

Doy-100, Mar 10 - 2025, folgun 26, 2081 BS.

- ① Bell Man Equation (12 min)
- ② Making Decisions (2 min)
- ③ State - Action function Definition & Example (15 min)
- ④ Learning the State-value function (16 min)
- ⑤ Algorithm Refinement ϵ -greedy Policy
- ⑥ The State of Reinforcement Learning.
- ⑦ Acknowledgements:
- ⑧ What's Daxt?

Bellman Equation, $Q(s, a) \rightarrow$ Help to compute Q of s of a .

Bellman Equation

$Q(s, a)$ = Return if you

- start in state s .
- take action a (once).
- then behave optimally after that.

s : current state

a : Current action

s' : State you get to after taking action a

a' : Action that you take in state s'

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

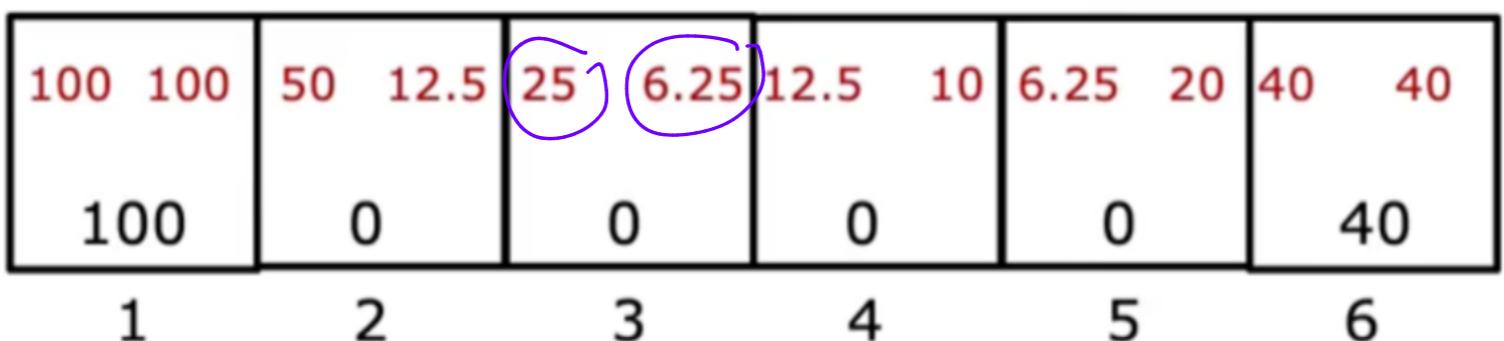


$$R(1) = 100 \quad R(2) = 0 \quad \dots \quad R(6) = 40$$

$R(s)$ = reward of current state

Bellman Equation

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



$$Q(2, \rightarrow) \Rightarrow R(2) + 0.5 \max_{a'} Q(3, a')$$
$$\Rightarrow 0 + 0.5(25)$$

$$[Q(2, \rightarrow) = 12.5]$$

$$S = 2 \\ a = \rightarrow \\ S' = 3$$

$$S = 4 \\ a \leftarrow \\ S' = 3$$

$$Q(4, \leftarrow) \Rightarrow R(4) + 0.5 \max_{a'} Q(3, a')$$
$$\Rightarrow 0 + (0.5) \times 25$$

$$[Q(4, \leftarrow) = 12.5]$$

Explanation of Bellman Equation

$Q(s, a) =$ Return if you
• start in state s .
• take action a (once).
• then behave optimally after that.

$s \rightarrow s'$

The best possible return from state s is $\max_a Q(s', a')$

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

$$R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$$

Broken down into Components

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

$R(s) \rightarrow$ Immediate Reward

$$\Rightarrow R_1 + \gamma [R_2 + \gamma R_3 + \gamma^2 R_4 + \dots]$$

$\max_{a'} Q(s', a')$ \rightarrow Return for behaving optimally starting from state s'

Explanation of Bellman Equation

3. Bellman Equation Explanation

The Bellman equation is:

$$Q(S, A) = R(S) + \gamma \max_{A'} Q(S', A')$$

Where:

- S' is the next state after taking action A from state S .
- A' is the next action to take from S' .
- $\max_{A'} Q(S', A')$ represents the best possible future return from state S' .

Breaking It Down:

1. First part: $R(S)$ is the immediate reward for being in state S .
2. Second part: $\gamma \max_{A'} Q(S', A')$ represents the discounted future rewards—what you will get starting from S' and behaving optimally.

$Q(s, a) = \text{Return if you}$

- start in state s .
- take action a (once).
- then behave optimally after that.

$s \rightarrow s'$

→ The best possible return from state s' is $\max_a Q(s', a')$

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

Reward you get right away Return from behaving optimally starting from state s' .

$$R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$$
$$Q(s, a) = R_1 + \gamma [R_2 + \gamma [R_3 + \gamma [R_4 + \dots]]]$$

4. Example Calculation

Consider the Markov Decision Process (MDP) example in the lecture:

- The reward values for states:
 - $R(1) = 100$
 - $R(2) = 0$
 - $R(3) = 0$
 - $R(4) = 0$
 - $R(5) = 0$
 - $R(6) = 40$
- Suppose the discount factor is $\gamma = 0.5$.

Computing $Q(2, \text{right})$:

DeepL

Explanation of Bellman Equation

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

100	100	50	12.5	25	6.25	12.5	10	6.25	20	40	40
100		0		0		0	0	0	40	40	

$$\begin{aligned} Q(4, \leftarrow) &= 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100 \\ &= R(4) + (0.5)[0 + (0.5)0 + (0.5)^2 100] \\ &= R(4) + (0.5) \max_{a'} Q(3, a') \end{aligned}$$

the reward you're getting right away.

5. The Intuition Behind the Bellman Equation

The Bellman equation breaks down the total expected return into two parts:

1. **Immediate reward $R(S)$** (what you get right now).
2. **Future rewards**, discounted by γ , assuming optimal actions thereafter.

Mathematically, this breakdown follows from the way rewards accumulate over time:

$$\text{Total Return} = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$$

Rearranging, we get:

$$\text{Total Return} = R(S) + \gamma(\text{Future Return})$$

This leads to the **Bellman equation**.

6. Key Takeaways

- Bellman equation recursively defines $Q(S, A)$.
- To find the best action in a state, choose the one that maximizes $Q(S, A)$.
- The equation helps compute optimal policies in RL.
- If you understand the basic intuition, you can apply the equation even if the details seem complex.

Computing $Q(2, \text{right})$:

- If you're in **State 2** and go **right**, you move to **State 3**.
- The Bellman equation gives:

$$Q(2, \text{right}) = R(2) + \gamma \max_{A'} Q(3, A')$$

- Since $R(2) = 0$, and the two possible $Q(3, A')$ values are **25** and **6.25**:

$$Q(2, \text{right}) = 0 + 0.5 \times 25 = 12.5$$

Computing $Q(4, \text{left})$:

- If you're in **State 4** and go **left**, you move to **State 3**.

$$Q(4, \text{left}) = R(4) + \gamma \max_{A'} Q(3, A')$$

- Since $R(4) = 0$ and the best possible value from **State 3** is **25**:

$$Q(4, \text{left}) = 0 + 0.5 \times 25 = 12.5$$

Terminal States:

For **terminal states**, the Bellman equation simplifies to:

$$Q(S, A) = R(S)$$

because there's no future state (S') to consider.

For example:

- If $S = 1$, then $Q(1, A) = 100$.
- If $S = 6$, then $Q(6, A) = 40$.

5. The Intuition Behind the Bellman Equation

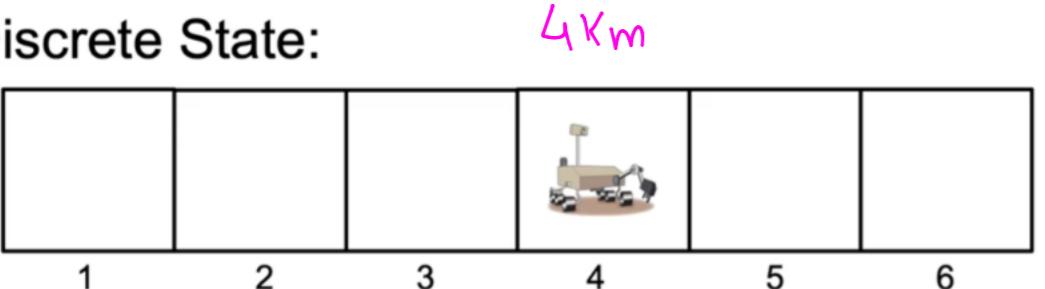
The Bellman equation breaks down the total expected return into two parts:



1. **Immediate reward $R(S)$** (what you get right now).

Discrete vs Continuous State

Discrete State:



Continuous State:

56.25 Km



2D
on on on
x
y
θ
2D

$s =$

$$S = \begin{bmatrix} x \\ y \\ \theta \\ x' \\ y' \\ \theta' \end{bmatrix}$$

$S = [6 \text{ numbers}]$

DeepLearning.AI Stanford instead, they can be in any of

Andrew Ng

Autonomous Helicopter (3D)



x
y
z
θ
φ
ω

State = ?
Angular. P.

$s =$

$S = [12 \text{ numbers}]$

$$S = \begin{bmatrix} x \\ y \\ z \\ \theta \\ \phi \\ \omega \\ x' \\ y' \\ z' \\ \theta' \\ \phi' \\ \omega' \end{bmatrix}$$

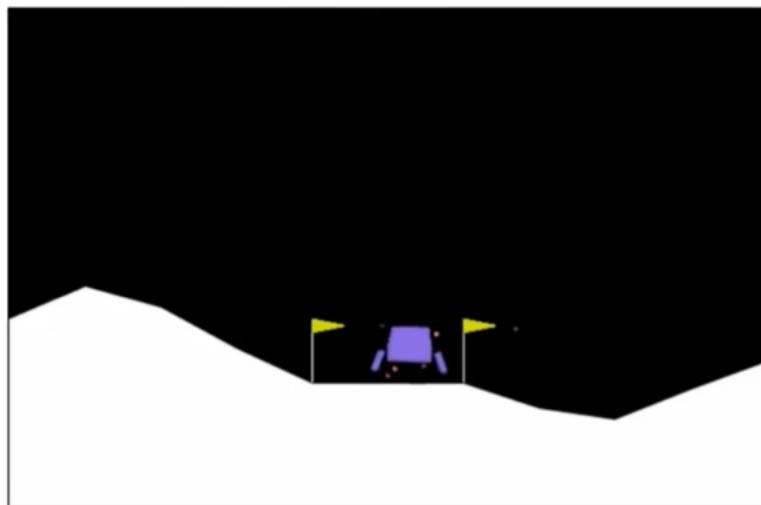
→ Vectors of numbers.

control an autonomous
helicopter,

Andrew Ng

#) lunar lander (Simulated).

Lunar Lander



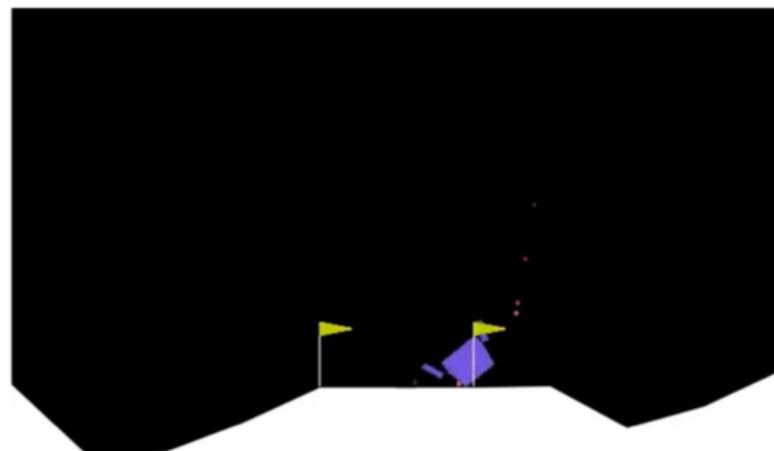
right to position itself to land
between these two yellow flags.

©DeepLearning.AI

Stanford ONLINE

Andrew Ng

Lunar Lander



it might look like where the lander
unfortunately has crashed

©DeepLearning.AI

Stanford ONLINE

Andrew Ng

→ Actions that make Soft landing.

→ Means Various Correct actions must be

triggered at a sequence of time

left \Rightarrow left leg.

\Rightarrow right leg

$x = y =$

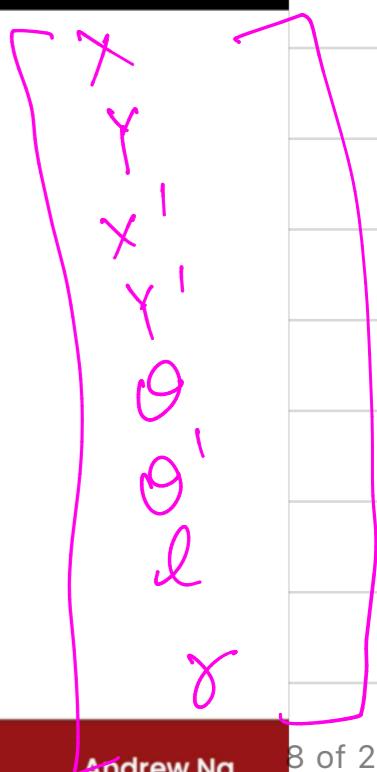
binary values

values $0 \text{ or } 1$

Lunar Lander



actions: $s =$
do nothing
left thruster
main thruster
right thruster



You could either do nothing,
in which case the forces of inertia and

©DeepLearning.AI

Stanford ONLINE

Andrew Ng

Lunar Lander and Reinforcement Learning

The Lunar Lander problem is a reinforcement learning environment where an agent (the lander) learns to land safely on a pad by optimizing its actions.

1. State Space:

- The lander's state is represented by:
 - Position: x, y
 - Velocity: \dot{x}, \dot{y} (horizontal and vertical speeds)
 - Angle: θ (tilt)
 - Angular velocity: $\dot{\theta}$
- Leg contact indicators: L and R (binary values indicating whether the left and right legs are touching the ground).

2. Action Space:

- The agent can take one of four discrete actions:
 - **Do nothing:** Gravity pulls the lander downward.
 - **Fire left thruster:** Pushes the lander to the right.
 - **Fire right thruster:** Pushes the lander to the left.
 - **Fire main thruster:** Slows the descent.

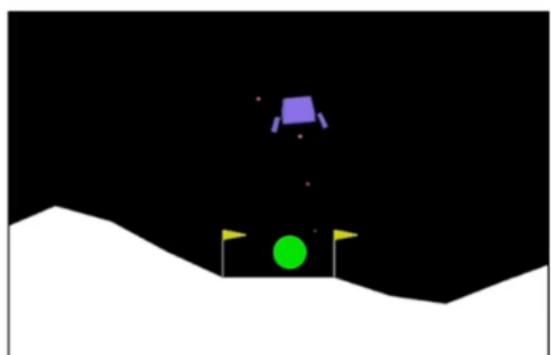
3. Reward Function:

- The goal is to reward good landings and penalize bad behavior:
 - Landing on the pad: **+100 to +140**.
 - Moving toward the pad: **positive reward**.
 - Moving away: **negative reward**.
 - Crashing: **-100**.
 - Landing softly: **+100 per leg**.
 - Landing with left or right leg: **+10 per leg**.
 - Fuel usage penalty:
 - Main engine: **-0.3 per use**.



Reward Function

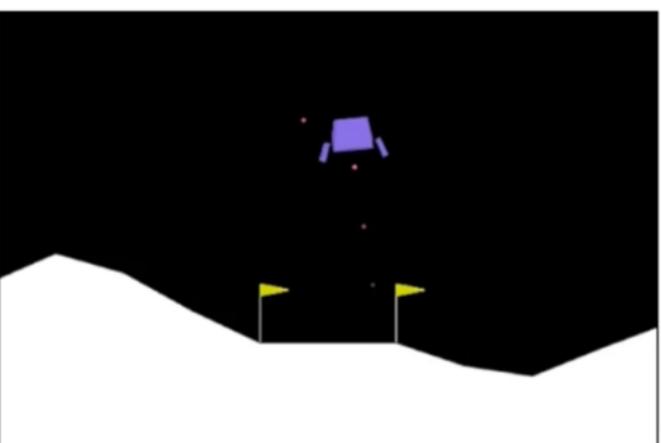
- Getting to landing pad: 100 – 140
- Additional reward for moving toward/away from pad.
- Crash: -100
- Soft landing: +100
- Leg grounded: +10
- Fire main engine: -0.3
- Fire side thruster: -0.03



Lunar Lander Problem

Learn a policy π that, given

$$s = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$

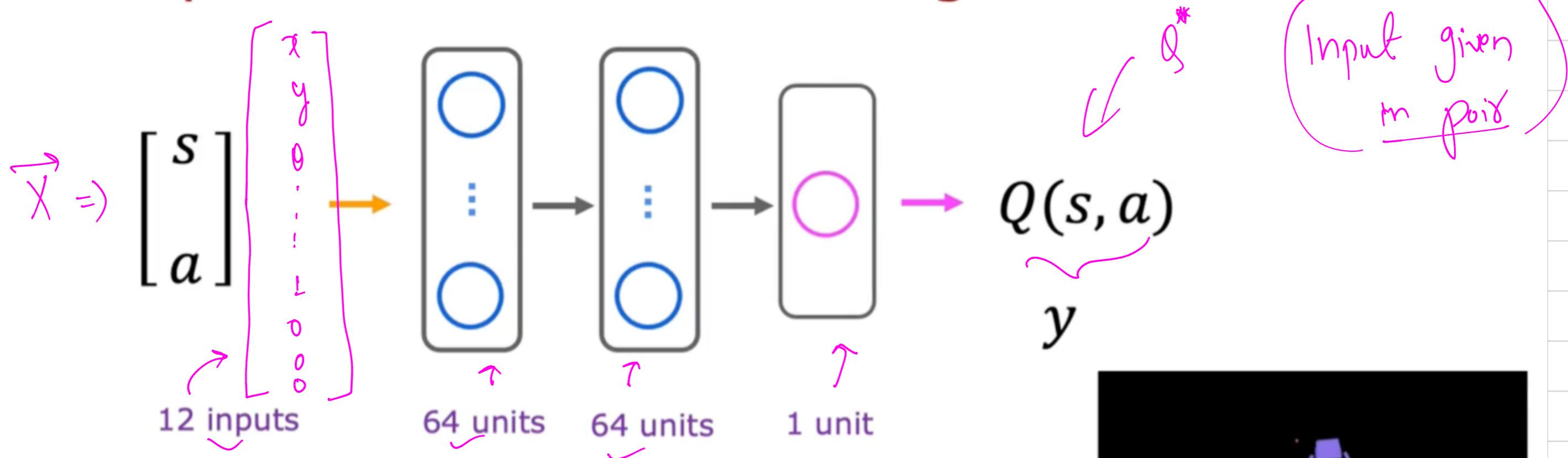


picks action $a = \pi(s)$ so as to maximize the return.

Use large value
 $\sqrt{= 0.985}$

Deep learning (State-Value Q^* function) \Rightarrow train NN to approximate $Q(s,a)$.

Deep Reinforcement Learning



In a state s , use neural network to compute

$Q(s, \text{nothing}), Q(s, \text{left}), Q(s, \text{main}), Q(s, \text{right})$

Pick the action a that maximizes $Q(s, a)$

any action a and
put them together.



One Hot feature vector:
[0 1 0 0]
↑
decides whole action

Bellman Equation

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

$$f_{w, b}(x) \approx y$$

$$(s, a, r(s), s') \leftarrow (s^{(1)}, a^{(1)}, r(s^{(1)}), s'^{(1)})$$

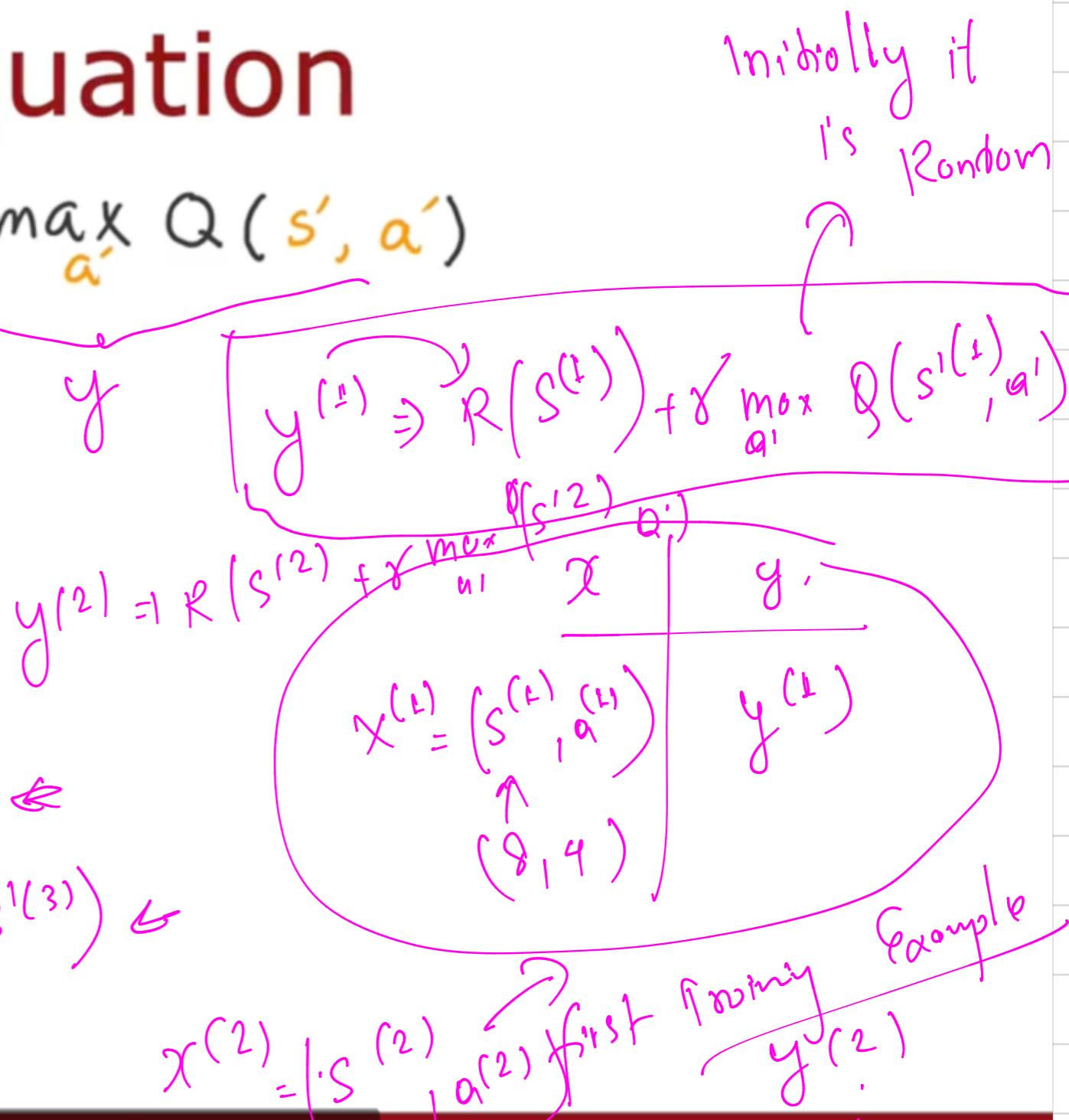
No policy then
Take Random Actions.

Result of Action gives, s'

$$(s^{(2)}, a^{(2)}, r(s^{(2)}), s'^{(2)}) \leftarrow$$

$$(s^{(3)}, a^{(3)}, r(s^{(3)}), s'^{(3)}) \leftarrow$$

Input



Bellman Equation

(A piece RL Algorithm)

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

$$f_{w, \beta}(x) \approx y$$

$$(s, a, R(s), s')$$

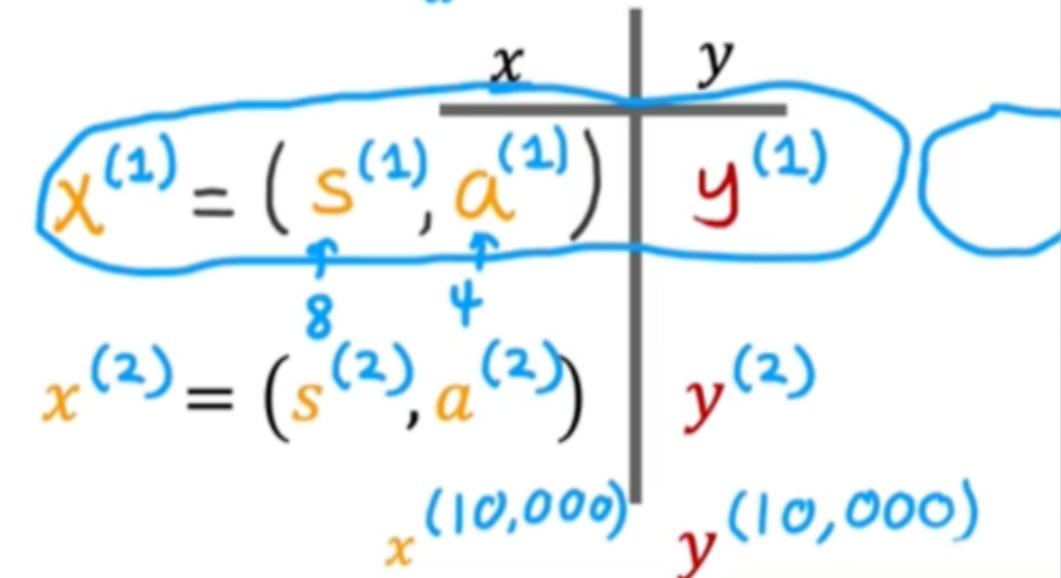
$$(s^{(1)}, a^{(1)}, R(s^{(1)}), s'^{(1)}) \leftarrow$$

$$(s^{(2)}, a^{(2)}, R(s^{(2)}), s'^{(2)}) \leftarrow$$

$$(s^{(3)}, a^{(3)}, R(s^{(3)}), s'^{(3)}) \leftarrow$$

$$y^{(1)} = R(s^{(1)}) + \gamma \max_{a'} Q(s'^{(1)}, a')$$

$$y^{(2)} = R(s^{(2)}) + \gamma \max_{a'} Q(s'^{(2)}, a')$$



What I describe here is

Learning Algorithm

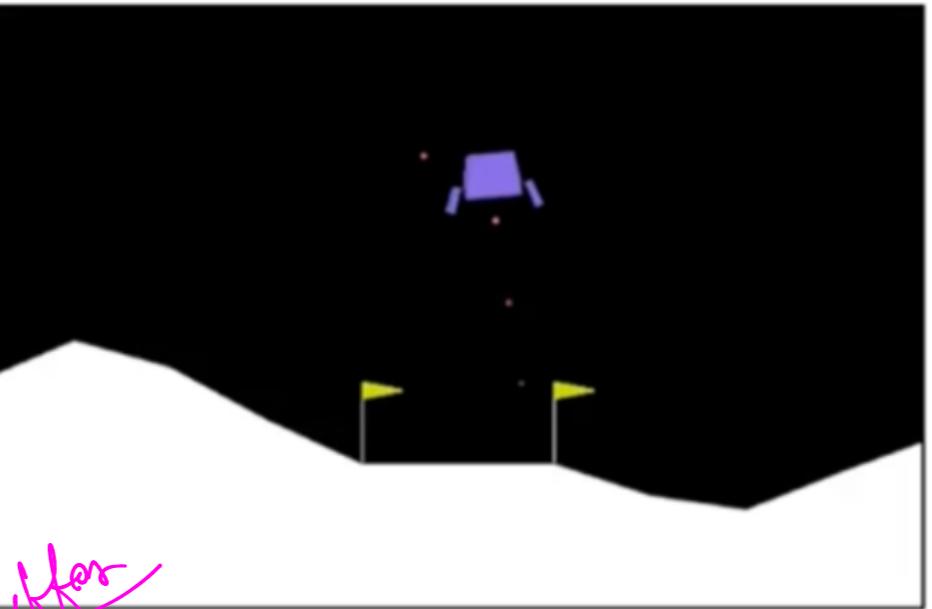
Initialize neural network randomly as guess of $\underline{Q(s,a)}$.

Repeat {

Take actions in the lunar lander. Get $(s, a, r(s), s')$.

Store 10,000 most recent $(s, a, r(s), s')$ tuples ↴

~~Replay Buffer~~



Train Neural Network.

Create training set of 10,000 examples using

$$X = (s, a) \text{ and } y = R(s) + \max_{a'} Q(s', a')$$

Twin Q s such that $Q_{\text{new}}(s, a) \approx y$

$$\delta_t + \beta = Q_{\text{new}}$$

training linear

$$f_{w, \beta}(x) \approx y$$

$$x_1, y_1, x'_1, y'_1 \\ \vdots$$

$$x^{1000}, y^{1000}$$

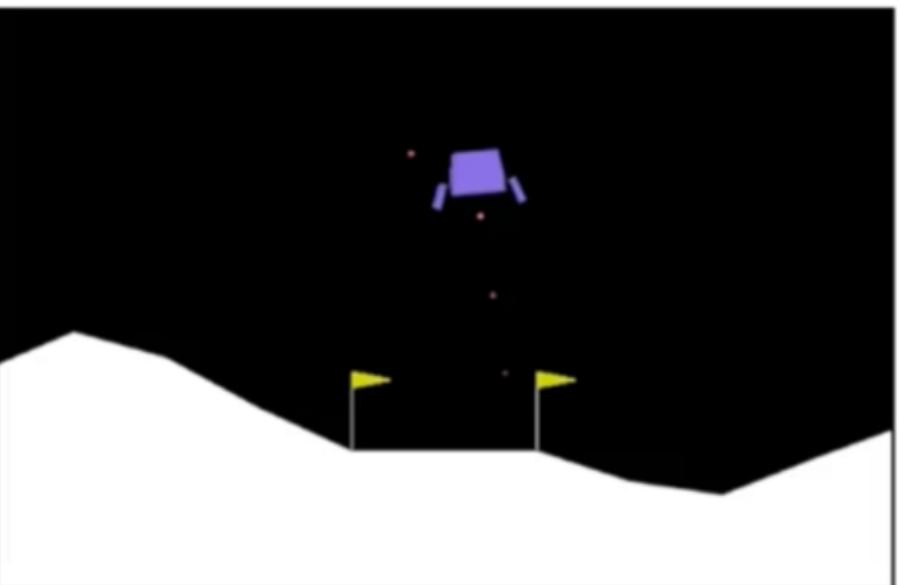
Learning Algorithm

Initialize neural network randomly as guess of $Q(s, a)$.

Repeat {

 Take actions in the lunar lander. Get $(s, a, R(s), s')$.

 Store 10,000 most recent $(s, a, R(s), s')$ tuples.



Replay Buffer

Train neural network:

 Create training set of 10,000 examples using

$$x = (s, a) \text{ and } y = R(s) + \gamma \max_{a'} Q(s', a')$$

 Train Q_{new} such that $Q_{new}(s, a) \approx y$.

 Set $Q = Q_{new}$.

$$f_{w, b}(x) \approx y$$

x, y

$x^{(1)}, y^{(1)}$

:

x^{10000}, y^{10000}

Then for the next time you train

1

Deep Q-Network (DQN) Algorithm for Learning the Q-Function

1. Problem Setup

The goal of reinforcement learning is to approximate the state-action value function $Q(s, a)$, which helps in selecting optimal actions.

For Lunar Lander, the state s consists of:

- (x, y) position
- (\dot{x}, \dot{y}) velocity
- $(\theta, \dot{\theta})$ orientation
- Leg contact indicators (Left, Right)

The actions a are one-hot encoded:

- Nothing: $[1, 0, 0, 0]$
- Left thruster: $[0, 1, 0, 0]$
- Main thruster: $[0, 0, 1, 0]$
- Right thruster: $[0, 0, 0, 1]$

Noed to
be
Reviewed

2. Q-Function Approximation with a Neural Network

A neural network is used to approximate $Q(s, a)$, where:

- Input: $X = [s, a]$ (12 features: 8 for state, 4 for action)
- Architecture:
 - Hidden layer 1: 64 neurons
 - Hidden layer 2: 64 neurons
 - Output: Single value $Q(s, a)$



2

2. Q-Function Approximation with a Neural Network

A neural network is used to approximate $Q(s, a)$, where:

- Input: $X = [s, a]$ (12 features: 8 for state, 4 for action)
- Architecture:
 - Hidden layer 1: 64 neurons
 - Hidden layer 2: 64 neurons
 - Output: Single value $Q(s, a)$

We train the network to predict $Q(s, a)$ using the Bellman equation:

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

where:

- $R(s)$ is the reward at state s
- s' is the next state after taking action a
- γ is the discount factor (typically 0.99)

3. Algorithm: Deep Q-Network (DQN)

Step 1: Initialize Neural Network

- Initialize all parameters randomly
- This serves as an initial guess for $Q(s, a)$

Step 2: Collect Experience Tuples

- Run the Lunar Lander simulator by taking random actions
- Store experience tuples:
- Maintain a Replay Buffer storing the last 10,000 tuples

$(s, a, R(s), s')$

↓

3

Step 4: Train Neural Network

- Use Mean Squared Error (MSE) Loss:

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

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N (Q(s^i, a^i) - y^i)^2$$

Step 5: Update Q-function

- Replace the old Q-function with the newly trained model

Step 6: Repeat

- Continue collecting new experiences, updating the neural network, and refining $Q(s, a)$
- Over time, $Q(s, a)$ improves, leading to better actions

4. Summary of Key Equations

1. Bellman Equation:

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

2. Target for Supervised Learning:

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

3. Loss Function (MSE):

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N (Q(s^i, a^i) - y^i)^2$$

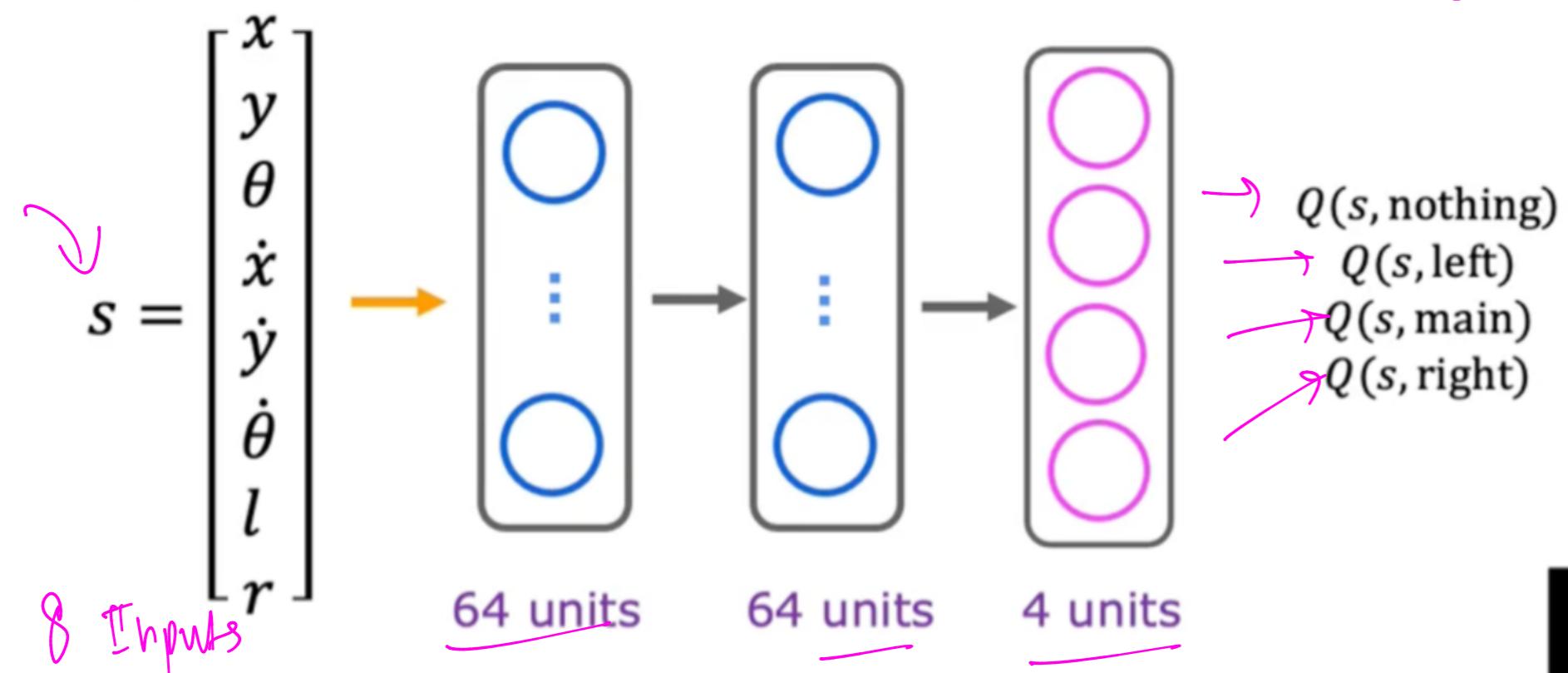
This DQN Algorithm allows reinforcement learning to efficiently learn optimal control strategies for the Lunar Lander problem. 🚀



Algorithm Refinement:

Deep Reinforcement Learning

(Modified / Better).

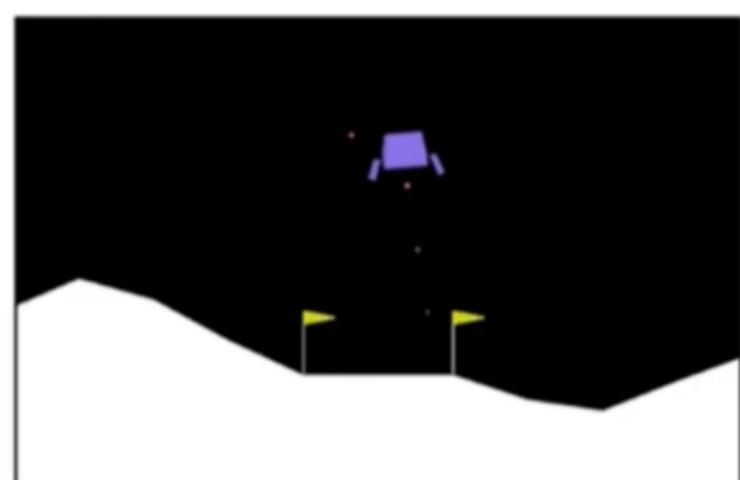


In a state s , input s to neural network.

Pick the action a that maximizes $Q(s, a)$.

Train S' at a time t .

$$R(s) + \gamma \max_a Q(s', a')$$



Here's a modified neural network

Algorithm Refinement ϵ -greedy Policy:

Learning Algorithm

Initialize neural network randomly as guess of $\underline{Q(s, a)}$.

Repeat {

Take actions in the lunar lander. Get $(s, a, R(s), s')$.

Store 10,000 most recent $(s, a, R(s), s')$ tuples.

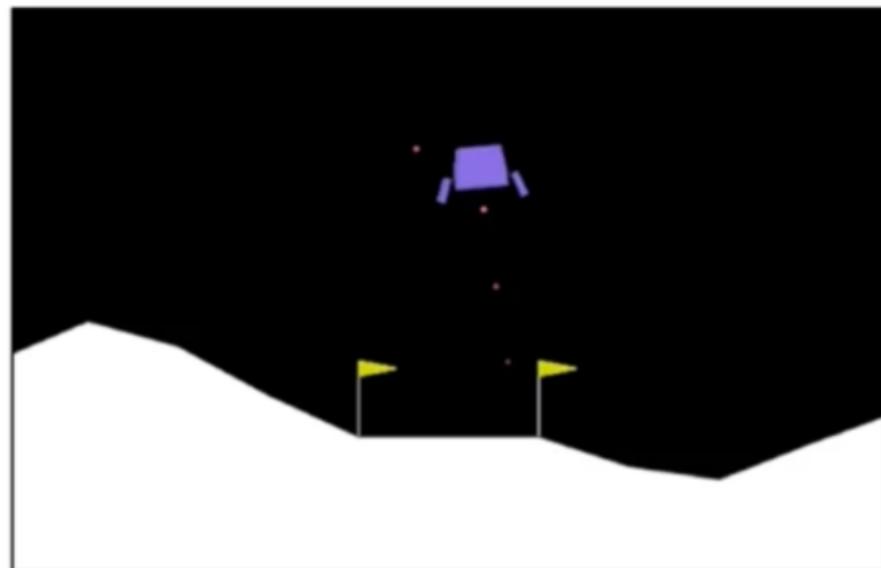
Train model:

Create training set of 10,000 examples using

$$x = (s, a) \text{ and } y = R(s) + \gamma \max_{a'} Q(s', a').$$

Train Q_{new} such that $Q_{new}(s, a) \approx y$. $f_{w.b}(x) \approx y$

Set $Q = Q_{new}$.



→ NN
Algorithm to
approximate
 $Q(s, a)$

Here's the algorithm
that you saw earlier.

How to choose actions while still learning?

In some state s

Option 1:

Pick the action a that maximizes $Q(s, a)$.

~~$\star \rightarrow \text{Best Optm}$~~

Maximize
Return

Option 2:

With probability 0.95, pick the action a that maximizes $Q(s, a)$.

Greedy

With probability 0.05, pick an action a randomly.

Small fraction of time t'
Exploration

$\rightarrow Q(s_{\text{main}})$ is low \rightarrow ever hover try to fire thruster
 $\rightarrow a \rightarrow$ just by chance it might not open thruster \Rightarrow we

keep 0.05 ' p ' for unexpected event.

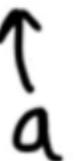
How to choose actions while still learning?

In some state s

Option 1:

Pick the action a that maximizes $\underline{Q(s,a)}$.

$Q(s,\text{main})$ is low



Option 2:

→ With probability 0.95, pick the action a that maximizes $\underline{Q(s,a)}$. Greedy, "Exploitation"

→ With probability 0.05, pick an action a randomly. "Exploration"

ϵ -greedy Policy ($\epsilon = 0.05$)

$0.95 \Rightarrow \epsilon$

$0.05 \Rightarrow$

→ Start ϵ high
1.0 ~ 0.01
→ Gradually decrease

versus trying to maximize
your return by say,

Mini-Batch and Soft Updates → (Helps RL Algorithm Efficiency)

How to choose actions while still learning?

x	y
2104	400
1416	232
1534	315
852	178
...	...
3210	870

Gradient Descent becomes slower for longer m .

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

repeat {
 $w = w - \alpha \frac{\partial}{\partial w} J(w, b)$
 $b = b - \alpha \frac{\partial}{\partial b} J(w, b)$

$m = 100,000,000$
↑ House sizes
Problem?
↳ Computing intensive.
 \sum , Average and
Do Backpropagation
and so on.

way back in the first course

How to choose actions while still learning?

x	y
2104	400
1416	232
1534	315
852	178
...	...
3210	870

100,000,000

Look at 1,000, rather than all.
(Mini-Batch -)

we would pick some subset of
1,000 or m prime examples.

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$m = 100,000,000$$

$$m' = 1,000$$

repeat {

$$w = w - \alpha \frac{\partial}{\partial w} \left[\frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right]$$

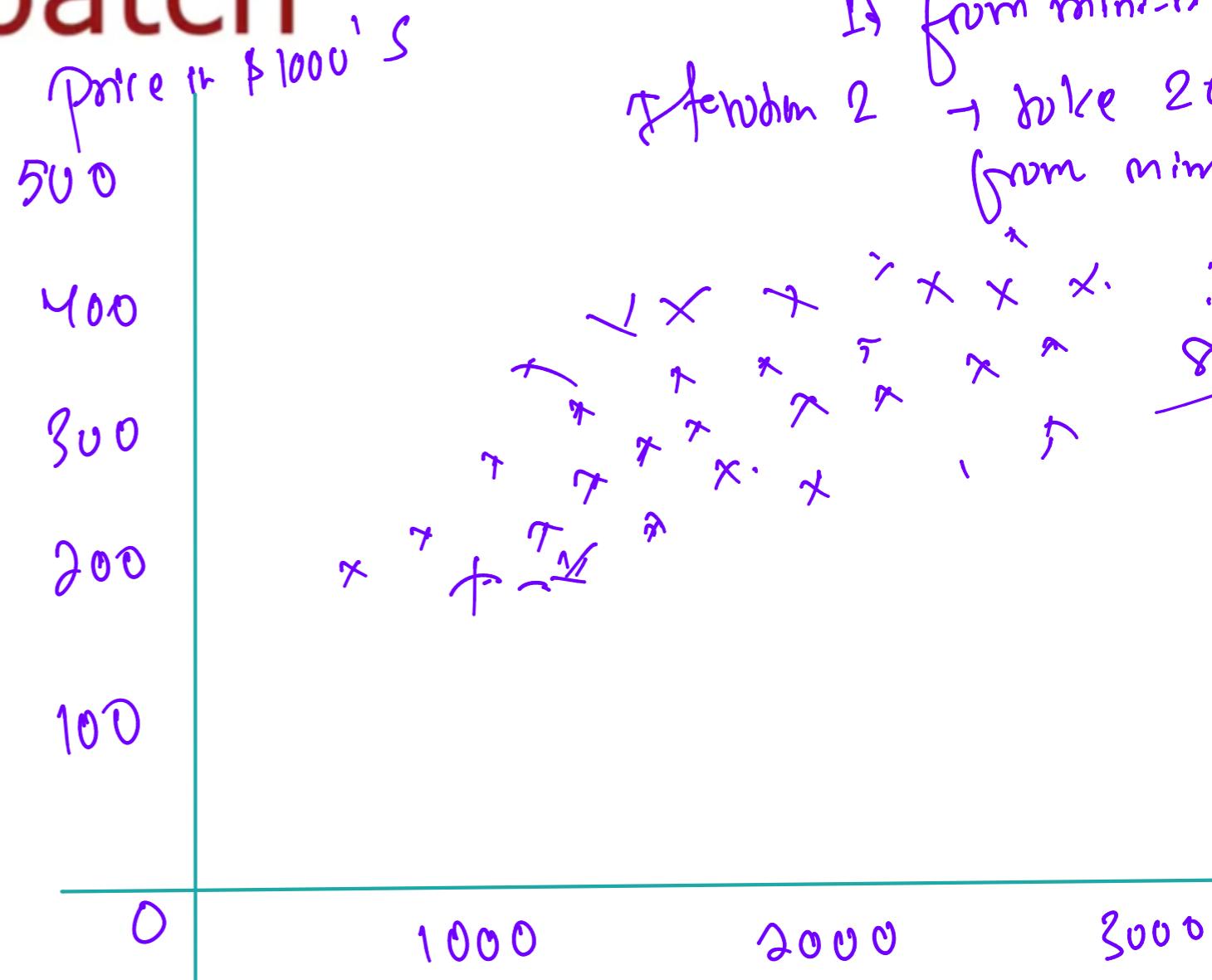
$$b = b - \alpha \frac{\partial}{\partial b} \left[\frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right]$$

}

Mini-batch

x	y
2104	400
1416	232
1534	315
852	178
...	...
3210	870

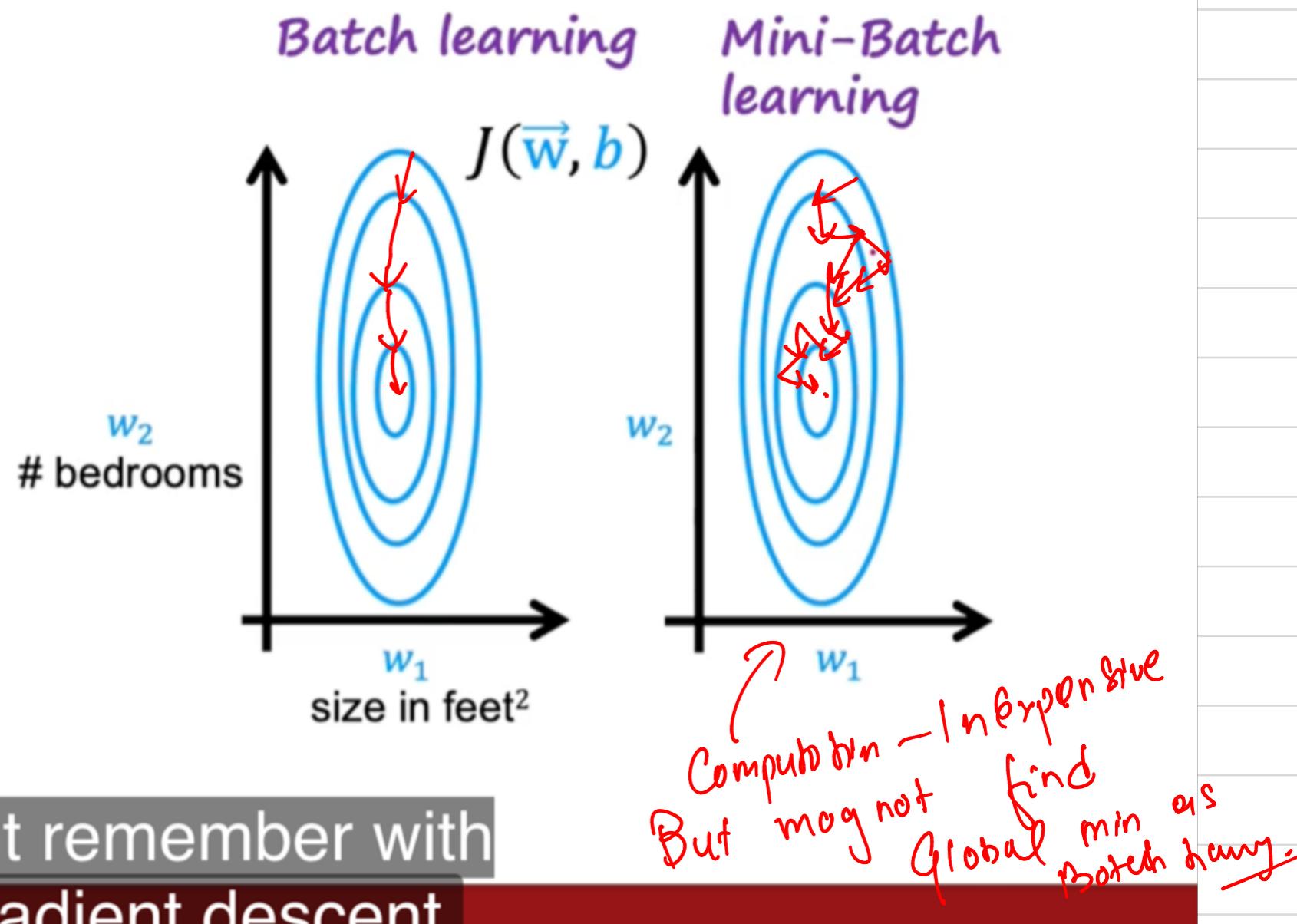
mini-batch 1
mini-batch 2
mini-batch 3



so that every iteration
is looking at

Mini-batch

x	y
2104	400
1416	232
1534	315
852	178
...	...
3210	870



Take Subset Not All (m) $r \in R$ Subsets $\emptyset \rightarrow$ May be Noisy

Learning Algorithm

Initialize neural network randomly as guess of $Q(s, a)$.

Repeat {

Take actions in the lunar lander. Get $(s, a, R(s), s')$.

Store 10,000 most recent $(s, a, R(s), s')$ tuples.

Train model:

Create training set of 10,000 examples using

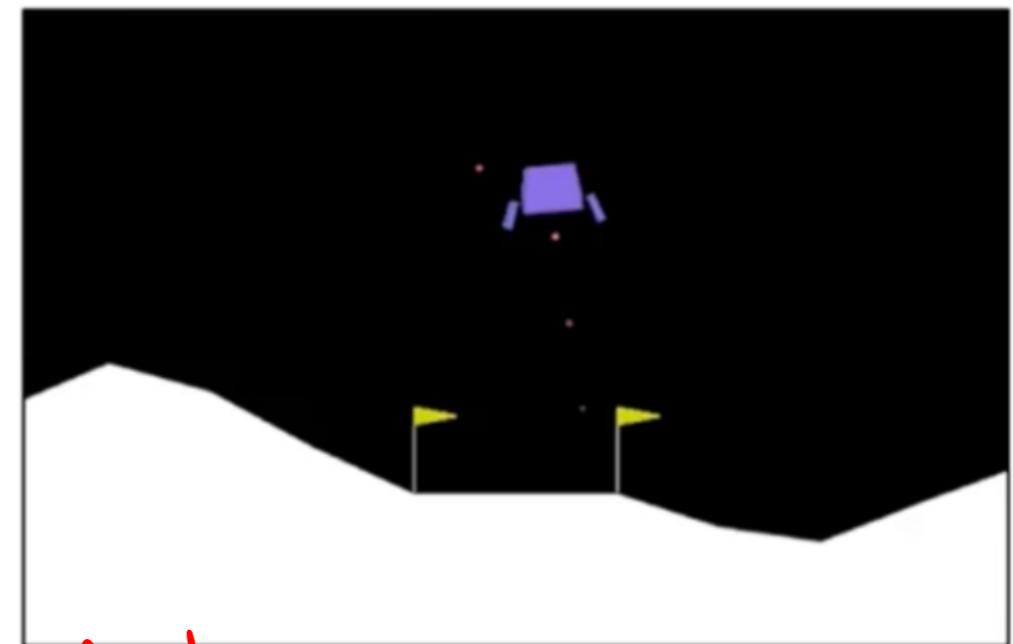
$$x = (s, a) \text{ and } y = R(s) + \gamma \max_{a'} Q(s', a').$$

Train Q_{new} such that $Q_{new}(s, a) \approx y$.

Set $Q = Q_{new}$.

recent tuples in
the replay buffer,

Moy
potentiel
Noisy



Replay Buffer $x^{(1)}, y^{(1)}$

$x^{(1000)}, y^{(1000)}$

Helps to Converge faster

Soft Update

Set $Q = Q_{\text{new}}$.

w, β

\uparrow

$w_{\text{now}}, \beta_{\text{now}}$

$$w = 0.01 w_{\text{new}} + 0.99 w$$

$$\beta = \beta_{\text{new}}$$

↳ Accept only a diff' bit value

$$\boxed{w = 0.01 w_{\text{now}} + 0.99 w}$$

$$\boxed{\beta = 0.01 \beta_{\text{now}} + 0.99 \beta}$$

$$\boxed{w = 2 w_{\text{new}} + 0 w}$$

$Q(S, a) \rightarrow$ works well applying Soft update for DNNs

Simba

When you train the
new neural network,

Steps - of RL

Limitations of Reinforcement Learning

- Much easier to get work in a simulation than a real robot!
- far fewer applications than SL and UCL. (Unsupervised learning)
- But -- exciting research direction with potential for future applications.

RL is sometimes
Hype but still
long way to go
(2025-03-03)

One of the reasons for some of the hype
about reinforcement learning is, it turns

This lecture introduces two key refinements to the reinforcement learning algorithm for Deep Q-Networks (DQN):

1. Mini-batch Gradient Descent – Speeds up training by using a subset of data.
2. Soft Updates – Stabilizes learning by gradually updating the Q-network.

1. Mini-Batch Gradient Descent in RL

Why Mini-Batches?

- Batch Gradient Descent: Computes gradients using all training examples.
- Issue: If the dataset is very large (e.g., 100M examples), every step is computationally expensive.
- Mini-Batch Gradient Descent: Uses a subset m' of training examples (e.g., 1,000).
 - Each step is much faster and updates the model with less computation.
 - Works well for supervised learning (e.g., linear regression, neural networks) and reinforcement learning.

How It Works in Reinforcement Learning?

- Store 10,000 tuples $(s, a, R(s), s')$ in a Replay Buffer.
- Instead of training on all 10,000 experiences, use a mini-batch of 1,000 samples.
- Update the Q-function using:

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

- This introduces some noise but significantly speeds up learning.

Nooks to
be reviewed

2. Soft Updates for Stable Learning

Why Soft Updates?

- In the original Q-update step, we set:

2. Soft Updates for Stable Learning

Why Soft Updates?

- In the original Q-update step, we set:

$$Q \leftarrow Q_{\text{new}}$$

- If the new network Q_{new} is worse due to a bad update, it can degrade learning.
- Instead, soft updates gradually adjust Q using:

$$W \leftarrow \tau W_{\text{new}} + (1 - \tau)W$$

$$B \leftarrow \tau B_{\text{new}} + (1 - \tau)B$$

where τ is a small value (e.g., 0.01).

Effect of Soft Updates

- Reduces oscillations and prevents instability.
- Allows gradual improvements instead of abrupt changes.
- Ensures smooth convergence in reinforcement learning.

Conclusion

- Mini-batch learning improves speed by training on small subsets.
- Soft updates make training more stable.
- These refinements enhance DQN performance and help the Lunar Lander land safely! 🚀

1. Mini-Batch Gradient Descent in Reinforcement Learning

Problem with Full-Batch Gradient Descent

- In supervised learning, batch gradient descent updates parameters using all m training examples:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

- Where $h(x^{(i)}) = wx^{(i)} + b$.
- The gradient descent updates are:

$$w := w - \alpha \frac{\partial J}{\partial w}$$

$$b := b - \alpha \frac{\partial J}{\partial b}$$

- Issue:** For large datasets (e.g., $m = 100M$), computing gradients for all examples in every step is computationally expensive.

Mini-Batch Gradient Descent Solution

- Instead of using all m samples, pick a mini-batch of m' (e.g., 1,000).
- The gradient updates become:

$$J_{\text{mini}}(w, b) = \frac{1}{2m'} \sum_{i=1}^{m'} (h(x^{(i)}) - y^{(i)})^2$$

- New update rules:

$$w := w - \alpha \frac{\partial J_{\text{mini}}}{\partial w}$$

$$b := b - \alpha \frac{\partial J_{\text{mini}}}{\partial b}$$

- Advantage:** Faster iterations with some noise, but still converges to the minimum.

Mini-Batch Learning in Reinforcement Learning

- In Deep Q-Networks (DQN), store past experiences (s, a, r, s') in a Replay Buffer.

- Advantage:** Faster iterations with some noise, but still converges to the minimum.

Mini-Batch Learning in Reinforcement Learning

- In Deep Q-Networks (DQN), store past experiences (s, a, r, s') in a **Replay Buffer**.
- Instead of training on the entire buffer, sample a **mini-batch** of m' experiences.
- Use **Q-learning update rule** for each mini-batch sample:

$$y^i = r^i + \gamma \max_{a'} Q(s'^i, a')$$

- Update the **Q-function** using gradient descent:

$$\theta := \theta - \alpha \nabla_{\theta} \sum_{i=1}^{m'} (y^i - Q(s^i, a^i; \theta))^2$$

- Reduces computational cost while still improving the policy.



2. Soft Updates for Stable Learning

Problem with Direct Q-Function Updates

- Original Q-learning update directly replaces old weights:

$$W \leftarrow W_{\text{new}}, \quad B \leftarrow B_{\text{new}}$$

- Issue:** If W_{new} is poorly trained, it may degrade performance.

Soft Update Solution

- Instead of full replacement, use a weighted update:

$$W \leftarrow \tau W_{\text{new}} + (1 - \tau)W$$

$$B \leftarrow \tau B_{\text{new}} + (1 - \tau)B$$

- Where τ (e.g., 0.01) controls how much of the new parameters are used.
- Ensures **smooth transitions** and prevents instability.

Acknowledgment: I want to thank Coursera (financial aid) for providing valuable materials.

Date of 100DaysOfMaths Completion: 2025|03|10 (2082|11|26 B.S)

Resources: ① B.Sc CSIT 1st and 2nd Semesters (KEC publication Maths Books) → 2078 B.S.

② Coursera - DeepLearning-AI Courses

(Statistics, Linear Algebra, Calculus, Machine Learning)

③ ChatGPT, DopeSook, Claude, Gemini (Coding & Additional Information)

(3 Courses)

Completed Courses: Coursera → Mathematics for Machine learning and Data Science

Coursera → Machine learning (3 Courses)

Comprehensive learning, Basic Strong Concepts, fundamental skills do works, and many more

More info? Visit: github.com/DilliF22/100DaysOfMaths