Learning

BITs Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Sarsa (on-policy TD control) for estimating $Q \approx q$

```

Initialize  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat {for each episode}
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Repeat {for each step of episode}
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal
  
```

Q_1	1.6	0	0	1.95	0
Q_2					
Q_3					

Episode 1: C, F, 2, W, S, 1, C, F

$$Q(W,S) = 0 + 0.6 [1 + (0.9 \cdot 1.6) - 0] = \sim 1.95$$

Sarsa (on-policy TD control) for estimating $Q \approx q$

```

Initialize  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat {for each episode}
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Repeat {for each step of episode}
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal
  
```



Terminating episode 1 after 2 steps. In practice, we can define an episode cutoff (e.g., after 2 or 10 steps or when Overheated occurs). If Overheated isn't reached, still the episode can be programmed to end artificially. That's how learning continues.

Deep Reinforcement Learning

BITs Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



Sarsa (on-policy TD control) for estimating $Q \approx q_s$.

```

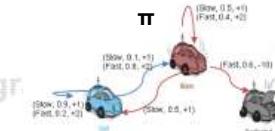
Initialize  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
  Initialize  $S$ 
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
  Repeat (for each step of episode):
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
     $S \leftarrow S'; A \leftarrow A'$ 
  until  $S$  is terminal
  
```

Deep Reinforcement Learning

Q₁	1.6	0	0	1.95	0
Q₂	3.62	0	-8	1.32	0
Q₃					

Episode 1: C, F, 2, W, S, 1, C, F
Episode 2: C, F, 2, C, F, 2, W, S, 1, W, S, 1, W, F, -10, O

t=0 : $Q(C,F) = 1.6 + 0.8 [2 + (0.9*1)*1.6] = 3.07$
t=1 : $Q(C,F) = 3.07 + 0.8 [2 + (0.9*1.95)*3.07] = \sim 3.62$
t=2 : $Q(W,S) = 1.95 + 0.8 [1 + (0.9*1.95)*1.95] = \sim -2.59$
t=3 : $Q(W,S) = 2.59 + 0.8 [1 + (0.9*0)*2.59] = \sim -1.32$
t=4 : $Q(W,F) = 0 + 0.8 [-10 + (0.9*0)*0] = \sim -8$

TT 

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



Temporal Difference Learning Methods

On Policy Control Algorithms

Work Integrated Learning Programmes



On Policy Control Algorithms

Work Integrated Learning Programmes



Sarsa (on-policy TD control) for estimating $Q \approx q_s$.

```

Initialize  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
  Initialize  $S$ 
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
  Repeat (for each step of episode):
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
     $S \leftarrow S'; A \leftarrow A'$ 
  until  $S$  is terminal
  
```

Continued for two more episodes....

Observe the difference in value updates w.r.t $Q(W,S)$ between successive time steps in episodes 3 & 4.

At the worst case on larger learning rate eg., 0.8 the pull is even larger

This large pull is an effect of the **VARIANCE** produced

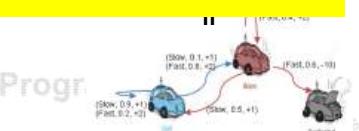
Solution to tackle : Expected SARSA

Deep Reinforcement Learning

Q₁	1.6	0	0	1.95	0
Q₂	3.62	0	-8	1.32	0
Q₃	3.62	0	6.48	-4.69	0
Q₄	3.62	0	6.48	+4.47	0

Episode 1: C, F, 2, W, S, 1, C, F
Episode 2: C, F, 2, C, F, 2, W, S, 1, W, S, 1, W, F, -10, O
Episode 3: W, S, +1, W, F, +10, O
t=0 : $Q(W,S) = 1.32 + 0.8 [1 + (0.9*-8) - 1.32] = \sim -4.69$
t=1 : $Q(W,F) = -8 + 0.8 [1 + (0.9*0) + 8] = \sim +6.4$

Episode 4: W, S, +1, W, F
t=0 : $Q(W,S) = -4.69 + 0.8 [1 + (0.9*6.4) + 4.69] = \sim +4.47$

TT 

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



Sarsa (on-policy TD control) for estimating $Q \approx q_s$.

```

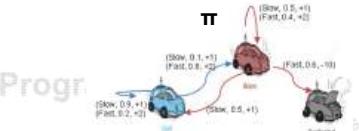
Initialize  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
  Initialize  $S$ 
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
  Repeat (for each step of episode):
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
     $S \leftarrow S'; A \leftarrow A'$ 
  until  $S$  is terminal
  
```

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \sum_{i=0}^{T-t-1} Q(S_{t+i+1}, A_{t+i+1}) | S_{t+i+1}] - Q(S_t, A_t)] \\ \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \sum_{i=0}^{T-t-1} \pi(a|s_{t+i})Q(S_{t+i}, a) - Q(S_t, A_t)] \quad (1)$$

Deep Reinforcement Learning

Q₀	0	0	0	0	0
Q₁					
Q₂					

Episode 1: C

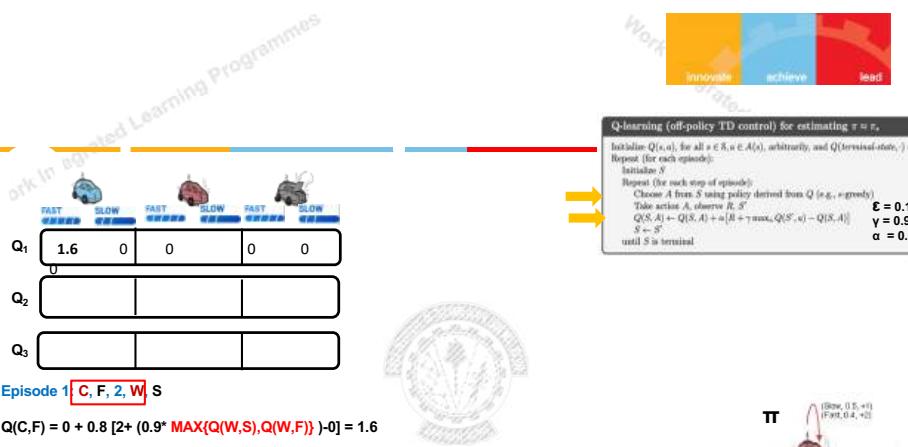
TT 

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Temporal Difference Learning Methods

Off Policy Control Algorithms

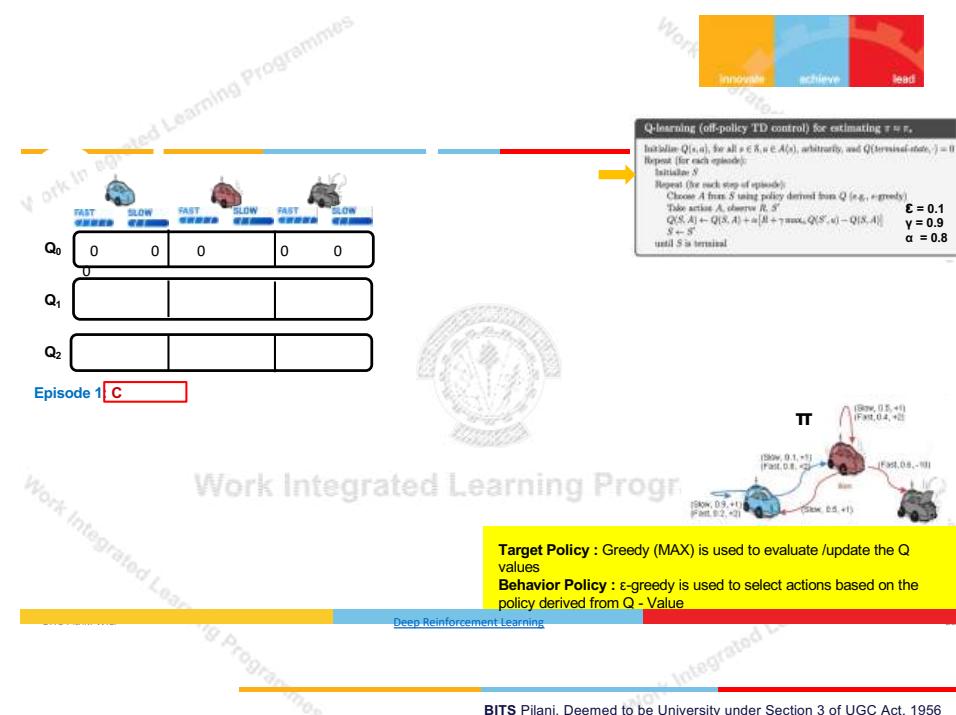
ITS Main, Deemed to be University under Section 3 of UGC Act, 1956



Behavior Policy : ϵ -greedy is used to select actions based on the policy derived from Q - Value . **Here action "F" for state "C"**
Target Policy : Greedy (MAX) is used to evaluate /update the Q values.

Deep Reinforcement Learning

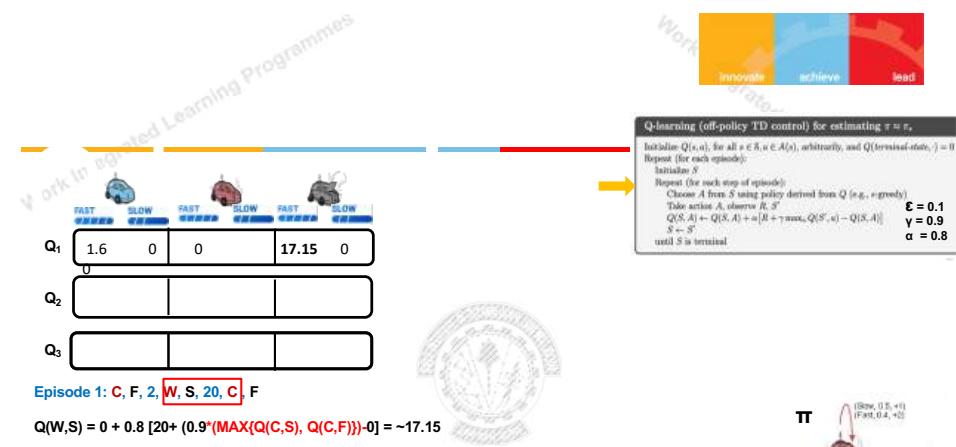
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



- Target Policy** : Greedy (MAX) is used to evaluate /update the Q values
- Behavior Policy** : ϵ -greedy is used to select actions based on the

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Terminating episode 1 after 2 steps. In practice, we can define an episode cutoff (e.g., after 2 or 10 steps or when Overheated occurs). If Overheated isn't reached, still the episode can be programmed to end artificially. That's how learning continues.

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$.

```

Initialize:  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat for each episode:
    Initialize:  $\pi$  (e.g.,  $\pi(a|s) = 1/\text{card}(A(s))$ )
    Initialize:  $R$  (for each step of episode)
    Initialize:  $\alpha$  (for each step of episode)
    Initialize:  $\gamma$  (for each step of episode)
    Initialize:  $\epsilon$  (for each step of episode)
    Initialize:  $\epsilon'$  (for each step of episode)
    Initialize:  $E = 0.1$ 
    Initialize:  $\gamma = 0.9$ 
    Initialize:  $\alpha = 0.8$ 
    Initialize:  $S$  (start state)
    Initialize:  $A$  (action chosen from  $\pi$ )
    Initialize:  $R, S'$  (rewards and next states)
    Initialize:  $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
    Initialize:  $S' \leftarrow S$ 
    until  $S$  is terminal
  
```

Q₁	1.6	0	0	17.15	0
0					
Q₂	16.16	0	-8	11.64	0
0					
Q₃					

Episode 1: C, F, 2, W, S, 20, C, F
 Episode 2: C, F, 2, C, F, 4, W, S, -5, W, S, 1, W, F, -10, O

t=0 : $Q(C,F) = 1.6 + 0.8(2+0.9(\text{MAX}(1.6,0))-1.6) \approx 3.07$
 t=1 : $Q(C,F) = 3.07 + 0.8(4+0.9(\text{MAX}(0,17.15))-3.07) \approx -16.16$
 t=2 : $Q(W,S) = 17.15 + 0.8(-5+0.9(\text{MAX}(0,17.15))-17.15) \approx -11.78$
 t=3 : $Q(W,S) = 11.78 + 0.8(1+0.9(\text{MAX}(0,11.78))-11.78) \approx -11.64$
 t=4 : $Q(W,F) = 0 + 0.8[-10 + (0.9*0)-0] = \approx -8$

Target Policy : Greedy (MAX) is used to evaluate /update the Q values
 Behavior Policy : ϵ -greedy is used to select actions based on the policy derived from Q - Value

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$.

```

Initialize:  $Q(s, a)$ , for all  $s \in S, a \in A(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat for each episode:
    Initialize:  $\pi$  (e.g.,  $\pi(a|s) = 1/\text{card}(A(s))$ )
    Initialize:  $R$  (for each step of episode)
    Initialize:  $\alpha$  (for each step of episode)
    Initialize:  $\gamma$  (for each step of episode)
    Initialize:  $\epsilon$  (for each step of episode)
    Initialize:  $\epsilon'$  (for each step of episode)
    Initialize:  $E = 0.1$ 
    Initialize:  $\gamma = 0.9$ 
    Initialize:  $\alpha = 0.8$ 
    Initialize:  $S$  (start state)
    Initialize:  $A$  (action chosen from  $\pi$ )
    Initialize:  $R, S'$  (rewards and next states)
    Initialize:  $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
    Initialize:  $S' \leftarrow S$ 
    until  $S$  is terminal
  
```

Q₁	1.6	0	0	17.15	0
0					
Q₂	16.16	0	-8	11.64	0
0					
Q₃	16.16	0	6.4	14.71	0
0					
Q₄	16.16	11.39	6.4	22.50	0
0					
Q₅	16.16	11.39	6.4	16.70	0
0					
Q₆	16.16	11.39	6.4	19.37	0

Episode 1: C, F, 2, W, S, 20, C
 Episode 2: C, F, 2, C, F, 4, W, S, -5, W, S, 1, W, F, -10, O

Observe the difference in value updates w.r.t $Q(W,S)$ between successive time steps in episodes 3 to 6.
 The true value is approx. $Q(W,S) = 15$
 Due to MAX bias the values are overestimated and this is more pronounced in the presence of NOISE
 Sometimes this diverges from true value.
 Solution to tackle : Double Q-Learning

** But in near-deterministic , less noisy environment , Q-learning helps in faster convergence!

Few more episodes are completed and the results of the update per Q(State, Action) is shown . Students are advised to practice this.

Deep Reinforcement Learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Temporal Difference Learning Methods

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Double Q-learning for estimating $Q_0 \geq Q_1 \geq Q_2$.

```

Algorithm (assumes: input dist  $\alpha \in [0,1]$ , and  $\beta > 0$ )
    Initialize  $Q_0(s, a)$  and  $Q_1(s, a)$ , for all  $s \in S, a \in A(s)$ , such that  $Q_0(\text{terminal}, \cdot) = 0$ 
    Loop for n-th episode:
        Initialize:  $\pi$  (e.g.,  $\pi(a|s) = 1/\text{card}(A(s))$ )
        Initialize:  $R$  (for each step of episode)
        Initialize:  $\alpha$  (for each step of episode)
        Initialize:  $\gamma$  (for each step of episode)
        Initialize:  $\epsilon$  (for each step of episode)
        Initialize:  $\epsilon'$  (for each step of episode)
        Initialize:  $E = 0.1$ 
        Initialize:  $\gamma = 0.9$ 
        Initialize:  $\alpha = 0.8$ 
        Initialize:  $S$  (start state)
        Initialize:  $A$  (action chosen from  $\pi$ )
        Initialize:  $R, S'$  (rewards and next states)
        Initialize:  $Q_0(S, A) \leftarrow Q_0(S, A) + \alpha[R + \gamma Q_2(S', \text{argmax}_a Q_2(S', a)) - Q_0(S, A)]$ 
        Initialize:  $Q_1(S, A) \leftarrow Q_1(S, A) + \alpha[R + \gamma Q_0(S', \text{argmax}_a Q_0(S', a)) - Q_1(S, A)]$ 
        Initialize:  $S' \leftarrow S$ 
        until  $S$  is terminal
    
```

Q₀	0	0	0	0	
5	1	5	15		
Q₁					
Q₂					

Episode 1: C

Target Policy : Q₂ - Table
 Behavior Policy : Q₁-Table is available from other sources/experience/logged data. Above algorithm learns to update this as well.

Work Integrated Learning Programmes

Innovate achieve lead

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_t$.

Algorithm (assumes input dist $\alpha \in [0, 1]$, and $\gamma > 0$)
 Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^*$, $a \in A(s)$, such that $Q(s/\text{terminal}, \cdot) = 0$
 Loop for each episode:
 Choose A from S using the policy π -greedy in $Q_1 + Q_2$
 With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha(R + \gamma Q_2(S', \text{argmax}_a Q_2(S', a)) - Q_1(S, A))$$
 else

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha(R + \gamma Q_1(S', \text{argmax}_a Q_1(S', a)) - Q_2(S, A))$$
 until S is terminal

Episode 1: C, F, 2, W, S

Q2(C,F)
 $= 0 + 0.8 [2 + (0.9^* \text{Q1}(W, \text{ARGMAX}(Q2(W,S), Q2(W,F))))] - 0$
 $= 0 + 0.8 [2 + (0.9^* \text{Q1}(W, F)) - 0]$ #Here both S & F had same values hence randomly "F" is chosen
 $= 0 + 0.8 [2 + (0.9^* 5) - 0]$
 $= 5.2$

Target Policy : Q2 - Table
Behavior Policy : Q1-Table is available from other sources/experience/logged data. Above algorithm learns to update this as well.

Deep Reinforcement Learning

Work Integrated Learning Programmes

Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_t$.

Algorithm (assumes input dist $\alpha \in [0, 1]$, and $\gamma > 0$)
 Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^*$, $a \in A(s)$, such that $Q(s/\text{terminal}, \cdot) = 0$
 Loop for each episode:
 Choose A from S using the policy π -greedy in $Q_1 + Q_2$
 With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha(R + \gamma Q_2(S', \text{argmax}_a Q_2(S', a)) - Q_1(S, A))$$
 else

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha(R + \gamma Q_1(S', \text{argmax}_a Q_1(S', a)) - Q_2(S, A))$$
 until S is terminal

Episode 1: C, F, 2, W, S, 20, C, F

Q2(W,S)
 $= 0 + 0.8 [20 + (0.9^* \text{Q1}(C, \text{ARGMAX}(Q2(C,S), Q2(C,F))))] - 0$
 $= 0 + 0.8 [20 + (0.9^* 5) - 0]$
 $= 0 + 0.8 [20 + (0.9^* 5) - 0]$
 $= -19.6$

Terminating episode 1 after 2 steps. In practice, we can define an episode cutoff (e.g., after 2 or 10 steps or when Overheated occurs). If Overheated isn't reached, still the episode can be programmed to end artificially. That's how learning continues.

Deep Reinforcement Learning

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_t$.

Algorithm (assumes input dist $\alpha \in [0, 1]$, and $\gamma > 0$)
 Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^*$, $a \in A(s)$, such that $Q(s/\text{terminal}, \cdot) = 0$
 Loop for each episode:
 Choose A from S using the policy π -greedy in $Q_1 + Q_2$
 With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha(R + \gamma Q_2(S', \text{argmax}_a Q_2(S', a)) - Q_1(S, A))$$
 else

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha(R + \gamma Q_1(S', \text{argmax}_a Q_1(S', a)) - Q_2(S, A))$$
 until S is terminal

Episode 1: C, F, 2, W, S, 20, C, F
Episode 2: C, F, 2, C, F, 4, W, S, -5, W, S, 1, W, F, -10, O
 $t=0 : Q2(C,F) = 5.2 + 0.8(2 + 0.9(\text{Q1}(C, \text{ARGMAX}(Q2(C,F), Q2(C,S)))) - 5) = -5.2 = -6.24$

For step 2 : The loop goes to Q1(S,A) where the Behavior policy is updated!
 $t=1 : \text{Q1}(C,F) = 5 + 0.8(4 + 0.9(\text{Q2}(W, \text{ARGMAX}(C(W,S), Q1(W,F)))) - 5) = 5 + 0.8(4 + 0.9(\text{Q2}(W, \text{ARGMAX}(5, 5))) - 5) = 5 + 0.8(4 + 0.9(5) - 5) = 5 + 0.8(4 + 0.9(19.6) - 5) = -18.3$

Only two time step is shown here . Other time steps updated target policy ie., Q2 table and the result is tabulated above.

Target Policy : Q2 - Table
Behavior Policy : Q1-Table is available from other sources/experience/logged data. Above algorithm learns to update this as well.

Deep Reinforcement Learning

Work Integrated Learning Programmes

Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_t$.

Algorithm (assumes input dist $\alpha \in [0, 1]$, and $\gamma > 0$)
 Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^*$, $a \in A(s)$, such that $Q(s/\text{terminal}, \cdot) = 0$
 Loop for each episode:
 Choose A from S using the policy π -greedy in $Q_1 + Q_2$
 With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha(R + \gamma Q_2(S', \text{argmax}_a Q_2(S', a)) - Q_1(S, A))$$
 else

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha(R + \gamma Q_1(S', \text{argmax}_a Q_1(S', a)) - Q_2(S, A))$$
 until S is terminal

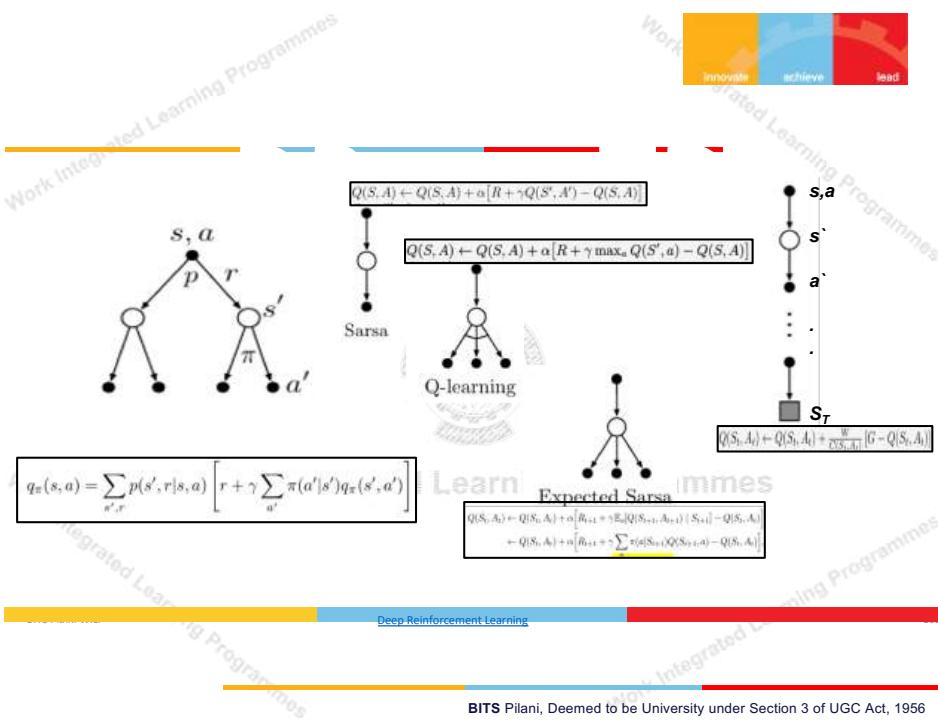
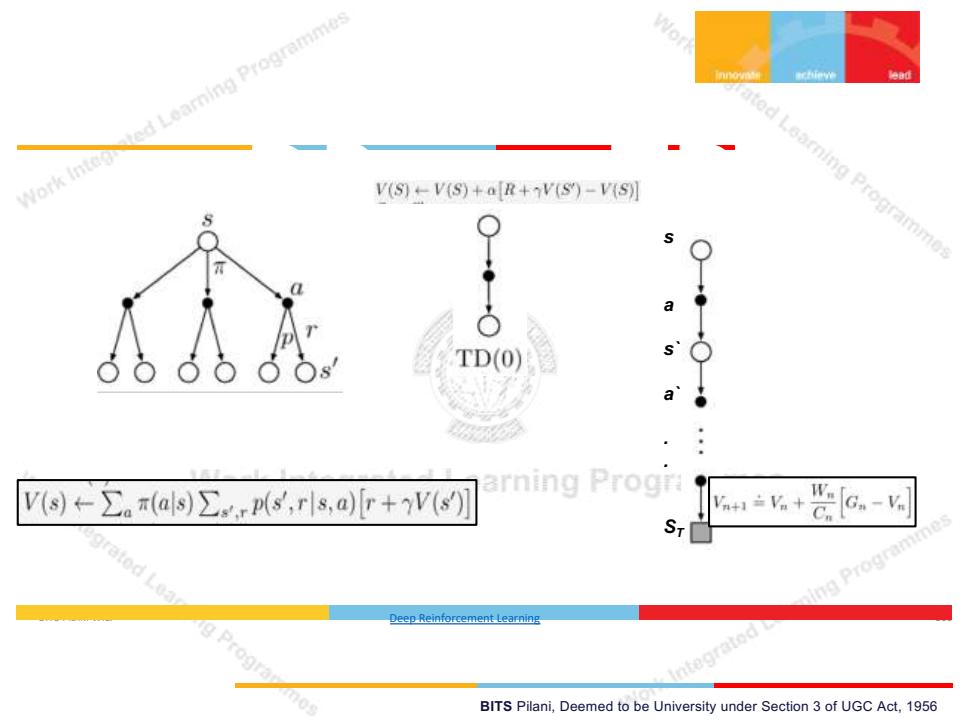
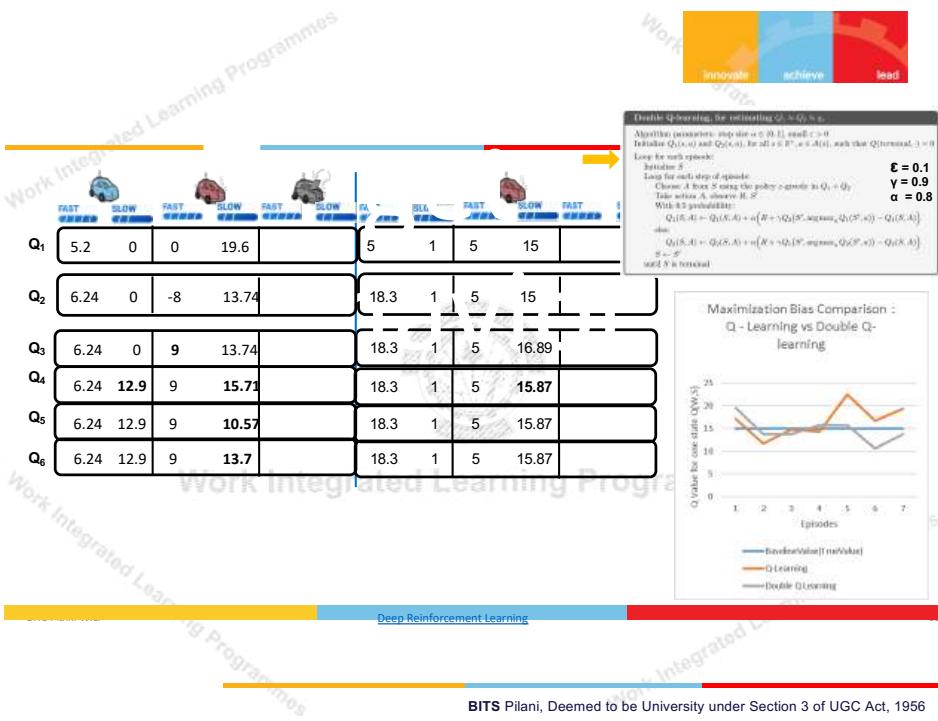
Episode 1: C, F, 2, W, S, 20, C
Episode 2: C, F, 2, C, F, 4, W, S, -5, W, S, 1, W, F, -10, O
Episode 3: C, F, 2, W, S, 20, C
Episode 4: C, S, +1, W, S, +10, C
Episode 5: W, S, -5, W
Episode 6: W, S, +5, W

Two more episodes are completed and the results of the update per Q(State, Action) is shown . Students are advised to practice this.

Deep Reinforcement Learning

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956





<AIMLCZG512, Deep Reinforcement Learning>

Lecture No. 10,11,12

Agenda for the classes

- Introduction
- Value Function Approximation
- Stochastic Gradient, Semi-Gradient Methods
- Role of Deep Learning for Function Approximation;
- Feature Construction Methods

Work Integrated Learning Programmes

Acknowledgements: Some of the slides were adopted *with permission* from the course [CSCE-689](#) (Texas A&M University) by Prof. Guni Sharon

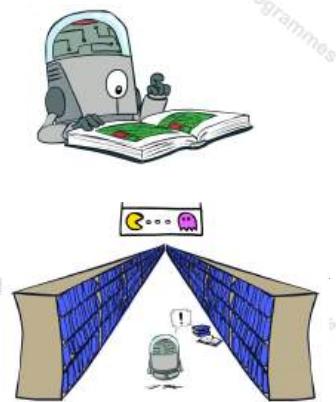
3

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Generalizing Across States

- Tabular Learning keeps a table of all state values
- In realistic situations, we cannot possibly learn about every single state!

- Too many states to visit them all during training
- Too many states to hold a value table in memory

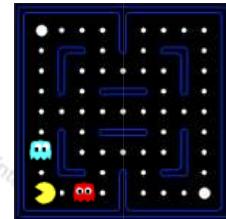


- Instead, we want to generalize:
 - Learn about some small number of training states from experience
 - Generalize that experience to new, similar situations

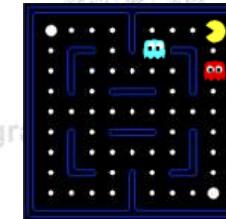
○ This is a fundamental idea in machine learning

Example: Pacman

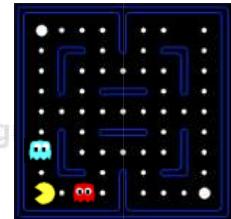
Let's say we discover through experience that this state is bad:



In naïve tabular-learning, we know nothing about this state:



Or even this one!



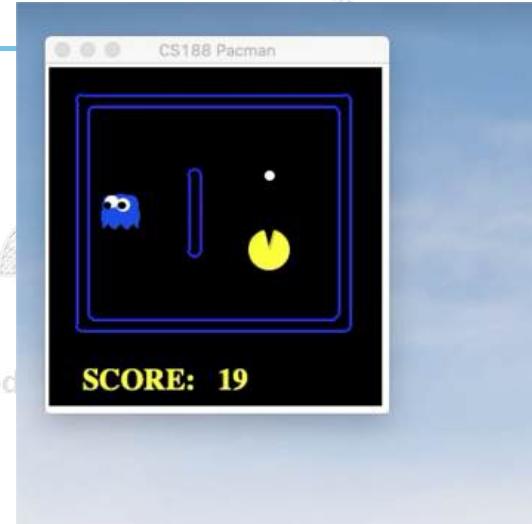
5

- Naïve Q-learning
- After 50 training episodes



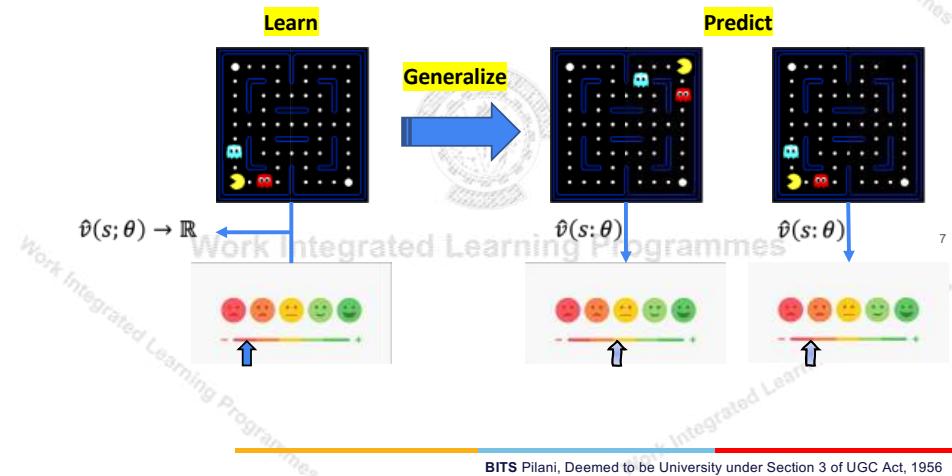
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

- Q-learning with function approximator



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

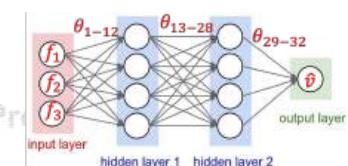
Learn an approximation function



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Parameterized function approximator

- Assume that each state is vector of features (f_1, f_2, \dots, f_n) , e.g.,
 - Packman location, Ghost1 location , Ghost2 location, food location
 - Or even screen pixels
- A parametrized value approximator $\hat{v}(s; \theta)$ might look like this:
 $\hat{v} = \sum_i \theta_i f_i$ or maybe like this $\hat{v} = \prod_i f_i^{\theta_i}$ or even this
- Assume we know the true value for a set of states:
 - $v(S_1) = 5, v(S_2) = 8, v(S_3) = 2$
 - How can we update θ to reflect this information?



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Gradient Descent

- Given: $v(S_1) = 5, v(S_2) = 8, v(S_3) = 2$
- We want to set θ such that $\forall s, \hat{v}(s; \theta) = v(s)$
 - Not possible in the general case, why?
 - Instead we'll try to minimize the errors: loss = $\sum_s |v(s) - \hat{v}(s; \theta)|$
 - Partial derivative of the loss with respect to θ_i = how to change θ_i such that loss will increase the most
 - Go the other way \rightarrow decrease loss
 - Ooops! Absolute value is not differentiable \rightarrow can't compute gradients
 - Simple fix: loss = $\frac{1}{2} \sum_s [v(s) - \hat{v}(s; \theta)]^2$ = squared loss function

Work Integrated Learning Programmes

10

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Gradient Descent

- $\text{loss} = \frac{1}{2} \sum_s [v(s) - \hat{v}(s; \theta)]^2$
- For each i
 - Push θ_i towards a direction that minimizes loss
 - $\theta_i = \theta_i - \frac{\partial \text{loss}}{\partial \theta_i}$
- More generally $\theta = \theta - \alpha \nabla \text{loss}$
- $\nabla \text{loss} = (\frac{\partial \text{loss}}{\partial \theta_1}, \frac{\partial \text{loss}}{\partial \theta_2}, \dots, \frac{\partial \text{loss}}{\partial \theta_n})$
- α is the learning rate, requires tuning per domain, too large learning diverges to small results in slow learning or even premature convergence



Work Integrated Learning Programmes

11

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



SGD for Monte Carlo estimation

Gradient Monte Carlo Algorithm for Estimating $\hat{v} \approx v_\pi$

Input: the policy π to be evaluated
Input: a differentiable function $\hat{v} : S \times \mathbb{R}^d \rightarrow \mathbb{R}$

Initialize value-function weights w as appropriate (e.g., $w = 0$)

Repeat forever:

Generate an episode $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$ using π
For $t = 0, 1, \dots, T - 1$:

$$w \leftarrow w + \alpha [G_t - \hat{v}(S_t, w)] \nabla \hat{v}(S_t, w)$$

w are the tunable parameters of the value approximation function

$$f(S) = [2, 2, 1]$$



- Guaranteed to converge to a local optimum because G_t is an unbiased estimate of $v_\pi(S_t)$

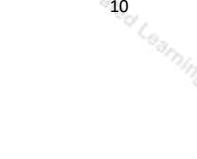
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Example

$$\mathcal{S} = \{f_1(S), f_2(S), f_3(S)\}$$

$$f(S') = [2, 4, 1]$$

- $f_{1,2}$ = distance to ghost 1, 2, f_3 = distance to food
- $\hat{v}(s) = \sum_i \theta_i f_i(s)$
- init: $\theta = [0, 0, 0]$
- $\theta = \theta + \alpha (G_t - \hat{v}(s; \theta)) \nabla \hat{v}(s; \theta)$
- $\theta = [0, 0, 0] + 0.1 (10 - [0, 0, 0] \cdot [2, 2, 1]) [2, 2, 1]$
- $\theta = [2, 2, 1]$
- $\hat{v}(S') = f(S') \cdot \theta = [2, 4, 1] \cdot [2, 2, 1] = 13$



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Learning approximation with bootstrapping

- Can we update the value approximation function at every step?
- Yes, define SGD as a function of the TD error
 - Tabular TD learning: $\hat{v}(s_t) = \hat{v}(s_t) + \alpha(r_t + \gamma\hat{v}(s_{t+1}) - \hat{v}(s_t))$
 - Approximation TD learning: $\theta = \theta + \alpha(r_t + \gamma\hat{v}(s_{t+1}; \theta) - \hat{v}(s_t; \theta))\nabla\hat{v}(s_t; \theta)$
- Known as Semi-gradient methods
- NOT guaranteed to converge to a local optimum because $\hat{v}(s_{t+1}; \theta)$ is a biased estimate of $v_\pi(s_{t+1})$
- Semi-gradient (bootstrapping) methods do not converge as robustly as (full) gradient methods

Semi-gradient methods

- They do converge reliably in important cases such as the linear approximation case
- They offer important advantages that make them often clearly preferred
- They typically enable significantly faster learning, as we have seen in Chapters 6 and 7
- They enable learning to be continual and online, without waiting for the end of an episode
- This enables them to be used on continuing problems and provides computational advantages

Semi-gradient TD(0)

Semi-gradient TD(0) for estimating $\hat{v} \approx v_\pi$

```

Input: the policy  $\pi$  to be evaluated
Input: a differentiable function  $\hat{v}: \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\hat{v}(\text{terminal}, \cdot) = 0$ 
Initialize value-function weights  $w$  arbitrarily (e.g.,  $w = 0$ )
Repeat (for each episode):
  Initialize  $S$ 
  Repeat (for each step of episode):
    Choose  $A \sim \pi(\cdot | S)$ 
    Take action  $A$ , observe  $R, S'$ 
     $w \leftarrow w + \alpha[R + \gamma\hat{v}(S', w) - \hat{v}(S, w)]\nabla\hat{v}(S, w)$ 
     $S \leftarrow S'$ 
  until  $S'$  is terminal
  
```

What's the difference from the tabular case?

Example

$$f(S) = [2,3,1]$$



$$R = +10$$

$$f(S') = [1,2,1]$$



$$\mathcal{S} = \{f_1(S), f_2(S), f_3(S)\}$$

- $f_{1,2}$ =distance to ghost 1,2, f_3 =distance to food
- $\hat{v}(s) = \sum_i \theta_i f_i(s)$
- init: $\theta = [0,0,0]$
- $\theta = \theta + \alpha(R + \gamma\hat{v}(S'; \theta) - \hat{v}(S; \theta))\nabla\hat{v}(S; \theta)$
- $\theta = [0,0,0] + 0.1(10 + [1,2,1] \cdot [0,0,0] - [2,3,1] \cdot [0,0,0])[2,3,1]$
- $\theta = [2,3,1]$
- $\hat{v}(U) = f(U) \cdot \theta = [2,4,1] \cdot [2,3,1] = 17$

$$f(U) = [2,4,1]$$



[Review] n-step TD Prediction

One-step return:

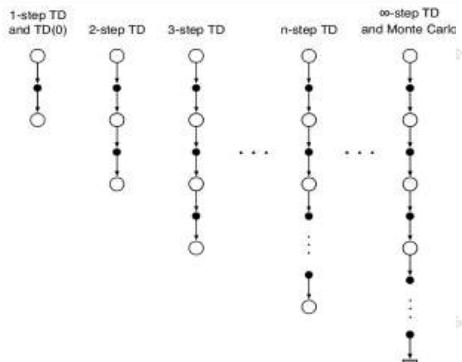
$$G_{t:t+1} \doteq R_{t+1} + \gamma V_t(S_{t+1})$$

Two-step return:

$$G_{t:t+2} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

Work Integrated Le

Ref Section 7.1 of TB



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

[Review] n-step TD Prediction

One-step return:

$$G_{t:t+1} \doteq R_{t+1} + \gamma V_t(S_{t+1})$$

Two-step return:

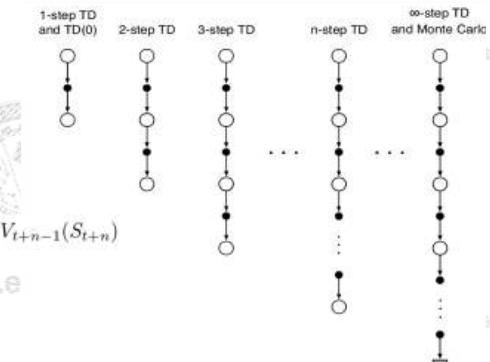
$$G_{t:t+2} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

n-step return:

$$G_{t:t+n} \doteq R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

Work Integrated Le

Ref Section 7.1 of TB



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



[Review] n-step TD Prediction

One-step return:

$$G_{t:t+1} \doteq R_{t+1} + \gamma V_t(S_{t+1})$$

Two-step return:

$$G_{t:t+2} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

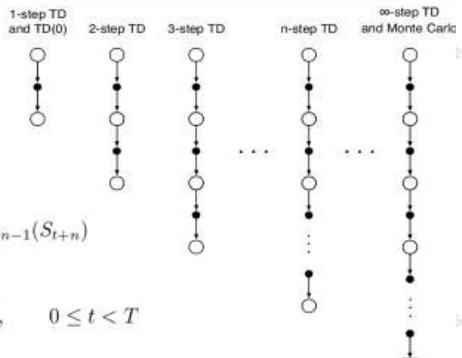
n-step return:

$$G_{t:t+n} \doteq R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

State-value learning algorithm for using n-step returns:

$$V_{t+n}(S_t) \doteq V_{t+n-1}(S_t) + \alpha [G_{t:t+n} - V_{t+n-1}(S_t)], \quad 0 \leq t < T$$

Ref Section 7.1 of TB



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

[Review] n-step TD Prediction

n-step TD for estimating $V \approx v_\pi$

Input: a policy π

Algorithm parameters: step size $\alpha \in (0, 1]$, a positive integer n

Initialize $V(s)$ arbitrarily, for all $s \in \mathcal{S}$

All store and access operations (for S_t and R_t) can take their index mod $n+1$

Loop for each episode:

 Initialize and store $S_0 \neq \text{terminal}$

$T \leftarrow \infty$

 Loop for $t = 0, 1, 2, \dots$:

 If $t < T$, then:

 Take an action according to $\pi(\cdot | S_t)$

 Observe and store the next reward as R_{t+1} and the next state as S_{t+1}

 If S_{t+1} is terminal, then $T \leftarrow t+1$

$\tau \leftarrow t - n + 1$ (τ is the time whose state's estimate is being updated)

 If $\tau \geq 0$:

$G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$

 If $\tau + n < T$, then: $G \leftarrow G + \gamma^n V(S_{\tau+n})$

$V(S_\tau) \leftarrow V(S_\tau) + \alpha [G - V(S_\tau)]$

($G_{\tau:T+n}$)

 Until $\tau = T - 1$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

n-step return

n-step semi-gradient TD for estimating $\hat{v} \approx v_\pi$

Input: the policy π to be evaluated
 Input: a differentiable function $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\hat{v}(\text{terminal}, \cdot) = 0$
 Parameters: step size $\alpha \in (0, 1]$, a positive integer n
 All store and access operations (S_t and R_t) can take their index mod n
 Initialize value-function weights \mathbf{w} arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)
 Repeat (for each episode):
 Initialize and store $S_0 \neq \text{terminal}$
 $T \leftarrow \infty$
 For $t = 0, 1, 2, \dots$:
 If $t < T$, then:
 Take an action according to $\pi(\cdot | S_t)$
 Observe and store the next reward as R_{t+1} and the next state as S_{t+1}
 If S_{t+1} is terminal, then $T \leftarrow t + 1$
 $\tau \leftarrow t - n + 1$ (τ is the time whose state's estimate is being updated)
 If $\tau \geq 0$:
 $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$
 If $\tau + n < T$, then: $G \leftarrow G + \gamma^n \hat{v}(S_{\tau+n}, \mathbf{w})$
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha [G - \hat{v}(S_\tau, \mathbf{w})] \nabla \hat{v}(S_\tau, \mathbf{w})$ (Gr:τ+n)
 Until $\tau = T - 1$

- Again, only a simple modification over the tabular setting

22

- Again, only a simple modification over the tabular setting

- Weight update instead of tabular entry update

23

n-step semi-gradient TD for estimating $\hat{v} \approx v_\pi$

Input: the policy π to be evaluated
 Input: a differentiable function $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\hat{v}(\text{terminal}, \cdot) = 0$
 Parameters: step size $\alpha \in (0, 1]$, a positive integer n
 All store and access operations (S_t and R_t) can take their index mod n
 Initialize value-function weights \mathbf{w} arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)
 Repeat (for each episode):
 Initialize and store $S_0 \neq \text{terminal}$
 $T \leftarrow \infty$
 For $t = 0, 1, 2, \dots$:
 If $t < T$, then:
 Take an action according to $\pi(\cdot | S_t)$
 Observe and store the next reward as R_{t+1} and the next state as S_{t+1}
 If S_{t+1} is terminal, then $T \leftarrow t + 1$
 $\tau \leftarrow t - n + 1$ (τ is the time whose state's estimate is being updated)
 If $\tau \geq 0$:
 $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$
 If $\tau + n < T$, then: $G \leftarrow G + \gamma^n \hat{v}(S_{\tau+n}, \mathbf{w})$ (Gr:τ+n)
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha [G - \hat{v}(S_\tau, \mathbf{w})] \nabla \hat{v}(S_\tau, \mathbf{w})$
 Until $\tau = T - 1$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [G - \hat{v}(S_\tau, \mathbf{w})] \nabla \hat{v}(S_\tau, \mathbf{w})$

Feature selection

- Assume a linear function approximator $\hat{v}(f(s); \theta) = f(s) \cdot \theta$
- What relevant features should represent states?

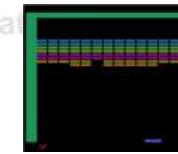
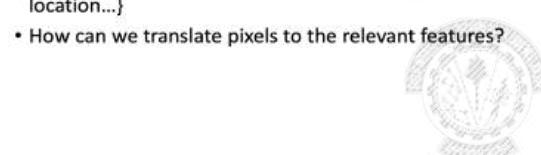


Features are domain depended requiring expert knowledge

24

Automatic features extraction

- Consider a game state that is given as a bit map
- Raw data of type $\text{pixel}(7, 3) = [0, 0, 0]$ (black)
- Desired features = {ball location, ball speed, ball direction, pan location...}
- How can we translate pixels to the relevant features?



25

Automatic features extraction for linear approximator

• Polynomials: $f_i(s) = \prod_{j=1}^k s_j^{c_{i,j}}$

- where each $c_{i,j}$ is an integer in the set $\{0, 1, \dots, n\}$ for an integer $n \geq 0$
- These features makeup the order- n polynomial basis for dimension k , which contains $(n + 1)^k$ different features

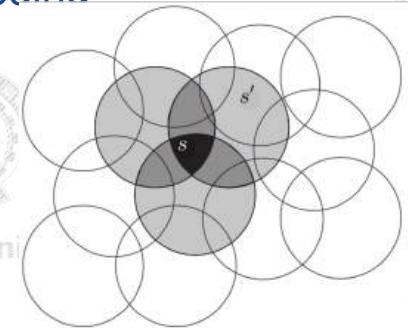
• Fourier Basis: $f_i(s) = \cos(\pi X^T c^i)$

- Where $c^i = (c_1^i, \dots, c_k^i)^T$, with $c_j^i \in \{0, \dots, n\}$ for $j = \{1, \dots, k\}$ and $i = \{0, \dots, (n + 1)^k\}$
- This defines a feature for each of the $(n + 1)^k$ possible integer vectors c^i
- The inner product $X^T c^i$ has the effect of assigning an integer in $\{0, \dots, n\}$ to each dimension of X
- This integer determines the feature's frequency along that dimension
- The features can be shifted and scaled to suit the bounded state space of a particular application

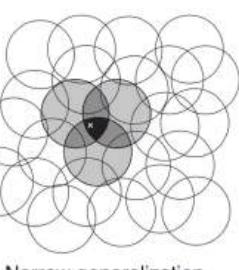
26

Automatic features extraction for linear approximator - Coarse Coding

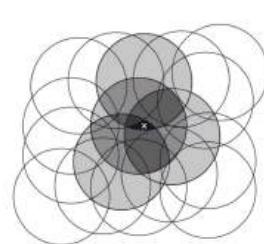
- Natural representation of the state set is continuous
- In 2-d, features corresponding to circles in state space
- Coding of a state:
 - If the state is inside a circle, then the corresponding feature has the value 1
 - otherwise the feature is 0
- Corresponding to each circle is a single weight (a component of w) that is learned
 - Training a state affects the weights of all the intersecting circles.



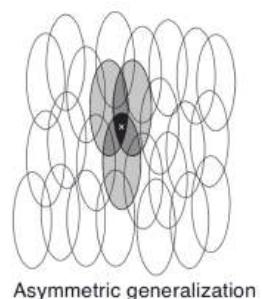
Automatic features extraction for linear approximator - Coarse Coding



Narrow generalization

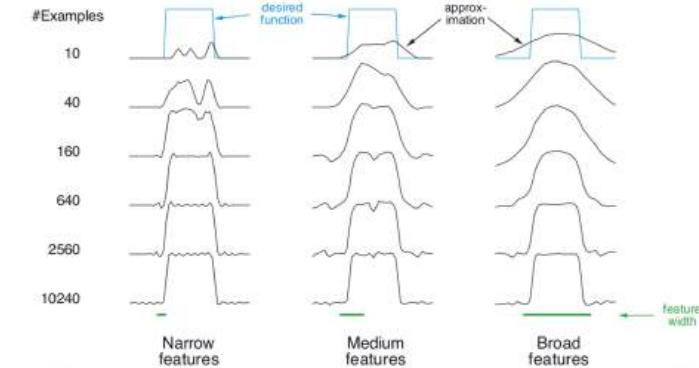


Broad generalization

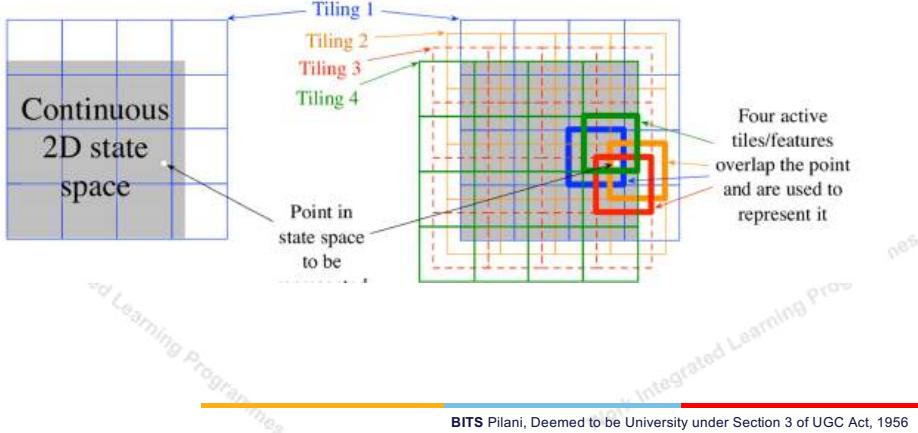


Asymmetric generalization

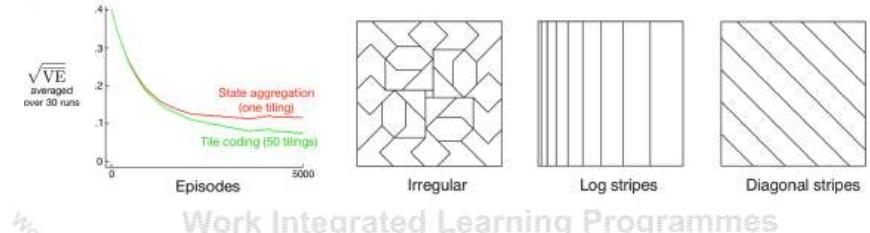
Automatic features extraction for linear approximator - Coarse Coding



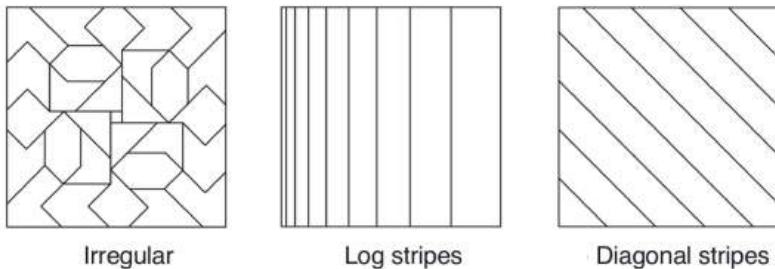
Automatic features extraction for linear approximator - Tile Coding



Automatic features extraction for linear approximator - Tile Coding

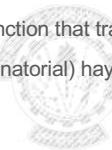


Automatic features extraction for linear approximator - Tile Coding



Automatic features extraction for linear approximator

- Other approaches include: Coarse Coding, Tile Coding, Radial Basis Functions (See chapter 9.5 in textbook)
- Each of these approaches defines a set of features, some useful yet most are not
 - E.g., is there a polynomial/Fourier function that translates pixels to pan location?
 - Probably but it's a needle in a (combinatorial) haystack
- Can we do better (generically)
 - Yes, using deep neural networks...



Work Integrated Learning Programmes



What did we learn?

- Reinforcement learning must generalize on observed experience if it is to be applicable to real world domains
- We can use parameterized function approximation to represent our knowledge about the domain state/action values
- Use stochastic gradient descend to update the tunable parameters such that the observed (TD, rollout) error is reduced
- When using a linear approximator, the Least squares TD method provides the most sample efficient approximation

Work Integrated Learning Programmes

34

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

T-Rex Chrome Dino Game

Feature selection

- Assume a linear function approximator $\hat{v}(f(s); \theta) = f(s) \cdot \theta$
- What relevant features should represent states?

→ State Space : Larger in Dimension

→ Feature Construction

Lets use coarse coding to construct feature for below sample states S1 & S2:

Features are domain dependent requiring expert knowledge

36

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

T-Rex Chrome Dino Game

→ State Space : Larger in Dimension

→ Action Space : { 0 - No Ops , 1 – Jump , 2 - Duck }

→ Choice : Value based function approximation

Require : Sample efficient learning with fast convergence

Choice : Off policy – Q learning

VS

Require : Online stable learning

Choice : On policy - SARSA

35

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

T-Rex Chrome Dino Game

Initialize Hyper parameters

→ State Space : Larger in Dimension

→ Lets use the assumed hyper parameters & value initializat double Q –learning!

→ By default let the policy select action based on ϵ -greedy. Assume that only actions {0,1} are

Initial weights:
 $\theta=[1.0, 0.5]$
 $\theta'=[0.5, 1.0]$
Discount $\gamma=0.2$
Learning rate $\alpha=0.9$
TD(n) : n=0

Features are domain dependent requiring expert knowledge

37

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

T-Rex Chrome Dino Game

Start the Episode

- State Space : Larger in Dimension
- Assume TD(0) 1st time step
- SARSA = (S1, action : 0, reward : 2points, S2, action : 0)

To estimate Q Value its desirable to encode the actions as a part of the feature vector! For simplicity lets append the action value with the state vector designed Eg., S1: $\phi(1,0,0)$ if no action is performed . S2: $\phi(1,1,0)$ if no action is performed S3: $\phi(1,1,1)$ if "jump" action is performed The parameters are modified to include the design

Initial weights:
 $\theta=[1,0,0.5,1,0]$
 $\theta'=[0.5,1,0,1,0]$

Discount $\gamma=0.2$
Learning rate $\alpha=0.9$
TD(n) : n=0

Lets use SARSA as well as Q-Learning to perform one iteration of GD on the parameters θ .

38

Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_\pi$

Input: a differentiable action-value function parameterization $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$
Algorithm parameters: step size $\alpha > 0$, small $\varepsilon > 0$
Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:
 $S, A \leftarrow$ initial state and action of episode (e.g., ε -greedy)
Loop for each step of episode:
Take action A , observe R, S'
If S' is terminal:
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$
Go to next episode
Choose A' as a function of $\hat{q}(S', \cdot, \mathbf{w})$ (e.g., ε -greedy)
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$
 $S \leftarrow S'$
 $A \leftarrow A'$

$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,a \sim p(\cdot), s' \sim \mathcal{E}} [(r + \gamma \hat{Q}(s', a'; \theta_i) - \hat{Q}(s, a; \theta_i)) \nabla_{\theta_i} \hat{Q}(s, a; \theta_i)]$

39

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Off Policy Approximation for Action Value

Original Q-learning

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

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'
 $\theta = \theta + \alpha (R + \gamma \max_a Q(S', a; \theta) - Q(S, A; \theta)) \nabla_{\theta} Q(S, A; \theta)$
 $S \leftarrow S'$
until S is terminal

40

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

On Policy Approximation for Action Value SARSA

Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_\pi$

Input: a differentiable action-value function parameterization $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$
Algorithm parameters: step size $\alpha > 0$, small $\varepsilon > 0$
Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:
 $S, A \leftarrow$ initial state and action of episode (e.g., ε -greedy)

Loop for each step of episode:
Take action A , observe R, S'

If S' is terminal:

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

Go to next episode

Choose A' as a function of $\hat{q}(S', \cdot, \mathbf{w})$ (e.g., ε -greedy)

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

$S \leftarrow S'$

$A \leftarrow A'$

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,a \sim p(\cdot), s' \sim \mathcal{E}} [(r + \gamma \hat{Q}(s', a'; \theta_i) - \hat{Q}(s, a; \theta_i)) \nabla_{\theta_i} \hat{Q}(s, a; \theta_i)]$$

mes

39

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

T-Rex Chrome Dino Game Estimate the VE: Value Error

→ SARSA = (S1, action : 0, reward : 2points, S2, action : 0)

→ Error = Target – Estimate

→ Current Estimate = $Q(s_1, a=0, \theta) = 1$

On policy – SARSA

Error = $r + \gamma Q(s', a'=0, \theta) - Q(s_1, a=0, \theta)$

$$= 2 + 0.2 \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1.0 \end{bmatrix} \begin{bmatrix} 1.0 \\ 0.5 \\ 0.5 \\ 1.0 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1.0 \end{bmatrix} \begin{bmatrix} 1.0 \\ 0.5 \\ 0.5 \\ 1.0 \end{bmatrix}$$

$$= 2.3 - 1 = 1.3$$

Vs Off policy – Q learning

Error = $r + \gamma \max \{$

$Q(s', a'=0, \theta),$

$Q(s', a'=1, \theta)$

$\}$

$- Q(s_1, a=0, \theta)$

$$= 2 + 0.2 \max \{ \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1.0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1.0 \end{bmatrix} \} - \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1.0 \end{bmatrix}$$

$$= 2 + 0.2 \max \{ 1.5, 2.5 \} - 1$$

$$= 2.5 - 1 = 1.5$$

$\}$

41

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

T-Rex Chrome Dino Game

Update the Model Parameters

- SARSA = (S1, action : 0, reward : 2points, S2, action : 0)
- Error = Target – Estimate
- Current Estimate = $Q(s_1, a=0, \theta) = 1$

On policy – SARSA

Error = 1.3

Update the parameter using SGD:

$$\theta \leftarrow \theta + \alpha \cdot \text{Error} \cdot \nabla \theta Q(s_1, a=0, \theta)$$

$$= \begin{bmatrix} 1.0 \\ 0.5 \\ 1.0 \end{bmatrix} + 0.9 \cdot 1.3 \cdot \begin{bmatrix} 1 \\ 0.5 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.17 \\ 0.5 \\ 1.0 \end{bmatrix}$$

Off policy – Q learning

Error = 1.5

Update the parameter using SGD:

$$\theta \leftarrow \theta + \alpha \cdot \text{Error} \cdot \nabla \theta Q(s_1, a=0, \theta)$$

$$= \begin{bmatrix} 1.0 \\ 0.5 \\ 1.0 \end{bmatrix} + 0.9 \cdot 1.5 \cdot \begin{bmatrix} 1 \\ 0.5 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.35 \\ 0.5 \\ 1.0 \end{bmatrix}$$

No internet

- Checking the networks cables, modem, and router
- Reconnecting to Wi-Fi
- Running Windows Network Diagnostics

Initial weights:
 $\theta = [1, 0, 0.5, 1, 0]$
 $\theta = [0, 1, 0, 3, 1, 0]$

Discount $\gamma=0.2$
Learning rate $\alpha=0.5$
 $TD(n) : n=0$
S1: $\phi(1, 0, 0)$
S2: $\phi(1, 1, 0)$

VS

42

T-Rex Chrome Dino Game

Optimize the Performance

- SARSA = (S1, action : 0, reward : 2points, S2, action : 0)
- Error = Target – Estimate
- Current Estimate = $Q(s_1, a=0, \theta) = 1$

Off policy – Q learning

Error = 1.5

Update the parameter using SGD:

$$\theta \leftarrow \theta + \alpha \cdot \text{Error} \cdot \nabla \theta Q(s_1, a=0, \theta)$$

$$= \begin{bmatrix} 1.0 \\ 0.5 \\ 1.0 \end{bmatrix} + 0.9 \cdot 1.5 \cdot \begin{bmatrix} 1 \\ 0.5 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.35 \\ 0.5 \\ 1.0 \end{bmatrix}$$

No internet

- Checking the networks cables, modem, and router
- Reconnecting to Wi-Fi
- Running Windows Network Diagnostics

Initial weights:
 $\theta = [1, 0, 0.5, 1, 0]$
 $\theta = [0, 1, 0, 3, 1, 0]$

Discount $\gamma=0.2$
Learning rate $\alpha=0.5$
 $TD(n) : n=0$
S1: $\phi(1, 0, 0)$
S2: $\phi(1, 1, 0)$

Understanding Maximization bias

Observation:
Though it may try to converge faster, it is biased due to MAX function resulting in Maximization bias

43

T-Rex Chrome Dino Game

Optimize the Performance

- SARSA = (S1, action : 0, reward : 2points, S2, action : 0)
- Error = Target – Estimate
- Current Estimate = $Q(s_1, a=0, \theta) = 1$

Understanding Maximization bias

True value of $Q(S_2, 0) = Q(S_2, 1) = 2.0$

The $\theta=[1, 0, 0.5, 1, 0]$ produces noisy estimates for :
 $Q(S_2, 0) = 1.5 & Q(S_2, 1) = 2.5$, resulting in 2.5 as target value which is an overestimation due to bias.

But given above true value target ~2.1 is expected.

Solution: Decouple the action selection from action evaluation using two independent models!!

No internet

- Checking the networks cables, modem, and router
- Reconnecting to Wi-Fi
- Running Windows Network Diagnostics

Initial weights:
 $\theta = [1, 0, 0.5, 1, 0]$
 $\theta = [0, 1, 0, 3, 1, 0]$

Discount $\gamma=0.2$
Learning rate $\alpha=0.5$
 $TD(n) : n=0$
S1: $\phi(1, 0, 0)$
S2: $\phi(1, 1, 0)$

Off policy – Q learning

Error = $r + \gamma \text{MAX} \{$

$$Q(s_2', a'=0, \theta),$$

$$Q(s_2', a'=1, \theta)$$

$$\}$$

$$- Q(s_1, a=0, \theta)$$

$$= 2 + 0.2 \text{MAX} \{ \begin{bmatrix} 1 \\ 1 \\ 0.5 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 0.5 \\ 1 \\ 1 \end{bmatrix} \} - \begin{bmatrix} 1 \\ 0 \\ 0.5 \\ 1 \\ 1 \end{bmatrix}$$

$$= 2 + 0.2 \text{MAX}[1.5, 2.5] - 1$$

$$= 2.5 - 1 = 1.5$$

VS

44

T-Rex Chrome Dino Game

Optimize the Performance

- SARSA = (S1, action : 0, reward : 2points, S2, action : 0)
- Error = Target – Estimate
- Current Estimate = $Q(s_1, a=0, \theta) = 1$

Understanding Maximization bias

True value of $Q(S_2, 0) = Q(S_2, 1) = 2.0$

The $\theta=[1, 0, 0.5, 1, 0]$ produces noisy estimates for :
 $Q(S_2, 0) = 1.5 & Q(S_2, 1) = 2.5$, resulting in 2.5 as target value which is an overestimation due to bias.

But given above true value target ~2.1 is expected.

Solution: Decouple the action selection from action evaluation using two independent models!!

No internet

- Checking the networks cables, modem, and router
- Reconnecting to Wi-Fi
- Running Windows Network Diagnostics

Initial weights:
 $\theta = [1, 0, 0.5, 1, 0]$
 $\theta = [0, 1, 0, 3, 1, 0]$

Discount $\gamma=0.2$
Learning rate $\alpha=0.5$
 $TD(n) : n=0$
S1: $\phi(1, 0, 0)$
S2: $\phi(1, 1, 0)$

Off policy – Double Q - Learning

Error = $r + \gamma \text{MAX} \{$

$$Q(s_2', a'=0, \theta), Q(s_2', a'=1, \theta)$$

$$\}$$

$$- Q(s_1, a=0, \theta)$$

Instead of **MAX** function in above target,
✓ Use **θ** for action a' selection in the next state
✓ Use **θ'** for estimating $Q(s_2', a', \theta')$
✓ Role of model can be reversed

45

T-Rex Chrome Dino Game

Optimize the Performance

→ SARSA = (S1, action : 0, reward : 2points, S2, action : 0)
 → Error = Target – Estimate
 → Current Estimate = Q(s1,a=0,θ) = 1

Effect of the Double Q - Learning
 True value of Q(S2,0) = Q(S2,1) = 2.0

Use θ for action a' selection in the next state S2
 : **ARGMAX** { Q(s2',a'=0,θ), Q(s2',a'=1,θ) }
 $= \text{ARGMAX} \left\{ \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix} \right\}$
 $= \text{ARGMAX} \{ 1.5, 2.5 \} \rightarrow \text{Action a1}$

No internet
 Tip:
 • Checking the networks cables, modem, and router
 • Reconnecting to Wi-Fi
 • Running Windows Network Diagnostics

Initial weights:
 $\theta = [1, 0, 0, 5, 1, 0]$
 $\theta' = [0, 1, 0, 3, 1, 0]$

Discount γ=0.2
 Learning rate α=0.5
 TD(n) : n=0
 S1: $\phi(1,0,0)$
 S2: $\phi(1,1,0)$

Off policy – Double Q - Learning

Error = $r + \gamma \max \{ Q(s',a',\theta), Q(s',a',\theta') \}$

Instead of MAX function in above target,
 Use θ for action a' selection in the next state S2
 ✓ Use θ' for estimating $Q(s',a',\theta')$
 ✓ Role of model can be reversed

T-Rex Chrome Dino Game

Optimize the Performance

→ SARSA = (S1, action : 0, reward : 2points, S2, action : 0)
 → Error = Target – Estimate
 → Current Estimate = Q(s1,a=0,θ) = 1

Effect of the Double Q - Learning
 True value of Q(S2,0) = Q(S2,1) = 2.0

Use θ for action a' selection in the next state S2:
 Action a1

Use θ for action a' selection in the next state S2:
 Instead of MAX function in above target,
 Use θ for action a' selection in the next state S2
 ✓ Use θ' for estimating $Q(s',a',\theta')$
 ✓ Role of model can be reversed

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

T-Rex Chrome Dino Game

Optimize the Performance

→ SARSA = (S1, action : 0, reward : 2points, S2, action : 0)
 → Error = Target – Estimate
 → Current Estimate = Q(s1,a=0,θ) = 1

Understanding Maximization bias
 True value of Q(S2,0) = Q(S2,1) = 0.5

Use θ for action a' selection in the next state S2:
 Action a1

Use θ' for estimating $Q(s',a',\theta') = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.3 \end{bmatrix} = 1.4$

No internet
 Tip:
 • Checking the networks cables, modem, and router
 • Reconnecting to Wi-Fi
 • Running Windows Network Diagnostics

Initial weights:
 $\theta = [1, 0, 0, 5, 1, 0]$
 $\theta' = [0, 1, 0, 3, 1, 0]$

Discount γ=0.2
 Learning rate α=0.5
 TD(n) : n=0
 S1: $\phi(1,0,0)$
 S2: $\phi(1,1,0)$

Off policy – Double Q - Learning

Error = $r + \gamma Q(s',a') - \arg\max(Q(s',a',\theta))$
 $- Q(s',a',\theta)$

$= 2 + 0.2 (1.4) - \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$
 $= 2.28 - 1$
 $= 1.28$

Update the parameter using SGD:
 $\theta \leftarrow \theta + \alpha \cdot \text{Error} \cdot \nabla \theta Q(s, a, \theta)$

$= \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix} + 0.9 * 1.28 * \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.152 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$

T-Rex Chrome Dino Game

Optimize the Performance

→ SARSA = (S1, action : 0, reward : 2points, S2, action : 0)
 → Error = Target – Estimate
 → Current Estimate = Q(s1,a=0,θ) = 1

Every consecutive sequence of states in this game is mostly highly correlated!! And if non-linear function approximated is used off policy learnt parameter this affects with diverging updates leading to unstable learning.

Solution : Use Experience Replay to store & use past experience (ie., SARSA). Instead of SGD use , mini Batch GD by randomly sampling from the buffer. This can be used in DQN as well as DDQN.

Off policy – Double Q - Learning

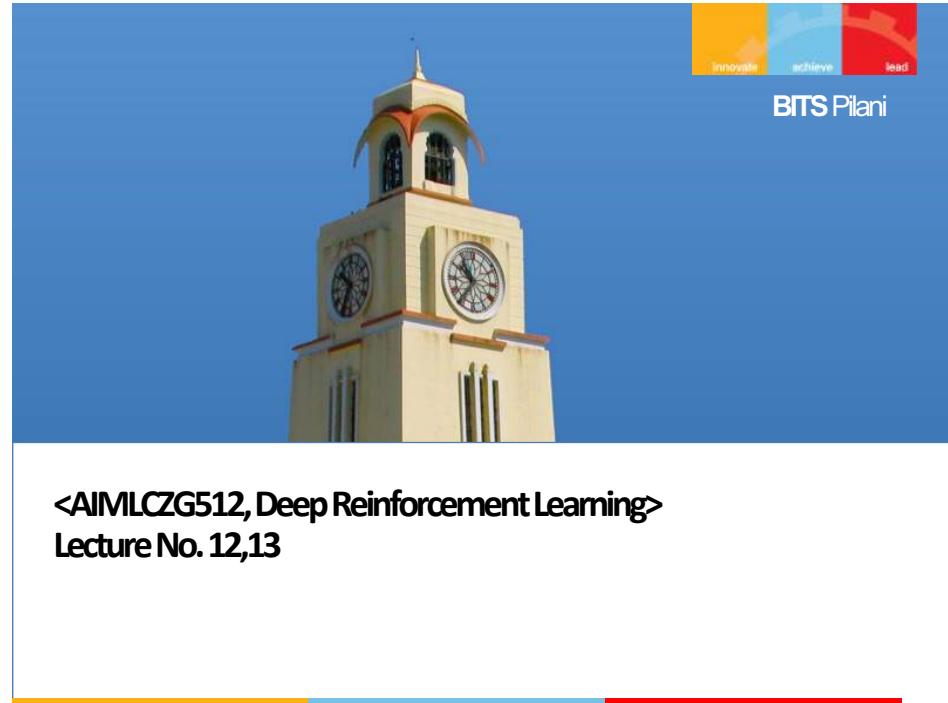
Error = $r + \gamma Q(s',a') - \arg\max(Q(s',a',\theta'))$
 $- Q(s',a',\theta)$

$= 2 + 0.2 (1.4) - \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1.0 \\ 1.0 \end{bmatrix}$
 $= 2.28 - 1$
 $= 1.28$

Update the parameter using SGD:
 $\theta \leftarrow \theta + \alpha \cdot \text{Error} \cdot \nabla \theta Q(s, a, \theta)$

$= \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix} + 0.9 * 1.28 * \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.152 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Mnih et al. 2015

- First deep learning model to successfully learn control policies directly from high-dimensional sensory input using reinforcement learning
- The model is a convolutional neural network, trained with a variant of Q-learning
- Input is raw pixels and output is an action-value function estimating future rewards
- Surpassed a human expert on various Atari video games

The age of deep learning

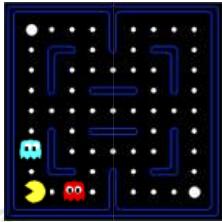
- Previous models relied on hand-crafted features combined with linear value functions
 - The performance of such systems heavily relies on the quality of the feature representation
- Advances in deep learning have made it possible to automatically extract high-level features

A row of five screenshots from classic Atari games, showing the transition from raw pixel input to processed feature maps. The games include Asteroids, Breakout, Ms. Pac-Man, Space Invaders, and Zaxxon.

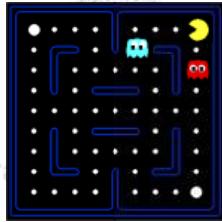


Example: Pacman

Let's say we discover through experience that this state is bad:



In naïve q-learning, we know nothing about this state:



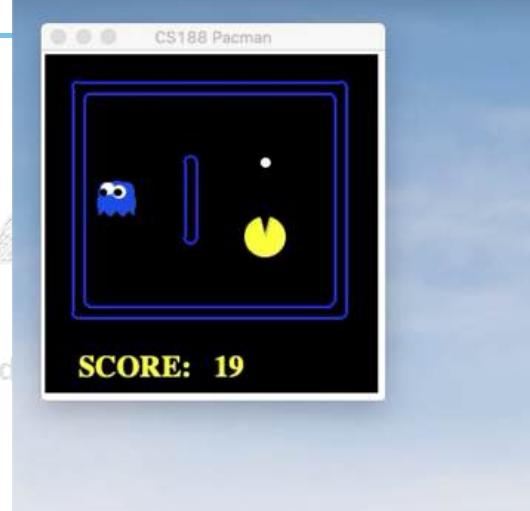
Or even this one!



We must generalize our knowledge!

5

- Generalize Q-learning with function approximator

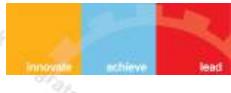


- Naïve Q-learning
- After 50 training episodes



Generalizing with Deep learning

- **Supervised:** Require large amounts of hand-labelled training data
 - In RL on the other hand, learns from a scalar reward signal that is frequently sparse, noisy, and delayed
- **Supervised:** Assume the data samples are independent
 - In RL one typically encounters sequences of highly correlated state
- **Supervised:** Assume a fixed underlying distribution
 - In RL the data distribution changes as the algorithm learns new behaviors
- DQN was first to demonstrate that a convolutional neural network can overcome these challenges to learn successful control policies from raw video data in complex RL environments



Deep Q learning [Mnih et al. 2015]

- Trains a generic neural network-based agent that successfully learns to operate in as many domains as possible
- The network is not provided with any domain-specific information or hand-designed features
- Must learn from nothing but the raw input (pixels), the reward, terminal signals, and the set of possible actions

9

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Deep Q learning [Mnih et al. 2015]

- DQN addresses problems of correlated data and non-stationary distributions
 - Use an experience replay mechanism
 - Randomly samples and trains on previous transitions
 - Results in a smoother training distribution over many past behaviors

11

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Original Q-learning

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

```

Repeat (for each episode):
  Initialize  $S$ 
  Repeat (for each step of episode):
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Take action  $A$ , observe  $R, S'$ 
     $\theta = \theta + \alpha \left( R + \gamma \max_a Q(S', a; \theta) - Q(S, A; \theta) \right) \nabla_{\theta} Q(S, A; \theta)$ 
     $S \leftarrow S'$ 
  until  $S$  is terminal
  
```

10

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Deep Q learning [Mnih et al. 2015]

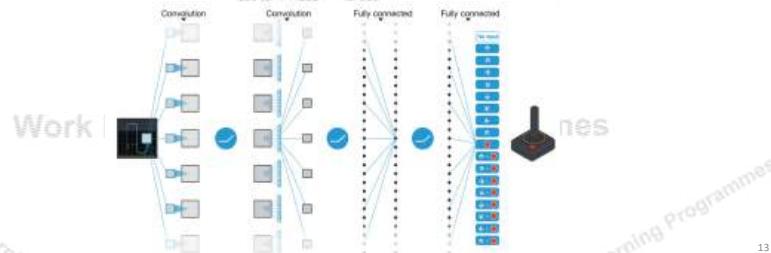
- Learn a function approximator (Q network), $\hat{Q}(s, a; \theta) \approx Q^*(s, a)$
- Value propagation: $Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a')]$
 - \mathcal{E} represents the environment (underlying MDP)
- Update $\hat{Q}(s, a; \theta)$ at each step i using SGD with squared loss:
- $L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} \left[(y_i - \hat{Q}(s, a; \theta_i))^2 \right]$
 - $y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} \hat{Q}(s', a'; \theta_{i-1})]$
 - $\rho(s, a)$ is the behavior distribution, e.g., epsilon greedy
 - θ_{i-1} is considered fix when optimizing the loss function (helps when taking the derivative with respect to θ and with stability)

12

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Deep Q learning [Mnih et al. 2015]

- $L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot)} \left[(y_i - \hat{Q}(s, a; \theta_i))^2 \right]$ Independent of θ_i (because θ_{i-1} is considered fix)
- $\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot), s' \sim \mathcal{E}} [(r + \gamma \max_{a'} \hat{Q}(s', a'; \theta_{i-1}) - \hat{Q}(s, a; \theta_i)) \nabla_{\theta_i} \hat{Q}(s, a; \theta_i)]$



13

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

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 Q̂ with weights θ⁻ = θ
For episode = 1, M do
    Initialize sequence s₁ = {x₁} and preprocessed sequence φ₁ = φ(s₁)
    For t = 1, T do
        With probability ε select a random action a,
        otherwise select a_t = argmax_a Q(φ(s_t), a; θ)
        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 φ_{t+1} = φ(s_{t+1})
        Store transition (φ_t, a_t, r_t, φ_{t+1}) in D
        Sample random minibatch of transitions (φ_j, a_j, r_j, φ_{j+1}) from D
        Set y_j = {
            r_j
            if episode terminates at step j+1
            r_j + γ max_a' Q(φ_{j+1}, a'; θ⁻) otherwise
        }
        Perform a gradient descent step on (y_j - Q(φ_j, a_j; θ))^2 with respect to the
        network parameters θ
        Every C steps reset Q̂ = Q
    End For
End For
Play m episodes (full games)
  
```

15

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

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 Q̂ with weights θ⁻ = θ
For episode = 1, M do
    Initialize sequence s₁ = {x₁} and preprocessed sequence φ₁ = φ(s₁)
    For t = 1, T do
        With probability ε select a random action a,
        otherwise select a_t = argmax_a Q(φ(s_t), a; θ)
        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 φ_{t+1} = φ(s_{t+1})
        Store transition (φ_t, a_t, r_t, φ_{t+1}) in D
        Sample random minibatch of transitions (φ_j, a_j, r_j, φ_{j+1}) from D
        Set y_j = {
            r_j
            if episode terminates at step j+1
            r_j + γ max_a' Q(φ_{j+1}, a'; θ⁻) otherwise
        }
        Perform a gradient descent step on (y_j - Q(φ_j, a_j; θ))^2 with respect to the
        network parameters θ
        Every C steps reset Q̂ = Q
    End For
End For
  
```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

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 Q̂ with weights θ⁻ = θ
For episode = 1, M do
    Initialize sequence s₁ = {x₁} and preprocessed sequence φ₁ = φ(s₁)
    For t = 1, T do
        With probability ε select a random action a,
        otherwise select a_t = argmax_a Q(φ(s_t), a; θ)
        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 φ_{t+1} = φ(s_{t+1})
        Store transition (φ_t, a_t, r_t, φ_{t+1}) in D
        Sample random minibatch of transitions (φ_j, a_j, r_j, φ_{j+1}) from D
        Set y_j = {
            r_j
            if episode terminates at step j+1
            r_j + γ max_a' Q(φ_{j+1}, a'; θ⁻) otherwise
        }
        Perform a gradient descent step on (y_j - Q(φ_j, a_j; θ))^2 with respect to the
        network parameters θ
        Every C steps reset Q̂ = Q
    End For
End For
  
```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Q-learning with experience replay

```

Initialize replay memory  $D$  to capacity  $N$ 
Initialize action-value function  $\hat{Q}$  with random weights  $\theta$ 
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  $\epsilon$  select a random action  $a_t$ ,
        otherwise select  $a_t = \operatorname{argmax}_a \hat{Q}(\phi(s_t), a; \theta)$ 
        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  $(y_j - \hat{Q}(\phi_j, a_j; \theta))^2$  with respect to the
        network parameters  $\theta$ 
        Every  $C$  steps reset  $\hat{Q} = Q$ 
    End For
End For

```

For each time step during the episode

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

17

Q-learning with experience replay

```

Initialize replay memory  $D$  to capacity  $N$ 
Initialize action-value function  $\hat{Q}$  with random weights  $\theta$ 
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  $\epsilon$  select a random action  $a_t$ ,
        otherwise select  $a_t = \operatorname{argmax}_a \hat{Q}(\phi(s_t), a; \theta)$ 
        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  $(y_j - \hat{Q}(\phi_j, a_j; \theta))^2$  with respect to the
        network parameters  $\theta$ 
        Every  $C$  steps reset  $\hat{Q} = Q$ 
    End For
End For

```

With small probability select a random action (explore), otherwise select the, currently known, best action (exploit).

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

18

Q-learning with experience replay

```

Initialize replay memory  $D$  to capacity  $N$ 
Initialize action-value function  $\hat{Q}$  with random weights  $\theta$ 
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  $\epsilon$  select a random action  $a_t$ ,
        otherwise select  $a_t = \operatorname{argmax}_a \hat{Q}(\phi(s_t), a; \theta)$ 
        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  $(y_j - \hat{Q}(\phi_j, a_j; \theta))^2$  with respect to the
        network parameters  $\theta$ 
        Every  $C$  steps reset  $\hat{Q} = Q$ 
    End For
End For

```

Execute the chosen action and store the (processed) observed transition in the replay memory

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

19

Q-learning with experience replay

```

Initialize replay memory  $D$  to capacity  $N$ 
Initialize action-value function  $\hat{Q}$  with random weights  $\theta$ 
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  $\epsilon$  select a random action  $a_t$ ,
        otherwise select  $a_t = \operatorname{argmax}_a \hat{Q}(\phi(s_t), a; \theta)$ 
        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  $(y_j - \hat{Q}(\phi_j, a_j; \theta))^2$  with respect to the
        network parameters  $\theta$ 
        Every  $C$  steps reset  $\hat{Q} = Q$ 
    End For
End For

```

Experience replay:
Sample a random minibatch of transitions from replay memory and perform gradient decent step on Q (not on \hat{Q})

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

20



Q-learning with experience replay

```

Initialize replay memory  $D$  to capacity  $N$ 
Initialize action-value function  $\hat{Q}$  with random weights  $\theta$ 
Initialize target action-value function  $\tilde{Q}$  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  $\epsilon$  select a random action  $a_t$ ,
        otherwise select  $a_t = \text{argmax}_a \hat{Q}(\phi(s_t), a; \theta)$ 
        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  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_a \tilde{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - \hat{Q}(\phi_j, a_j; \theta))^2$  with respect to the
        network parameters  $\theta$ 
        Every  $C$  steps reset  $\tilde{Q} = \hat{Q}$ 
    End For
End For

```

Once every several steps set the
target function, \tilde{Q} , to equal \hat{Q}

21

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Q-learning with experience replay

```

Initialize replay memory  $D$  to capacity  $N$ 
Initialize action-value function  $\hat{Q}$  with random weights  $\theta$ 
Initialize target action-value function  $\tilde{Q}$  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  $\epsilon$  select a random action  $a_t$ ,
        otherwise select  $a_t = \text{argmax}_a \hat{Q}(\phi(s_t), a; \theta)$ 
        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  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_a \tilde{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - \hat{Q}(\phi_j, a_j; \theta))^2$  with respect to the
        network parameters  $\theta$ 
        Every  $C$  steps reset  $\tilde{Q} = \hat{Q}$ 
    End For
End For

```

Such delayed online learning helps in practice:
 "This modification makes the algorithm more stable compared to standard online Q-learning, where an update that increases $Q(s_t, a_t)$ often also increases $Q(s_{t+1}, a)$ for all a and hence also increases the target y_j , possibly leading to oscillations or divergence of the policy" (Human-level control through deep reinforcement learning. Nature 518:7540 (2015): 529.)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Deep Q learning

- model-free:** solves the reinforcement learning task directly using samples from the emulator \mathcal{E} , without explicitly learning an estimate of E
- off-policy:** learns the optimal policy, $a = \text{argmax}_a Q(s, a; \theta)$, while following a different behavior policy
 - One that ensures adequate exploration of the state space



Work Integrated Learning Programmes

23

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Experience replay

- Utilizing experience replay has several advantages
 - Each step of experience is potentially used in many weight updates, which allows for greater data efficiency
 - Learning directly from consecutive samples is inefficient, due to the strong correlations between the samples; randomizing the samples breaks these correlations and therefore reduces the variance of the updates
 - The behavior distribution is averaged over many of its previous states, smoothing out learning and avoiding oscillations or divergence in the parameters
- Note that when learning by experience replay, it is necessary to learn off-policy (because our current parameters are different to those used to generate the sample), which motivates the choice of Q-learning



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

24

Experience replay

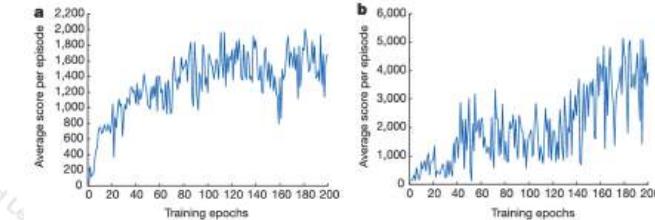
- DQN only stores the last N experience tuples in the replay memory
 - Old transitions are overwritten
- Samples uniformly at random from the buffer when performing updates
- Is there room for improvement?
 - Important transitions?
 - Prioritized sweeping
 - Prioritize deletions from the replay memory
 - see prioritized experience reply, <https://arxiv.org/abs/1511.05952>

25

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Results: DQN

- Each point is the average score achieved per episode after the agent is run with an epsilon-greedy policy ($\epsilon = 0.05$) for 520k frames on SpaceInvaders (a) and Seaquest (b)

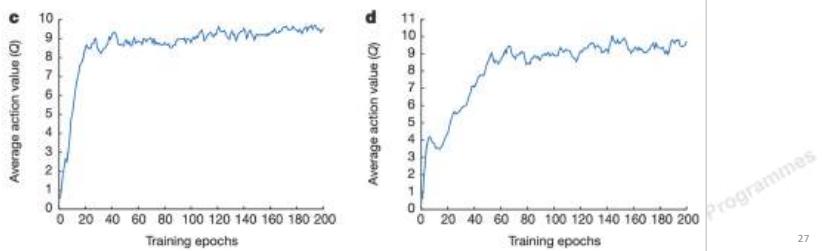


26

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

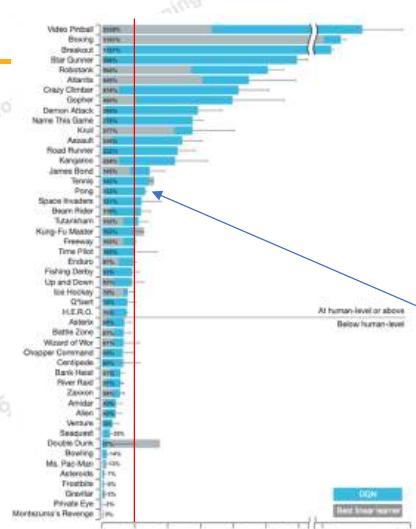
Results: DQN

- Average predicted action-value on a held-out set of states on Space Invaders (c) and Seaquest (d)



27

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



28

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

- Normalized between a professional human games tester (100%) and random play (0%)
- E.g., in Pong, DQN achieved a factor of 1.32 higher score on average when compared to a professional human player

Maximization bias

- See lecture 5TD_learning.pptx, slide 41
- The max operator in Q-learning uses the same values both to select and to evaluate an action
 - $y_i = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_a \hat{Q}(s', a; \theta_{i-1}) \right]$
 - Makes it more likely to select overestimated values, resulting in overoptimistic value estimates
- We can use a technique similar to the previously discussed Double Q-learning (lecture 5TD_learning.pptx, slide 43)
 - Two value functions are trained. Update only one of the two value functions at each step while considering the max from the other

29

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Double Deep Q networks (DDQN)

```

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  $a_t = \arg\max_a Q(\phi(s_t), a; \theta)$ 
    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 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  $\theta$ 
    Every C steps reset  $\hat{Q} = Q$ 
  End For
End For

```

We have two available Q networks. Let's use them for double learning

30

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Double Deep Q networks (DDQN)

```

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  $a_t = \arg\max_a Q(\phi(s_t), a; \theta)$ 
    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 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  $\theta$ 
    Every C steps reset  $\hat{Q} = Q$ 
  End For
End For

```

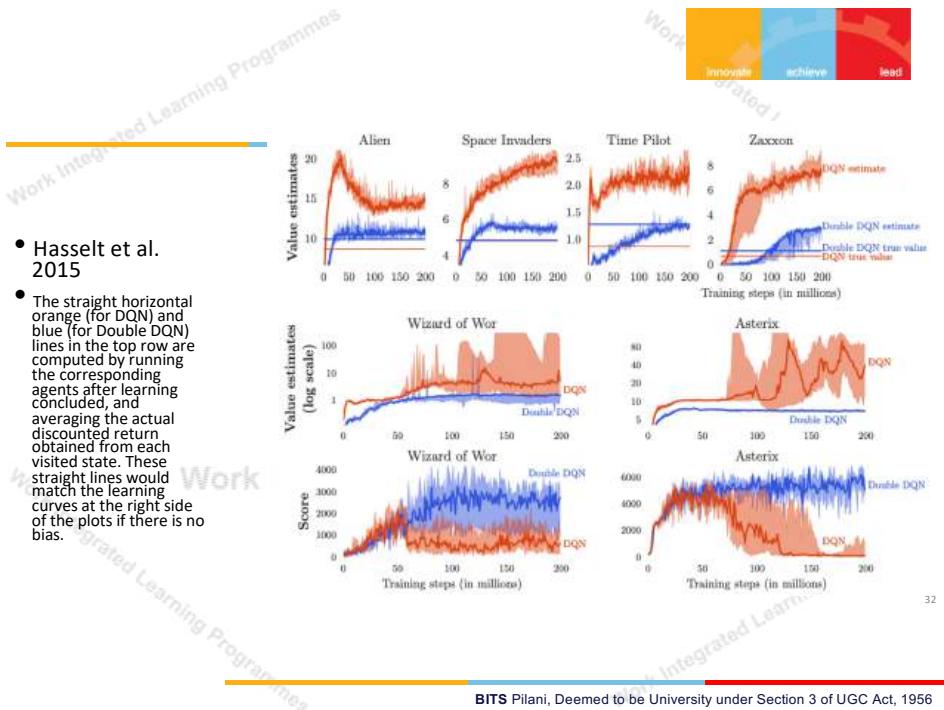
We have two available Q networks. Let's use them for double learning

31

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

• Hasselt et al. 2015

- The straight horizontal orange (for DQN) and blue (for Double DQN) lines in the top row are computed by running the corresponding agents after learning concluded, and averaging the actual discounted return obtained from each visited state. These straight lines would match the learning curves at the right side of the plots if there is no bias.



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

32

What did we learn?

- Using deep neural networks as function approximators in RL is tricky
 - Sparse samples
 - Correlated samples
 - Evolving policy (nonstationary sample distribution)
- DQN attempts to address these issues
 - Reuse previous transitions at each training (SGD) step
 - Randomly sample previous transitions to break correlation
 - Use off-policy, TD(0) learning to allow convergence to the true target values (Q^*)
 - No guarantees for non-linear (DNN) approximators

33

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Deep Reinforcement Learning

Prof. Vimal SP
CSIS

BITS Pilani



<AIMLCZG512, Deep Reinforcement Learning>
Lecture No. 13,14

Agenda for the classes

- Introduction
- Policy gradients
- REINFORCE algorithm
- Actor-critic methods
- REINFORCE - example

Work Integrated Learning Programmes

Acknowledgements: Some of the slides were adopted *with permission* from the course CSCE-689 (Texas A&M University) by Prof. Gunil Sharon

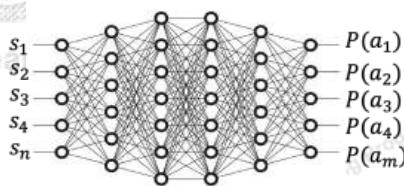
3

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Notation

- The policy is a parametrized function: $\pi_\theta(a|s)$
 - For policy gradient we need a continuous, differentiable policy... (Softmax activation can be useful for discrete action spaces)
 - $\pi(a|s)$ is assumed to be differentiable with respect to θ
 - E.g., a DNN where θ is the set of weights and biases
- $J(\theta)$ is a scalar policy performance measure (expected sum of discounted rewards) with respect to the policy params



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Improving the policy

- SGD: $\theta_{t+1} = \theta_t + \alpha \nabla J(\theta_t)$
 - Where $\nabla J(\theta_t)$ is a stochastic estimate, whose expectation approximates the gradient of the performance measure
- All methods that follow this general schema we call policy gradient methods
 - Might also learn an approximate value function for normalization (baseline)
- Methods that use approximations to both policy and value functions for computing the policy's gradient are called actor-critic methods (more on this later)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

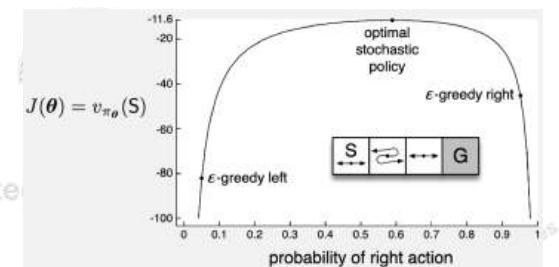
Advantages of PG

- The policy converges over time as opposed to an epsilon-greedy value-based approach
- Naturally applies to continuous action space as opposed to a Q learning approach
- In many domains the policy is a simpler function to approximate
 - Though this is not always the case
- Choice of policy parameterization is sometimes a good way of injecting prior knowledge
 - E.g., in phase assignment by a traffic controller
- Can converge on stochastic optimal policies, as opposed to value-based approaches
 - Useful in games with imperfect information where the optimal play is often to do two different things with specific probabilities, e.g., bluffing in Poker

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Stochastic policy - Example – Short corridor with switched actions

- Reward = -1 per step
- $\mathcal{A} = \{\text{left, right}\}$
- Left goes left and right goes right except in the middle state where they are reversed
- States are identical from the policy's perspective
- $\pi^* = [0.41 \quad 0.59]$



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Evaluate the gradient in performance

- $\widehat{\nabla_{\theta} J(\theta)} = ?$
- J depends on both the action selections and the distribution of states in which those selections are made
 - Both of these are affected by the policy parameter θ
- Seems impossible to solve without knowing the transition function (or the distribution of visited states)
 - $p(s'|s, a)$ is unknown in model free RL
- The PG theorem allows us to evaluate $\nabla_{\theta} J(\theta)$ without the need for $p(s'|s, a)$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Monte-Carlo Policy Gradient

- Sample-based gradient estimation
- $\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla \pi(a|s)$
- $\sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla \pi(a|s) = \mathbb{E}_{\pi}[\sum_a q_{\pi}(S_t, a) \nabla_{\theta} \pi(a|S_t)]$
- We can now use this gradient estimation for PG steps
 - $\theta_{t+1} = \theta_t + \alpha \sum_a \widehat{q}_{\pi}(S_t, a) \nabla_{\theta} \pi(a|S_t; \theta)$
 - Requires an approximation to q_{π} , denoted \widehat{q} , for every possible action
 - Can we avoid this approximation?

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



REINFORCE [Williams, 1992]

- $\theta_{t+1} = \theta_t + \alpha \sum_a \widehat{q}_{\pi}(S_t, a) \nabla_{\theta} \pi(a|S_t)$
 - Can we avoid this approximation? YES!
 - $\nabla J(\theta) \propto \mathbb{E}_{\pi}[\sum_a q_{\pi}(S_t, a) \nabla_{\theta} \pi(a|S_t)]$
 - $= \mathbb{E}_{\pi} \left[\sum_a \pi(a|S_t; \theta) q_{\pi}(S_t, a) \frac{\nabla_{\theta} \pi(a|S_t)}{\pi(a|S_t)} \right]$ (multiplied and divided by the same number)
 - $= \mathbb{E}_{\pi} \left[q_{\pi}(S_t, A_t) \frac{\nabla_{\theta} \pi(A_t|S_t)}{\pi(A_t|S_t)} \right] (\sum_x \Pr(x) f(x) = \mathbb{E}_{x \sim \Pr(x)}[f(x)])$
 - $= \mathbb{E}_{\pi} \left[G_t \frac{\nabla_{\theta} \pi(A_t|S_t)}{\pi(A_t|S_t)} \right] (\mathbb{E}_{\pi}[G_t|S_t, A_t] = q_{\pi}(S_t, A_t))$
- We can now use this gradient estimation for PG steps
 - $\theta_{t+1} = \theta_t + \alpha G_t \frac{\nabla_{\theta} \pi(A_t|S_t)}{\pi(A_t|S_t)}$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



REINFORCE [Williams, 1992]

- $\theta_{t+1} = \theta_t + \alpha G_t \frac{\nabla_{\theta} \pi(A_t|S_t; \theta)}{\pi(A_t|S_t; \theta)}$
- Each increment is proportional to the product of a return G_t and a vector
 - The gradient of the probability of taking the action actually taken over the (current) probability of taking that action
 - The vector is the direction in parameter space that most increases the probability of repeating the action A_t on future visits to state S_t
- The update increases the parameter vector (*) proportional to the return, and (**) inversely proportional to the action probability
 - (*) makes sense because it causes the parameter to move most in the directions that favor actions that yield the highest return
 - (**) makes sense because otherwise actions that are selected frequently are at an advantage (the updates will be more often in their direction) and might win out even if they do not yield the highest return

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

REINFORCE [Williams, 1992]

$$\theta_{t+1} = \theta_t + \alpha G_t \frac{\nabla_\theta \pi(A_t|S_t; \theta)}{\pi(A_t|S_t; \theta)}$$

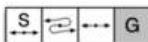
- REINFORCE uses G_t : the complete return from time t until the end of the episode
- In this sense REINFORCE is a Monte Carlo algorithm and is well defined only for the episodic case with all updates made in retrospect after the episode is completed

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

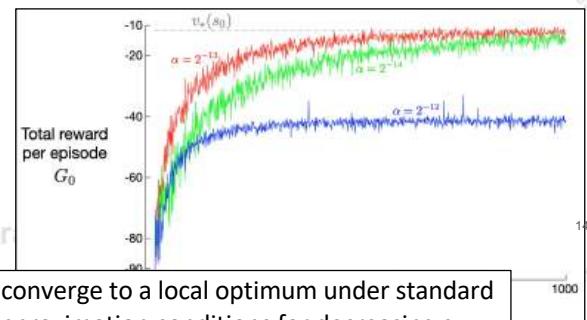
REINFORCE [Williams, 1992]

- The crazy corridor domain



- With the right step size, the total reward per episode approaches the optimal value of the start state

Guaranteed to converge to a local optimum under standard stochastic approximation conditions for decreasing α



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

REINFORCE [Williams, 1992]

* In the literature you might encounter "grad log pi" instead of "grad ln pi". That's not a problem since: $\nabla \ln(f(x)) \propto \nabla \log(f(x))$
Specifically: $\nabla \ln(f(x)) = \nabla \log_a(f(x)) \ln(a)$

REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$

Repeat forever:

Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot| \cdot, \theta)$

For each step of the episode $t = 0, \dots, T-1$:

$G \leftarrow$ return from step t

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla_\theta \ln \pi(A_t|S_t; \theta)$$

How did we get here
from $\frac{\nabla_\theta \pi(A_t|S_t; \theta)}{\pi(A_t|S_t; \theta)}$?

$$\nabla \ln(f(x)) = \frac{\nabla f(x)}{f(x)}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Calibrating REINFORCE

- Imagine a domain with two possible episode trajectories (rollouts)

- $\tau_1 = \{S_0, A_0, R_1, S_1, \dots, A_{T-1}, R_T, S_T\}$ with $G = 1001$
- $\tau_2 = \{S'_0, A'_0, R'_1, S'_1, \dots, A'_{T-1}, R'_T, S'_T\}$ with $G' = 1000$
- Further assume that $R_{i < T} = R'_i = 0$ and $R_T = G, R'_T = G'$

- Following REINFORCE: $\theta_{t+1} = \theta_t + \alpha G_t \nabla_\theta \ln \pi(A_t|S_t; \theta)$

- How will the policy be updated after sampling each of the trajectories?

- $\tau_1 ++, \tau_2 ++$

- How should the policy be updated after sampling each of the trajectories?

- $\tau_1 +, \tau_2 -$

- How can we address this gap?



PG with baseline

- Normalize the returns with a baseline!

- $\nabla J(\theta) \propto \sum_s \mu(s) \sum_a (q_\pi(s, a) - b(s)) \nabla \pi(a|s)$
- Are we allowed to do that???
- Can we subtract $\sum_a b(s) \nabla \pi(a|s)$ from $\nabla J(\theta)$?
- Yes!** As long as $b(s)$ isn't a function of a
- $\sum_a b(s) \nabla \pi(a|s) = b(s) \nabla \sum_a \pi(a|s) = b(s) \nabla 1 = 0$ (actions' probabilities sum to 1)
- REINFORCE with baseline: $\theta_{t+1} = \theta_t + \alpha(G_t - b(S_t)) \nabla_\theta \ln \pi_\theta(A_t | S_t)$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

PG with baseline

- In the bandit algorithms the baseline was the average of the rewards seen so far
- For MDPs the baseline should be different for each states
 - In some states all actions have high (q) values, and we need a high baseline to differentiate the higher valued actions from the less highly valued ones
 - In other states all actions will have low values and a low baseline is appropriate
- What would be the equivalent of average reward (bandits) for a given state (MDP)?
 - $v_\pi(s)$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



PG with baseline

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a (q_\pi(s, a) - v_\pi(s)) \nabla \pi(a|s)$$

- Actions that improve on the current policy are encouraged
 - $q_\pi(s, a) - v_\pi(s) > 0$
- Actions that damage the current policy are discouraged
 - $q_\pi(s, a) - v_\pi(s) < 0$
- Also known as the **advantage** of action a in state s
- How can we learn $v_\pi(s)$?
- Another function approximator $\hat{v}(s; w)$
 - On top of the policy

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

REINFORCE with baseline

REINFORCE with Baseline (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Input: a differentiable state-value parameterization $\hat{v}(s, w)$

Parameters: step sizes $\alpha^\theta > 0$, $\alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$

Repeat forever:

Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot | \cdot, \theta)$

For each step of the episode $t = 0, \dots, T-1$:

$G_t \leftarrow$ return from step t

$\delta \leftarrow G_t - \hat{v}(S_t, w)$

$w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S_t, w)$

$\theta \leftarrow \theta + \alpha^\theta \delta \nabla_\theta \ln \pi(A_t | S_t, \theta)$

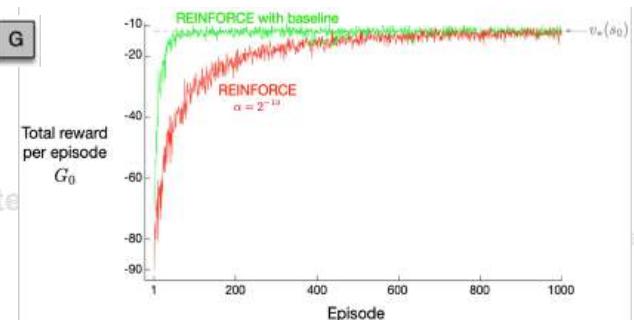
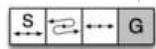
Each approximator has its unique learning rate

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

REINFORCE with baseline

The crazy corridor domain

- $\alpha^\theta = 2^{-9}$
- $\alpha^w = 2^{-6}$



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Add a critic

$$\widehat{\nabla J(\theta_t)} \propto (q_\pi(S_t, A_t) - b(S_t)) \frac{\nabla_\theta \pi(A_t | S_t; \theta)}{\pi(A_t | S_t; \theta)}$$

- Define a new estimator: $\hat{q}_\pi(s, a; \theta) \approx q_\pi(s, a)$
 - To be used instead of G_t in REINFORCE
- How should we train $\hat{q}_\pi(s, a; \theta)$?
 - Monte-Carlo updates
 - High variance samples
 - Requires a full episode
 - Bootstrapping e.g., Q-learning
 - Lower variance (though introduces bias)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Policy Gradient

- $\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$
- PG theorem: $\widehat{\nabla J(\theta_t)} \propto (q_\pi(S_t, A_t) - b(S_t)) \nabla_\theta \log \pi(A_t | S_t; \theta)$
 - REINFORCE+baseline: $\theta_{t+1} = \theta_t + \alpha(G_t - b(S_t)) \nabla_\theta \log \pi(A_t | S_t; \theta)$
- Pros: approximating the optimal policy directly is more accurate and faster to converge in many domains when compared to value-based approaches
- Cons: using G_t as an estimator for $q_\pi(S_t, A_t)$ is noisy (high variance)²¹
 - Unstable learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Critic's duties

- Approximate state or action or advantage ($q(s, a) - v(s)$) values
 - Trained via bootstrapping, criticizes the action chosen by the actor adjusting the actor's (policy) gradient
- Is REINFORCE (+ state value baseline) an actor-critic framework?
 - $\theta_{t+1} = \theta_t + \alpha(G_t - \hat{v}_\pi(S_t)) \nabla_\theta \ln \pi(A_t | S_t; \theta)$
 - NO! The state-value function is used only as a baseline, not as a critic. Moreover, it doesn't utilize bootstrapping (uses high-variance MC updates)
 - The bias introduced through bootstrapping is often worthwhile because it reduces variance and accelerates learning

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Benefits from a critic

- REINFORCE with baseline is unbiased* and will converge asymptotically to a local optimum
 - With a linear state value approximator, and when b is not a function of a
 - Like all Monte-Carlo methods it tends to learn slowly (produce estimates of high variance)
 - Not suitable for online or for continuing problems
- Temporal-difference methods can eliminate these inconveniences
- In order to gain the TD advantages in the case of policy gradient methods we use actor–critic methods

Actor+critic

- Actor-critic algorithms are a derivative of policy iteration, which alternates between policy evaluation—computing the value function for a policy—and policy improvement—using the value function to obtain a better policy
- In large-scale reinforcement learning problems, it is typically impractical to run either of these steps to convergence, and instead the value function and policy are optimized jointly
- The policy is referred to as the actor, and the value function as the critic

Advantage function

- Eventually we would like to shift the policy towards actions that result in higher return
- what is the benefit from taking action a at state s while following policy π ?
 - $A_\pi(s, a) = q_\pi(s, a) - v_\pi(s)$
- This resembles PG with baseline but not the same as q_π is approximated:
 - $\widehat{V}(\theta_t) = A_\pi(s_t, a_t) \nabla_\theta \ln \pi(a_t | s_t; \theta)$
- Actions that improve on the current policy are encouraged
 - $A_\pi(s, a) > 0$
- Actions that damage the current policy are discouraged
 - $A_\pi(s, a) < 0$

One-step actor-critic

- $A_\pi(s, a) = q_\pi(s, a) - v_\pi(s)$
 - Does that mean that approximating the advantage function requires two function approximators (\hat{q} and \hat{v})?
 - No since q values can be derived from state values + one step transition
 - $\hat{q}_\pi(s_t, a_t) - \hat{v}_\pi(s_t; w) = \hat{\mathbb{E}}[r_{t+1} + \gamma \hat{v}(s_{t+1}; w)] - \hat{v}(s_t; w)$
- $\theta_{t+1} = \theta_t + \alpha(r_{t+1} + \gamma \hat{v}(s_{t+1}; w) - \hat{v}(s_t; w)) \nabla_\theta \ln \pi(a_t | s_t; \theta)$

δ – TD error

- One-step Actor-Critic is a fully online, incremental algorithm, with states, actions, and rewards processed as they occur and then never revisited

One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
 Input: a differentiable state-value parameterization $\hat{v}(s, w)$
 Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$

Repeat forever:

 Initialize S (first state of episode)

$I \leftarrow 1$

 While S is not terminal:

$A \sim \pi(\cdot|S, \theta)$

 Take action A , observe S', R

$\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ (if S' is terminal, then $\hat{v}(S', w) \doteq 0$)

$w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S, w)$

$\theta \leftarrow \theta + \alpha^\theta I \delta \nabla_\theta \ln \pi(A|S, \theta)$

$I \leftarrow \gamma I$

$S \leftarrow S'$

Keep track of accumulated discount

28

One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
 Input: a differentiable state-value parameterization $\hat{v}(s, w)$
 Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$

Repeat forever:

 Initialize S (first state of episode)

$I \leftarrow 1$

 While S is not terminal:

$A \sim \pi(\cdot|S, \theta)$

 Take action A , observe S', R

$\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ (if S' is terminal, then $\hat{v}(S', w) \doteq 0$)

$w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S, w)$

$\theta \leftarrow \theta + \alpha^\theta I \delta \nabla_\theta \ln \pi(A|S, \theta)$

$I \leftarrow \gamma I$

$S \leftarrow S'$

Follow the current policy

29

Work Integrated Learning Programmes



One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
 Input: a differentiable state-value parameterization $\hat{v}(s, w)$
 Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$
 Repeat forever:

- Initialize S (first state of episode)
 $I \leftarrow 1$
- While S is not terminal:
 - $A \sim \pi(\cdot|S, \theta)$
 - Take action A , observe S', R
 - $\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ (if S' is terminal, then $\hat{v}(S', w) \doteq 0$)
 - $w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S, w)$
 - $\theta \leftarrow \theta + \alpha^\theta I \delta \nabla_\theta \ln \pi(A|S, \theta)$
 - $I \leftarrow \gamma I$
 - $S \leftarrow S'$

Update the actor with discounting. Early actions matter more.

32

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
 Input: a differentiable state-value parameterization $\hat{v}(s, w)$
 Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$
 Repeat forever:

- Initialize S (first state of episode)
 $I \leftarrow 1$
- While S is not terminal:
 - $A \sim \pi(\cdot|S, \theta)$
 - Take action A , observe S', R
 - $\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ (if S' is terminal, then $\hat{v}(S', w) \doteq 0$)
 - $w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S, w)$
 - $\theta \leftarrow \theta + \alpha^\theta I \delta \nabla_\theta \ln \pi(A|S, \theta)$
 - $I \leftarrow \gamma I$
 - $S \leftarrow S'$

In practice: training the network at every step on a single observation is inefficient (slow and correlated)

34

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
 Input: a differentiable state-value parameterization $\hat{v}(s, w)$
 Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$
 Repeat forever:

- Initialize S (first state of episode)
 $I \leftarrow 1$
- While S is not terminal:
 - $A \sim \pi(\cdot|S, \theta)$
 - Take action A , observe S', R
 - $\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ (if S' is terminal, then $\hat{v}(S', w) \doteq 0$)
 - $w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S, w)$
 - $\theta \leftarrow \theta + \alpha^\theta I \delta \nabla_\theta \ln \pi(A|S, \theta)$
 - $I \leftarrow \gamma I$
 - $S \leftarrow S'$

Update accumulated discount and progress to the next state

33

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
 Input: a differentiable state-value parameterization $\hat{v}(s, w)$
 Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$
 Repeat forever:

- Initialize S (first state of episode)
 $I \leftarrow 1$
- While S is not terminal:
 - $A \sim \pi(\cdot|S, \theta)$
 - Take action A , observe S', R
 - $\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ (if S' is terminal, then $\hat{v}(S', w) \doteq 0$)
 - $w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S, w)$
 - $\theta \leftarrow \theta + \alpha^\theta I \delta \nabla_\theta \ln \pi(A|S, \theta)$
 - $I \leftarrow \gamma I$
 - $S \leftarrow S'$

In practice: training the network at every step on a single observation is inefficient (slow and correlated)

34

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes



One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
 Input: a differentiable state-value parameterization $\hat{v}(s, w)$
 Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$
 Repeat forever:

- Initialize S (first state of episode)
 $I \leftarrow 1$
- While S is not terminal:
 - $A \sim \pi(\cdot|S, \theta)$
 - Take action A , observe S', R
 - $\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ (if S' is terminal, then $\hat{v}(S', w) \doteq 0$)
 - $w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S, w)$
 - $\theta \leftarrow \theta + \alpha^\theta I \delta \nabla_\theta \ln \pi(A|S, \theta)$
 - $I \leftarrow \gamma I$
 - $S \leftarrow S'$

Instead: store all state approx values, log probabilities, and rewards along the episode. Train once at the end of the episode.

35

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

One-step Actor-Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
 Input: a differentiable state-value parameterization $\hat{v}(s, w)$
 Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$

Repeat forever:

 Initialize S (first state of episode)

$I \leftarrow 1$

 While S is not terminal:

$A \sim \pi(\cdot|S, \theta)$

 Take action A , observe S', R

$\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$

$w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S, w)$

$\theta \leftarrow \theta + \alpha^\theta I \delta \nabla_\theta \ln \pi(A|S, \theta)$

$I \leftarrow \gamma I$

$S \leftarrow S'$

Instead: store all state approx. values, log probabilities, and rewards along the episode. Train once at the end of the episode.

```
values = torch.FloatTensor(values)
Qvals = torch.FloatTensor(Qvals)
log_probs = torch.stack(log_probs)
advantage = Qvals - values
```

36

What did we learn?

- In many domains it's more efficient to learn the policy directly
 - Instead of deriving one from state or action values
- Applicable for continuous action spaces and stochastic policies
- Approximate the change to the policy that results in the highest increase in the expected return: $\nabla_\theta J(\theta)$
- Such approximation is made possible when by the policy gradient theorem
- Should we reinforce policies that result in high return?
 - Not always, a different policy might yield higher return
 - We need to use a baseline to determine if a policy is relatively good

Advantage Actor-Critic (A2C)• Initialize the actor π_θ and the critic \hat{V}_w

- For each episode:
 - Init empty episode memory
 - $s = \text{init env}$
 - For each time step:
 - $a \sim \pi(s)$
 - $s', r_t \sim \text{env}(s, a)$
 - Store (s, a, r) in episode memory
 - $s = s'$
 - Reset $d\theta = 0, dw = 0$
- Backwards iteration (i from length to 0) over the episode
 - Compute advantage $\delta_t = r_i + \gamma \hat{V}_w(s_{i+1}) - \hat{V}_w(s_i)$
 - Accumulate the policy gradient using the critic: $d\theta \leftarrow d\theta + \delta \nabla_\theta \ln \pi_\theta(s_i, a_i)$
 - Accumulate the critic gradient: $dw \leftarrow dw + \delta \nabla_w \hat{V}_w(s_i)$
 - Update the actor and the critic with the accumulated gradients using gradient descent

What did we learn?

• REINFORCE:

$$\nabla J(\theta_t) = G_t \nabla_\theta \ln \pi(A_t|S_t; \theta)$$

• Q Actor-Critic:

$$\nabla J(\theta_t) = \hat{q}(S_t, A_t; w) \nabla_\theta \ln \pi(A_t|S_t; \theta)$$

• REINFORCE + baseline:

$$\nabla J(\theta_t) = (G_t - \hat{v}(S_t; w)) \nabla_\theta \ln \pi(A_t|S_t; \theta)$$

• Advantage Actor-Critic:

$$\nabla J(\theta_t) = (R_{t+1} + \gamma \hat{v}(S_{t+1}; w) - \hat{v}(S_t; w)) \nabla_\theta \ln \pi(A_t|S_t; \theta)$$

Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

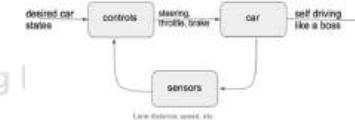
Eg.,

State = function (Distance of Lane edges, Speed of Vehicle, etc.,)

Action = function (Steering Wheel Turn Angle, Amount(Force) of Acceleration/Deceleration, etc.,)

Challenge 1:

Continuous Large State Space & Action Space



Solution : Policy based function approximation

[Source Reference: Example inspired from AI Researcher & Scientist : Eder Santana](#)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

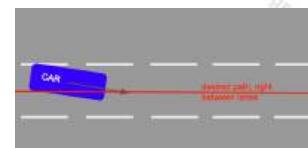
Eg.,

State $\in \{S_0: \text{Left}, S_1: \text{Center}, S_2: \text{Right}\}$

Action $\in [-1.0, +1.0]$ (Steer Angle is any value between)

-1=> Full Left , 0=> Straight(no turn) , +1=> Full right

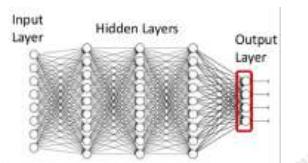
Reward $\in \{+1 \text{ if vehicle is in the center} \atop 0 \text{ if its in left or right}\}$



Output Layer : Parameters of the action distribution

Steer Angle: let it be a $\sim N(\mu, \sigma)$

For simplicity, let's Assume $\sigma = 0.2$ (a constant) & only $\mu = \theta^T X$, a linear approximation is a function the network parameter θ



Note : Output can also be designed with $P(a|s)$ for discrete larger action space!. Apply soft max in this case

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

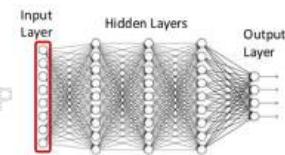
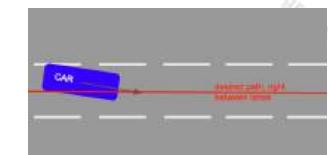
Eg.,

State $\in \{S_0: \text{Left}, S_1: \text{Center}, S_2: \text{Right}\}$

Action $\in [-1.0, +1.0]$ (Steer Angle is any value between)

-1=> Full Left , 0=> Straight(no turn) , +1=> Full right

Reward $\in \{+1 \text{ if vehicle is in the center} \atop 0 \text{ if its in left or right}\}$



Work Integrated Learning P

Input Layer : One hot encoding of the states.

For $S_1 = [X_1 \ X_2 \ X_3]^T = [0 \ 1 \ 0]^T$

[Source Reference: Example inspired from AI Researcher & Scientist : Eder Santana](#)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)



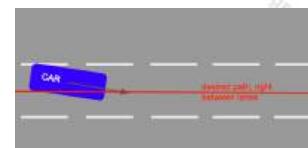
Eg.,

State $\in \{S_0: \text{Left}, S_1: \text{Center}, S_2: \text{Right}\}$

Action $\in [-1.0, +1.0]$ (Steer Angle is any value between)

-1=> Full Left , 0=> Straight(no turn) , +1=> Full right

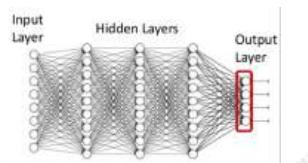
Reward $\in \{+1 \text{ if vehicle is in the center} \atop 0 \text{ if its in left or right}\}$



Output Layer : Parameters of the action distribution

Steer Angle: let it be a $\sim N(\mu, \sigma)$

For simplicity, let's Assume $\sigma = 0.2$ (a constant) & only $\mu = \theta^T X$, a linear approximation is a function the network parameter θ



Note : Output can also be designed with $P(a|s)$ for discrete larger action space!. Apply soft max in this case

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane)

How the Policy Network Parameter θ 's are learnt?

> Update by Gradient Ascent

> Refer to the below gradient of the Gaussian

> Initialize the parameter for $\mu = \theta^T X = \theta_1 X_1 + \theta_2 X_2 + \dots$

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for π .

Input: a differentiable policy parametrization $\pi(a|s, \theta)$
 Algorithm parameter: step size $\alpha > 0$
 Initialize policy parameter $\theta \in \mathbb{R}^d$ (e.g., to 0)
 Loop forever (for each episode):
 Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot | \cdot, \theta)$
 Loop for each step of the episode $t = 0, 1, \dots, T-1$:

$$G \leftarrow \sum_{k=t+1}^T k^{t-1} \pi(a_k | s_k, \theta) R_k$$

$$\theta \leftarrow \theta + \alpha \nabla G \nabla \ln \pi(a_t | s_t, \theta)$$

Work Integrated Learning Programmes

Refer to TB-1 Page section 13.7 →
$$\nabla \ln \pi(a|s, \theta_\mu) = \frac{\nabla \pi(a|s, \theta_\mu)}{\pi(a|s, \theta)} = \frac{1}{\pi(s, \theta)^2} (a - \mu(s, \theta)) \mathbf{x}_\mu(s)$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane)

How the Policy Network Parameter θ 's are learnt?

- > Update by Gradient Ascent
- > Refer to the below gradient of the Gaussian
- > Initialize the parameter for $\mu = \theta^T X = \theta_1 X_1 + \theta_2 X_2 + \dots$

$$\text{Let } \theta^T = [-0.5 \ 0 \ 0.5]$$

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for π .

Input: a differentiable policy parameterization $\pi(a|s, \theta)$
Algorithm parameter: step size $\alpha > 0$
Initialize policy parameter $\theta \in \mathbb{R}^d$ (e.g., to 0)

Loop forever (for each episode):
Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot | \cdot, \theta)$
Loop for each step of the episode $t = 0, 1, \dots, T-1$:
 $G \leftarrow \sum_{i=t+1}^T r^{k-i-1} R_i$
 $\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t | S_t, \theta)$

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

- > Update by Gradient Ascent
- > Refer to the below gradient of the Gaussian
- > Initialize the parameter for $\mu = \theta^T X = \theta_1 X_1 + \theta_2 X_2 + \theta_3 X_3$:

$$\text{Let } \theta^T = [-0.5 \ 0 \ 0.5]$$

Trajectory 1: $(S_0, a_0, r_1, S_1) = (\text{Left}, -0.3, +1, \text{Center})$

Start with $S_0 = [1 \ 0 \ 0]^T$

Apply $\mu = \theta^T X = -0.5$ derived in Gaussian to derive Action from above policy

$a_0 = -0.3$ (Steer left slightly)

Next State: $S_1 = [0 \ 1 \ 0]^T$

Reward = +1

(Assume episode terminates) $\rightarrow G = +1$ with $\gamma = 1$ [No discount]

Refer to TB-1 Page section 13.7 $\rightarrow \nabla \ln \pi(a|s, \theta_\mu) = \frac{\nabla \pi(a|s, \theta_\mu)}{\pi(a|s, \theta)} = \frac{1}{\pi(s, \theta)^2} (a - \mu(s, \theta)) \mathbf{x}_\mu(s)$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

- > Update by Gradient Ascent
- > Refer to the below gradient of the Gaussian
- > Initialize the parameter for $\mu = \theta^T X = \theta_1 X_1 + \theta_2 X_2 + \theta_3 X_3$

$$\text{Let } \theta^T = [-0.5 \ 0 \ 0.5]$$

Trajectory 1: $(S_0, a_0, r_1, S_1) = (\text{Left}, -0.3, +1, \text{Center})$

Start with $S_0 = [1 \ 0 \ 0]^T$

Apply $\mu = \theta^T X = -0.5$ derived in Gaussian to derive Action from above policy

$a_0 = -0.3$ (Steer left slightly)

Next State: $S_1 = [0 \ 1 \ 0]^T$

Reward = +1

Refer to TB-1 Page section 13.7 $\rightarrow \nabla \ln \pi(a|s, \theta_\mu) = \frac{\nabla \pi(a|s, \theta_\mu)}{\pi(a|s, \theta)} = \frac{1}{\pi(s, \theta)^2} (a - \mu(s, \theta)) \mathbf{x}_\mu(s)$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

- > Compute the gradient, return
- > Update θ with e.g., $\alpha = 0.01$

$$\theta_{\text{old}}^T = [-0.5 \ 0 \ 0.5]$$

Trajectory 1:

$(S_0, a_0, r_1, S_1) = (\text{Left}, -0.3, +1, \text{Center})$

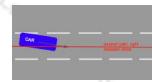
$G = +1$

$$\begin{aligned} \theta_{\text{new}}^T &= \theta_{\text{old}}^T + \alpha \mathbf{G} \nabla \ln \pi(A_t | S_t, \theta_{\text{old}}) \\ &= [-0.5 \ 0 \ 0.5]^T + 0.01 \begin{pmatrix} 1 \\ \frac{(-0.3 - (-0.5))}{(0.2)^2} \end{pmatrix} [1 \ 0 \ 0]^T \\ &= [-0.5 \ 0 \ 0.5]^T + 0.01 \begin{pmatrix} 5 \\ 0.2 \end{pmatrix} [1 \ 0 \ 0]^T \\ &= [-0.45 \ 0 \ 0.5] \end{aligned}$$

Refer to TB-1 Page section 13.7 $\rightarrow \nabla \ln \pi(a|s, \theta_\mu) = \frac{\nabla \pi(a|s, \theta_\mu)}{\pi(a|s, \theta)} = \frac{1}{\pi(s, \theta)^2} (a - \mu(s, \theta)) \mathbf{x}_\mu(s)$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

$$\text{With } \theta_{\text{old}}^T = [-0.45 \ 0 \ 0.5],$$

after repeating the exercise for one more say ,

$$\text{Trajectory 2: } (S_0, a_0, r_1, S_2) = (\text{Left}, -0.9, 0, \text{Right})$$

Update is as follows :

Trajectory 2: (assume that the episode terminates after S2)

$$(S_0, a_0, r_1, S_2) = (\text{Left}, -0.9, 0, \text{Right})$$

$$G = 0$$

$$\begin{aligned} \theta_{\text{new}}^T &= \theta_{\text{old}}^T + \alpha \mathbf{G} \nabla \ln \pi(A_t | S_t, \theta_{\text{old}}) \\ &= [-0.45 \ 0 \ 0.5]^T + 0.01 \begin{pmatrix} 0 \\ (-0.9 - (-0.45))^2 \\ 0.2 \end{pmatrix} [1 \ 0 \ 0]^T \\ &= [-0.45 \ 0.5 \ 0.5]^T + 0.01 \begin{pmatrix} 0 \\ 11.25 \\ 0.2 \end{pmatrix} [1 \ 0 \ 0]^T \\ &= \boxed{[-0.45 \ 0.5 \ 0.5]} \end{aligned}$$

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for π .

Input: a differentiable policy parameterizing $\pi(a|s, \theta)$
 Algorithm parameter: step size $\alpha > 0$
 Initialize policy parameter $\theta \in \mathbb{R}^d$ (e.g., to 0)

Loop forever (for each episode):
 Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot | \theta)$
 Loop for each step of the episode $t = 0, 1, \dots, T-1$:

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

$$\theta \leftarrow \theta + \alpha^\pi G \nabla \ln \pi(A_t | S_t, \theta)$$

$$\theta \leftarrow \theta + \alpha^\pi \nabla^2 \ln \pi(A_t | S_t, \theta)$$

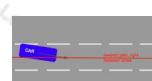
Observation :

Gradient update vanishes or tend to vary due to full dependence on reward

Solution : Induce bias using baseline. Use value network w

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

> Compute the gradient, return

> Update θ with eg., $\alpha_0 = 0.01$

$$\theta_{\text{old}}^T = [-0.5 \ 0 \ 0.5]$$

Typically the value network w will also be learnt as per the algorithm.
 For simplicity assume that it estimated values for $S_0: V(\text{left})=0.4$,

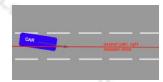
$$V(\text{center})=0.9, V(\text{right})=0.2$$

Trajectory 1:

$$(S_0, a_0, r_1, S_1) = (\text{Left}, -0.3, +1, \text{Center})$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

> Compute the gradient, return

> Update θ with eg., $\alpha_0 = 0.01$

$$\theta_{\text{old}}^T = [-0.5 \ 0 \ 0.5]$$

Typically the value network w will also be learnt as per the algorithm.
 For simplicity assume that it estimated values for $S_0: V(\text{left})=0.4$,

$$V(\text{center})=0.9, V(\text{right})=0.2$$

REINFORCE with Baseline (episodic): for estimating $v_\theta = v$.

Input: a differentiable policy parameterization $\pi(a|s, \theta)$.
 Input: a differentiable state-value function parameterization $V(s, w)$.
 Algorithm parameter: step size $\alpha > 0$, $\alpha^\pi > 0$.
 Initialize policy parameter $\theta \in \mathbb{R}^d$ and state-value weights $w \in \mathbb{R}^d$ (e.g., to 0)

Loop forever (for each episode):
 Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot | \theta)$
 Loop for each step of the episode $t = 0, 1, \dots, T-1$:

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

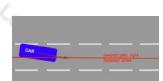
$$\theta \leftarrow \theta + \alpha^\pi G \nabla \ln \pi(A_t | S_t, \theta)$$

$$w \leftarrow w + \alpha^\pi \nabla^2 \ln \pi(A_t | S_t, \theta)$$

$$\theta \leftarrow \theta + \alpha^\pi \nabla^2 \ln \pi(A_t | S_t, \theta)$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

> Compute the gradient, return

> Update θ with eg., $\alpha_0 = 0.01$

$$\theta_{\text{old}}^T = [-0.5 \ 0 \ 0.5]$$

Typically the value network w will also be learnt as per the algorithm. For simplicity assume that it estimated values for $S_0: V(\text{left})=0.4, V(\text{center})=0.9, V(\text{right})=0.2$

Trajectory 1:

$$(S_0, a_0, r_1, S_1) = (\text{Left}, -0.3, +1, \text{Center})$$

$$G = +1$$

$$\begin{aligned} \theta_{\text{new}}^T &= \theta_{\text{old}}^T + \alpha (G - V(S_1)) \nabla \ln \pi(A_t | S_t, \theta_{\text{old}}) \\ &= [-0.5 \ 0 \ 0.5]^T + 0.01 \begin{pmatrix} 1-0.4 \\ (-0.3 - (-0.5))^2 \\ 0.2 \end{pmatrix} [1 \ 0 \ 0]^T \end{aligned}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant

Objective : Policy to keep in the vehicle in line (center of the lane!)

$$\text{With } \theta_{\text{old}}^T = [-0.47 \ 0 \ 0.5],$$

after repeating the exercise for one more say ,

$$\text{Trajectory 2: } (S_0, a_0, r_1, S_2) = (\text{Left}, -0.9, 0, \text{Right})$$

Update is as follows :

Trajectory 2: assume that the episode terminates after S2

$$(S_0, a_0, r_1, S_2) = (\text{Left}, -0.9, 0, \text{Right})$$

$$G = 0$$

$$\begin{aligned} \theta_{\text{new}}^T &= \theta_{\text{old}}^T + \alpha \mathbf{G} \nabla \ln \pi(A_t | S_t, \theta_{\text{old}}) \\ &= [-0.47 \ 0 \ 0.5]^T + 0.01 \quad \text{(0-0.4)} \quad \frac{(-0.9 - (-0.45))}{(0.2)^2} [1 \ 0 \ 0]^T \\ &= [-0.47 \ 0 \ 0.5]^T + 0.01 \quad \text{(4.3)} \quad [1 \ 0 \ 0]^T \\ &= \boxed{[-0.427 \ 0 \ 0.5]} \end{aligned}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



(Simplified) Lane Keeping Assistant

Objective : Policy to keep in the vehicle in line (center of the lane!)

$$\text{With } \theta_{\text{old}}^T = [-0.47 \ 0 \ 0.5],$$

after repeating the exercise for one more say ,

$$\text{Trajectory 2: } (S_0, a_0, r_1, S_2) = (\text{Left}, -0.9, 0, \text{Right})$$

Update is as follows :

Trajectory 2: assume that the episode terminates after S2

$$(S_0, a_0, r_1, S_2) = (\text{Left}, -0.9, 0, \text{Right})$$

$$G = 0$$

$$\begin{aligned} \theta_{\text{new}}^T &= \theta_{\text{old}}^T + \alpha \mathbf{G} \nabla \ln \pi(A_t | S_t, \theta_{\text{old}}) \\ &= [-0.47 \ 0 \ 0.5]^T + 0.01 \quad \text{(0-0.4)} \quad \frac{(-0.9 - (-0.45))}{(0.2)^2} [1 \ 0 \ 0]^T \\ &= [-0.47 \ 0 \ 0.5]^T + 0.01 \quad \text{(4.3)} \quad [1 \ 0 \ 0]^T \\ &= \boxed{[-0.427 \ 0 \ 0.5]} \end{aligned}$$

REINFORCE with Baseline (episode), for estimating $\pi_\theta = \pi$.
 Input: a differentiable policy parameterization $\pi(a|s, \theta)$.
 Input: a differentiable state-value function parameterization $V(s, w)$.
 Parameters: step size $\alpha^{\theta} > 0$, $\alpha^w > 0$.
 Initialize: policy parameter $\theta \in \mathbb{R}^d$ and state-value weights $w \in \mathbb{R}^d$ (e.g., to 0).
 Long loop (for each episode):
 Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot | \theta)$.
 Loop for each step of the episode $t = 0, 1, \dots, T-1$:

$$\begin{aligned} G &\leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \\ \theta &\leftarrow \theta + \alpha^{\theta} \nabla \ln \pi(A_t | S_t, \theta) \\ w &\leftarrow w + \alpha^w \delta \nabla V(S_t, w) \\ \theta &\leftarrow \theta + \alpha^{\theta} \delta \nabla \ln \pi(A_t | S_t, \theta) \end{aligned} \quad \text{(G.)}$$

Observation :

- > Relatively less variance
- > But Monte Carlo based updates may not be suited for lengthier episodes/continuing tasks
- > Need online learning with faster updates

Solution : Use Temporal Difference .

Use value network w as critic
 Policy net θ as actor for policy improvement
 Gradient Descent with TD error based advantage function

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



(Simplified) Lane Keeping Assistant

Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

> Compute the gradient, return

> Update θ with eg., $\alpha_\theta = 0.01$

$$\theta_{\text{old}}^T = [-0.5 \ 0 \ 0.5]$$

Typically the value network w will also be learnt as per the algorithm.
 For simplicity assume that it estimated values for $S_0: V(\text{left})=0.2$, $V(\text{center})=0.9$, $V(\text{right})=0.2$

One-step Actor-Critic (episode), for estimating $\pi_\theta = \pi$.
 Input: a differentiable policy parameterization $\pi(a|s, \theta)$.
 Input: a differentiable state-value function parameterization $V(s, w)$.
 Parameters: step size $\alpha^{\theta} > 0$, $\alpha^w > 0$.
 Initialize: policy parameter $\theta \in \mathbb{R}^d$ and state-value weights $w \in \mathbb{R}^d$ (e.g., to 0).
 Long loop (for each episode):
 Initialize S (first state of episode)
 $t \leftarrow 1$
 Loop while S is not terminal (for each time step):
 $A \sim \pi(\cdot | S, \theta)$
 Take A , observe S' , R
 $\delta \leftarrow R + \gamma V(S', w) - V(S, w)$ $\text{if } S' \text{ is terminal, then } \delta(S', w) = 0$
 $w \leftarrow w + \alpha^w \delta \nabla V(S, w)$
 $\theta \leftarrow \theta + \alpha^{\theta} \delta \nabla \ln \pi(A | S, \theta)$
 $t \leftarrow t + 1$
 $S \leftarrow S'$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



(Simplified) Lane Keeping Assistant

Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

> Compute the gradient, return

> Update θ with eg., $\alpha_\theta = 0.01$

$$\theta_{\text{old}}^T = [-0.5 \ 0 \ 0.5]$$

Typically the value network w will also be learnt as per the algorithm.
 For simplicity assume that it estimated values for $S_0: V(\text{left})=0.3$, $V(\text{center})=0.9$, $V(\text{right})=0.2$

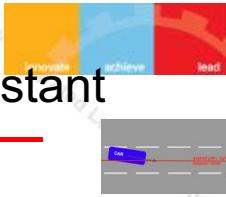
Trajectory 1:

$$(S_0, a_0, r_1, S_1) = (\text{Left}, -0.3, +1, \text{Center})$$

One-step Actor-Critic (episode), for estimating $\pi_\theta = \pi$.
 Input: a differentiable policy parameterization $\pi(a|s, \theta)$.
 Input: a differentiable state-value function parameterization $V(s, w)$.
 Parameters: step size $\alpha^{\theta} > 0$, $\alpha^w > 0$.
 Initialize: policy parameter $\theta \in \mathbb{R}^d$ and state-value weights $w \in \mathbb{R}^d$ (e.g., to 0).
 Long loop (for each episode):
 $t \leftarrow 1$
 Loop while S is not terminal (for each time step):
 Take A , observe S' , R
 $\delta \leftarrow R + \gamma V(S', w) - V(S, w)$ $\text{if } S' \text{ is terminal, then } \delta(S', w) = 0$
 $w \leftarrow w + \alpha^w \delta \nabla V(S, w)$
 $\theta \leftarrow \theta + \alpha^{\theta} \delta \nabla \ln \pi(A | S, \theta)$
 $t \leftarrow t + 1$
 $S \leftarrow S'$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

> Compute the gradient, return

> Update θ with eg., $\alpha_\theta = 0.01$

$$\theta_{\text{old}}^T = [-0.5 \ 0 \ 0.5]$$

Typically the value network w will also be learnt as per the algorithm. For simplicity assume that it estimated values for S0: V(left)=0.3, V(center)=0.5, V(right)=0.2

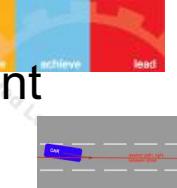
Trajectory 1:

(S0, a0, r1, S1) = (Left, -0.3, +1, Center)
 $\theta_{\text{new}}^T = \theta_{\text{old}}^T + \alpha \left(\mathbf{r} + \mathbf{V(S_{t+1})} - \mathbf{V(S_t)} \right) \nabla \ln \pi(A_t | S_t, \theta_{\text{old}})$
 $= [-0.5 \ 0.0]T + 0.01 \left(\frac{-0.3 - (-0.5)}{(0.2)^2} \right) [1 \ 0 \ 0]^T$

One-step Actor-Critic (episodic), for estimating $\pi_\theta \approx \pi_c$.
Input: a differentiable policy parameterization $\pi(a|s)$,
Input: a differentiable state-action function parameterization $v(s, w)$
Parameters: step size $\alpha^P > 0, \alpha^W > 0$
Initial policy parameter $\theta \in \mathbb{R}^d$ and state-value weights $w \in \mathbb{R}^d$ (e.g., to 0)
Long Keeve (for each episode):
Initialise S (first state of episode)
 $I \leftarrow 1$
Loop while S is not terminal (for each time step):
 $A \sim \pi(\cdot | S, \theta)$
 $R \sim P(A, S)$ if A observes S' , R
 $\delta \leftarrow R + \gamma v(S', w) - v(S, w)$ (If S' is terminal, then $v(S', w) \leftarrow 0$)
 $w \leftarrow w + \alpha^W \delta \nabla v(S, w)$
 $\theta \leftarrow \theta + \alpha^P (\delta \nabla \ln \pi(A | S, \theta))$
 $I \leftarrow I + 1$
 $S \leftarrow S'$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

With $\theta_{\text{old}}^T = [-0.42 \ 0 \ 0.5]$,

after repeating the exercise for one more say ,

Trajectory 2: (S0, a0, r1, S2) = (Left, -0.9, 0, Right)

Update is as follows :

Trajectory 2: (assume that the episode terminates after S2)

(S0, a0, r1, S2) = (Left, -0.9, 0, Right)

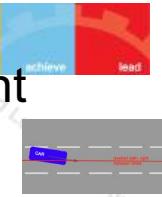
G = 0

$$\theta_{\text{new}}^T = \theta_{\text{old}}^T + \alpha \left(\mathbf{r} + \mathbf{V(S_{t+1})} - \mathbf{V(S_t)} \right) \nabla \ln \pi(A_t | S_t, \theta_{\text{old}})$$
 $= [-0.42 \ 0.0]T + 0.01 \left(\frac{-0.9 - (-0.42)}{(0.2)^2} \right) [1 \ 0 \ 0]^T$
 $= [-0.42 \ 0.05]T + 0.01 \left(\frac{1.2}{(0.2)^2} \right) [1 \ 0 \ 0]^T$
 $= \boxed{[-0.408 \ 0 \ 0.5]}$

This iteration is added only for additional practice!

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

(Simplified) Lane Keeping Assistant



Objective : Policy to keep in the vehicle in line (center of the lane!)

How the Policy Network Parameter θ 's are learnt?

> Compute the gradient, return

> Update θ with eg., $\alpha_\theta = 0.01$

$$\theta_{\text{old}}^T = [-0.5 \ 0 \ 0.5]$$

Typically the value network w will also be learnt as per the algorithm. For simplicity assume that it estimated values for S0: V(left)=0.3, V(center)=0.5, V(right)=0.2

Trajectory 1:

(S0, a0, r1, S1) = (Left, -0.3, +1, Center)
 $\theta_{\text{new}}^T = \theta_{\text{old}}^T + \alpha \left(\mathbf{r} + \mathbf{V(S_{t+1})} - \mathbf{V(S_t)} \right) \nabla \ln \pi(A_t | S_t, \theta_{\text{old}})$
 $= [-0.5 \ 0.5]T + 0.01 \left(\frac{1+0.9-0.3}{(0.2)^2} \right) [1 \ 0 \ 0]^T$
 $= [-0.5 \ 0.5]T + 0.01 \left(\frac{8}{(0.2)^2} \right) [1 \ 0 \ 0]^T$
 $= \boxed{[-0.42 \ 0 \ 0.5]}$

One-step Actor-Critic (episodic), for estimating $\pi_\theta \approx \pi_c$.
Input: a differentiable policy parameterization $\pi(a|s)$,
Input: a differentiable state-action function parameterization $v(s, w)$
Parameters: step size $\alpha^P > 0, \alpha^W > 0$
Initial policy parameter $\theta \in \mathbb{R}^d$ and state-value weights $w \in \mathbb{R}^d$ (e.g., to 0)
Long Keeve (for each episode):
Initialise S (first state of episode)
 $I \leftarrow 1$
Loop while S is not terminal (for each time step):
 $A \sim \pi(\cdot | S, \theta)$
 $R \sim P(A, S)$ if A observes S' , R
 $\delta \leftarrow R + \gamma v(S', w) - v(S, w)$ (If S' is terminal, then $v(S', w) \leftarrow 0$)
 $w \leftarrow w + \alpha^W \delta \nabla v(S, w)$
 $\theta \leftarrow \theta + \alpha^P (\delta \nabla \ln \pi(A | S, \theta))$
 $I \leftarrow I + 1$
 $S \leftarrow S'$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Deep Reinforcement Learning
Prof. Vimal SP
CSIS



BITS Pilani





<AIMLCZG512, Deep Reinforcement Learning> Lecture No. 14,15

Agenda for the classes

- Introduction
- Upper-Confidence-bound Action Selection
- Monte-Carlo Tree Search
- AlphaGo Zero
- MuZero, PlaNet



Work Integrated Learning Programmes

Acknowledgements: Some of the slides were adopted with permission from the course [CSCE-689](#) (Texas A&M University) by Prof. Guni Sharon

3

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Model Based Algorithms

- Learn the model of an environment's transition dynamics or make use of a known dynamic model
- Once an agent has a model of the environment, $P(s'|s,a)$, it can "imagine" what will happen in the future by predicting the trajectory for a few time steps.
- If the environment is in state s , an agent can estimate how the state will change if it makes a sequence of actions a_1, a_2, \dots, a_n by repeatedly applying $P(s'|s,a)$ all without actually producing an action to change the environment.
- Hence, the predicted trajectory occurs in the agent's "head" using a model.

Model Based Algorithms

- Most commonly applied to games with a target state, such as winning or losing in a game of chess, or navigation tasks with a goal state s^* .
- Do not model any rewards

Work Integrated Learning Programmes

Model based algorithms - Advantages

- Very appealing : it can play out scenarios and understand the consequences of its actions without having to actually act in an environment.
- Require many fewer samples of data to learn good policies since having a model enables an agent to supplement its actual experiences with imagined ones.

Upper Confidence bound action selection

UCB is a deterministic algorithm for Reinforcement Learning that focuses on exploration and exploitation based on a confidence boundary that the algorithm assigns to each machine on each round of exploration. These boundary decreases when a machine is used more in comparison to other machines.

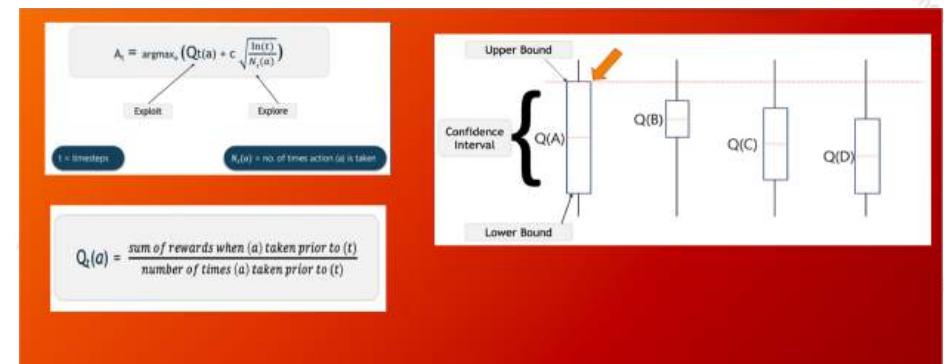
The **Upper Confidence Bound** follows the principle of optimism in the face of uncertainty which implies that if we are uncertain about an action, we should optimistically assume that it is the correct action.

Initially, UCB explores more to systematically reduce uncertainty but its exploration reduces over time. Thus we can say that UCB obtains greater reward on average than other algorithms such as Epsilon-greedy, Optimistic Initial Values, etc.

Model based algorithms - challenges

- Models are hard to come by.
- An environment with a large state space and action space can be very difficult to model; doing so may even be intractable, especially if the transitions are extremely complex
- Models are only useful when they can accurately predict the transitions of an environment many steps into the future - prediction errors

UCB Algorithm





Upper Confidence bound action selection

Steps followed in Upper Confidence Bound

1. At each round n , we consider two numbers for machine m .
-> $N_m(n)$ = number of times the machine m was selected up to round n .
-> $R_m(n)$ = number of rewards of the machine m up to round n .
2. From these two numbers we have to calculate,
a. The average reward of machine m up to round n , $r_m(n) = R_m(n) / N_m(n)$.
b. The confidence interval $[r_m(n) - \Delta_m(n), r_m(n) + \Delta_m(n)]$ at round n with, $\Delta_m(n) = \sqrt{1.5 * \log(n) / N_m(n)}$
3. We select the machine m that has the maximum UCB, $(r_m(n) + \Delta_m(n))$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Monte Carlo Tree Search

- Monte Carlo Tree Search (MCTS) is a heuristic search set of rules that has won big attention and reputation within the discipline of synthetic intelligence, specially in the area of choice-making and game playing.
- MCTS combines the standards of Monte Carlo strategies, which rely upon random sampling and statistical evaluation, with tree-primarily based search techniques.
- **Idea :**
 - build a seek tree incrementally by using simulating more than one random performs (regularly known as rollouts or playouts) from the current recreation nation.
 - These simulations are carried out until a terminal state or a predefined intensity is reached.
 - The results of these simulations are then backpropagated up the tree, updating the records of the nodes visited at some stage in the play, which includes the wide variety of visits and the win ratios.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Monte Carlo Tree Search

- It is a probabilistic and heuristic driven search algorithm that combines the classic tree search implementations alongside machine learning principles of reinforcement learning.
- **exploration-exploitation trade-off :** exploits the actions and strategies that is found to be the best till now but also must continue to explore the local space of alternative decisions and find out if they could replace the current best.
- exploration expands the tree's breadth more than its depth.But it quickly becomes inefficient in situations with large number of steps or repetitions.
- Exploitation sticks to a single path that has the greatest estimated value. This is a greedy approach and this will extend the tree's depth more than its breadth.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Monte Carlo Tree Search

Why use Monte Carlo Tree Search (MCTS) ?

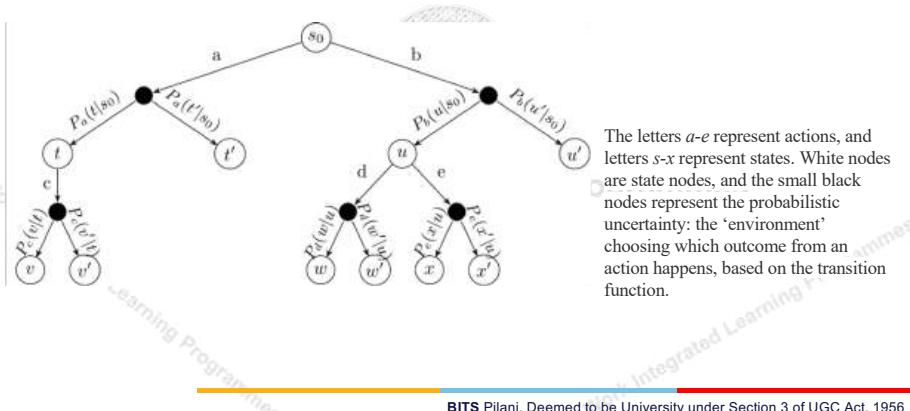
1. Handling Complex and Strategic Games
2. Unknown or Imperfect Information
3. Learning from Simulations
4. Optimizing Exploration and Exploitation
5. Scalability and Parallelization
6. Applicability Beyond Games
7. Domain Independence

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Monte Carlo Tree Search

Nodes are the building blocks of the search tree and are formed based on the outcome of a number of simulations.

MDPs can be represented as trees (or graphs), called *ExpectiMax* trees:



Monte Carlo Tree Search – Overview

The algorithm is online, which means the action selection is interleaved with action execution. Thus, MCTS is invoked every time an agent visits a new state.

Fundamental features:

1. The Q-value $Q(s,a)$ for each is approximated using *random simulation*.
2. For a single-agent problem, an *ExpectiMax search tree* is built incrementally
3. The search terminates when some pre-defined computational budget is used up, such as a time limit or a number of expanded nodes. Therefore, it is an *anytime* algorithm, as it can be terminated at any time and still give an answer.
4. The best performing action is returned.
 - o This is complete if there are *no dead-ends*.
 - o This is optimal if an entire search can be performed (which is unusual – if the problem is that small we should just use a dynamic programming technique such as value iteration).

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

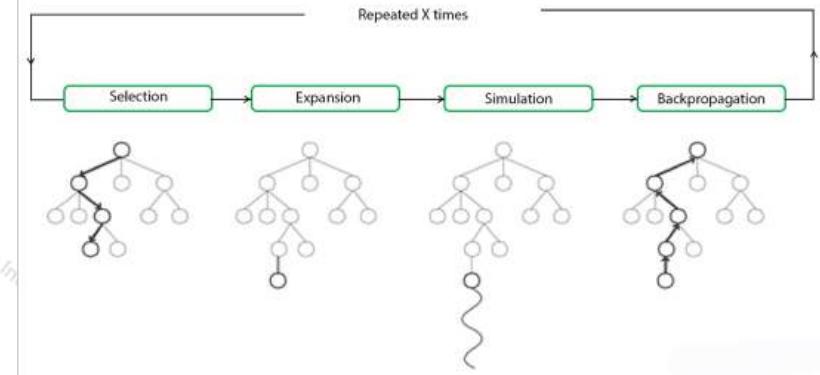
The Framework of MCTS

The basic framework is to build up a tree using simulation. The states that have been evaluated are stored in a search tree. The set of evaluated states is *incrementally* built by iterating over the following four steps:

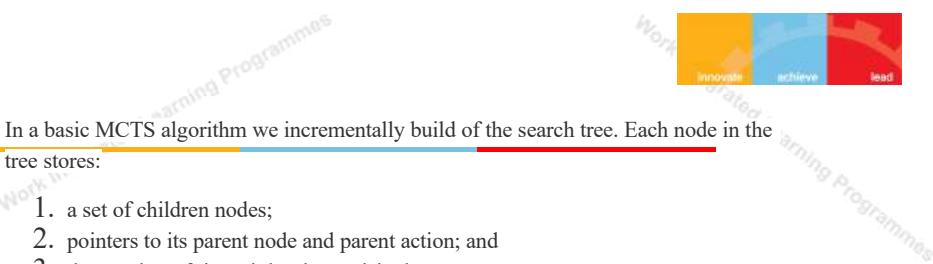
- **Select:** Select a single node in the tree that is *not fully expanded*. By this, we mean at least one of its children is not yet explored.
- **Expand:** Expand this node by applying one available action (as defined by the MDP) from the node.
- **Simulation:** From one of the outcomes of the expanded, perform a complete random simulation of the MDP to a terminating state. This therefore assumes that the simulation is finite, but versions of MCTS exist in which we just execute for some time and then estimate the outcome.
- **Backpropagate:** Finally, the value of the node is *backpropagated* to the root node, updating the value of each ancestor node on the way using expected value

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

MCTS Framework



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Algorithm – Monte-Carlo Tree Search

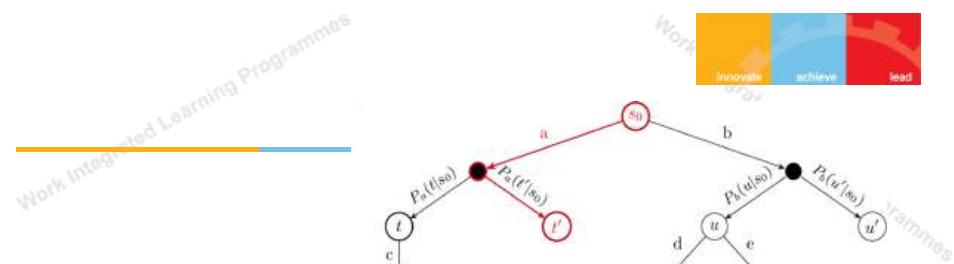
```

Input: MDP  $M = \langle S, s_0, A, P_a(s' | s), r(s, a, s') \rangle$ , base value function  $Q$ , time limit  $T$ .
Output: updated Q-function  $Q$ 

while  $currentTime < T$ 
     $selected\_node \leftarrow Select(s_0)$ 
     $child \leftarrow Expand(selected\_node)$  – expand and choose a child to simulate
     $G \leftarrow Simulate(child)$  – simulate from  $child$ 
    Backpropagate( $selected\_node, child, G$ )
return  $Q$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Selection: The first loop progressively selects a branch in the tree using a multi-armed bandit algorithm using $Q(s|a)$. The outcome that occurs from an action is chosen according to $P(s'|s,a)$ defined in the MDP.

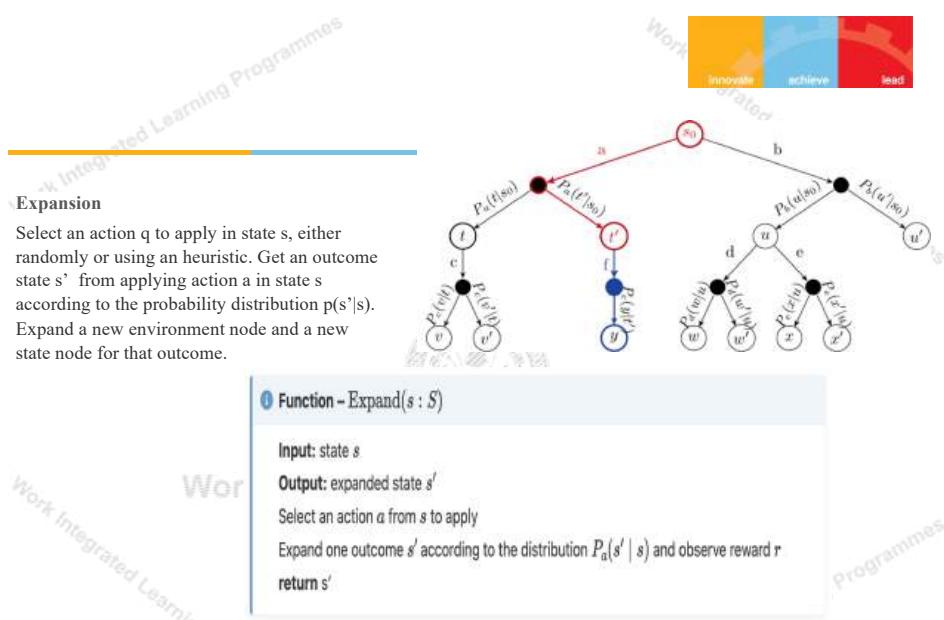
Function – Select($s : S$)

```

Input: state  $s$ 
Output: unexpanded state
while  $s$  is fully expanded
    Select action  $a$  to apply in  $s$  using a multi-armed bandit algorithm
    Choose one outcome  $s'$  according to  $P_a(s' | s)$ 
     $s \leftarrow s'$ 
return  $s$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Expansion

Select an action q to apply in state s , either randomly or using an heuristic. Get an outcome state s' from applying action a in state s according to the probability distribution $p(s'|s)$. Expand a new environment node and a new state node for that outcome.

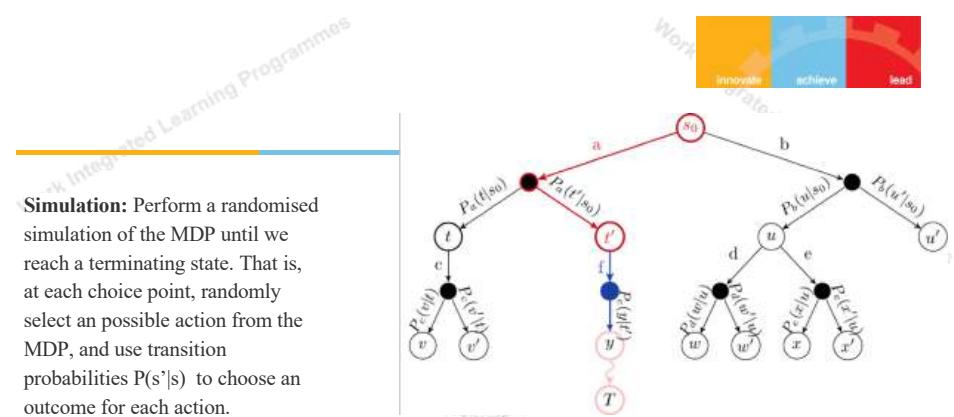
Function – Expand($s : S$)

```

Input: state  $s$ 
Output: expanded state  $s'$ 
Select an action  $a$  from  $s$  to apply
Expand one outcome  $s'$  according to the distribution  $P_a(s' | s)$  and observe reward  $r$ 
return  $s'$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Simulation: Perform a randomised simulation of the MDP until we reach a terminating state. That is, at each choice point, randomly select an possible action from the MDP, and use transition probabilities $P(s'|s)$ to choose an outcome for each action.

Heuristics can be used to improve the random simulation by guiding it towards more promising states. G is the cumulative discounted reward received from the simulation starting at s^* until the simulation terminates.

To avoid memory explosion, we discard all nodes generated from the simulation. In any non-trivial search, we are unlikely to ever need them again.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Backpropagation: The reward from the simulation is backpropagated from the selected node to its ancestors recursively. We must not forget the *discount factor!* For each state s and action a selected in the Select step, update the cumulative reward of that state.

```

Procedure – Backpropagation( $s : S; a : A$ )
Input: state-action pair  $(s, a)$ 
Output: none
do
     $N(s, a) \leftarrow N(s, a) + 1$ 
     $G \leftarrow r + \gamma G$ 
     $Q(s, a) \leftarrow Q(s, a) + \frac{1}{N(s, a)}[G - Q(s, a)]$ 
     $s \leftarrow \text{parent of } s$ 
     $a \leftarrow \text{parent action of } s$ 
while  $s \neq s_0$ 

```

Example : Backpropagation

Consider the following ExpectiMax tree that has been expanded several times. Assume $\gamma = 0.8$, $r = X$ represents reward X received at a state, V represents the value of the state (the value $\max_{a' \in \text{children}} Q(s, a')$) and the length of the simulation is 14. After the simulation step, but before backpropagation, our tree would look like this:

Work Integrated Learning Programmes

Innovate achieve lead

In the next iteration, the red actions are selected, and the blue node is expanded to its three children nodes. A simulation is run from y , which terminates after 3 steps. A reward of 31.25 is received in the terminal state. This would mean that the discounted reward returned at node y would be $\gamma \times 31.25$, which is 20.

Before backpropagation, we have the following:

$$Q(s, a) = 18$$

$$Q(t, f) = 0$$

We leave out the remainder of the Q-function as it is not relevant to the example.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

The backpropagation step is then calculated for the nodes y , t , and s as follows:

$$Q(y, g) = \gamma^2 \times 31.25 \quad (\text{simulation is 3 steps long and receives reward of 31.25})$$

$$= 20$$

$$Q(t, f) = Q(t, f) + \frac{1}{N(t, f)}[r + \gamma G - Q(t, f)]$$

$$= 0 + \frac{1}{2}[0 + 0.8 \cdot 20 - 0]$$

$$= 8$$

$$Q(s, a) = Q(s, a) + \frac{1}{N(s, a)}[r + \gamma G - Q(s, a)]$$

$$= 18 + \frac{1}{5}[6 + 0.8 \cdot (0.8 \cdot 20) - 18]$$

$$= 18 + \frac{1}{5}[6 + 12.8 - 18]$$

$$= 18.16$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Action selection

Once we have run out of computational time, we select the action that maximises our expected return, which is simply the one with the highest Q-value from our simulations:

$$\text{argmax}_{a \in A(s)} Q(s, a)$$

We execute that action and wait to see which outcome occurs for the action.

Once we see the outcome state, which we will call s' , we start the process all over again, except with $s_0 \leftarrow s'$.

However, importantly, we can keep the sub-tree from state s' , as we already have done simulations from that state. We discard the rest of the tree (all child of s_0 other than the chosen action) and incrementally build from s' .

AlphaZero

- Combining MCTS and TD learning: Alpha Zero
- Alpha Zero (or more accurately its predecessor AlphaGo) made headlines when it beat Go world champion Lee Sedol in 2016.
- It uses a combination of MCTS and (deep) reinforcement learning to learn a policy.

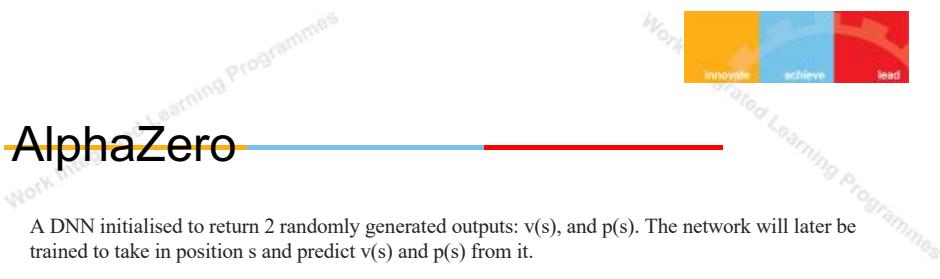
Issues in Monte Carlo Tree Search:

1. Exploration-Exploitation Trade-off
2. Sample Efficiency
3. High Variance
4. Heuristic Design
5. Computation and Memory Requirements
6. Overfitting
7. Domain-specific Challenges

	Value iteration	MCTS
Cost	Higher cost (exhaustive)	Lower cost (does not solve for entire state space)
Coverage/ Robustness	Higher (works from any state)	Lower (works only from initial state or state reachable from initial state)

AlphaZero - Overview

1. AlphaZero uses a deep neural network to estimate the Q-function. More accurately, it gives an estimate of the probability of selecting action a in state s , $P(a|s)$ and the *value* of the state ($V(s)$), which represents the probability of the player winning from s .
2. It is trained via *self-play*. Self-play is when the same policy is used to generate the moves of both the learning agent and any of its opponents. In AlphaZero, this means that initially, both players make random moves, but both also learn the same policy and use it to select subsequent moves.
3. In short, AlphaZero is a **game-playing program** that, through a combination of self-play and neural network reinforcement learning (more on that later), is able to learn to play games such as chess and Go from scratch — that is, after being fed nothing more than the rules of said games.

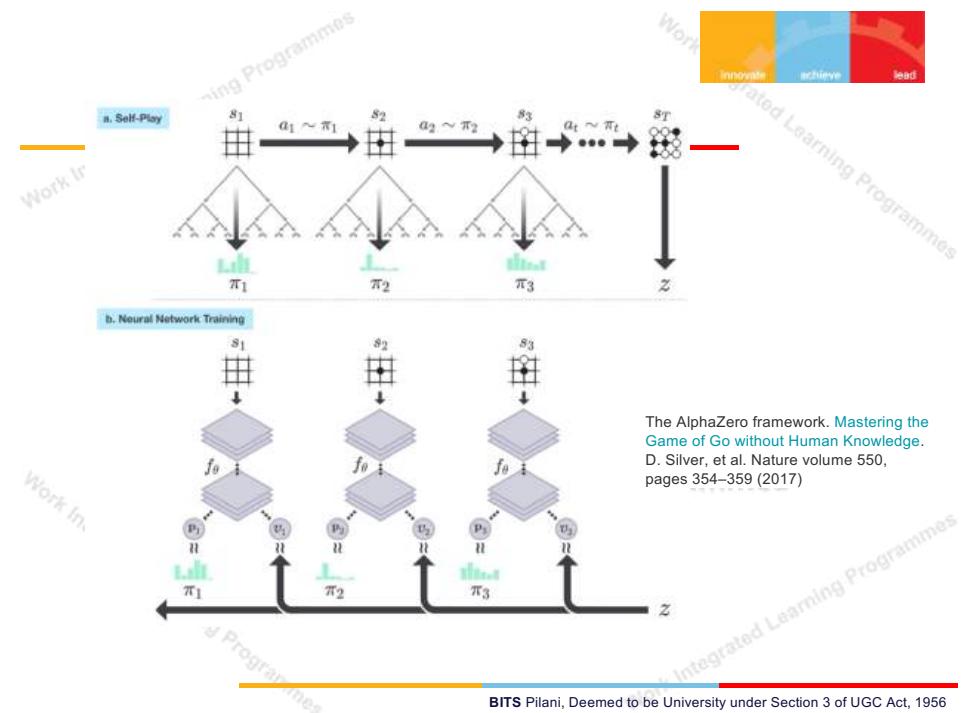


AlphaZero

A DNN initialised to return 2 randomly generated outputs: $v(s)$, and $p(s)$. The network will later be trained to take in position s and predict $v(s)$ and $p(s)$ from it.

At each move, AlphaZero:

1. Executes an MCTS search using UCB-like selection: $Q(s, a) + P(s, a)/1 + N(s, a)$, which returns the probabilities of playing each move.
2. The neural network is used to guide the MCTS by influencing $Q(s, a)$.
3. The final result of a simulated game is used as the reward for each simulation.
4. After a set number of MCTS simulations, the best move is chosen for self-play.
5. Repeat steps 1-4 for each move until the self-play game ends.
6. Then, feedback the result of the self-play game to update the Q function for each move.



Implications of AlphaZero

- Increase the probability of a win come the end of the game
- AlphaZero's evaluation function is the product of a highly sophisticated neural network's experience accumulated over millions of games
- AlphaZero does not make use of any human knowledge. we can expect it to come up with brand-new ideas prev

Tree Search: Evaluate the MCTS value of a move by playing out many games from that move.

Heuristic Evaluation: Use DCNN to choose moves with highest predicted value within simulated playouts.

Learning: Update DCNN using new MCTS value for state.

MuZero

- AlphaZero (2017) - learning to play on its own, without any human data or domain knowledge, or even by mastering three different games (Go, Chess, Shogi) with only one algorithm being only provided with the known set of rules.
- MuZero (2020) - learnt to master these games (Go, Chess, Shogi, Atari) without being provided known rules
- Big step forward towards general-purpose algorithms

MuZero

MuZero models three elements that are key to the planner:

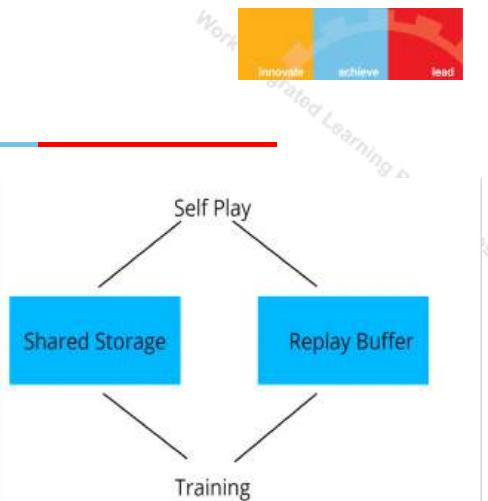
- **Value:** how good is the current state?
- **Policy:** which action should be the next one?
- **Reward:** how good was the last action taken?



MuZero - overview

- The MuZero algorithm aims to accurately predict characteristics and details of the future which it deems important for planning.
- The algorithm initially receives an input, for instance an image of a chess board, which is translated into a hidden state.
- The hidden state then undergoes iterations based on the previous hidden state and a proposed subsequent plan of action.
- Each time the hidden state is updated the model predicts three variables: policy, value function and immediate reward.
- The policy is the next move to be played, the value function is the predicted winner, and the immediate reward is the strength of the move (if it improves the player's position).
- The model is then trained to accurately predict the values of the three aforementioned variables.

- The MuZero algorithm does not receive any rules of the game, for instance in a chess game, the algorithm is unaware of the legal moves or what constitutes as win, draw or a loss.
- MuZero only receives the rewards of its actions when a game is terminated, either by winning, drawing, or losing. However, during its learning, MuZero receives rewards periodically based on how well it is doing in its own mini-games.



Training and Loss function

- function `train_network` repeatedly trains the neural network by making use of the `replayBuffer`.
- The `train_network` function works by looping the total number of training steps, which is set to one million by default.
- The function then samples a batch at every step and uses that data to update the neural network.
- The tuples within a batch are: the current state of the game, a list of actions taken from the current position and lastly the targets used to train the neural networks.
- The targets used to train the neural networks are calculated by using Temporal Difference (TD) Learning .
- The loss function of MuZero is responsible for how the weights of the neural network are updated.

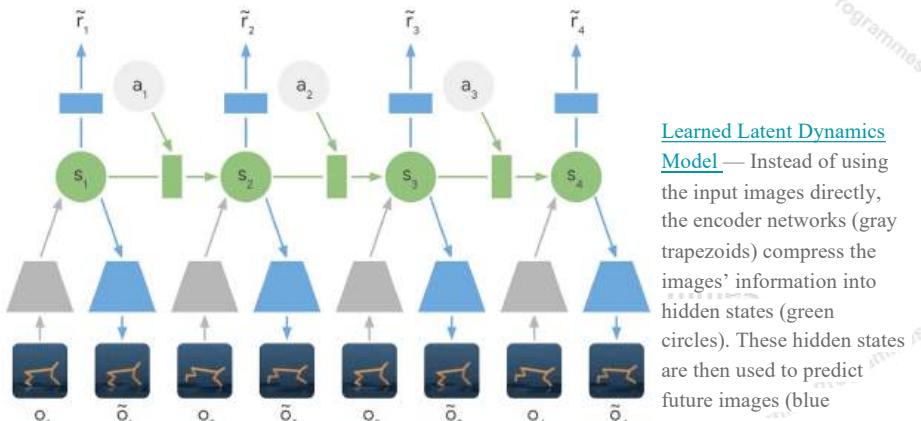
PlaNet

- Deep Planning Network - Google AI & DeepMind Initiative
- PlaNet agent was tasked with ‘planning’ a sequence of actions to achieve a goal like pole balancing, teaching a virtual entity (human or cheetah) to walk, or keeping a box rotating by hitting it in a specific location.
- Common goals between these tasks that the PlaNet needed to achieve:
 1. The Agent needs to predict a variety of possible futures (for robust planning)
 2. The Agent needs to update the plan based on the outcomes/rewards of a recent action.
 3. The Agent needs to retain information over many time steps

So how did the Google AI team achieve these goals?

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

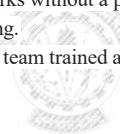
- 1) Latent Dynamics Model** - Key benefits to using compact latent state spaces are that it allows the agent to learn more abstract representations like the objects' positions and velocities and



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

PlaNet Overview

- 1. Learning with a latent dynamics model** — PlaNet learns from a series of hidden or latent states *instead of images* to predict the latent state moving forward.
- 2. Model-based planning** — PlaNet works without a policy network and instead makes decisions based on continuous planning.
- 3. Transfer learning** — The Google AI team trained a single PlaNet agent to solve all six different tasks.



Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

2. Model based planning

Better
Sample Efficient

Model-based
(100 time steps)

- one agent for all tasks

Off-policy
Q-learning
(1 M time steps)



Actor-critic

On-policy
Policy Gradient
(10 M time steps)

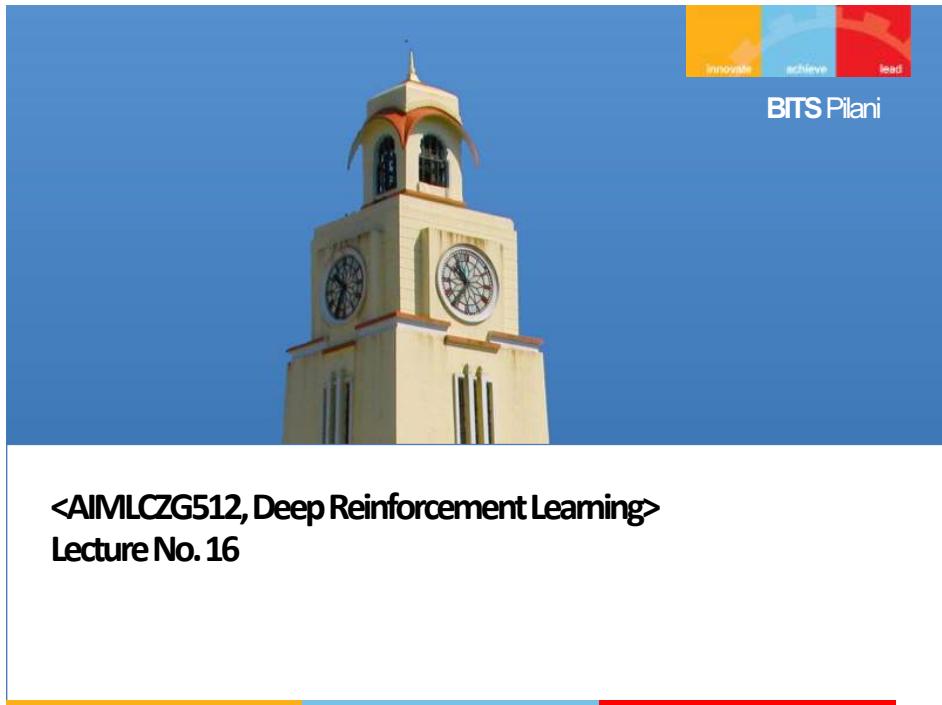
Less
Sample Efficient

Evolutionary/
gradient-free
(100 M time steps)

Work Integrated Learning Programmes

The agent is randomly placed into different environments without knowing the task, so it needs to infer the task from its image observations. Without changes to the hyper parameters, the multi-task agent achieves the same mean performance as individual agents.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Agenda for the classes

- Introduction
- Imitation Learning Via Supervised Learning
- Behavior Cloning
- DAGGER
- Inverse Reinforcement Learning, GAIL
- Applications



3

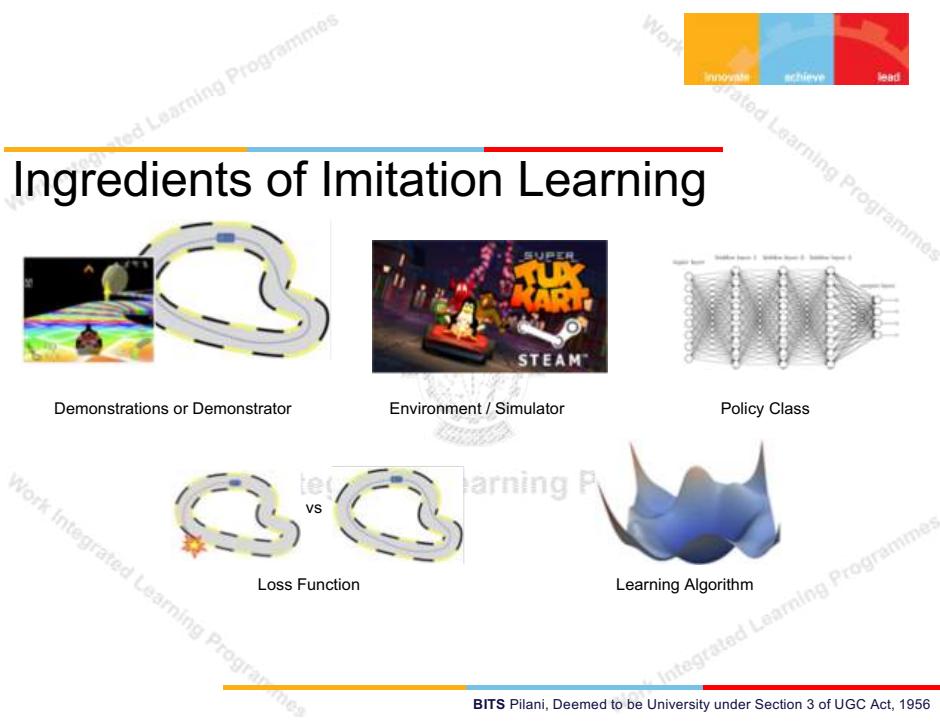
Imitation Learning in a Nutshell

Given: demonstrations or demonstrator

Goal: train a policy to mimic demonstrations



Images from Stephane Ross



Some Interesting Examples

- **ALVINN** <https://www.ri.cmu.edu/publications/alvinn-an-autonomous-land-vehicle-in-a-neural-network/>
Dean Pomerleau et al., 1989-1999 <https://www.youtube.com/watch?v=ilP4aPDTBPE>
- **Helicopter Acrobatics**
Learning for Control from Multiple Demonstrations - Adam Coates, Pieter Abbeel, Andrew Ng, ICML 2008
An Application of Reinforcement Learning to Aerobatic Helicopter Flight - Pieter Abbeel, Adam Coates, Morgan Quigley, Andrew Y. Ng, NIPS 2006
<https://www.youtube.com/watch?v=0JL04JJjocc>
- **Ghosting (Sports Analytics)** - Next Slide.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Notation & Set-up

State: s (sometimes x) (**state may only be partially observed)

Action: a (sometimes y)

Policy: π_θ (sometimes h)

- Policy maps states to actions: $\pi_\theta(s) \rightarrow a$
- ...or distributions over actions: $\pi_\theta(s) \rightarrow P(a)$

State Dynamics: $P(s'|s,a)$

- Typically not known to policy.
- Essentially the simulator/environment

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Notation & Set-up

Rollout: sequentially execute $\pi(s_0)$ on an initial state

- Produce trajectory $\tau = (s_0, a_0, s_1, a_1, \dots)$

P($\tau|\pi$): distribution of trajectories induced by a policy

- Sample s_0 from P_0 (distribution over initial states), initialize $t = 1$.
- Sample action a_t from $\pi(s_{t-1})$
- Sample next state s_t from applying a_t to s_{t-1} (requires access to environment)
- Repeat from Step 2 with $t=t+1$

P($s|\pi$): distribution of states induced by a policy

- Let $P_t(s|\pi)$ denote distribution over t -th state
- $P(s|\pi) = (1/T)\sum_t P_t(s|\pi)$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Example #1: Racing Game (Super Tux Kart)

s = game screen
a = turning angle

Training set: $D = \{r := (s, a)\}$ from π^*

- s = sequence of s
- a = sequence of a

Goal: learn $\pi_\theta(s) \rightarrow a$

Images from Stephane Ross

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Example #2: Basketball Trajectories

s = location of players & ball
a = next location of player

Training set: $D = \{r := (s, a)\}$ from π^*

- s = sequence of s
- a = sequence of a

Goal: learn $\pi_\theta(s) \rightarrow a$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Innovate achieve lead

Behavioral Cloning

= Reduction to Supervised Learning (Ignoring regularization for brevity.)

Define $P^* = P(s|\pi^*)$ (distribution of states visited by expert)

(recall $P(s|\pi^*) = (1/T)\sum_t P_t(s|\pi^*)$)

(sometimes abuse notation: $P^* = P(s, a^* = \pi^*(s)|\pi^*)$)

Learning objective:

$$\text{argmin}_\theta E_{(s,a^*) \sim P^*} L(a^*, \pi_\theta(s))$$

Interpretations:

- Assuming perfect imitation so far, learn to continue imitating perfectly
- Minimize 1-step deviation error along the expert trajectories

Images from Stephane Ross

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Behavioral Cloning vs. Imitation Learning

Behavioral Cloning (Supervised Learning):

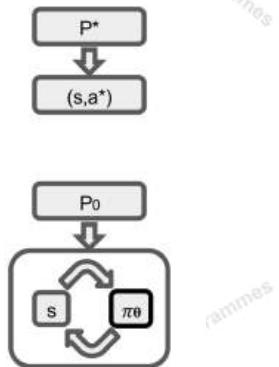
$$\text{argmax}_{\theta} E_{(s,a^*) \sim P^*} L(a^*, \pi_\theta(s))$$

Distribution provided exogenously

(General) Imitation Learning:

$$\text{argmax}_{\theta} E_{s \sim P(s|\theta)} L(\pi^*(s), \pi_\theta(s))$$

Distribution depends on rollout.
 $P(s|\theta)$ = state distribution of π_θ

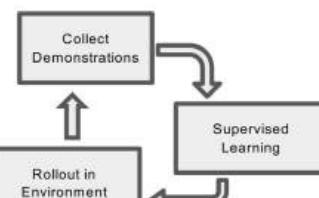


Types of Imitation Learning

Behavioral Cloning
 $\text{argmax}_{\theta} E_{(s,a) \sim P^*} L(a^*, \pi_\theta(s))$
 Works well when P^* close to P_θ

Inverse RL
 Learn r such that:
 $\pi^* = \text{argmax}_{\theta} E_{s \sim P(s|\theta)} r(s, \pi_\theta(s))$
 RL problem
 Assumes learning r is statistically easier than directly learning π^*

Direct Policy Learning
 via Interactive Demonstrator



Requires Interactive Demonstrator
 (BC is 1-step special case)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Types of Imitation Learning

	Direct Policy Learning	Reward Learning	Access to Environment	Interactive Demonstrator	Pre-collected Demonstrations
Behavioral Cloning	Yes	No	No	No	Yes
Direct Policy Learning (Interactive IL)	Yes	No	Yes	Yes	Optional
Inverse Reinforcement Learning	No	Yes	Yes	No	Yes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Interactive Expert

Can query expert at any state

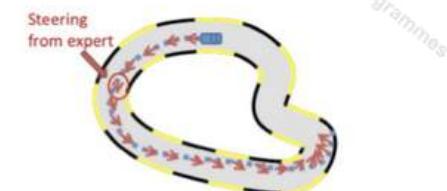
Construct loss function

- $L(\pi^*(s), \pi(s))$

Typically applied to rollout trajectories

- $s \sim P(s|\pi)$

Driving example: $L(\pi^*(s), \pi(s)) = (\pi^*(s) - \pi(s))^2$



Example from Super Tux Kart
 (Image courtesy of Stephane Ross)

Expert provides feedback on
 state visited by policy

Images from Stephane Ross

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Interactive Expert

Can query expert at any state

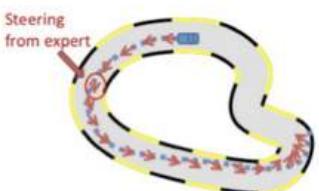
Construct loss function

- $L(\pi^*(s), \pi(s))$

Typically applied to rollout trajectories

- $s \sim P(s|\pi)$

Driving example: $L(\pi^*(s), \pi(s)) = (\pi^*(s) - \pi(s))^2$



Example from Super Tux Kart
 (Image courtesy of Stephane Ross)

Expert provides feedback on
 state visited by policy

Images from Stephane Ross

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



DAGGER: Dataset Aggregation

```

Initialize  $\mathcal{D} \leftarrow \emptyset$ .
Initialize  $\hat{\pi}_1$  to any policy in  $\Pi$ .
for  $i = 1$  to  $N$  do
    Let  $\pi_i = \beta_i\pi^* + (1 - \beta_i)\hat{\pi}_i$ .
    Sample  $T$ -step trajectories using  $\pi_i$ .
    Get dataset  $\mathcal{D}_i = \{(s, \pi^*(s))\}$  of visited states by  $\pi_i$ 
    and actions given by expert.
    Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ .
    Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ .
end for
Return best  $\hat{\pi}_i$  on validation.

```

- Idea: Get more labels of the expert action along the path taken by the policy computed by behavior cloning
- Obtains a stationary deterministic policy with good performance under its induced state distribution

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Recycling Robot

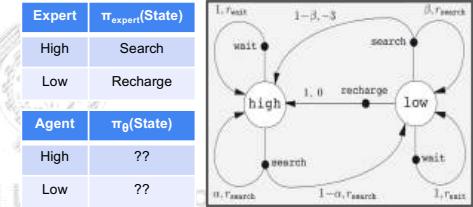
Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}
- A(low) = {search, wait, recharge}

$$\text{Obj: } \hat{\pi}^* = \arg \min_{\pi} \sum_{\zeta \in \Sigma} \sum_{x \in \zeta} L(\pi(x), \pi^*(x))$$



Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Recycling Robot

Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}
- A(low) = {search, wait, recharge}

Behavioral Cloning Exp-1:

Train the SL model using training data

$$P(\text{Search} | \text{High}) = 0.7, P(\text{Wait} | \text{High}) = 0.3$$

$$P(\text{Search} | \text{Low}) = 0.3, P(\text{Wait} | \text{Low}) = 0.6, P(\text{Recharge} | \text{Low}) = 0.3$$

Evaluate model Loss

Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1



Recycling Robot

Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

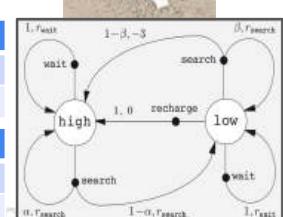
Action = function (Charge level based reactions)

- A(high) = {search, wait}
- A(low) = {search, wait, recharge}

Behavioral Cloning Exp-1:

Behave (Rollout) using π_{θ_n}

$$\tau: (\text{High}, \text{Search}, +1, \text{Low}) \rightarrow (\text{Low}, \text{Wait}, 0.5, \text{Low}) \rightarrow (\text{Low}, \text{Wait}, 0.5, \text{Low}), \dots$$



- $r_{\text{search}} = 1$
- $r_{\text{wait}} = 0.5$
- No Charge = -3
- Recharge = 0

Loss

Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1

$$L(\theta_n) = \sum_i 1 (\pi_{\text{expert}}(\text{State}) \neq \pi_{\theta_n}(\text{State})) = 0+1+1+\dots^{(\infty)} = \text{High Error}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}
- A(low) = {search, wait, recharge}

Behavioral Cloning Exp-1:

Behave (Rollout) using π_{θ_n}

τ : (High, Search, Low) \rightarrow (Low, Wait, Low) \rightarrow (Low, Wait, Low).....

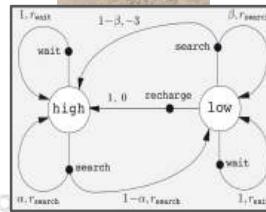
Observation: Compounding Error due to Distributional Shift

Loss

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

$$L(\theta_n) = \sum_i 1 (\pi_{\text{expert}}(\text{State}) \neq \pi_{\theta_n}(\text{State})) = 0 + 1 + 1 + \dots + (\infty) = \text{High Error}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}
- A(low) = {search, wait, recharge}

DAGGER

Behavioral Cloning with Data augmentation Exp-2:

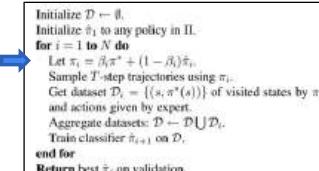
Behave (Rollout) using π_{θ_n}

$$\beta = 0.6 \rightarrow \pi_{\theta_i} = 0.6 \pi_{\text{expert}} + 0.4 \pi_{\text{agent}}$$

E \rightarrow A \rightarrow E

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

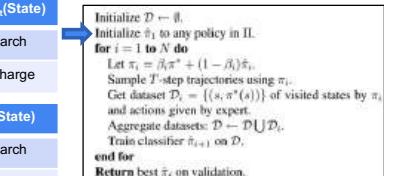
Action = function (Charge level based reactions)

- A(high) = {search, wait}
- A(low) = {search, wait, recharge}

DAGGER

Behavioral Cloning with Data augmentation Exp-2:

SL Training initializes the agent policy



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}
- A(low) = {search, wait, recharge}

DAGGER

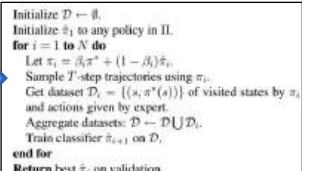
Behavioral Cloning with Data augmentation Exp-2:

Behave (Rollout) using π_{θ_n}

$$\beta = 0.6 \rightarrow \pi_{\theta_i} = 0.6 \pi_{\text{expert}} + 0.4 \pi_{\text{agent}}$$

E \rightarrow A \rightarrow E

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}

- A(low) = {search, wait, recharge}

DAGGER

Behavioral Cloning with Data augmentation Exp-2:

Behave (Rollout) using π_{θ_n}

τ : (High, Search, Low) \rightarrow (Low, Wait, Low) \rightarrow (Low, Recharge, High)

E: Search Recharge Recharge

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

Expert $\pi_{\text{expert}}(\text{State})$

High Search

Low Recharge

Agent $\pi_{\theta}(\text{State})$

High Search

Low Wait

```
Initialize  $\mathcal{D} \leftarrow \emptyset$ .
Initialize  $\hat{\pi}_1$  to any policy in IL.
for  $i = 1$  to  $N$  do
    Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i)\hat{\pi}_i$ .
    Sample  $T$ -step trajectories using  $\pi_i$ .
    Get dataset  $\mathcal{D}_i = \{(s, \pi^*(a))\}$  of visited states by  $\pi_i$ 
    and actions given by expert.
    Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ .
    Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ .
end for
Return best  $\hat{\pi}_n$  on validation.
```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}

- A(low) = {search, wait, recharge}

DAGGER

Behavioral Cloning with Data augmentation Exp-2:

Behave (Rollout) using π_{θ_n}

τ : (High, Search, Low) \rightarrow (Low, Wait, Low) \rightarrow (Low, Recharge, High)

E: Search Recharge Recharge

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

Expert $\pi_{\text{expert}}(\text{State})$

High Search

Low Recharge

Agent $\pi_{\theta(n+1)}(\text{State})$

High Search

Low Recharge

```
Initialize  $\mathcal{D} \leftarrow \emptyset$ .
Initialize  $\hat{\pi}_1$  to any policy in IL.
for  $i = 1$  to  $N$  do
    Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i)\hat{\pi}_i$ .
    Sample  $T$ -step trajectories using  $\pi_i$ .
    Get dataset  $\mathcal{D}_i = \{(s, \pi^*(a))\}$  of visited states by  $\pi_i$ 
    and actions given by expert.
    Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ .
    Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ .
end for
Return best  $\hat{\pi}_{n+1}$  on validation.
```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}

- A(low) = {search, wait, recharge}

DAGGER

Behavioral Cloning with Data augmentation Exp-2:

Behave (Rollout) using π_{θ_n}

τ : (High, Search, Low) \rightarrow (Low, Wait, Low) \rightarrow (Low, Recharge, High)

E: Search Recharge Recharge

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

Expert $\pi_{\text{expert}}(\text{State})$

High Search

Low Recharge

Agent $\pi_{\theta}(\text{State})$

High Search

Low Wait

```
Initialize  $\mathcal{D} \leftarrow \emptyset$ .
Initialize  $\hat{\pi}_1$  to any policy in IL.
for  $i = 1$  to  $N$  do
    Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i)\hat{\pi}_i$ .
    Sample  $T$ -step trajectories using  $\pi_i$ .
    Get dataset  $\mathcal{D}_i = \{(s, \pi^*(a))\}$  of visited states by  $\pi_i$ 
    and actions given by expert.
    Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ .
    Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ .
end for
Return best  $\hat{\pi}_n$  on validation.
```

x	y
.....
High	Search
Low	Recharge
Low	Recharge

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}

- A(low) = {search, wait, recharge}

DAGGER

Behavioral Cloning with Data augmentation Exp-2:

Behave (Rollout) using $\pi_{\theta_{n+1}}$

τ : (High, Search, Low) \rightarrow (Low, Recharge, High) \rightarrow (Low, Recharge, High)

E: Search Recharge Recharge

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

Expert $\pi_{\text{expert}}(\text{State})$

High Search

Low Recharge

Agent $\pi_{\theta(n+1)}(\text{State})$

High Search

Low Recharge

```
Initialize  $\mathcal{D} \leftarrow \emptyset$ .
Initialize  $\hat{\pi}_1$  to any policy in IL.
for  $i = 1$  to  $N$  do
    Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i)\hat{\pi}_i$ .
    Sample  $T$ -step trajectories using  $\pi_i$ .
    Get dataset  $\mathcal{D}_i = \{(s, \pi^*(a))\}$  of visited states by  $\pi_i$ 
    and actions given by expert.
    Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ .
    Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ .
end for
Return best  $\hat{\pi}_{n+1}$  on validation.
```

x	y
.....
High	Search
Low	Recharge
Low	Recharge
High	Search
Low	Recharge
Low	Recharge

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}

- A(low) = {search, wait, recharge}

DAGGER

Behavioral Cloning with Data augmentation Exp-2:

Behave (Rollout) using $\pi_{\theta_{n+1}}$

τ : (High, Search, Low) \rightarrow (Low, Recharge, High) \rightarrow (Low, Recharge, High)

Observation: What if the π_{expert} is suboptimal? Trusts that

E: expert behavior is reward maximized!

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

Expert	$\pi_{\text{expert}}(\text{State})$
High	Search
Low	Recharge

Agent	$\pi_{\theta(n+1)}(\text{State})$
High	Search
Low	Recharge

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}

- A(low) = {search, wait, recharge}

Variations of DAGGER (Few Ideas!)

Behavioral Cloning with Data augmentation Exp-2:

Behave (Rollout) using π_{θ_n}

E \rightarrow A \rightarrow A

τ : (High, Search, +1, Low) \rightarrow (Low, Wait, 0, Low) \rightarrow (Low, Wait, 0, Low)

E: Search [Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

Expert	$\pi_{\text{expert}}(\text{State})$
High	Search
Low	Recharge

Agent	$\pi_{\theta(n+1)}(\text{State})$
High	Search
Low	Wait

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}

- A(low) = {search, wait, recharge}

Variations of DAGGER (Few Ideas!)

Behavioral Cloning with Data augmentation Exp-2:

Behave (Rollout) using π_{θ_n}

τ : (High, Search, +1, Low) \rightarrow (Low, Wait, 0.5, Low) \rightarrow (Low, Recharge, 0, High)

E: Search [Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)

Expert	$\pi_{\text{expert}}(\text{State})$
High	Search
Low	Recharge

Agent	$\pi_{\theta(n+1)}(\text{State})$
High	Search
Low	Wait

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Recycling Robot



Objective : Policy to search for cans efficiently based on current charge level of the battery

Eg.,

State = function (Charge Levels: High, Low)

Action = function (Charge level based reactions)

- A(high) = {search, wait}

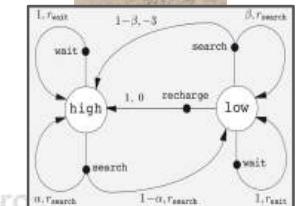
- A(low) = {search, wait, recharge}

Observation:

What is the π_{expert} intent? What if the agent's works in dynamic environment?

Solution : Learn Reward functions using Feature Expectation by parameterizing the $R(s,a)$

[Source Reference: Example inspired from Recycle Robot example in Chapter 3 of Text Book-1](#)



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL

- What if we don't have an online demonstrator?
 - We only have access to an offline set of demonstrated trajectories
- Behavioral cloning is not robust
 - Suffers from overfitting
 - We know what to do in observed states but can't generalize well to other states
- How can we learn to mimic the demonstrator in a general why?
 - Learn the demonstrator's objective (reward) function
 - Apply RL

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL



$$\begin{aligned} \mathbf{a}_1, \dots, \mathbf{a}_T &= \arg \max_{\mathbf{a}_1, \dots, \mathbf{a}_T} \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \\ \mathbf{s}_{t+1} &= f(\mathbf{s}_t, \mathbf{a}_t) \end{aligned}$$

optimize this to explain the data

$$\begin{aligned} \pi &= \arg \max_{\pi} E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t), \mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)} [r(\mathbf{s}_t, \mathbf{a}_t)] \\ \mathbf{a}_t &\sim \pi(\mathbf{a}_t | \mathbf{s}_t) \end{aligned}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL

- What if we don't have an online demonstrator?
 - We only have access to an offline set of demonstrated trajectories
- Behavioral cloning is not robust
 - Suffers from overfitting
 - We know what to do in observed states but can't generalize well to other states
- How can we learn to mimic the demonstrator in a general why?
 - Learn the demonstrator's objective (reward) function
 - Apply RL

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL



$$\begin{aligned} \mathbf{a}_1, \dots, \mathbf{a}_T &= \arg \max_{\mathbf{a}_1, \dots, \mathbf{a}_T} \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \\ \mathbf{s}_{t+1} &= f(\mathbf{s}_t, \mathbf{a}_t) \end{aligned}$$

optimize this to explain the data

$$\begin{aligned} \pi &= \arg \max_{\pi} E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t), \mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)} [r(\mathbf{s}_t, \mathbf{a}_t)] \\ \mathbf{a}_t &\sim \pi(\mathbf{a}_t | \mathbf{s}_t) \end{aligned}$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Work Integrated Learning Programmes

Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Feature Representation:

F1 : 1 if action = Search

F2 : 1 if battery status = Low

F3 : 1 if action = Recharge

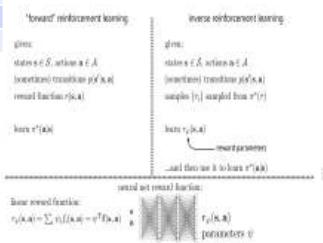
Imitation Learning Exp-3:

$\tau_1: (H,S) \rightarrow (L,R) \rightarrow (H,S)$

$\tau_2: (H,S) \rightarrow (L,R) \rightarrow (H,S)$



	F1	F2	F3
(H,S)	1	0	0
(H,W)	0	0	0
(H,R)	0	0	1
(L,S)	1	1	0
(L,W)	0	1	0
(L,R)	0	1	1



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Feature Representation:

F1 : 1 if action = Search

F2 : 1 if battery status = Low

F3 : 1 if action = Recharge

Imitation Learning Exp-3: Standard Apprenticeship algorithm

$\tau_1: (H,S) \rightarrow (L,R) \rightarrow (H,S) \quad [1,0,0] \rightarrow [0,1,1] \rightarrow [1,0,0]$

$\tau_2: (H,S) \rightarrow (L,R) \quad [1,0,0] \rightarrow [0,1,1]$

Note: For simplicity only two demo is shown . But in reality large dataset of demo is used to compute μ_{π^*}

	F1	F2	F3
(H,S)	1	0	0
(H,W)	0	0	0
(H,R)	0	0	1
(L,S)	1	1	0
(L,W)	0	1	0
(L,R)	0	1	1

Algorithm 2: Apprenticeship Learning
 Inputs: π^*, τ_1, τ_2
 Results:
 Initialize policy μ
 $\text{for } t = 1 \text{ to } T \text{ do}$
 Compute $v(\pi_{t-1})$ (or approximate via Monte Carlo)
 Solve problem (10.5) with policies $\{\pi_0, \dots, \pi_{t-1}\}$ to compute w_t and δ_t
 $(w_t, b_t) = \arg \max_{(w,b)}$
 s.t. $w^\top \mu(\pi^*) \geq w^\top \mu(\pi) + 1, \quad \forall \pi \in \{\pi_0, \dots, \pi_{t-1}\},$
 $\|w\|_2 \leq 1$
 $(w_t)^\top \mu(\pi^*) \geq (w_t)^\top \mu(\pi) + 1, \quad \forall \pi \in \{\pi_0, \dots, \pi_{t-1}\},$
 $\|w_t\|_2 \leq 1$
 end
 if $t \leq T_{\text{init}}$
 $\mu \leftarrow \text{best feature matching policy from } \{\pi_0, \dots, \pi_{t-1}\}$
 else
 $\mu \leftarrow \mu$
 end
 Use RL to find an optimal policy π^* for reward function defined by w_t

$$\mu_{\pi^*} = [0.6, 0.4, 0.4]$$

$$\epsilon = 0.1$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Feature Representation:

F1 : 1 if action = Search

F2 : 1 if battery status = Low

F3 : 1 if action = Recharge

	F1	F2	F3
(H,S)	1	0	0
(H,W)	0	0	0
(H,R)	0	0	1
(L,S)	1	1	0
(L,W)	0	1	0
(L,R)	0	1	1



Imitation Learning Exp-3: Standard Apprenticeship algorithm

$\tau_1: (H,W) \rightarrow (L,S) \quad [0,0,0] \rightarrow [1,1,0]$

$\tau_2: (H,R) \rightarrow (L,W) \quad [0,0,1] \rightarrow [0,1,0]$

$\tau_3: (H,S) \rightarrow (L,R) \quad [1,0,0] \rightarrow [0,1,1]$

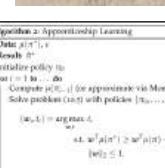
$$\mu_{\pi^*} = [0.6, 0.4, 0.4]$$

$$\mu_{\pi^0} = [0.33, 0.5, 0.33]$$

$$\theta = \frac{(\mu_{\pi^*} - \mu_{\pi^0})}{\text{Norm}((\mu_{\pi^*} - \mu_{\pi^0}))} = \frac{[0.267, -0.1, 0.067]}{0.292}$$

$$\epsilon = 0.1 < t = 0.292$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Feature Representation:

F1 : 1 if action = Search

F2 : 1 if battery status = Low

F3 : 1 if action = Recharge

Imitation Learning Exp-3: Standard Apprenticeship algorithm

$\tau_1: (H,W) \rightarrow (L,S) \quad [0,0,0] \rightarrow [1,1,0]$

$\tau_2: (H,R) \rightarrow (L,W) \quad [0,0,1] \rightarrow [0,1,0]$

$\tau_3: (H,S) \rightarrow (L,R) \quad [1,0,0] \rightarrow [0,1,1]$

$$\text{Reward} = \theta^\top f = [0.914, -0.342, 0.229]^\top * [1,0,0] = \textcolor{red}{-0.914}$$

Algorithm 2: Apprenticeship Learning
 Inputs: π^*, τ_1, τ_2
 Results:
 Initialize policy μ
 $\text{for } t = 1 \text{ to } T \text{ do}$
 Compute $v(\pi_{t-1})$ (or approximate via Monte Carlo)
 Solve problem (10.5) with policies $\{\pi_0, \dots, \pi_{t-1}\}$ to compute w_t and δ_t
 $(w_t, b_t) = \arg \max_{(w,b)}$
 s.t. $w^\top \mu(\pi^*) \geq w^\top \mu(\pi) + 1, \quad \forall \pi \in \{\pi_0, \dots, \pi_{t-1}\},$
 $\|w\|_2 \leq 1$
 $(w_t)^\top \mu(\pi^*) \geq (w_t)^\top \mu(\pi) + 1, \quad \forall \pi \in \{\pi_0, \dots, \pi_{t-1}\},$
 $\|w_t\|_2 \leq 1$
 end
 if $t \leq T_{\text{init}}$
 $\mu \leftarrow \text{best feature matching policy from } \{\pi_0, \dots, \pi_{t-1}\}$
 else
 $\mu \leftarrow \mu$
 end
 Use RL to find an optimal policy π^* for reward function defined by w_t

$$\mu_{\pi^*} = [0.6, 0.4, 0.4]$$

$$\mu_{\pi^0} = [0.33, 0.5, 0.33]$$

$$\theta = \frac{(\mu_{\pi^*} - \mu_{\pi^0})}{\text{Norm}((\mu_{\pi^*} - \mu_{\pi^0}))} = \frac{[0.267, -0.1, 0.067]}{0.292}$$

$$= [0.914, -0.342, 0.229]$$

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Agent	$\pi_{\theta(n+1)}(\text{State})$ Greedy
High	Search
Low	Search

Feature Representation:

- F1 : 1 if action = Search
- F2 : 1 if battery status = Low
- F3 : 1 if action = Recharge

Imitation Learning Exp-3: Standard Apprenticeship algorithm

$\tau_1: (H,W) \rightarrow (L,S)$ $[0,0,0] \rightarrow [1,1,0]$

$\tau_2: (H,R) \rightarrow (L,W)$ $[0,0,1] \rightarrow [0,1,0]$

$\tau_3: (H,S) \rightarrow (L,R)$ $[1,0,0] \rightarrow [0,1,1]$

Reward = $\theta_0^T f = [0.914, -0.342, 0.229]^T \cdot [1,0,0] = \text{~0.914}$

Algorithm 2: Apprenticeship Learning

```

    Inputs:  $\pi^*, \pi_{\theta}$ ,  $\mathcal{B}$ 
    Initialize policy  $\pi$ 
    for  $t = 1$  to ... do
        Compute  $\pi^*(\cdot|s_t)$  (or approximate via Monte Carlo)
        Solve problem  $(\pi, \pi^*)$  with policies  $\{\pi_0, \dots, \pi_{t-1}\}$  to compute  $\pi_t$  and  $\theta_t$ 
         $(\pi_t, \theta_t) = \arg\max_{(\pi, \theta)}$ 
            s.t.  $\pi^T \mu(\pi^*) \geq \pi^T \mu(\pi) + 1$ ,  $\forall \pi \in \{\pi_0, \dots, \pi_{t-1}\}$ 
             $\|\theta\|_2 \leq L$ 
    end
    if  $t \leq t_{\text{stop}}$ 
         $\pi^* \leftarrow \text{best function matching policy from } \{\pi_0, \dots, \pi_{t-1}\}$ 
    return  $\pi^*$ 
    Use RL to find an optimal policy  $\pi^*$  for reward function defined by  $\theta_t$ 

```

Algorithm 2: Apprenticeship Learning

```

    Inputs:  $\pi^*, \pi_{\theta}$ ,  $\mathcal{B}$ 
    Initialize policy  $\pi$ 
    for  $t = 1$  to ... do
        Compute  $\pi^*(\cdot|s_t)$  (or approximate via Monte Carlo)
        Solve problem  $(\pi, \pi^*)$  with policies  $\{\pi_0, \dots, \pi_{t-1}\}$  to compute  $\pi_t$  and  $\theta_t$ 
         $(\pi_t, \theta_t) = \arg\max_{(\pi, \theta)}$ 
            s.t.  $\pi^T \mu(\pi^*) \geq \pi^T \mu(\pi) + 1$ ,  $\forall \pi \in \{\pi_0, \dots, \pi_{t-1}\}$ 
             $\|\theta\|_2 \leq L$ 
    end
    if  $t \leq t_{\text{stop}}$ 
         $\pi^* \leftarrow \text{best function matching policy from } \{\pi_0, \dots, \pi_{t-1}\}$ 
    return  $\pi^*$ 
    Use RL to find an optimal policy  $\pi^*$  for reward function defined by  $\theta_t$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Agent	$\pi_{\theta(n+1)}(\text{State})$ Greedy
High	Search
Low	Search

Feature Representation:

- F1 : 1 if action = Search
- F2 : 1 if battery status = Low
- F3 : 1 if action = Recharge

Imitation Learning Exp-3: Standard Apprenticeship algorithm

$\tau_1: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\tau_2: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\tau_3: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\mu_{\pi^*} = [0.6, 0.4, 0.4]$
 $\mu_{\pi_1} = [0.33, 0.5, 0.33]$
 $\theta = \frac{(\mu_{\pi^*} - \mu_{\pi_1})}{\text{Norm}((\mu_{\pi^*} - \mu_{\pi_1}))} = \frac{[0.267, -0.1, 0.067]}{0.292}$
 $= [0.914, -0.342, 0.229]$
 $\epsilon = 0.1$ $t = 1$ $\rightarrow [0.292]$ continue

Reward = $\theta_0^T f = [-0.697, -0.174, 0.697]^T \cdot [1,0,0] = \text{~-0.697}$

Algorithm 2: Apprenticeship Learning

```

    Inputs:  $\pi^*, \pi_{\theta}$ ,  $\mathcal{B}$ 
    Initialize policy  $\pi$ 
    for  $t = 1$  to ... do
        Compute  $\pi^*(\cdot|s_t)$  (or approximate via Monte Carlo)
        Solve problem  $(\pi, \pi^*)$  with policies  $\{\pi_0, \dots, \pi_{t-1}\}$  to compute  $\pi_t$  and  $\theta_t$ 
         $(\pi_t, \theta_t) = \arg\max_{(\pi, \theta)}$ 
            s.t.  $\pi^T \mu(\pi^*) \geq \pi^T \mu(\pi) + 1$ ,  $\forall \pi \in \{\pi_0, \dots, \pi_{t-1}\}$ 
             $\|\theta\|_2 \leq L$ 
    end
    if  $t \leq t_{\text{stop}}$ 
         $\pi^* \leftarrow \text{best function matching policy from } \{\pi_0, \dots, \pi_{t-1}\}$ 
    return  $\pi^*$ 
    Use RL to find an optimal policy  $\pi^*$  for reward function defined by  $\theta_t$ 

```

Algorithm 2: Apprenticeship Learning

```

    Inputs:  $\pi^*, \pi_{\theta}$ ,  $\mathcal{B}$ 
    Initialize policy  $\pi$ 
    for  $t = 1$  to ... do
        Compute  $\pi^*(\cdot|s_t)$  (or approximate via Monte Carlo)
        Solve problem  $(\pi, \pi^*)$  with policies  $\{\pi_0, \dots, \pi_{t-1}\}$  to compute  $\pi_t$  and  $\theta_t$ 
         $(\pi_t, \theta_t) = \arg\max_{(\pi, \theta)}$ 
            s.t.  $\pi^T \mu(\pi^*) \geq \pi^T \mu(\pi) + 1$ ,  $\forall \pi \in \{\pi_0, \dots, \pi_{t-1}\}$ 
             $\|\theta\|_2 \leq L$ 
    end
    if  $t \leq t_{\text{stop}}$ 
         $\pi^* \leftarrow \text{best function matching policy from } \{\pi_0, \dots, \pi_{t-1}\}$ 
    return  $\pi^*$ 
    Use RL to find an optimal policy  $\pi^*$  for reward function defined by  $\theta_t$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Agent	$\pi_{\theta(n+1)}(\text{State})$ Greedy
High	Search
Low	Search

Feature Representation:

- F1 : 1 if action = Search
- F2 : 1 if battery status = Low
- F3 : 1 if action = Recharge

Imitation Learning Exp-3: Standard Apprenticeship algorithm

$\tau_1: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\tau_2: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\tau_3: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\mu_{\pi^*} = [0.6, 0.4, 0.4]$
 $\mu_{\pi_1} = [1, 0.5, 0]$
 $\theta = \frac{(\mu_{\pi^*} - \mu_{\pi_1})}{\text{Norm}((\mu_{\pi^*} - \mu_{\pi_1}))} = \frac{[-0.4, -0.1, 0.4]}{0.574}$
 $= [-0.697, -0.174, 0.697]$
 $\epsilon = 0.1$ $t = 1$ $\rightarrow [0.574]$ continue

Algorithm 2: Apprenticeship Learning

```

    Inputs:  $\pi^*, \pi_{\theta}$ ,  $\mathcal{B}$ 
    Initialize policy  $\pi$ 
    for  $t = 1$  to ... do
        Compute  $\pi^*(\cdot|s_t)$  (or approximate via Monte Carlo)
        Solve problem  $(\pi, \pi^*)$  with policies  $\{\pi_0, \dots, \pi_{t-1}\}$  to compute  $\pi_t$  and  $\theta_t$ 
         $(\pi_t, \theta_t) = \arg\max_{(\pi, \theta)}$ 
            s.t.  $\pi^T \mu(\pi^*) \geq \pi^T \mu(\pi) + 1$ ,  $\forall \pi \in \{\pi_0, \dots, \pi_{t-1}\}$ 
             $\|\theta\|_2 \leq L$ 
    end
    if  $t \leq t_{\text{stop}}$ 
         $\pi^* \leftarrow \text{best function matching policy from } \{\pi_0, \dots, \pi_{t-1}\}$ 
    return  $\pi^*$ 
    Use RL to find an optimal policy  $\pi^*$  for reward function defined by  $\theta_t$ 

```

Algorithm 2: Apprenticeship Learning

```

    Inputs:  $\pi^*, \pi_{\theta}$ ,  $\mathcal{B}$ 
    Initialize policy  $\pi$ 
    for  $t = 1$  to ... do
        Compute  $\pi^*(\cdot|s_t)$  (or approximate via Monte Carlo)
        Solve problem  $(\pi, \pi^*)$  with policies  $\{\pi_0, \dots, \pi_{t-1}\}$  to compute  $\pi_t$  and  $\theta_t$ 
         $(\pi_t, \theta_t) = \arg\max_{(\pi, \theta)}$ 
            s.t.  $\pi^T \mu(\pi^*) \geq \pi^T \mu(\pi) + 1$ ,  $\forall \pi \in \{\pi_0, \dots, \pi_{t-1}\}$ 
             $\|\theta\|_2 \leq L$ 
    end
    if  $t \leq t_{\text{stop}}$ 
         $\pi^* \leftarrow \text{best function matching policy from } \{\pi_0, \dots, \pi_{t-1}\}$ 
    return  $\pi^*$ 
    Use RL to find an optimal policy  $\pi^*$  for reward function defined by  $\theta_t$ 

```

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Agent	$\pi_{\theta(n+1)}(\text{State})$ Greedy
High	Search
Low	Search

Feature Representation:

- F1 : 1 if action = Search
- F2 : 1 if battery status = Low
- F3 : 1 if action = Recharge

Imitation Learning Exp-3: Standard Apprenticeship algorithm

$\tau_1: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\tau_2: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\tau_3: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\mu_{\pi^*} = [0.6, 0.4, 0.4]$
 $\mu_{\pi_1} = [1, 0.5, 0]$
 $\theta = \frac{(\mu_{\pi^*} - \mu_{\pi_1})}{\text{Norm}((\mu_{\pi^*} - \mu_{\pi_1}))} = \frac{[-0.4, -0.1, 0.4]}{0.574}$
 $= [-0.697, -0.174, 0.697]$
 $\epsilon = 0.1$ $t = 1$ $\rightarrow [0.574]$ continue

But for (H,R) ie., [0,0,1] Reward = 0.697

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL

Objective : Policy to search for cans efficiently based on current charge level of the battery

Feature Representation:

F1 : 1 if action = Search

F2 : 1 if battery status = Low

F3 : 1 if action = Recharge

Imitation Learning Exp-3: Standard Apprenticeship algorithm

$\tau_1: (H,S) \rightarrow (L,S)$

$[1,0,0] \rightarrow [1,1,0]$

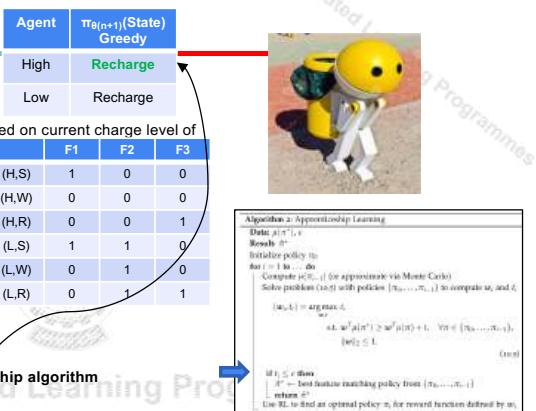
$\tau_2: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$\tau_3: (H,S) \rightarrow (L,S)$ $[1,0,0] \rightarrow [1,1,0]$

$$\text{Reward} = \theta_0^\top f = [-0.697, -0.174, 0.697]^\top \cdot [1,0,0] = \mathbf{-0.697}$$

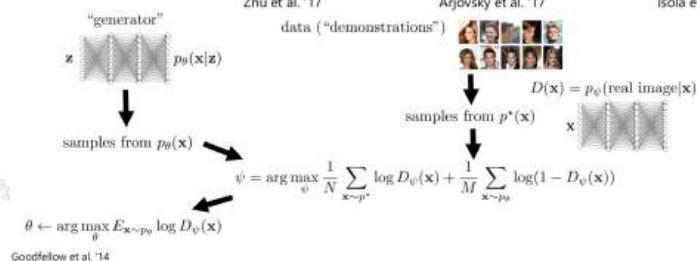
But for (H,R) ie., $[0,0,1]$ Reward = 0.697

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



$$\begin{aligned} \mu_{\pi^*} &= [0.6, 0.4, 0.4] \\ \mu_{\pi^1} &= [1, 0.5, 0] \\ \theta &= \frac{(H_m - L_m)}{\text{Norm}((\mu_{\pi^*} - \mu_{\pi^1}))} = \frac{[-0.4, -0.1, 0.4]}{0.574} \\ &= [-0.697, -0.174, 0.697] \\ \epsilon &= 0.1 \quad t = 0.574 \quad \dots \text{continue} \end{aligned}$$

GAN



BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Inverse RL as GAN

generator/policy
 $\pi_\theta(\tau)$

data/demonstrations
samples from $p^*(\tau)$

$$\nabla_\theta \mathcal{L} \approx \frac{1}{M} \sum_{j=1}^M \nabla_\theta \log \pi_\theta(\tau_j) r_\psi(\tau_j)$$

$$D_\psi(\tau) = \frac{\frac{1}{Z} \exp(r(\tau))}{\prod_t \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) + \frac{1}{Z} \exp(r(\tau))}$$

policy changed to make it *harder* to distinguish from demos

Inverse RL as GAN

initial policy π

samples from $\pi_\theta(\tau)$

$$\nabla_\theta \mathcal{L} \approx \frac{1}{M} \sum_{j=1}^M \nabla_\theta \log \pi_\theta(\tau_j) r_\psi(\tau_j)$$

policy changed to make it *harder* to distinguish from demos

human demonstrations

samples from $\pi^*(\tau)$

$$\nabla_\psi \mathcal{L} \approx \frac{1}{N} \sum_{i=1}^N \nabla_\psi \log \pi_\psi(\tau_i) - \frac{1}{\sum_j w_j} \sum_{j=1}^M w_j \nabla_\psi \log \pi_\psi(\tau_j)$$

demos are made more likely, samples less likely



Deep Reinforcement Learning

Prof. Vimal SP
CSIS

BITS Pilani



BITS Pilani



<AIMLCZG512, Deep Reinforcement Learning> AlphaGo_Zero Research paper



Mastering the Game of Go without Human Knowledge

David Silver*, Julian Schrittwieser*, Karen Simonyan*, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrián Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, Demis Hassabis.

DeepMind, 5 New Street Square, London EC4A 3TW.

*These authors contributed equally to this work.

A long-standing goal of artificial intelligence is an algorithm that learns, *tabula rasa*, superhuman proficiency in challenging domains. Recently, *AlphaGo* became the first program to defeat a world champion in the game of Go. The tree search in *AlphaGo* evaluated positions and selected moves using deep neural networks. These neural networks were trained by supervised learning from human expert moves, and by reinforcement learning from self-play. Here, we introduce an algorithm based solely on reinforcement learning, without human data, guidance, or domain knowledge beyond game rules. *AlphaGo* becomes its own teacher: a neural network is trained to predict *AlphaGo*'s own move selections and also the winner of *AlphaGo*'s games. This neural network improves the strength of tree search, resulting in higher quality move selection and stronger self-play in the next iteration. Starting *tabula rasa*, our new program *AlphaGo Zero* achieved superhuman performance, winning 100-0 against the previously published, champion-defeating *AlphaGo*.

Much progress towards artificial intelligence has been made using supervised learning systems that are trained to replicate the decisions of human experts^{1–4}. However, expert data is often expensive, unreliable, or simply unavailable. Even when reliable data is available it may impose a ceiling on the performance of systems trained in this manner⁵. In contrast, reinforcement learning systems are trained from their own experience, in principle allowing them to exceed human capabilities, and to operate in domains where human expertise is lacking. Recently, there has been rapid progress towards this goal, using deep neural networks trained by reinforcement learning. These systems have outperformed humans in computer games such as Atari^{6–7} and 3D virtual environments^{8–10}. However, the most challenging domains in terms of human intellect – such as the

game of Go, widely viewed as a grand challenge for artificial intelligence¹¹ – require precise and sophisticated lookahead in vast search spaces. Fully general methods have not previously achieved human-level performance in these domains.

AlphaGo was the first program to achieve superhuman performance in Go. The published version¹², which we refer to as *AlphaGo Fan*, defeated the European champion Fan Hui in October 2015. *AlphaGo Fan* utilised two deep neural networks: a policy network that outputs move probabilities, and a value network that outputs a position evaluation. The policy network was trained initially by supervised learning to accurately predict human expert moves, and was subsequently refined by policy-gradient reinforcement learning. The value network was trained to predict the winner of games played by the policy network against itself. Once trained, these networks were combined with a Monte-Carlo Tree Search (MCTS)^{13–15} to provide a lookahead search, using the policy network to narrow down the search to high-probability moves, and using the value network (in conjunction with Monte-Carlo rollouts using a fast rollout policy) to evaluate positions in the tree. A subsequent version, which we refer to as *AlphaGo Lee*, used a similar approach (see Methods), and defeated Lee Sedol, the winner of 18 international titles, in March 2016.

Our program, *AlphaGo Zero*, differs from *AlphaGo Fan* and *AlphaGo Lee*¹² in several important aspects. First and foremost, it is trained solely by self-play reinforcement learning, starting from random play, without any supervision or use of human data. Second, it only uses the black and white stones from the board as input features. Third, it uses a single neural network, rather than separate policy and value networks. Finally, it uses a simpler tree search that relies upon this single neural network to evaluate positions and sample moves, without performing any Monte-Carlo rollouts. To achieve these results, we introduce a new reinforcement learning algorithm that incorporates lookahead search *inside* the training loop, resulting in rapid improvement and precise and stable learning. Further technical differences in the search algorithm, training procedure and network architecture are described in Methods.



1 Reinforcement Learning in AlphaGo Zero

Our new method uses a deep neural network f_θ with parameters θ . This neural network takes as an input the raw board representation s of the position and its history, and outputs both move probabilities and a value, $(\mathbf{p}, v) = f_\theta(s)$. The vector of move probabilities \mathbf{p} represents the probability of selecting each move (including pass), $p_a = Pr(a|s)$. The value v is a scalar evaluation, estimating the probability of the current player winning from position s . This neural network combines the roles of both policy network and value network¹² into a single architecture. The neural network consists of many residual blocks⁴ of convolutional layers^{16,17} with batch normalisation¹⁸ and rectifier non-linearities¹⁹ (see Methods).

The neural network in *AlphaGo Zero* is trained from games of self-play by a novel reinforcement learning algorithm. In each position s_t , an MCTS search is executed, guided by the neural network f_θ . The MCTS search outputs probabilities π_t of playing each move. These search probabilities usually select much stronger moves than the raw move probabilities \mathbf{p} of the neural network $f_\theta(s)$; MCTS may therefore be viewed as a powerful *policy improvement* operator^{20,21}. Self-play with search – using the improved MCTS-based policy to select each move, then using the game winner z as a sample of the value – may be viewed as a powerful *policy evaluation* operator. The main idea of our reinforcement learning algorithm is to use these search operators repeatedly in a policy iteration procedure^{22,23}: the neural network's parameters are updated to make the move probabilities and value $(\mathbf{p}, v) = f_\theta(s)$ more closely match the improved search probabilities and self-play winner (π, z) ; these new parameters are used in the next iteration of self-play to make the search even stronger. Figure 1 illustrates the self-play training pipeline.

The Monte-Carlo tree search uses the neural network f_θ to guide its simulations (see Figure 2). Each edge (s, a) in the search tree stores a prior probability $P(s, a)$, a visit count $N(s, a)$, and an action-value $Q(s, a)$. Each simulation starts from the root state and iteratively selects moves that maximise an upper confidence bound $Q(s, a) + U(s, a)$, where $U(s, a) \propto P(s, a)/(1 + N(s, a))$ ^{12,24}, until a leaf node s' is encountered. This leaf position is expanded and evaluated just

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

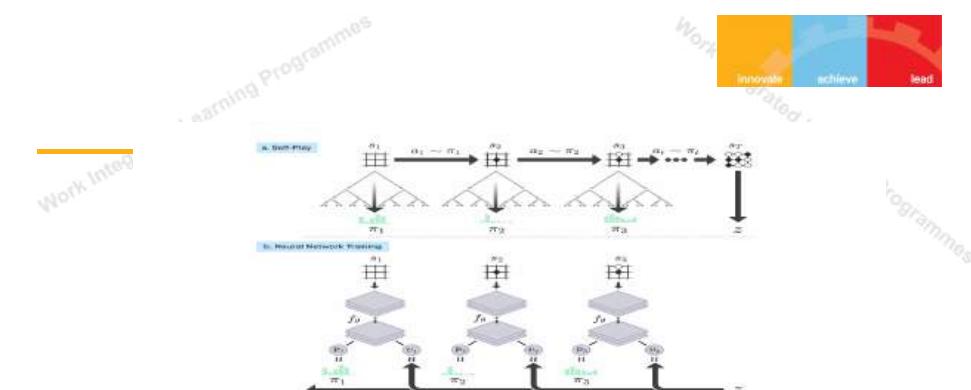


Figure 1: Self-play reinforcement learning in *AlphaGo Zero*. **a** The program plays a game s_1, \dots, s_T against itself. In each position s_t , a Monte-Carlo tree search (MCTS) α_{s_t} is executed (see Figure 2) using the latest neural network f_θ . Moves are selected according to the search probabilities computed by the MCTS, $a_t \sim \pi_t$. The terminal position s_T is scored according to the rules of the game to compute the game winner z . **b** Neural network training in *AlphaGo Zero*. The neural network takes the raw board position s_t as its input, passes it through many convolutional layers with parameters θ , and outputs both a vector \mathbf{p}_t , representing a probability distribution over moves, and a scalar value v_t , representing the probability of the current player winning in position s_t . The neural network parameters θ are updated so as to maximise the similarity of the policy vector \mathbf{p}_t to the search probabilities π_t , and to minimise the error between the predicted winner v_t and the game winner z (see Equation 1). The new parameters are used in the next iteration of self-play **a**.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

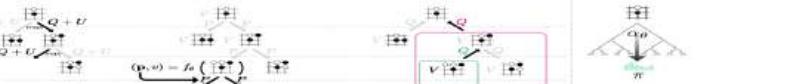


Figure 2: Monte-Carlo tree search in *AlphaGo Zero*. **a** Each simulation traverses the tree by selecting the edge with maximum action-value Q , plus an upper confidence bound U that depends on a stored prior probability P and visit count N for that edge (which is incremented once traversed). **b** The leaf node is expanded and the associated position s is evaluated by the neural network ($P(s, \cdot), V(s) = f_\theta(s)$); the vector of P values are stored in the outgoing edges from s . **c** Action-values Q are updated to track the mean of all evaluations V in the subtree below that action. **d** Once the search is complete, search probabilities π are returned, proportional to $N^{1/\tau}$, where N is the visit count of each move from the root state and τ is a parameter controlling temperature.

once by the network to generate both prior probabilities and evaluation, $(P(s', \cdot), V(s')) = f_\theta(s')$. Each edge (s, a) traversed in the simulation is updated to increment its visit count $N(s, a)$, and to update its action-value to the mean evaluation over these simulations, $Q(s, a) = 1/N(s, a) \sum_{s' \mid s, a \rightarrow s'} V(s')$, where $s, a \rightarrow s'$ indicates that a simulation eventually reached s' after taking move a from position s .

MCTS may be viewed as a self-play algorithm that, given neural network parameters θ and a root position s , computes a vector of search probabilities recommending moves to play, $\pi = \pi_\theta(s)$, proportional to the exponentiated visit count for each move, $\pi_a \propto N(s, a)^{1/\tau}$, where τ is a temperature parameter.

The neural network is trained by a self-play reinforcement learning algorithm that uses MCTS to play each move. First, the neural network is initialised to random weights θ_0 . At each subsequent iteration $i \geq 1$, games of self-play are generated (Figure 1a). At each time-step t , an MCTS search $\pi_t = \alpha_{t-1}(s_t)$ is executed using the previous iteration of neural network $f_{\theta_{i-1}}$, and a move is played by sampling the search probabilities π_t . A game terminates at step T when

both players pass, when the search value drops below a resignation threshold, or when the game exceeds a maximum length; the game is then scored to give a final reward of $r_T \in \{-1, +1\}$ (see Methods for details). The data for each time-step t is stored as (s_t, π_t, z_t) where $z_t = \pm r_T$ is the game winner from the perspective of the current player at step t . In parallel (Figure 1b), new network parameters θ_i are trained from data (s, π, z) sampled uniformly among all time-steps of the last iteration(s) of self-play. The neural network $(\mathbf{p}, v) = f_{\theta_i}(s)$ is adjusted to minimise the error between the predicted value v and the self-play winner z , and to maximise the similarity of the neural network move probabilities \mathbf{p} to the search probabilities π . Specifically, the parameters θ are adjusted by gradient descent on a loss function l that sums over mean-squared error and cross-entropy losses respectively,

$$(\mathbf{p}, v) = f_\theta(s), \quad l = (z - v)^2 - \pi^\top \log \mathbf{p} + c \|\theta\|^2 \quad (1)$$

where c is a parameter controlling the level of L2 weight regularisation (to prevent overfitting).

2 Empirical Analysis of AlphaGo Zero Training

We applied our reinforcement learning pipeline to train our program *AlphaGo Zero*. Training started from completely random behaviour and continued without human intervention for approximately 3 days.

Over the course of training, 4.9 million games of self-play were generated, using 1,600 simulations for each MCTS, which corresponds to approximately 0.4s thinking time per move. Parameters were updated from 700,000 mini-batches of 2,048 positions. The neural network contained 20 residual blocks (see Methods for further details).

Figure 3a shows the performance of *AlphaGo Zero* during self-play reinforcement learning, as a function of training time, on an Elo scale²⁵. Learning progressed smoothly throughout training, and did not suffer from the oscillations or catastrophic forgetting suggested in prior literature

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Figure 3: Empirical evaluation of AlphaGo Zero. **a**, Performance of self-play reinforcement learning. The plot shows the performance of each MCTS player α_i , from each iteration i of reinforcement learning in AlphaGo Zero. Elo ratings were computed from evaluation games between different players, using 0.4 seconds of thinking time per move (see Methods). For comparison, a similar player trained by supervised learning from human data, using the KGS data-set, is also shown. **b**, Prediction accuracy on human professional moves. The plot shows the accuracy of the neural network f_{θ_i} at each iteration of self-play i , in predicting human professional moves from the GoKifu data-set. The accuracy measures the percentage of positions in which the neural network assigns the highest probability to the human move. The accuracy of a neural network trained by supervised learning is also shown. **c**, Mean-squared error (MSE) on human professional game outcomes. The plot shows the MSE of the neural network f_{θ_i} at each iteration of self-play i , in predicting the outcome of human professional games from the GoKifu data-set. The MSE is between the actual outcome $z \in \{-1, +1\}$ and the neural network value v_i , scaled by a factor of $\frac{1}{3}$ to the range [0, 1]. The MSE of a neural network trained by supervised learning is also shown.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Surprisingly, AlphaGo Zero outperformed AlphaGo Lee after just 36 hours; for comparison, AlphaGo Lee was trained over several months. After 72 hours, we evaluated AlphaGo Zero against the exact version of AlphaGo Lee that defeated Lee Sedol, under the 2 hour time controls and match conditions as were used in the man-machine match in Seoul (see Methods). AlphaGo Zero used a single machine with 4 Tensor Processing Units (TPUs)²⁹, while AlphaGo Lee was distributed over many machines and used 48 TPUs. AlphaGo Zero defeated AlphaGo Lee by 100 games to 0 (see Extended Data Figure 5 and Supplementary Information).

To assess the merits of self-play reinforcement learning, compared to learning from human data, we trained a second neural network (using the same architecture) to predict expert moves in the KGS data-set; this achieved state-of-the-art prediction accuracy compared to prior work^{13,30–33} (see Extended Data Table 1 and 2 respectively). Supervised learning achieved better initial performance, and was better at predicting the outcome of human professional games (Figure 3). Notably, although supervised learning achieved higher move prediction accuracy, the self-trained player performed much better overall, defeating the human-trained player within the first 24 hours of training. This suggests that AlphaGo Zero may be learning a strategy that is qualitatively different to human play.

To separate the contributions of architecture and algorithm, we compared the performance of the neural network architecture in AlphaGo Zero with the previous neural network architecture used in AlphaGo Lee (see Figure 4). Four neural networks were created, using either separate policy and value networks, as in AlphaGo Lee, or combined policy and value networks, as in AlphaGo Zero; and using either the convolutional network architecture from AlphaGo Lee or the residual network architecture from AlphaGo Zero. Each network was trained to minimise the same loss function (Equation 1) using a fixed data-set of self-play games generated by AlphaGo Zero after 72 hours of self-play training. Using a residual network was more accurate, achieved lower error, and improved performance in AlphaGo by over 600 Elo. Combining policy and value together into a single network slightly reduced the move prediction accuracy, but reduced the value error and boosted playing performance in AlphaGo by around another 600 Elo. This is partly due to

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Figure 4: Comparison of neural network architectures in AlphaGo Zero and AlphaGo Lee. Comparison of neural network architectures using either separate ("sep") or combined policy and value networks ("dual"), and using either convolutions ("conv") or residual networks ("res"). The combinations "dual-res" and "sep-conv" correspond to the neural network architectures used in AlphaGo Zero and AlphaGo Lee respectively. Each network was trained on a fixed data-set generated by a previous run of AlphaGo Zero. **a**, Each trained network was combined with AlphaGo Zero's search to obtain a different player. Elo ratings were computed from evaluation games between these different players, using 5 seconds of thinking time per move. **b**, Prediction accuracy on human professional moves (from the GoKifu data-set) for each network architecture. **c**, Mean-squared error on human professional game outcomes (from the GoKifu data-set) for each network architecture.

improved computational efficiency, but more importantly the dual objective regularises the network to a common representation that supports multiple use cases.

3 Knowledge Learned by AlphaGo Zero

AlphaGo Zero discovered a remarkable level of Go knowledge during its self-play training process. This included fundamental elements of human Go knowledge, and also non-standard strategies beyond the scope of traditional Go knowledge.

Figure 5 shows a timeline indicating when professional *joseki* (corner sequences) were discovered (Figure 5a, Extended Data Figure 1); ultimately AlphaGo Zero preferred new *joseki* variants that were previously unknown (Figure 5b, Extended Data Figure 2). Figure 5c and the Supplementary Information show several fast self-play games played at different stages of training. Tournament-like games played at regular intervals throughout training are shown in Extended Data Figure 3 and Supplementary Information. AlphaGo Zero rapidly progressed from entirely random moves towards a sophisticated understanding of Go concepts including *fuseki* (opening), *tesuji* (tactics), life-and-death, *ko* (repeated board situations), *yose* (endgame), capturing races, *sente* (initiative), shape, influence and territory, all discovered from first principles. Surprisingly, *shicho* ("ladder") capture sequences that may span the whole board – one of the first elements of Go knowledge learned by humans – were only understood by AlphaGo Zero much later in training.

4 Final Performance of AlphaGo Zero

We subsequently applied our reinforcement learning pipeline to a second instance of AlphaGo Zero using a larger neural network and over a longer duration. Training again started from completely random behaviour and continued for approximately 40 days.

Over the course of training, 29 million games of self-play were generated. Parameters were updated from 3.1 million mini-batches of 2,048 positions each. The neural network contained

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

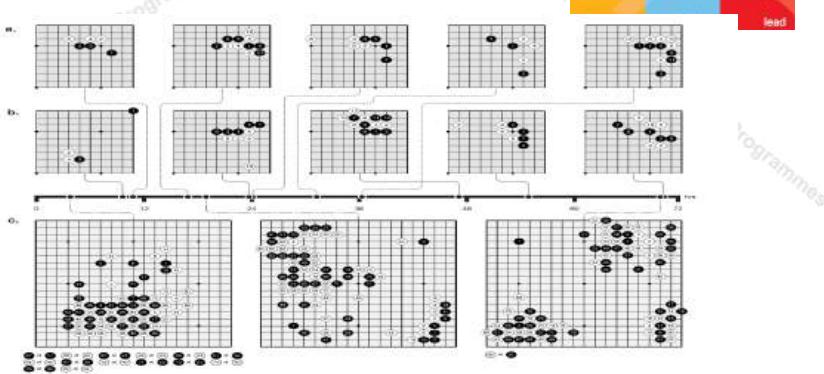


Figure 5: Go knowledge learned by AlphaGo Zero. **a** Five human *joseki* (common corner sequences) discovered during *AlphaGo Zero* training. The associated timestamps indicate the first time each sequence occurred (taking account of rotation and reflection) during self-play training. Extended Data Figure 1 provides the frequency of occurrence over training for each sequence. **b** Five *joseki* favoured at different stages of self-play training. Each displayed corner sequence was played with the greatest frequency, among all corner sequences, during an iteration of self-play training. The timestamp of that iteration is indicated on the timeline. At 10 hours a weak corner move was preferred. At 47 hours the 3-3 invasion was most frequently played. This *joseki* is also common in human professional play; however *AlphaGo Zero* later discovered and preferred a new variation. Extended Data Figure 2 provides the frequency of occurrence over time for all five sequences and the new variation. **c** The first 80 moves of three self-play games that were played at different stages of training, using 1,600 simulations (around 0.4s per search). At 3 hours, the game focuses greedily on capturing stones, much like a human beginner. At 19 hours, the game exhibits the fundamentals of life-and-death, influence and territory. At 70 hours, the game is beautifully balanced, involving multiple battles and a complicated *ko* fight, eventually resolving into a half-point win for white. See Supplementary Information for the full games.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

40 residual blocks. The learning curve is shown in Figure 6a. Games played at regular intervals throughout training are shown in Extended Data Figure 4 and Supplementary Information.

We evaluated the fully trained *AlphaGo Zero* using an internal tournament against *AlphaGo Fan*, *AlphaGo Lee*, and several previous Go programs. We also played games against the strongest existing program, *AlphaGo Master* – a program based on the algorithm and architecture presented in this paper but utilising human data and features (see Methods) – which defeated the strongest human professional players 60–0 in online games²⁴ in January 2017. In our evaluation, all programs were allowed 5 seconds of thinking time per move. *AlphaGo Zero* and *AlphaGo Master* each played on a single machine with 4 TPUs; *AlphaGo Fan* and *AlphaGo Lee* were distributed over 176 GPUs and 48 TPUs respectively. We also included a player based solely on the raw neural network of *AlphaGo Zero*; this player simply selected the move with maximum probability.

Figure 6b shows the performance of each program on an Elo scale, without using any lookahead, achieved an Elo rating of 3,055. *AlphaGo Zero* achieved a rating of 5,185, compared to 4,858 for *AlphaGo Master*, 3,739 for *AlphaGo Lee* and 3,144 for *AlphaGo Fan*.

Finally, we evaluated *AlphaGo Zero* head to head against *AlphaGo Master* in a 100 game match with 2 hour time controls. *AlphaGo Zero* won by 89 games to 11 (see Extended Data Figure 6) and Supplementary Information.

5 Conclusion

Our results comprehensively demonstrate that a pure reinforcement learning approach is fully feasible, even in the most challenging of domains: it is possible to train to superhuman level, without human examples or guidance, given no knowledge of the domain beyond basic rules. Furthermore, a pure reinforcement learning approach requires just a few more hours to train, and achieves much better asymptotic performance, compared to training on human expert data. Using this ap-

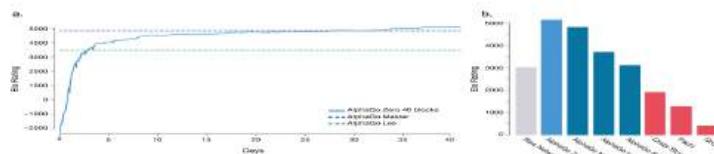


Figure 6: Performance of AlphaGo Zero. **a** Learning curve for *AlphaGo Zero* using larger 40 block residual network over 40 days. The plot shows the performance of each player α_i , from each iteration i of our reinforcement learning algorithm. Elo ratings were computed from evaluation games between different players, using 0.4 seconds per search (see Methods). **b** Final performance of *AlphaGo Zero*. *AlphaGo Zero* was trained for 40 days using a 40 residual block neural network. The plot shows the results of a tournament between: *AlphaGo Zero*, *AlphaGo Master* (defeated top human professionals 60–0 in online games), *AlphaGo Lee* (defeated Lee Sedol), *AlphaGo Fan* (defeated Fan Hui), as well as previous Go programs *Crazy Stone*, *Pachi* and *GnuGo*. Each program was given 5 seconds of thinking time per move. *AlphaGo Zero* and *AlphaGo Master* played on a single machine on the Google Cloud; *AlphaGo Fan* and *AlphaGo Lee* were distributed over many machines. The raw neural network from *AlphaGo Zero* is also included, which directly selects the move a with maximum probability p_a , without using MCTS. Programs were evaluated on an Elo scale²⁵; a 200 point gap corresponds to a 75% probability of winning.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

proach, *AlphaGo Zero* defeated the strongest previous versions of *AlphaGo*, which were trained from human data using handcrafted features, by a large margin.

Humankind has accumulated Go knowledge from millions of games played over thousands of years, collectively distilled into patterns, proverbs and books. In the space of a few days, starting *tabula rasa*, *AlphaGo Zero* was able to rediscover much of this Go knowledge, as well as novel strategies that provide new insights into the oldest of games.

Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Methods

Reinforcement learning Policy iteration^{20,21} is a classic algorithm that generates a sequence of improving policies, by alternating between *policy evaluation* – estimating the value function of the current policy – and *policy improvement* – using the current value function to generate a better policy. A simple approach to policy evaluation is to estimate the value function from the outcomes of sampled trajectories^{22,23}. A simple approach to policy improvement is to select actions greedily with respect to the value function²⁰. In large state spaces, approximations are necessary to evaluate each policy and to represent its improvement^{22,23}.

Classification-based reinforcement learning²⁷ improves the policy using a simple Monte-Carlo search. Many rollouts are executed for each action; the action with the maximum mean value provides a positive training example, while all other actions provide negative training examples; a policy is then trained to classify actions as positive or negative, and used in subsequent rollouts. This may be viewed as a precursor to the policy component of *AlphaGo Zero*'s training algorithm when $\tau \rightarrow 0$.

A more recent instantiation, classification-based modified policy iteration (CBMPI), also performs policy evaluation by regressing a value function towards truncated rollout values, similar to the value component of *AlphaGo Zero*; this achieved state-of-the-art results in the game of Tetris³⁸. However, this prior work was limited to simple rollouts and linear function approximation using handcrafted features.

The *AlphaGo Zero* self-play algorithm can similarly be understood as an approximate policy iteration scheme in which MCTS is used for both policy improvement and policy evaluation. Policy improvement starts with a neural network policy, executes an MCTS based on that policy's recommendations, and then projects the (much stronger) search policy back into the function space of the neural network. Policy evaluation is applied to the (much stronger) search policy: the outcomes of self-play games are also projected back into the function space of the neural network. These projection steps are achieved by training the neural network parameters to match the search

probabilities and self-play game outcome respectively.

Guo et al.⁷ also project the output of MCTS into a neural network, either by regressing a value network towards the search value, or by classifying the action selected by MCTS. This approach was used to train a neural network for playing Atari games; however, the MCTS was fixed — there was no policy iteration — and did not make any use of the trained networks.

Self-play reinforcement learning in games Our approach is most directly applicable to zero-sum games of perfect information. We follow the formalism of alternating Markov games described in previous work¹², noting that algorithms based on value or policy iteration extend naturally to this setting²⁰.

Self-play reinforcement learning has previously been applied to the game of Go. *NeuroGo*^{40,41} used a neural network to represent a value function, using a sophisticated architecture based on Go knowledge regarding connectivity, territory and eyes. This neural network was trained by temporal-difference learning⁴² to predict territory in games of self-play, building on prior work⁴³. A related approach, *RLGO*⁴⁴, represented the value function instead by a linear combination of features, exhaustively enumerating all 3×3 patterns of stones; it was trained by temporal-difference learning to predict the winner in games of self-play. Both *NeuroGo* and *RLGO* achieved a weak amateur level of play.

Monte-Carlo tree search (MCTS) may also be viewed as a form of self-play reinforcement learning⁴⁵. The nodes of the search tree contain the value function for the positions encountered during search; these values are updated to predict the winner of simulated games of self-play. MCTS programs have previously achieved strong amateur level in Go^{46,47}, but used substantial domain expertise: a fast *rollout policy*, based on handcrafted features^{48,49}, that evaluates positions by running simulations until the end of the game; and a *tree policy*, also based on handcrafted features, that selects moves within the search tree⁴⁷.

Self-play reinforcement learning approaches have achieved high levels of performance in other games: chess^{49–51}, checkers⁵², backgammon⁵³, othello⁵⁴, Scrabble⁵⁵ and most recently

poker⁵⁶. In all of these examples, a value function was trained by regression^{54–56} or temporal-difference learning^{49–53} from training data generated by self-play. The trained value function was used as an evaluation function in an alpha-beta search^{49–54}, a simple Monte-Carlo search^{55,57}, or counterfactual regret minimisation⁵⁶. However, these methods utilised handcrafted input features^{49–53,56} or handcrafted feature templates^{54,55}. In addition, the learning process used supervised learning to initialise weights⁵⁸, hand-selected weights for piece values^{49,51,52}, handcrafted restrictions on the action space⁵⁶, or used pre-existing computer programs as training opponents^{49,50} or to generate game records⁵¹.

Many of the most successful and widely used reinforcement learning methods were first introduced in the context of zero-sum games: temporal-difference learning was first introduced for a checkers-playing program⁵⁹, while MCTS was introduced for the game of Go¹². However, very similar algorithms have subsequently proven highly effective in video games^{6–8,10}, robotics⁶⁰, industrial control^{61–63}, and online recommendation systems^{64,65}.

AlphaGo versions We compare three distinct versions of *AlphaGo*:

1. *AlphaGo Fan* is the previously published program¹² that played against Pan Hui in October 2015. This program was distributed over many machines using 176 GPUs.
2. *AlphaGo Lee* is the program that defeated Lee Sedol 4–1 in March, 2016. It was previously unpublished but is similar in most regards to *AlphaGo Fan*¹². However, we highlight several key differences to facilitate a fair comparison. First, the value network was trained from the outcomes of fast games of self-play by *AlphaGo*, rather than games of self-play by the policy network; this procedure was iterated several times – an initial step towards the tabula rasa algorithm presented in this paper. Second, the policy and value networks were larger than those described in the original paper – using 12 convolutional layers of 256 planes respectively – and were trained for more iterations. This player was also distributed over many machines using 48 TPUs, rather than GPUs, enabling it to evaluate neural networks faster during search.

3. *AlphaGo Master* is the program that defeated top human players by 60–0 in January, 2017³⁴. It was previously unpublished but uses the same neural network architecture, reinforcement learning algorithm, and MCTS algorithm as described in this paper. However, it uses the same handcrafted features and rollouts as *AlphaGo Lee*¹² and training was initialised by supervised learning from human data.

4. *AlphaGo Zero* is the program described in this paper. It learns from self-play reinforcement learning, starting from random initial weights, without using rollouts, with no human supervision, and using only the raw board history as input features. It uses just a single machine in the Google Cloud with 4 TPUs (*AlphaGo Zero* could also be distributed but we chose to use the simplest possible search algorithm).

Domain Knowledge Our primary contribution is to demonstrate that superhuman performance can be achieved without human domain knowledge. To clarify this contribution, we enumerate the domain knowledge that *AlphaGo Zero* uses, explicitly or implicitly, either in its training procedure or its Monte-Carlo tree search; these are the items of knowledge that would need to be replaced for *AlphaGo Zero* to learn a different (alternating Markov) game.

1. *AlphaGo Zero* is provided with perfect knowledge of the game rules. These are used during MCTS, to simulate the positions resulting from a sequence of moves, and to score any simulations that reach a terminal state. Games terminate when both players pass, or after $19 \times 19 \cdot 2 = 722$ moves. In addition, the player is provided with the set of legal moves in each position.
2. *AlphaGo Zero* uses Tromp-Taylor scoring⁶⁶ during MCTS simulations and self-play training. This is because human scores (Chinese, Japanese or Korean rules) are not well-defined if the game terminates before territorial boundaries are resolved. However, all tournament and evaluation games were scored using Chinese rules.
3. The input features describing the position are structured as a 19×19 image; i.e. the neural network architecture is matched to the grid-structure of the board.

4. The rules of Go are invariant under rotation and reflection; this knowledge has been utilised in *AlphaGo Zero* both by augmenting the data set during training to include rotations and reflections of each position, and to sample random rotations or reflections of the position during MCTS (see Search Algorithm). Aside from komi, the rules of Go are also invariant to colour transposition; this knowledge is exploited by representing the board from the perspective of the current player (see Neural network architecture).

AlphaGo Zero does not use any form of domain knowledge beyond the points listed above. It only uses its deep neural network to evaluate leaf nodes and to select moves (see section below). It does not use any rollout policy or tree policy, and the MCTS is not augmented by any other heuristics or domain-specific rules. No legal moves are excluded – even those filling in the player’s own eyes (a standard heuristic used in all previous programs⁶⁷).

The algorithm was started with random initial parameters for the neural network. The neural network architecture (see Neural Network Architecture) is based on the current state of the art in image recognition^{4,18}, and hyperparameters for training were chosen accordingly (see Self-Play Training Pipeline). MCTS search parameters were selected by Gaussian process optimisation⁶⁸, so as to optimise self-play performance of *AlphaGo Zero* using a neural network trained in a preliminary run. For the larger run (40 block, 40 days), MCTS search parameters were re-optimised using the neural network trained in the smaller run (20 block, 3 days). The training algorithm was executed autonomously without human intervention.

Self-Play Training Pipeline *AlphaGo Zero*’s self-play training pipeline consists of three main components, all executed asynchronously in parallel. Neural network parameters θ , are continually optimised from recent self-play data; *AlphaGo Zero* players α_{θ} , are continually evaluated; and the best performing player so far, $\alpha_{\theta_{\text{best}}}$, is used to generate new self-play data.

Optimisation Each neural network f_{θ_t} is optimised on the Google Cloud using TensorFlow, with 64 GPU workers and 19 CPU parameter servers. The batch-size is 32 per worker, for a total mini-batch size of 2,048. Each mini-batch of data is sampled uniformly at random from

all positions from the most recent 500,000 games of self-play. Neural network parameters are optimised by stochastic gradient descent with momentum and learning rate annealing, using the loss in Equation 1. The learning rate is annealed according to the standard schedule in Extended Data Table 3. The momentum parameter is set to 0.9. The cross-entropy and mean-squared error losses are weighted equally (this is reasonable because rewards are unit scaled, $r \in \{-1, +1\}$) and the L2 regularisation parameter is set to $c = 10^{-4}$. The optimisation process produces a new checkpoint every 1,000 training steps. This checkpoint is evaluated by the evaluator and it may be used for generating the next batch of self-play games, as we explain next.

Evaluator To ensure we always generate the best quality data, we evaluate each new neural network checkpoint against the current best network f_{θ_t} before using it for data generation. The neural network f_{θ_t} is evaluated by the performance of an MCTS search α_{θ_t} that uses f_{θ_t} to evaluate leaf positions and prior probabilities (see Search Algorithm). Each evaluation consists of 400 games, using an MCTS with 1,600 simulations to select each move, using an infinitesimal temperature $\tau \rightarrow 0$ (i.e. we deterministically select the move with maximum visit count, to give the strongest possible play). If the new player wins by a margin of > 55% (to avoid selecting on noise alone) then it becomes the best player $\alpha_{\theta_{\text{best}}}$, and is subsequently used for self-play generation, and also becomes the baseline for subsequent comparisons.

Self-Play The best current player $\alpha_{\theta_{\text{best}}}$, as selected by the evaluator, is used to generate data. In each iteration, $\alpha_{\theta_{\text{best}}}$ plays 25,000 games of self-play, using 1,600 simulations of MCTS to select each move (this requires approximately 0.4s per search). For the first 30 moves of each game, the temperature is set to $\tau = 1$; this selects moves proportionally to their visit count in MCTS, and ensures a diverse set of positions are encountered. For the remainder of the game, an infinitesimal temperature is used, $\tau \rightarrow 0$. Additional exploration is achieved by adding Dirichlet noise to the prior probabilities in the root node s_0 , specifically $P(s, a) = (1 - \epsilon)p_a + \epsilon\eta_a$, where $\eta \sim \text{Dir}(0.03)$ and $\epsilon = 0.25$; this noise ensures that all moves may be tried, but the search may still overrule bad moves. In order to save computation, clearly lost games are resigned. The resignation threshold v_{resign} is selected automatically to keep the fraction of false positives (games that could have been

won if *AlphaGo* had not resigned) below 5%. To measure false positives, we disable resignation in 10% of self-play games and play until termination.

Supervised Learning For comparison, we also trained neural network parameters θ_{SL} by supervised learning. The neural network architecture was identical to *AlphaGo Zero*. Mini-batches of data (s, π, z) were sampled at random from the KGS data-set, setting $\pi_a = 1$ for the human expert move a . Parameters were optimised by stochastic gradient descent with momentum and learning rate annealing, using the same loss as in Equation 1, but weighting the mean-squared error component by a factor of 0.01. The learning rate was annealed according to the standard schedule in Extended Data Table 3. The momentum parameter was set to 0.9, and the L2 regularisation parameter was set to $c = 10^{-4}$.

By using a combined policy and value network architecture, and by using a low weight on the value component, it was possible to avoid overfitting to the values (a problem described in prior work¹²). After 72 hours the move prediction accuracy exceeded the state of the art reported in previous work^{12,30-33}, reaching 60.4% on the KGS test set; the value prediction error was also substantially better than previously reported¹². The validation set was composed of professional games from *GoKifu*. Accuracies and mean squared errors are reported in Extended Data Table 1 and Extended Data Table 2 respectively.

Search Algorithm *AlphaGo Zero* uses a much simpler variant of the asynchronous policy and value MCTS algorithm (APV-MCTS) used in *AlphaGo Fan* and *AlphaGo Lee*.

Each node s in the search tree contains edges (s, a) for all legal actions $a \in \mathcal{A}(s)$. Each edge stores a set of statistics:

$$\{N(s, a), W(s, a), Q(s, a), P(s, a)\},$$

where $N(s, a)$ is the visit count, $W(s, a)$ is the total action-value, $Q(s, a)$ is the mean action-value, and $P(s, a)$ is the prior probability of selecting that edge. Multiple simulations are executed in parallel on separate search threads. The algorithm proceeds by iterating over three phases (a–c in Figure 2), and then selects a move to play (d in Figure 2).

Select (Figure 2a). The selection phase is almost identical to *AlphaGo Fan*¹²; we recapitulate here for completeness. The first *in-tree phase* of each simulation begins at the root node of the search tree, s_0 , and finishes when the simulation reaches a leaf node s_L at time-step L . At each of these time-steps, $t < L$, an action is selected according to the statistics in the search tree, $a_t = \underset{a}{\operatorname{argmax}}(Q(s_t, a) + U(s_t, a))$, using a variant of the PUCT algorithm²⁴,

$$U(s, a) = c_{\text{puct}} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}$$

where c_{puct} is a constant determining the level of exploration; this search control strategy initially prefers actions with high prior probability and low visit count, but asymptotically prefers actions with high action-value.

Expand and evaluate (Figure 2b). The leaf node s_L is added to a queue for neural network evaluation, $(d_i(p), v) = f_{\theta}(d_i(s_L))$, where d_i is a dihedral reflection or rotation selected uniformly at random from $i \in [1..8]$.

Positions in the queue are evaluated by the neural network using a mini-batch size of 8; the search thread is locked until evaluation completes. The leaf node is expanded and each edge (s_L, a) is initialised to $\{N(s_L, a) = 0, W(s_L, a) = 0, Q(s_L, a) = 0, P(s_L, a) = p_a\}$; the value v is then backed up.

Backup (Figure 2c). The edge statistics are updated in a backward pass through each step $t \leq L$. The visit counts are incremented, $N(s_t, a) = N(s_t, a_t) + 1$, and the action-value is updated to the mean value, $W(s_t, a_t) = W(s_t, a_t) + v$, $Q(s_t, a_t) = \frac{W(s_t, a_t)}{N(s_t, a_t)}$. We use virtual loss to ensure each thread evaluates different nodes⁶⁹.

Play (Figure 2d). At the end of the search *AlphaGo Zero* selects a move a to play in the root position s_0 , proportional to its exponentiated visit count, $\pi(a|s_0) = N(s_0, a)^{1/\tau} / \sum_b N(s_0, b)^{1/\tau}$, where τ is a temperature parameter that controls the level of exploration. The search tree is reused at subsequent time-steps: the child node corresponding to the played action becomes the new root node; the subtree below this child is retained along with all its statistics, while the remainder of

the tree is discarded. *AlphaGo Zero* resigns if its root value and best child value are lower than a threshold value v_{resign} .

Compared to the MCTS in *AlphaGo Fan* and *AlphaGo Lee*, the principal differences are that *AlphaGo Zero* does not use any rollouts; it uses a single neural network instead of separate policy and value networks; leaf nodes are always expanded, rather than using dynamic expansion; each search thread simply waits for the neural network evaluation, rather than performing evaluation and backup asynchronously; and there is no tree policy. A transposition table was also used in the large (40 block, 40 day) instance of *AlphaGo Zero*.

Neural Network Architecture The input to the neural network is a $19 \times 19 \times 17$ image stack comprising 17 binary feature planes. 8 feature planes X_t consist of binary values indicating the presence of the current player's stones ($X_t^i = 1$ if intersection i contains a stone of the player's colour at time-step t ; 0 if the intersection is empty, contains an opponent stone, or if $t < 0$). A further 8 feature planes, Y_t , represent the corresponding features for the opponent's stones. The final feature plane, C , represents the colour to play, and has a constant value of either 1 if black is to play or 0 if white is to play. These planes are concatenated together to give input features $s_t = [X_t, Y_t, X_{t-1}, Y_{t-1}, \dots, X_{t-7}, Y_{t-7}, C]$. History features X_t, Y_t are necessary because Go is not fully observable solely from the current stones, as repetitions are forbidden; similarly, the colour feature C is necessary because the *komi* is not observable.

The input features s_t are processed by a residual tower that consists of a single convolutional block followed by either 19 or 39 residual blocks¹⁸.

The convolutional block applies the following modules:

1. A convolution of 256 filters of kernel size 3×3 with stride 1
2. Batch normalisation¹⁹
3. A rectifier non-linearity

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

4. A fully connected linear layer to a hidden layer of size 256
5. A rectifier non-linearity
6. A fully connected linear layer to a scalar
7. A tanh non-linearity outputting a scalar in the range $[-1, 1]$

The overall network depth, in the 20 or 40 block network, is 39 or 79 parameterised layers respectively for the residual tower, plus an additional 2 layers for the policy head and 3 layers for the value head.

We note that a different variant of residual networks was simultaneously applied to computer Go²⁰ and achieved amateur dan-level performance; however this was restricted to a single-headed policy network trained solely by supervised learning.

Neural Network Architecture Comparison Figure 4 shows the results of a comparison between network architectures. Specifically, we compared four different neural networks:

1. *dual-res*: The network contains a 20-block residual tower, as described above, followed by both a policy head and a value head. This is the architecture used in *AlphaGo Zero*.
2. *sep-res*: The network contains two 20-block residual towers. The first tower is followed by a policy head and the second tower is followed by a value head.
3. *dual-conv*: The network contains a non-residual tower of 12 convolutional blocks, followed by both a policy head and a value head.
4. *sep-conv*: The network contains two non-residual towers of 12 convolutional blocks. The first tower is followed by a policy head and the second tower is followed by a value head. This is the architecture used in *AlphaGo Lee*.

Each network was trained on a fixed data-set containing the final 2 million games of self-play data generated by a previous run of *AlphaGo Zero*, using stochastic gradient descent with

Each residual block applies the following modules sequentially to its input:

1. A convolution of 256 filters of kernel size 3×3 with stride 1
2. Batch normalisation
3. A rectifier non-linearity
4. A convolution of 256 filters of kernel size 3×3 with stride 1
5. Batch normalisation
6. A skip connection that adds the input to the block
7. A rectifier non-linearity

The output of the residual tower is passed into two separate "heads" for computing the policy and value respectively. The policy head applies the following modules:

1. A convolution of 2 filters of kernel size 1×1 with stride 1
2. Batch normalisation
3. A rectifier non-linearity
4. A fully connected linear layer that outputs a vector of size $19^2 + 1 = 362$ corresponding to logit probabilities for all intersections and the pass move

The value head applies the following modules:

1. A convolution of 1 filter of kernel size 1×1 with stride 1
2. Batch normalisation
3. A rectifier non-linearity

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

the annealing rate, momentum, and regularisation hyperparameters described for the supervised learning experiment; however, cross-entropy and mean-squared error components were weighted equally, since more data was available.

Evaluation We evaluated the relative strength of *AlphaGo Zero* (Figure 3a and 6) by measuring the Elo rating of each player. We estimate the probability that player a will defeat player b by a logistic function $p(a \text{ defeats } b) = \frac{1}{1 + \exp(\alpha_{ab} \text{elo}_b - \text{elo}_a)}$, and estimate the ratings $\text{elo}(.)$ by Bayesian logistic regression, computed by the *BayesElo* program²¹ using the standard constant $c_{\text{elo}} = 1/400$.

Elo ratings were computed from the results of a 5 second per move tournament between *AlphaGo Zero*, *AlphaGo Master*, *AlphaGo Lee*, and *AlphaGo Fan*. The raw neural network from *AlphaGo Zero* was also included in the tournament. The Elo ratings of *AlphaGo Fan*, *Crazy Stone*, *Pachi* and *GinuGo* were anchored to the tournament values from prior work¹², and correspond to the players reported in that work. The results of the matches of *AlphaGo Fan* against Fan Hui and *AlphaGo Lee* against Lee Sedol were also included to ground the scale to human references, as otherwise the Elo ratings of *AlphaGo* are unrealistically high due to self-play bias.

The Elo ratings in Figure 3a, 4a and 6a were computed from the results of evaluation games between each iteration of player α_{elo} during self-play training. Further evaluations were also performed against baseline players with Elo ratings anchored to the previously published values¹².

We measured the head-to-head performance of *AlphaGo Zero* against *AlphaGo Lee*, and the 40 block instance of *AlphaGo Zero* against *AlphaGo Master*, using the same player and match conditions as were used against Lee Sedol in Seoul, 2016. Each player received 2 hours of thinking time plus 3 byoyomi periods of 60 seconds per move. All games were scored using Chinese rules with a *komi* of 7.5 points.

Data Availability The datasets used for validation and testing are the *GoKifu* dataset (available from <http://gokifu.com/>) and the *KGS* dataset (available from <https://u-go.net/gamerecords/>).

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

	<i>KGS</i> train	<i>KGS</i> test	<i>GoKifu</i> validation
Supervised learning (20 block)	62.0	60.4	54.3
Supervised learning (12 layer ¹²)	59.1	55.9	-
Reinforcement learning (20 block)	-	-	49.0
Reinforcement learning (40 block)	-	-	51.3

Extended Data Table 1: Move prediction accuracy. Percentage accuracy of move prediction for neural networks trained by reinforcement learning (i.e., *AlphaGo Zero*) or supervised learning respectively. For supervised learning, the network was trained for 3 days on *KGS* data (amateur games); comparative results are also shown from Silver et al.¹². For reinforcement learning, the 20 block network was trained for 3 days and the 40 block network was trained for 40 days. Networks were also evaluated on a validation set based on professional games from the *GoKifu* data set.

	<i>KGS</i> train	<i>KGS</i> test	<i>GoKifu</i> validation
Supervised learning (20 block)	0.177	0.185	0.207
Supervised learning (12 layer ¹²)	0.19	0.37	-
Reinforcement learning (20 block)	-	-	0.177
Reinforcement learning (40 block)	-	-	0.180

Extended Data Table 2: Game outcome prediction error. Mean squared error on game outcome predictions for neural networks trained by reinforcement learning (i.e., *AlphaGo Zero*) or supervised learning respectively. For supervised learning, the network was trained for 3 days on *KGS* data (amateur games); comparative results are also shown from Silver et al.¹². For reinforcement learning, the 20 block network was trained for 3 days and the 40 block network was trained for 40 days. Networks were also evaluated on a validation set based on professional games from the *GoKifu* data set.

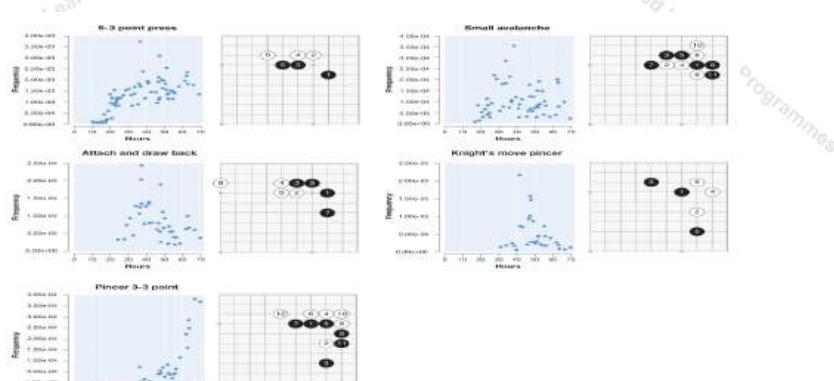
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

Thousands of steps	Reinforcement learning	Supervised learning
0–200	10^{-2}	10^{-1}
200–400	10^{-2}	10^{-2}
400–600	10^{-3}	10^{-3}
600–700	10^{-4}	10^{-4}
700–800	10^{-4}	10^{-5}
>800	10^{-4}	-

Extended Data Table 3: Learning rate schedule. Learning rate used during reinforcement learning and supervised learning experiments, measured in thousands of steps (mini-batch updates).

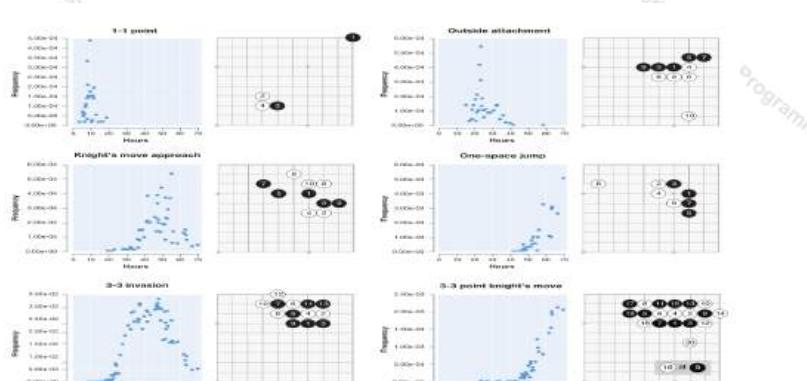
Work Integrated Learning Programmes

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Extended Data Figure 1: Frequency of occurrence over time during training, for each *joseki* from Figure 5a (corner sequences common in professional play that were discovered by *AlphaGo Zero*). The corresponding *joseki* are reproduced in this figure as insets.

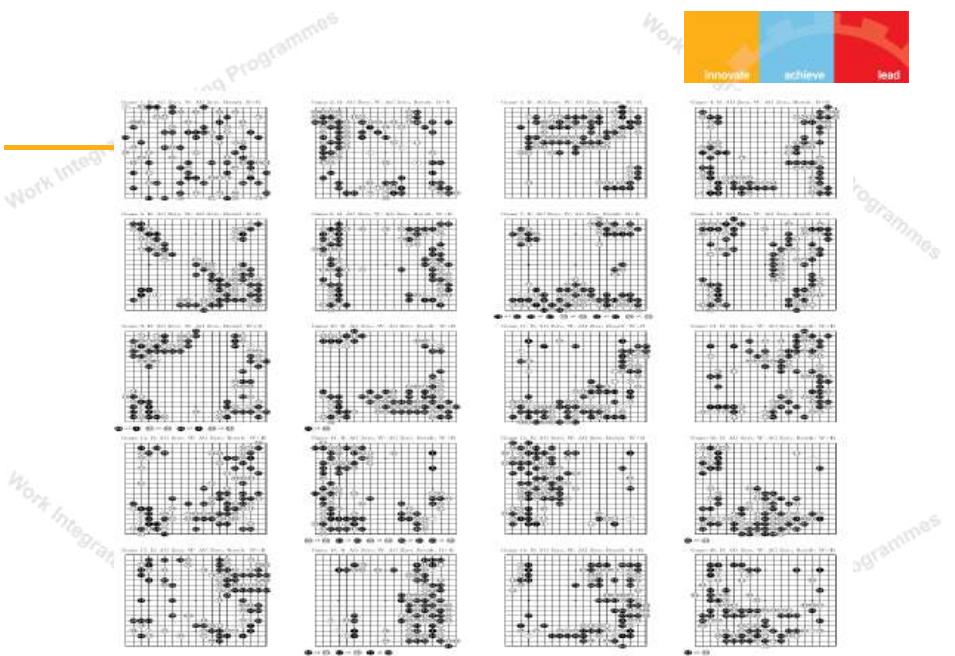
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Extended Data Figure 2: Frequency of occurrence over time during training, for each *joseki* from Figure 5b (corner sequences that *AlphaGo Zero* favoured for at least one iteration), and one additional variation. The corresponding *joseki* are reproduced in this figure as insets.

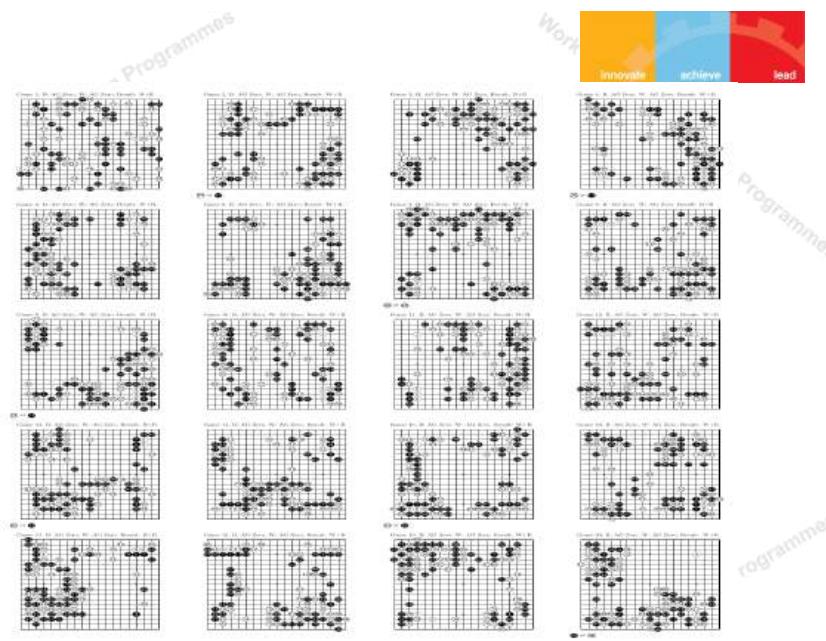
BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



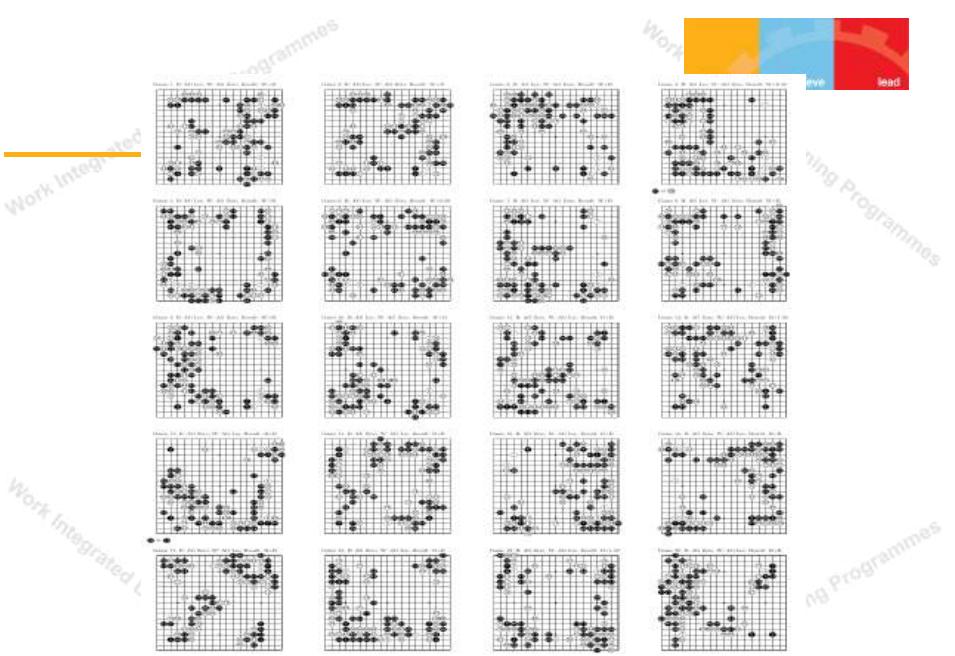
Extended Data Figure 3: *AlphaGo Zero* (20 block) self-play games. The 3 day training run was subdivided into 20 periods. The best player from each period (as selected by the evaluator) played a single game against itself, with 2 hour time controls. 100 moves are shown for each game; full games are provided in Supplementary Information.

UGC Act, 1956



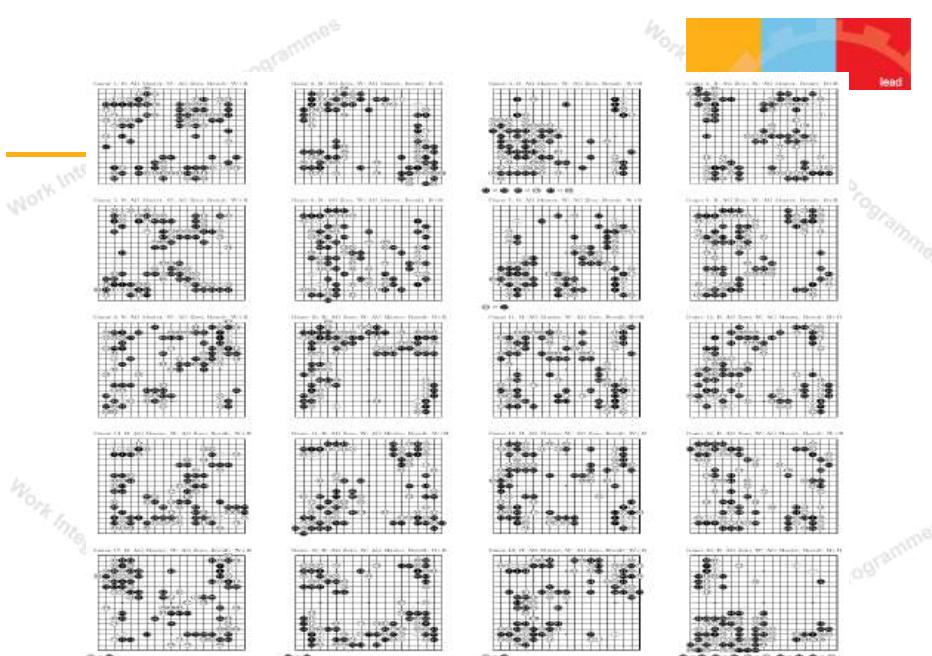
Extended Data Figure 4: *AlphaGo Zero* (40 block) self-play games. The 40 day training run was subdivided into 20 periods. The best player from each period (as selected by the evaluator) played a single game against itself, with 2 hour time controls. 100 moves are shown for each game; full games are provided in Supplementary Information.

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956



Extended Data Figure 5: Tournament games between *AlphaGo Zero* (20 block, 3 day) versus *AlphaGo Lee* using 2 hour time controls. 100 moves of the first 20 games are shown; full games are provided in Supplementary Information.

on 3 of UGC Act, 1956



Extended Data Figure 6: *AlphaGo Zero* (40 block, 40 day) versus *AlphaGo Master* tournament games using 2 hour time controls. 100 moves of the first 20 games are shown; full games are provided in Supplementary Information.

UGC Act, 1956