

Reinforcement Learning in Control

Dr. Saeed Shamaghdari

**Electrical Engineering Department
Control Group**

Fall 2025 | 4041

Deep Reinforcement Learning

Playing Atari with Deep Reinforcement Learning

Volodymyr Mnih Koray Kavukcuoglu David Silver Alex Graves Ioannis Antonoglou

Daan Wierstra Martin Riedmiller

DeepMind Technologies

`{vlad, koray, david, alex.graves, ioannis, daan, martin.riedmiller}` @ deepmind.com

☆ Save  Cite Cited by 19338 Related articles All 37 versions 

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

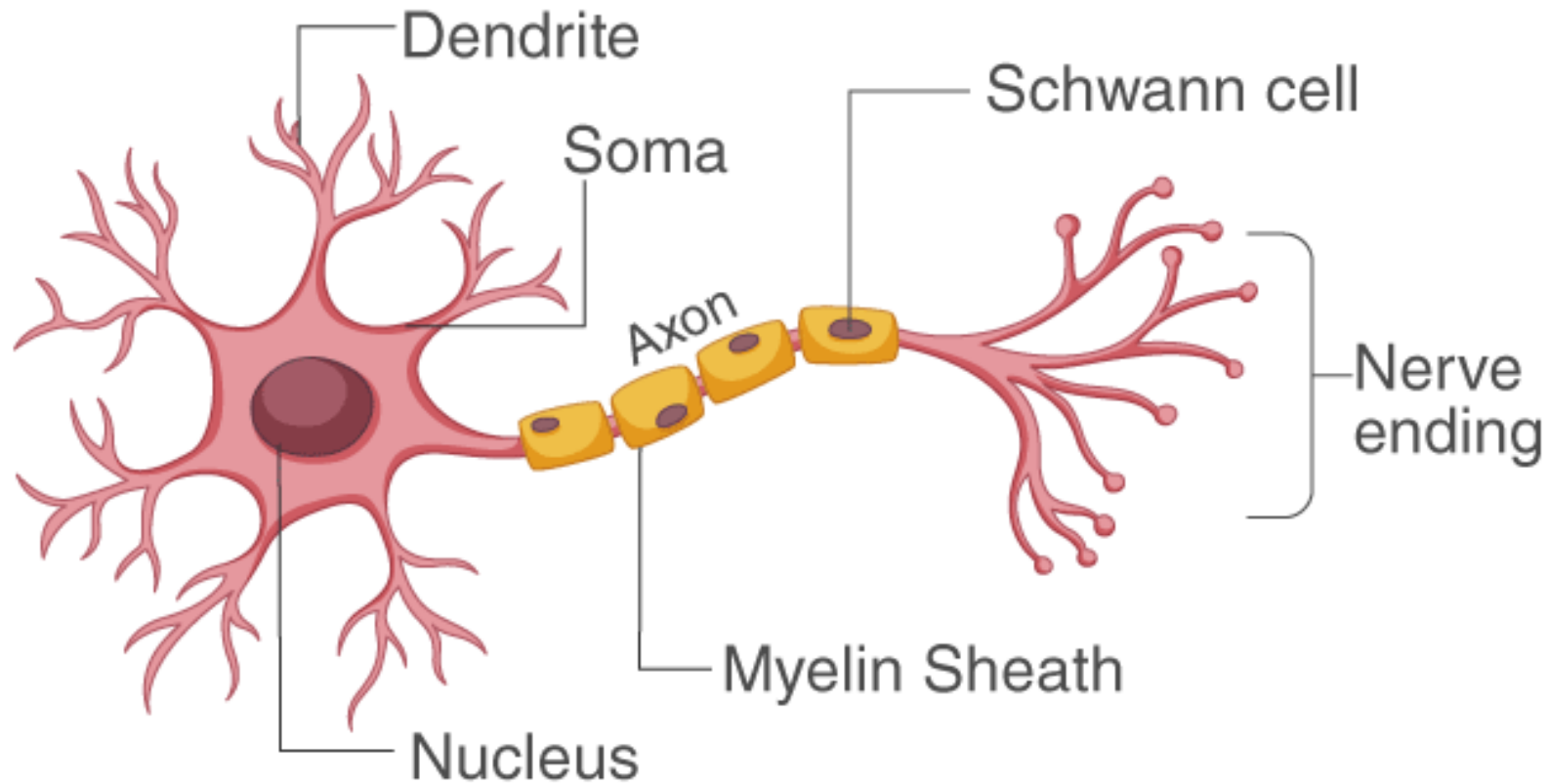
¹Google DeepMind, 5 New Street Square, London EC4A 3TW, UK.

*These authors contributed equally to this work.

26 FEBRUARY 2015 | VOL 518 | NATURE | 529

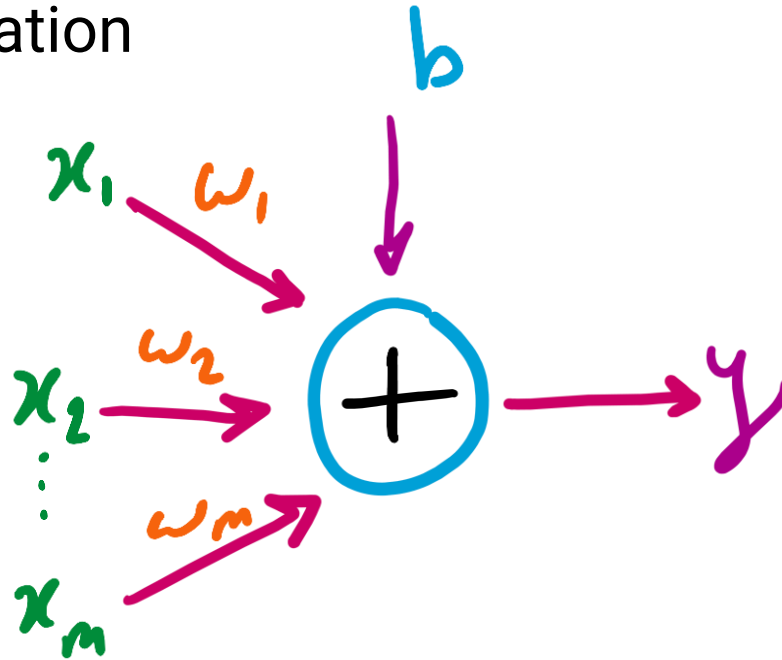
Neural Networks and Deep Learning: A Simple Review

Neuron



Neural Networks and Deep Learning: A Simple Review

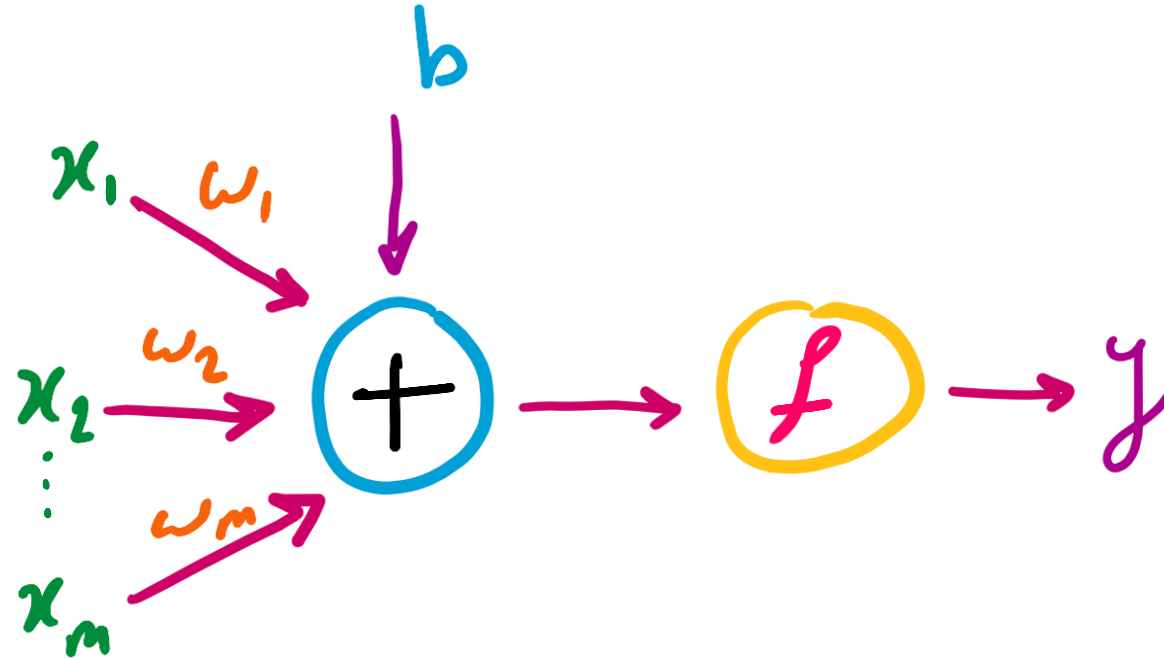
Artificial Neuron Formulation



$$y = w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b$$

Neural Networks and Deep Learning: A Simple Review

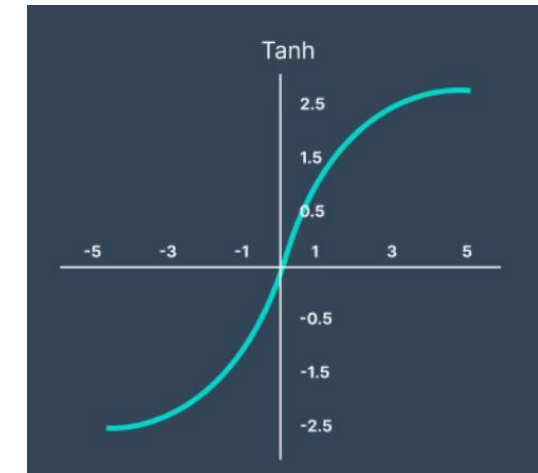
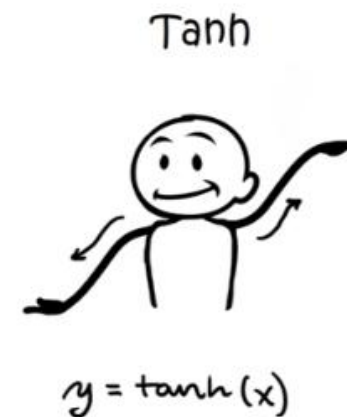
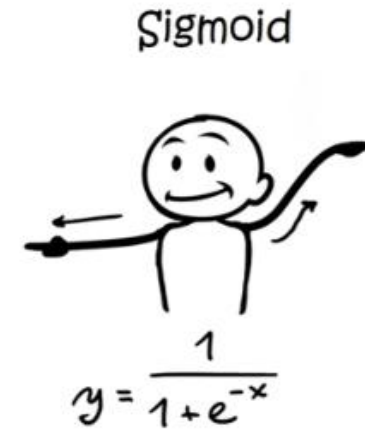
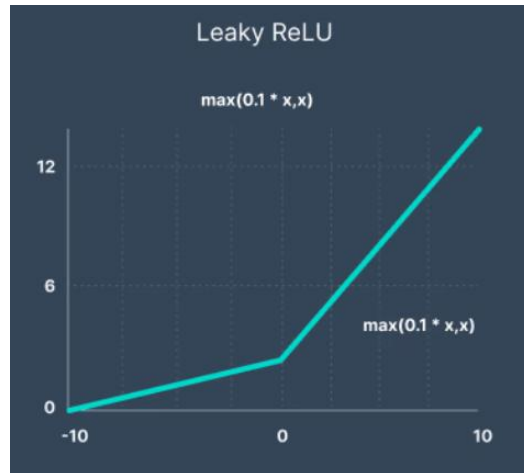
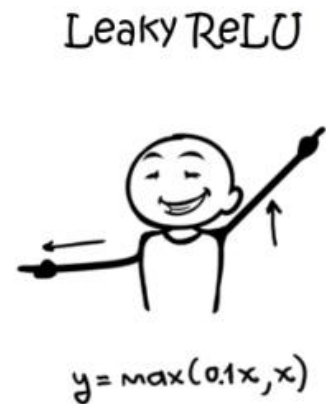
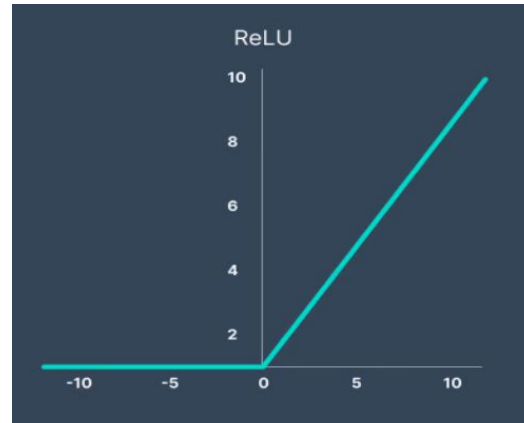
Q: Problems?



$$y = f(w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b) = f(\omega^T x + b)$$

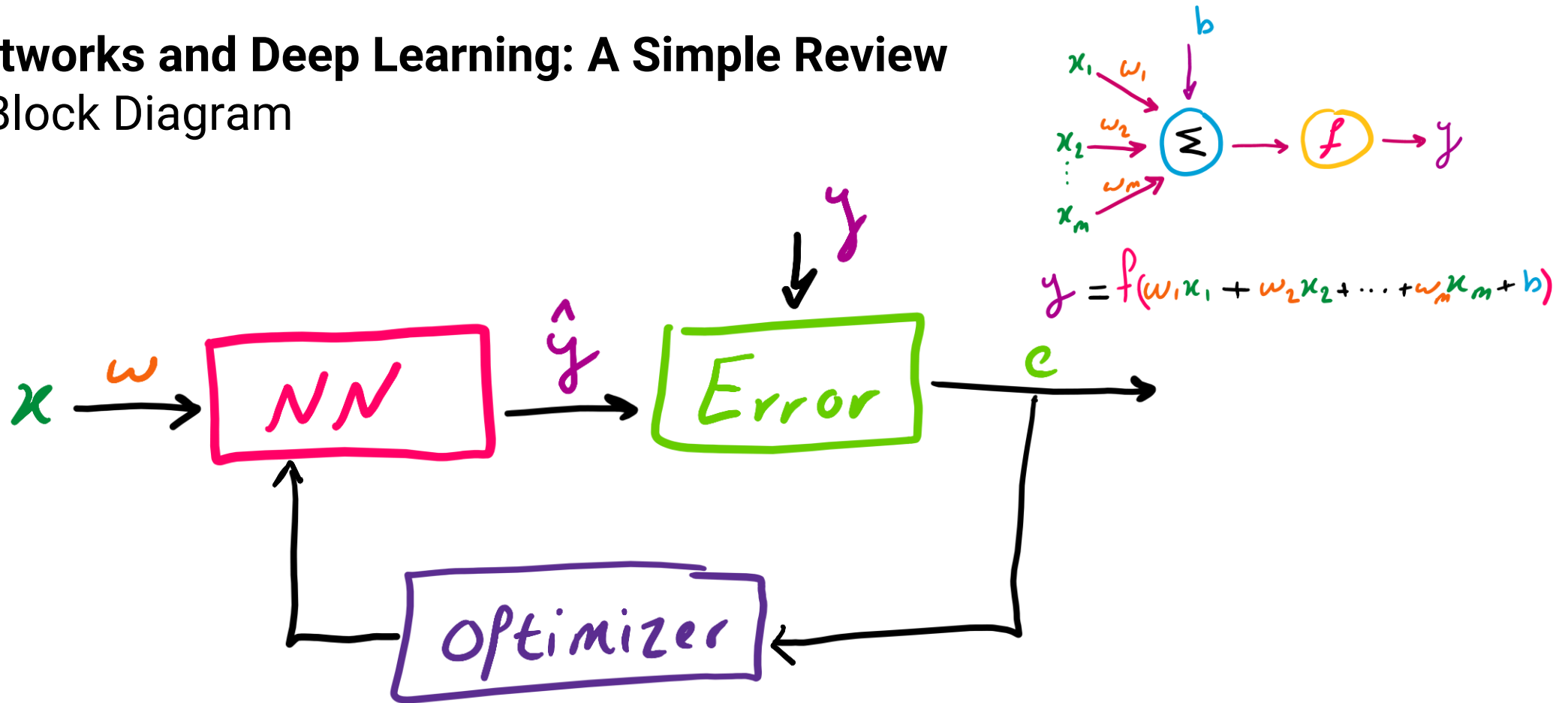
Neural Networks and Deep Learning: A Simple Review

Activation Functions



Neural Networks and Deep Learning: A Simple Review

Learning Block Diagram



Neural Networks and Deep Learning: A Simple Review

Learning Block Diagram

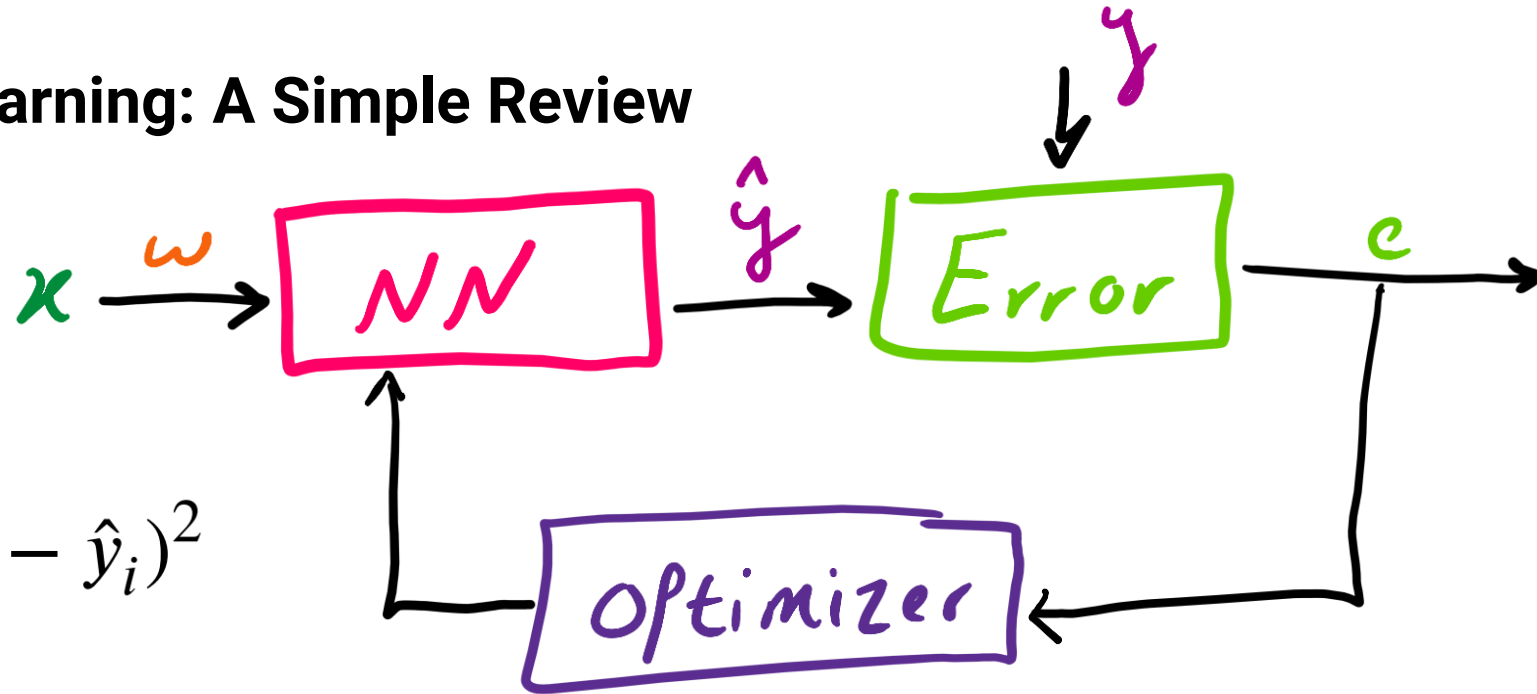
Loss Function:

MSE:

$$L(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

MAE:

$$L(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$



Neural Networks and Deep Learning: A Simple Review

Learning Block Diagram

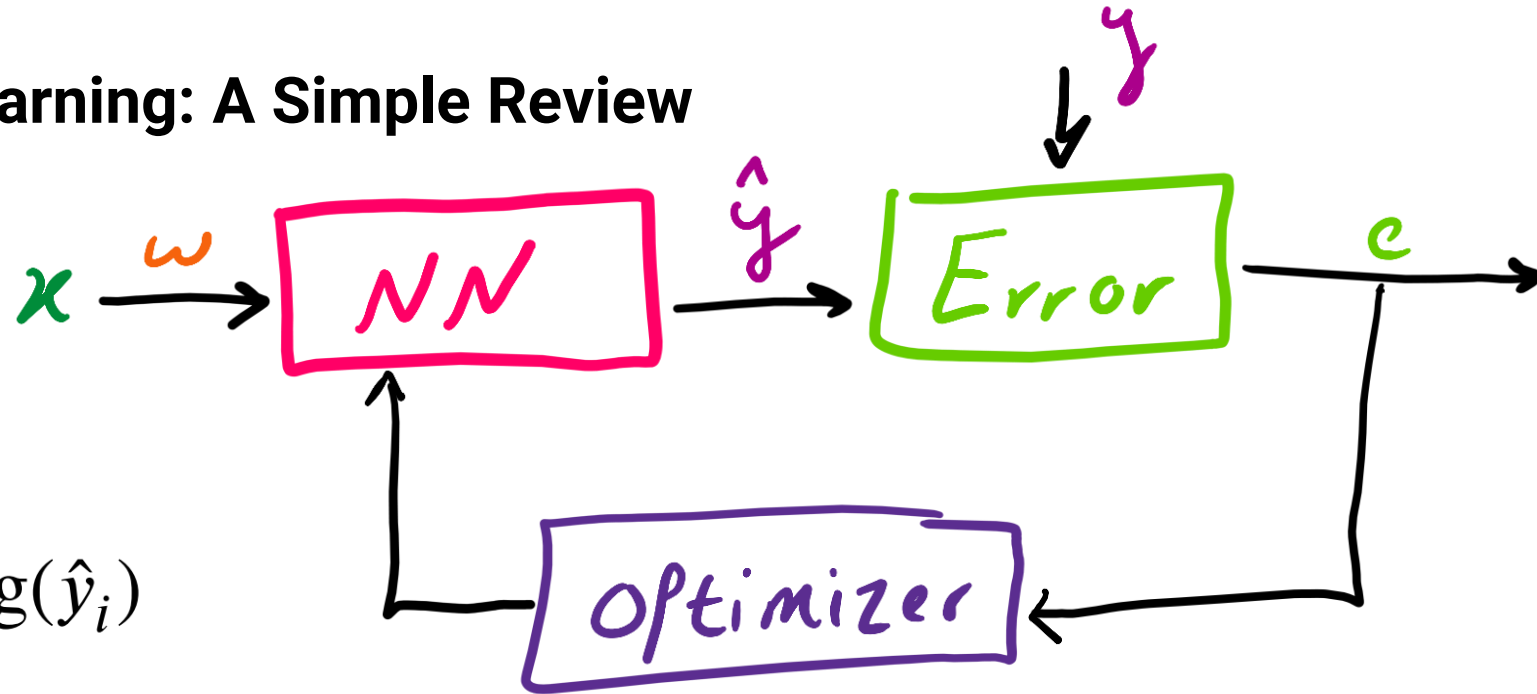
Loss Function:

Cross Entropy:

$$L(y, \hat{y}) = - \sum_{i=1}^N y_i \log(\hat{y}_i)$$

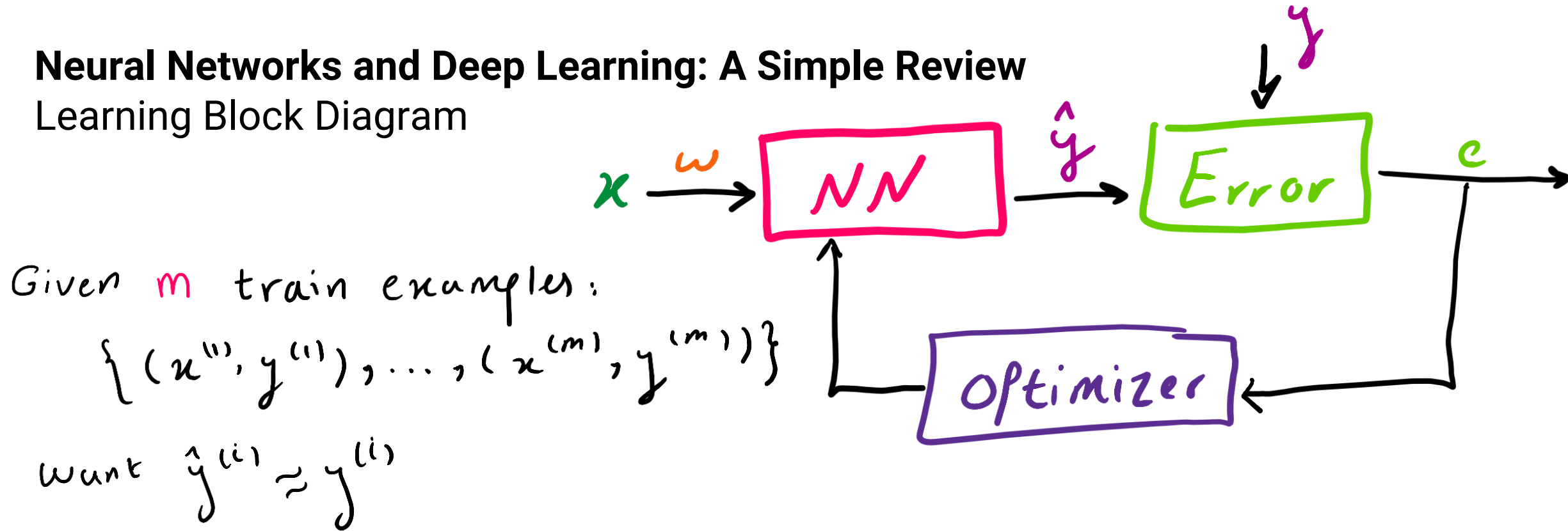
Binary Cross Entropy:

$$L(y, \hat{y}) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$$



Neural Networks and Deep Learning: A Simple Review

Learning Block Diagram



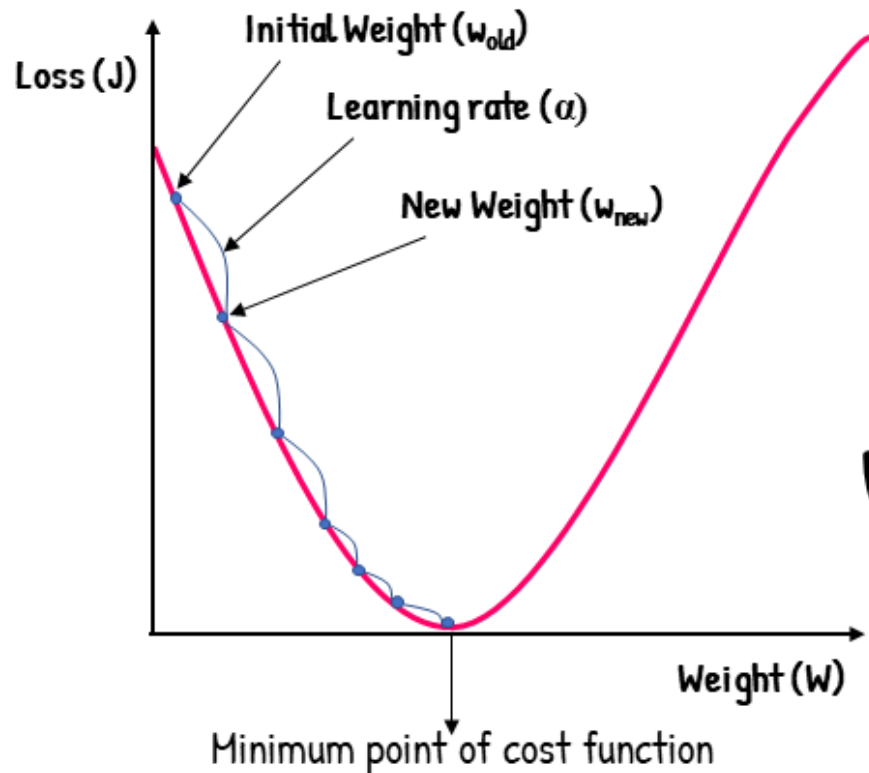
So the cost function is:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(y, \hat{y})$$

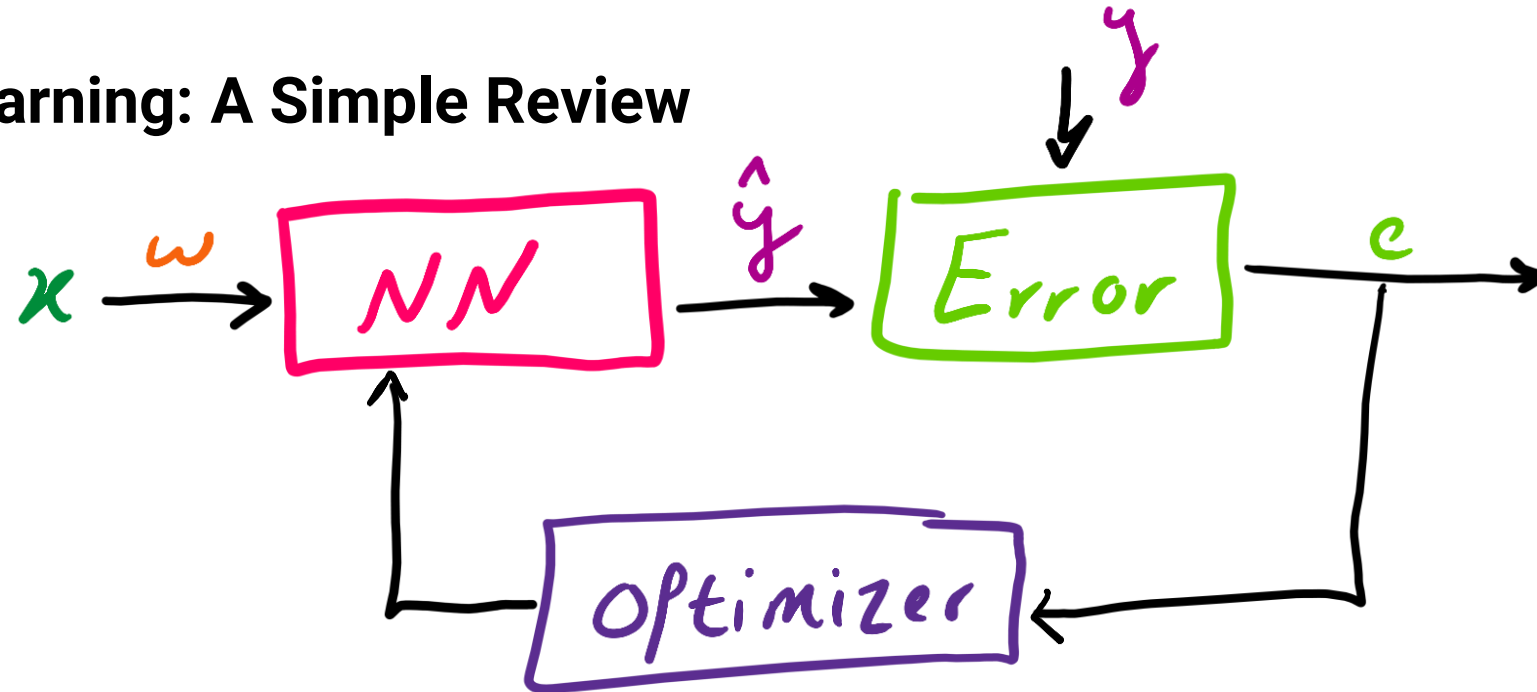
Neural Networks and Deep Learning: A Simple Review

Learning Block Diagram

Gradient Descent:

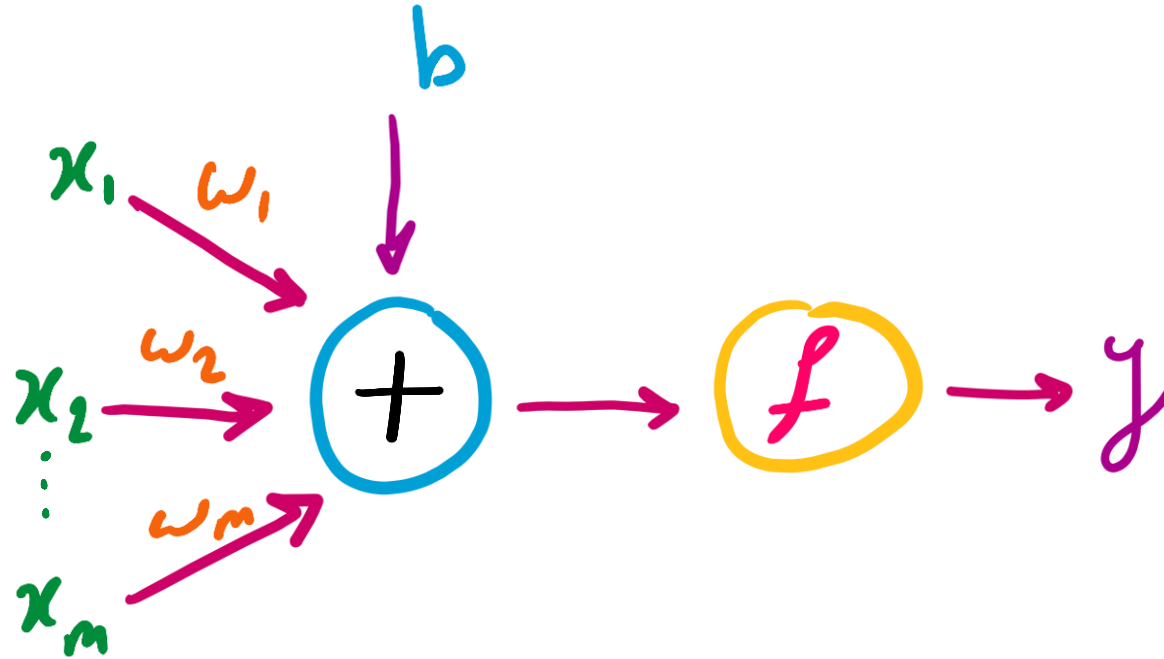


$$w_{new} = w_{old} - \alpha \frac{\delta J}{\delta w}$$



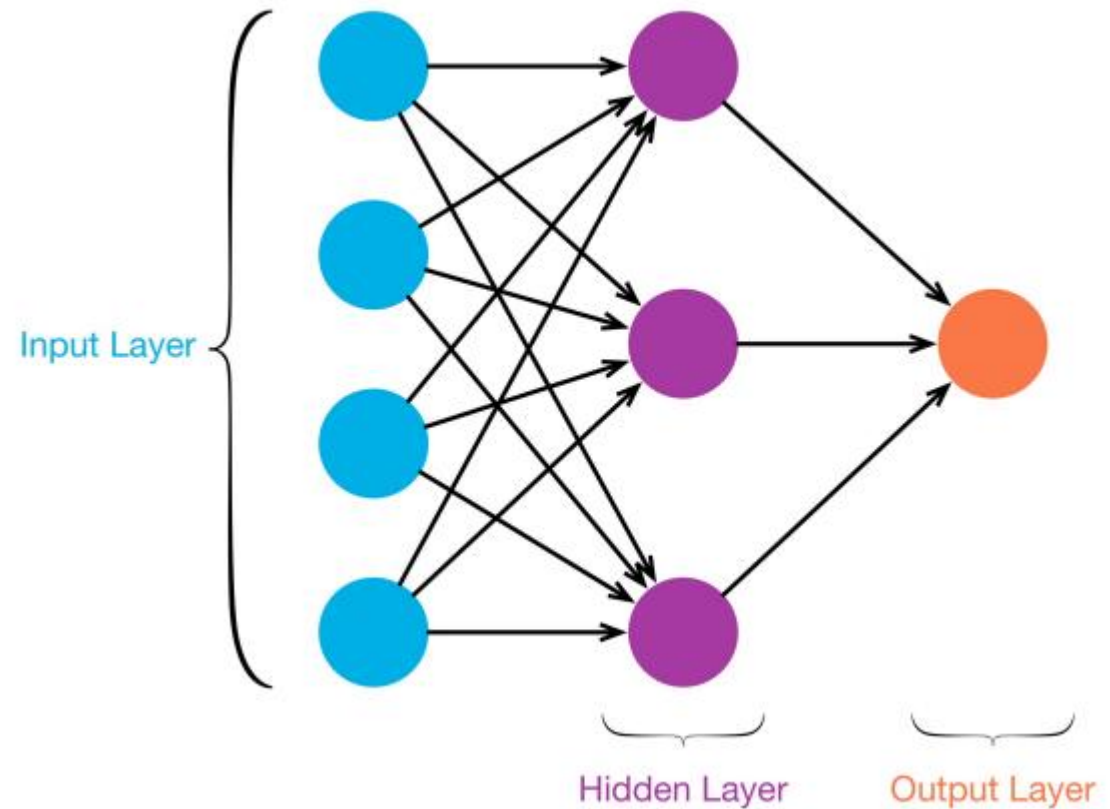
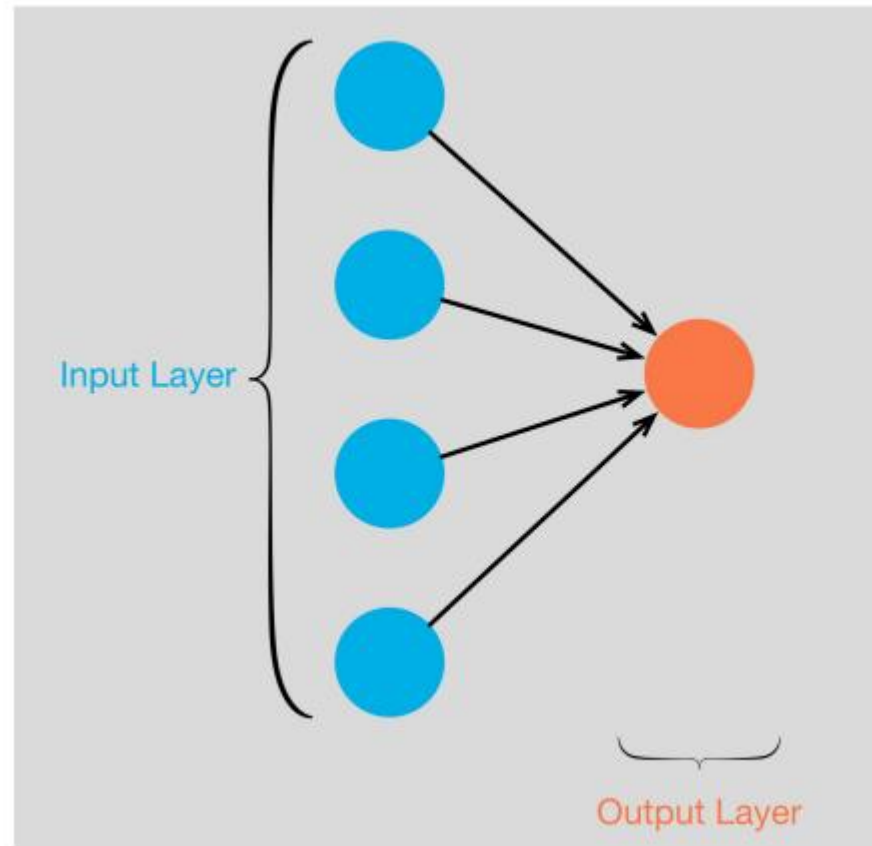
Neural Networks and Deep Learning: A Simple Review

Single-Layer Perceptron

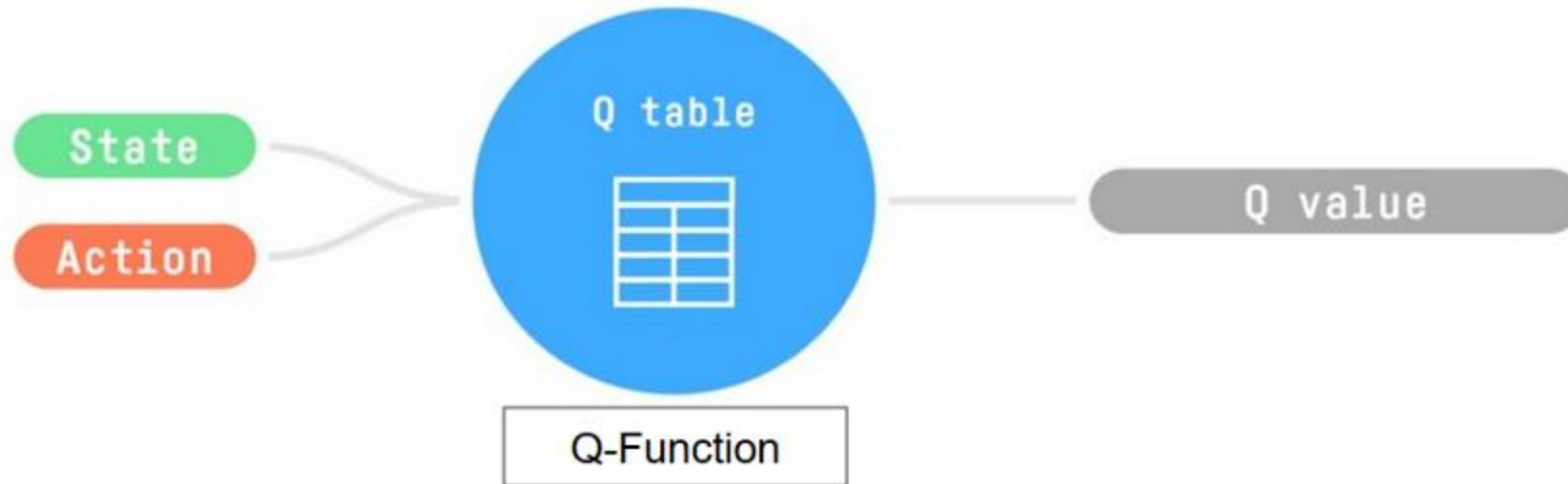


Neural Networks and Deep Learning: A Simple Review

Multi-Layer Perceptron ([MLP](#)) ([Play!](#))



Q-Learning Recap ...



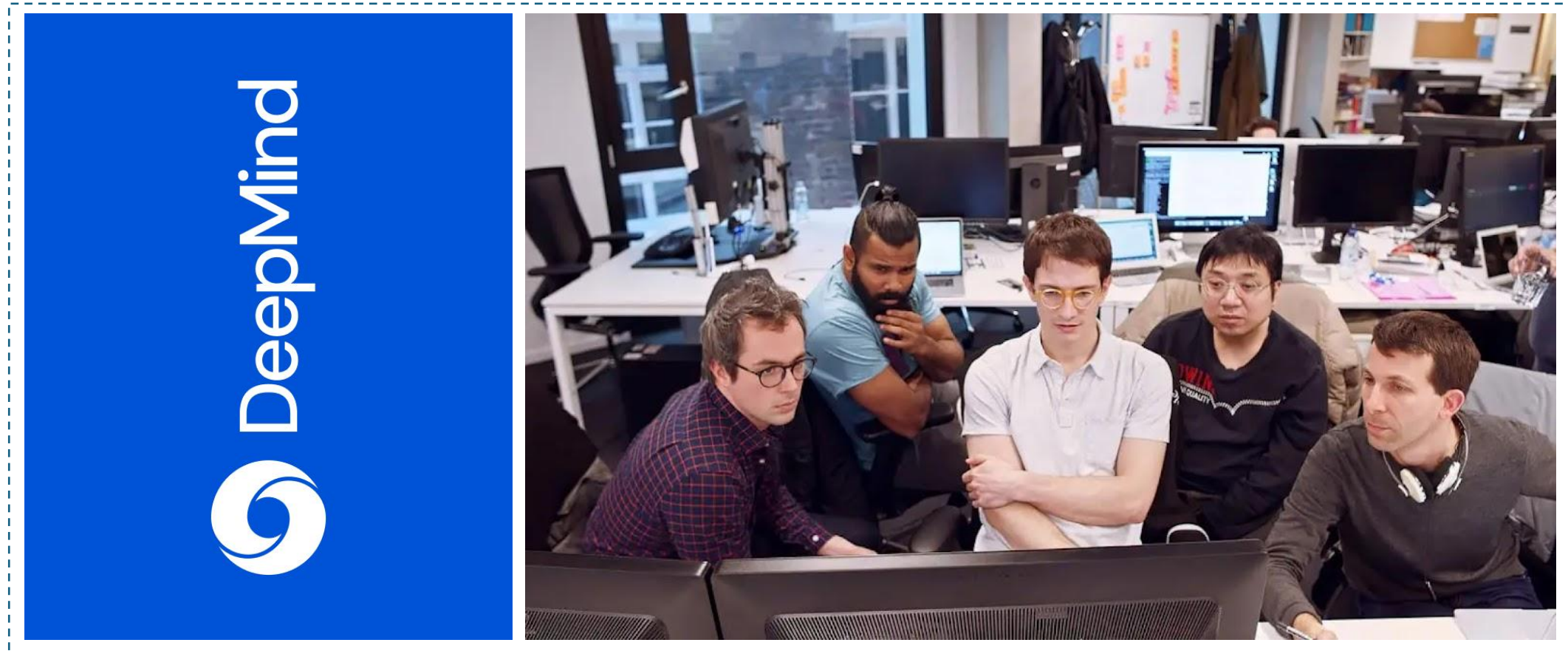
Pseudocode

Q-Learning

Algorithm 14: Sarsamax (Q-Learning)

Input: policy π , positive integer $num_episodes$, small positive fraction α , GLIE $\{\epsilon_i\}$
Output: value function Q ($\approx q_\pi$ if $num_episodes$ is large enough)
Initialize Q arbitrarily (e.g., $Q(s, a) = 0$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$, and $Q(\text{terminal-state}, \cdot) = 0$)
for $i \leftarrow 1$ **to** $num_episodes$ **do** ↖ Step 1
 $\epsilon \leftarrow \epsilon_i$
 Observe S_0
 $t \leftarrow 0$
 repeat
 Choose action A_t using policy derived from Q (e.g., ϵ -greedy) Step 2
 Take action A_t and observe R_{t+1}, S_{t+1} Step 3
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$ Step 4
 $t \leftarrow t + 1$
 until S_t is terminal;
end
return Q

Why **Deep** Reinforcement Learning?



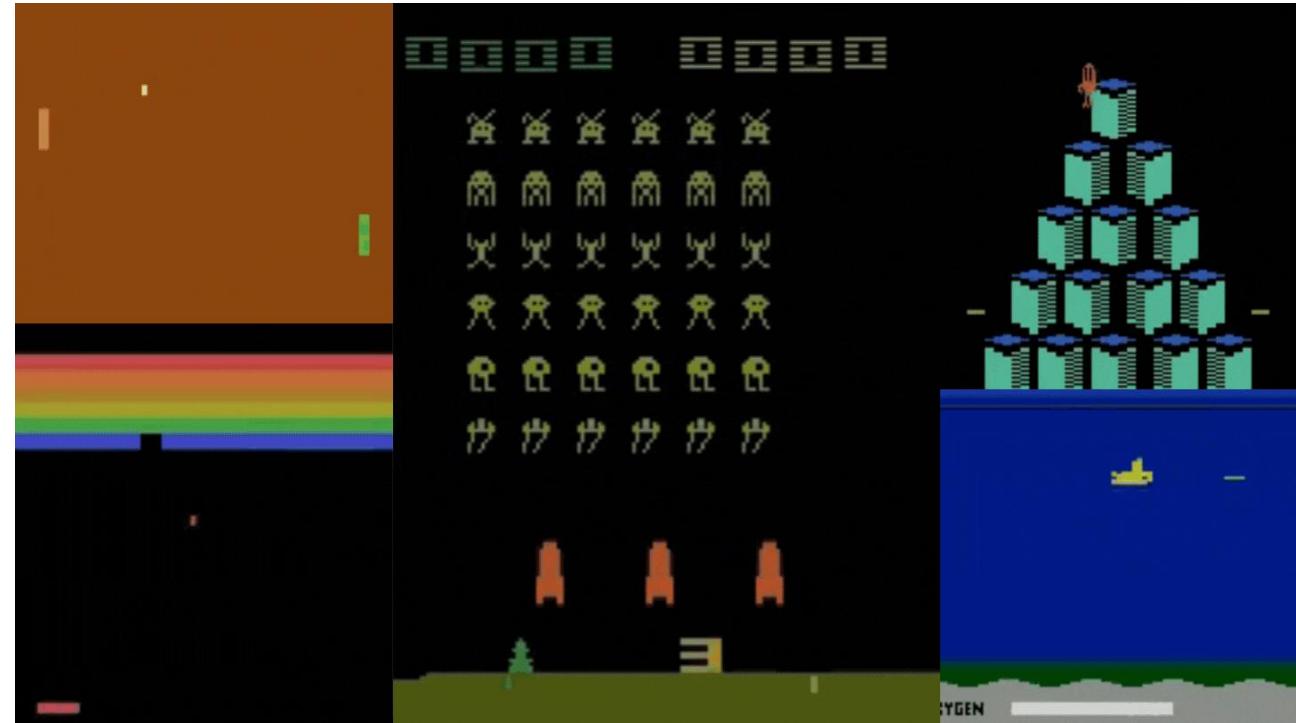
Playing Atari with Deep Reinforcement Learning

([Paper](#))

Why **Deep** Reinforcement Learning?

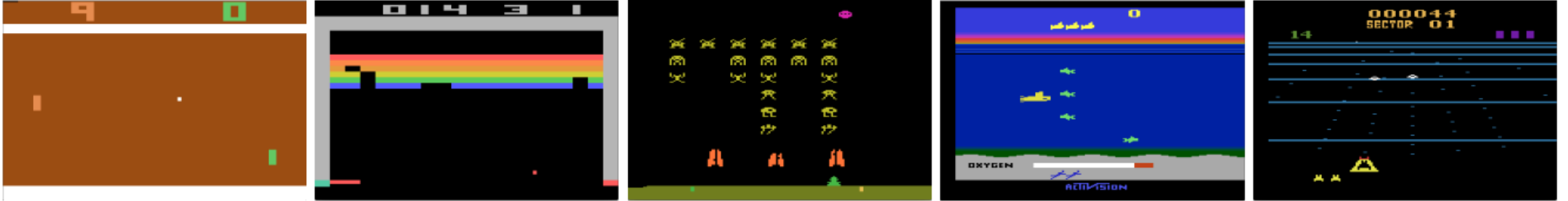
چالش جدی RL: کنترل مبتنی بر یادگیری مستقیماً از دیتای با ابعاد بزرگ (تصویر یا صوت یا سنسورها)

Q: Number of states in an 8*8
Gridworld?
What about an **Atari** game?



Playing Atari with Deep Reinforcement Learning ([Paper](#))

Why **Deep** Reinforcement Learning?

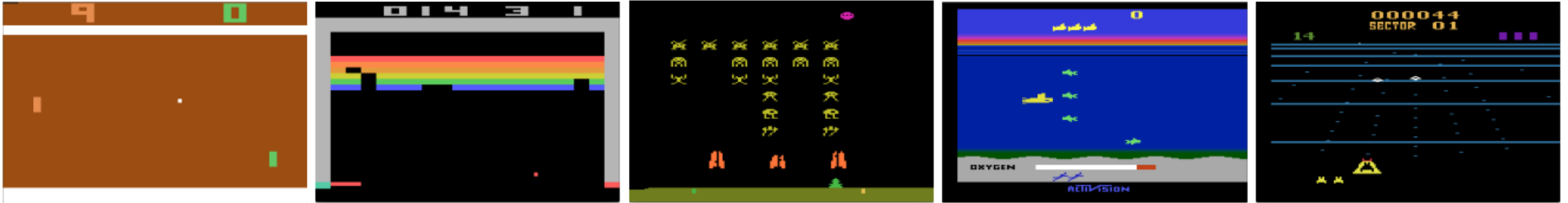


(Left-to-right) Pong, Breakout, Space Invaders, Seaquest, Beam Rider

*“Atari 2600 :
visual input (210×160 RGB video at 60Hz)*

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Why **Deep** Reinforcement Learning?



(Left-to-right) Pong, Breakout, Space Invaders, Seaquest, Beam Rider

Each Frame: (210, 160, 3) containing values ranging from 0 to 255

Q: Number of states?

A: $256^{210 \times 160 \times 3} = 256^{100800}$

Idea: approximate Q-values

Using parametrized Q-function $Q_{\theta}(s, a)$

Playing Atari with Deep Reinforcement Learning ([Paper](#))

DeepRL: Why Is It Still Hard?

- روش های موفق Deep Learning: دیتاهای عظیم لیبل دار
- یادگیری RL از پاداش اسکالربه صورت sparse و نویزی و تاخیر دار (برعکس DL)

Playing Atari with Deep Reinforcement Learning ([Paper](#))

DeepRL: Why Is It Still Hard?

DL data samples: IID

RL: Highly correlated states.

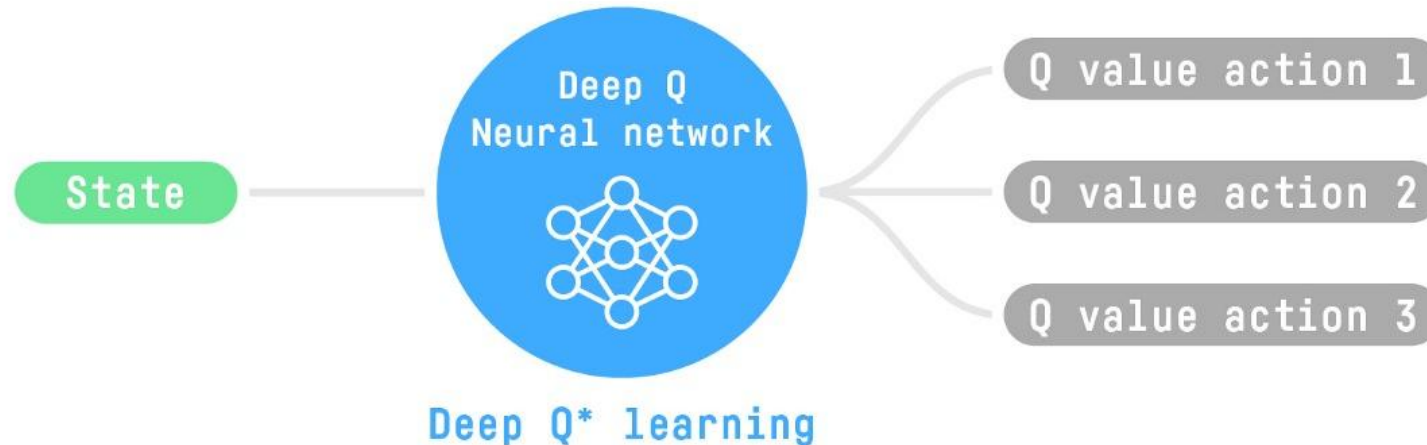
تغییر توزیع دیتا در RL در فرایند یادگیری

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Deep Q-Learning (DQN)

The best idea is to approximate the Q-values using a parametrized Q-function $Q_{\theta}(s, a)$.

مشابه DL:



Playing Atari with Deep Reinforcement Learning ([Paper](#))

Deep Q-Learning (DQN)

Emulator: agent interacts with an environment \mathcal{E}

Goal: maximizes future rewards.

The future discounted *return* at time t :

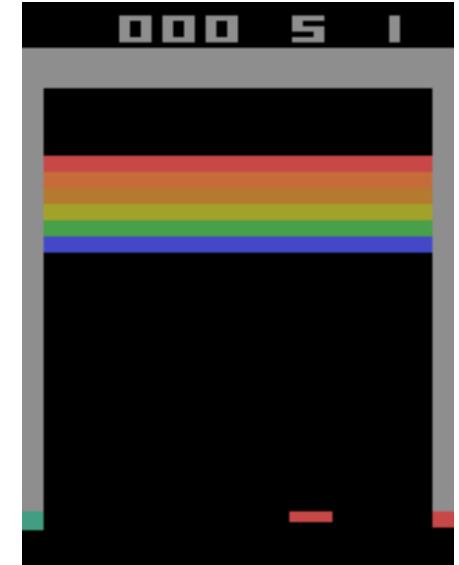
$$R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$$

r_t : reward

The optimal action-value function:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} [R_t | s_t = s, a_t = a, \pi]$$

Playing Atari with Deep Reinforcement Learning ([Paper](#))



Deep Q-Learning (DQN)

optimal strategy :

select a' :

maximizing the expected value of $r + \gamma Q^*(s', a')$

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

Reminder Box: Value Iteration

$$Q_{i+1}(s, a) = \mathbb{E} [r + \gamma \max_{a'} Q_i(s', a') \mid s, a]$$
$$Q_i \rightarrow Q^* \text{ as } i \rightarrow \infty$$

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Deep Q-Learning (DQN)

A function approximator to estimate the action-value function:

$$Q(s, a)$$

A Neural Network!

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Deep Q-Learning (DQN)

*Neural network function approximator with weights θ : **Q-network***

Training:

minimizing

$$L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot)} \left[(y_i - Q(s, a; \theta_i))^2 \right]$$

“Where

target:

$$y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$$

$\rho(s,a)$: *behaviour*

تارگت وابسته به وزنهای شبکه است برعکس Supervised DL که ثابت است

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Deep Q-Learning (DQN)

مشتق تابع Loss نسبت به وزن‌ها:

Reminder Box

$$L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot)} \left[(y_i - Q(s, a; \theta_i))^2 \right]$$

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right) \nabla_{\theta_i} Q(s, a; \theta_i) \right]$$

stochastic gradient descent

model-free and off-policy.

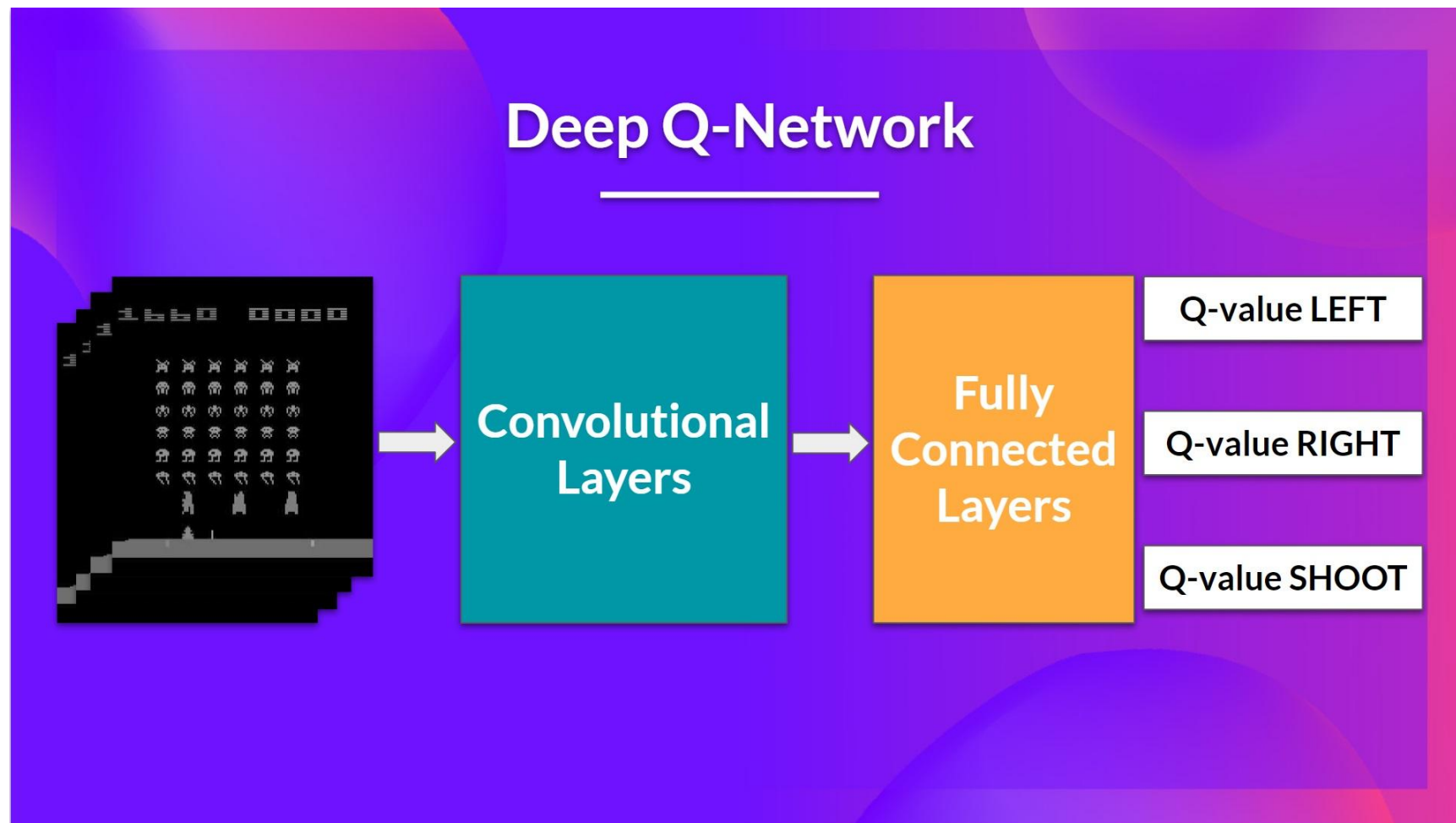
Playing Atari with Deep Reinforcement Learning ([Paper](#))

Architecture

Input: Stack of 4 gray-scale frames

Q: Why?

Output: A vector of Q-values for each possible action at that state



Playing Atari with Deep Reinforcement Learning ([Paper](#))

Architecture: What Is a Convolutional Layer?



Input

Convolutional
Layers

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Architecture: What Is a Convolutional Layer?

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Input Image

1	0	1
0	1	0
1	0	1

Kernel/Filter

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

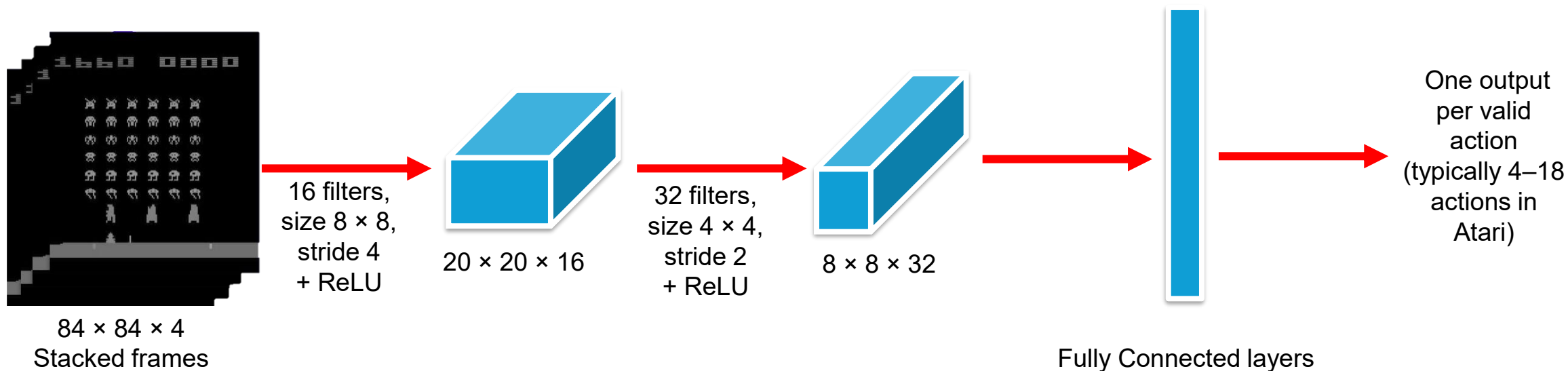
4		

Convolved
Feature

Convolutional
Layers

Playing Atari with Deep Reinforcement Learning ([Paper](#))

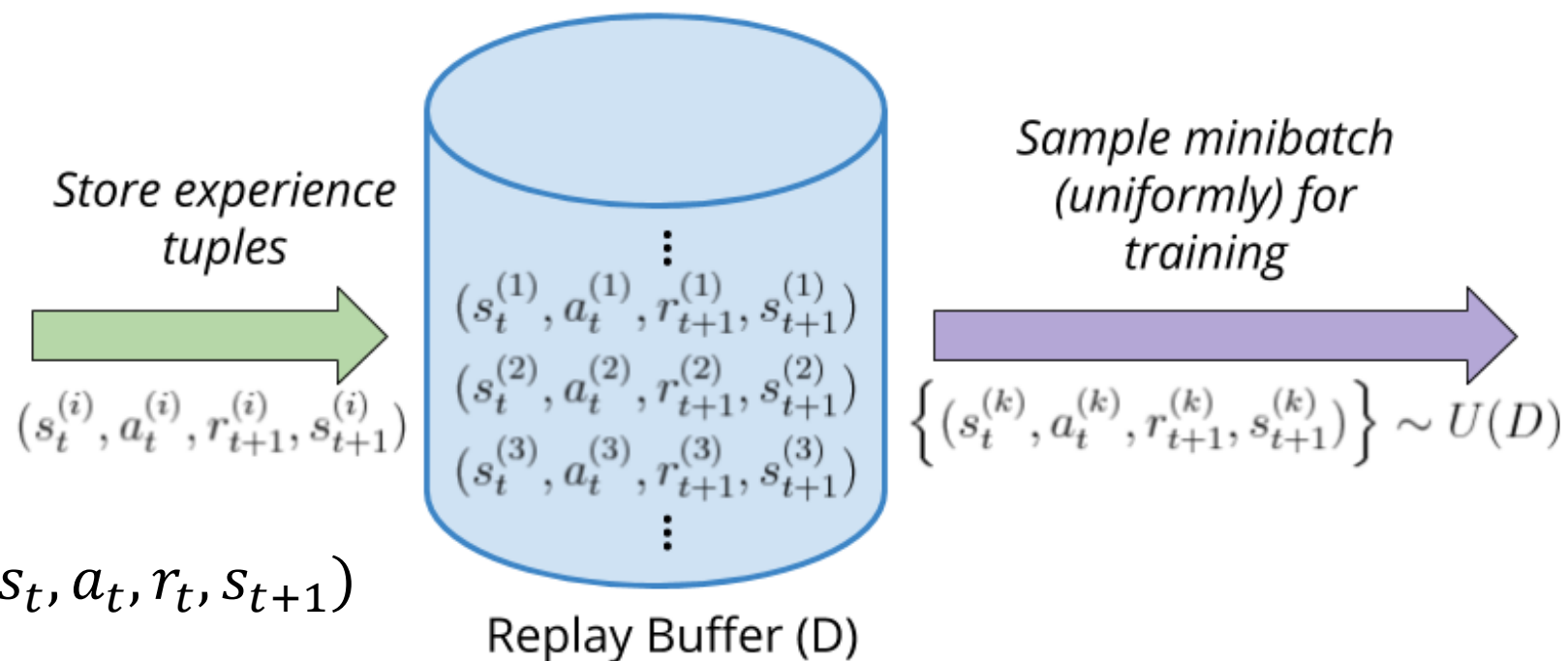
Architecture



This architecture is known as a **Deep Q-Network (DQN)**

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Experience Replay Replay Buffer



Agent's experiences : $e_t = (s_t, a_t, r_t, s_{t+1})$

Data-set $D = e_1, \dots, e_N$

Q -learning updates, or minibatch updates :

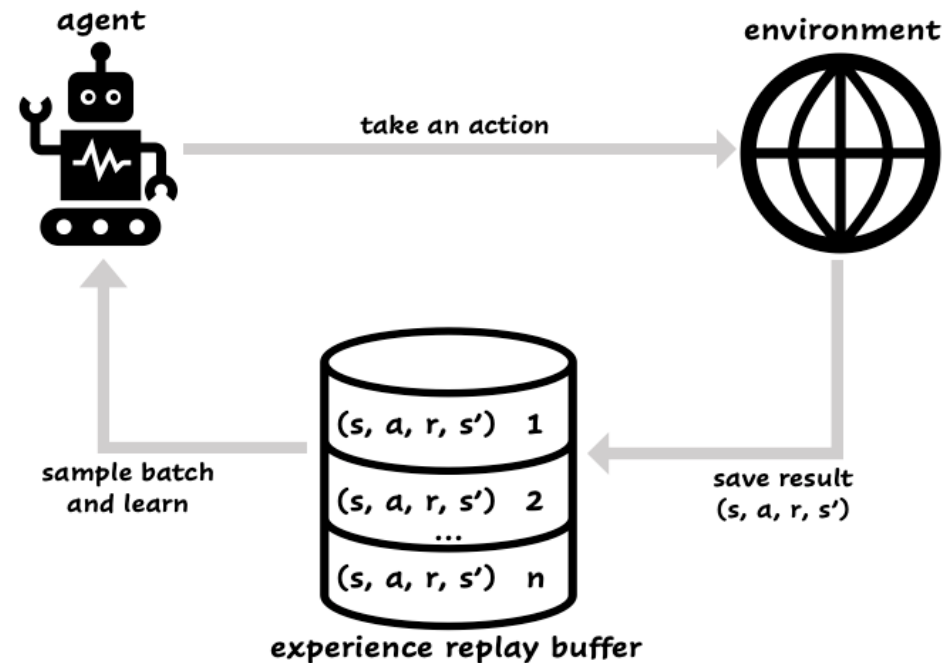
$$\text{Random : } e \sim D$$

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Algorithm: Replay Buffer

Stores the last N experience

Samples uniformly at random from D



Q: Why Replay Buffer?

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Algorithm

Reminder Box

$$\underbrace{Q(S_t, A_t)}_{\text{New Q-value estimation}} \leftarrow \underbrace{Q(S_t, A_t)}_{\text{Former Q-value estimation}} + \underbrace{\alpha}_{\text{Learning Rate}} [\underbrace{R_{t+1}}_{\text{Immediate Reward}} + \underbrace{\gamma \max_a Q(S_{t+1}, a)}_{\text{Discounted Estimate optimal Q-value of next state}} - \underbrace{Q(S_t, A_t)}_{\text{Former Q-value estimation}}]$$

New
Q-value
estimation

Former
Q-value
estimation

Learning
Rate

Immediate
Reward

Discounted Estimate
optimal Q-value
of next state

Former
Q-value
estimation

TD Target

TD Error

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Algorithm

Intuition

Q-Target

$$y_j = r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-)$$

$$\overbrace{R_{t+1}}^{\text{Immediate Reward}} + \overbrace{\gamma \max_a Q(S_{t+1}, a)}^{\text{Discounted Estimate optimal Q-value of next state}}$$

Immediate Reward Discounted Estimate optimal Q-value of next state

TD Target

Q-Loss

$$y_j - Q(\phi_j, a_j; \theta)$$

$$\overbrace{[R_{t+1} + \gamma \max_a Q(S_{t+1}, a)]}^{\text{TD Target}} - \overbrace{Q(S_t, A_t)}^{\text{Former Q-value estimation}}$$

Immediate Reward Discounted Estimate optimal Q-value of next state Former Q-value estimation

TD Target

TD Error

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Algorithm

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

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

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

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

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

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

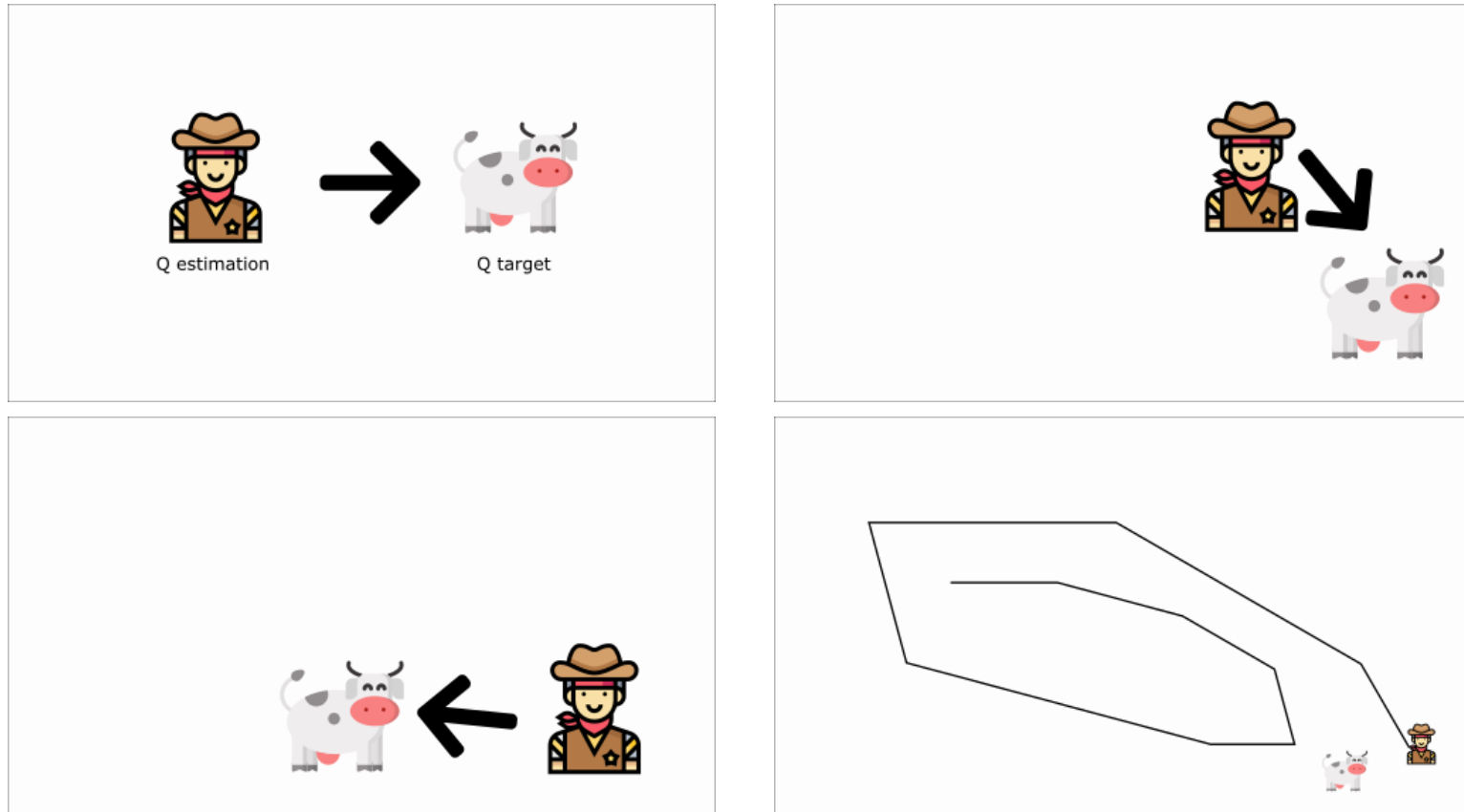
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

Playing Atari with Deep Reinforcement Learning ([Paper](#))

Algorithm: A Challenge!



Playing Atari with Deep Reinforcement Learning ([Paper](#))

Algorithm: A Challenge!

Solution:

- Use a **separate network** with fixed parameters for estimating the TD Target
- Copy the parameters from our Deep Q-Network **every C steps** to update the target network.

Playing Atari with Deep Reinforcement Learning ([Paper](#), [Paper](#))

Algorithm

Algorithm 1: deep Q-learning with experience replay.

 Initialize replay memory D to capacity N

 Initialize action-value function Q with random weights θ

 Initialize target action-value function \hat{Q} with weights $\theta^- = \theta$
For episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
For $t = 1, T$ **do**

 With probability ε select a random action a_t

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

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

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

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

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

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

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

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

 Playing Atari with Deep Reinforcement Learning ([Paper](#), [Paper](#))

Conclusion

“We apply our method to seven Atari 2600 games from the Arcade Learning Environment, with no adjustment of the architecture or learning algorithm. We find that it outperforms all previous approaches on six of the games and surpasses a human expert on three of them, with no adjustment of the architecture or hyperparameters.”

Playing Atari with Deep Reinforcement Learning ([Paper](#), [Paper](#))