

The background features a collection of abstract, colorful geometric shapes. There are large circles in shades of pink, purple, and blue. A bright yellow circle is positioned in the upper right. A green curved line arches over a large purple circle on the right. A pink curved line is located below the green one. In the bottom right, there are overlapping circles in red and pink. A blue circle is partially visible at the bottom center. The overall aesthetic is modern and vibrant.

DQN Reinforcement Learning

Advancing Deep Q-Networks for Optimal Navigation in Frozen Lake:
Evaluating Adaptability and Performance in Dynamic Environments

Aviv Salomon
Ben Gornizky

TABLE OF CONTENTS

01

Introduction

02

Architecture

03

Articles

04

Environment

05

Methodology

06

Project Plans

The background features several large, soft-edged, organic shapes in various colors: a large orange shape in the top left, a blue shape in the top right, a pink shape on the left, a purple shape in the bottom left, and a blue shape in the bottom right. A small pink sphere and a green curved line are also visible near the bottom left.

01

Introduction

History

- 1992 Q-learning algorithm introduced by Watkins
- 2013 DQN introduced by DeepMind in the paper 'Playing Atari with Deep Reinforcement Learning' at NIPS
- 2015 Comprehensive DQN paper 'Human-level control through deep reinforcement learning' published in Nature
- 2018 Rainbow DQN algorithm introduced, combining several improvements



Google DeepMind is a research company owned by Google that focuses on artificial intelligence (AI) and machine learning.

DQN In A Nutshell

DQN is a reinforcement learning algorithm that combines Q-Learning with deep neural networks to learn how to act in an environment.

- **Key Concepts:**
 - States: Environment configurations.
 - Actions: Possible moves.
 - Rewards: Feedback signals.
- **Applications:** Video games (Atari), robotics, finance (trading strategies).
- **Advantages:** Handles high-dimensional state spaces, effective for complex policy learning.
- **Challenges:** Requires extensive training data, tuning for stability.
- **Purpose:** It aims to enable an agent to make decisions that maximize cumulative rewards over time.

The background features several large, soft-edged, organic shapes in various colors: a large orange shape in the top left, a blue shape in the top right, a pink shape on the left, a purple shape in the bottom left, and a blue shape in the bottom right. A small red shape is also visible near the top center.

02

Architecture

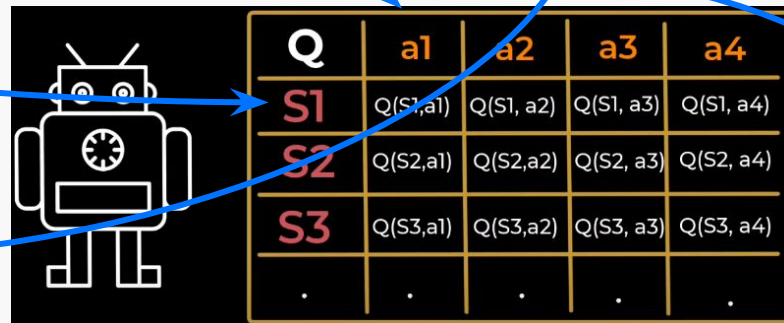
Architecture

RL agents policy (π):

In each **State** the Agent know which **Action** to do, based on **Q-Table** to achieve the maximum **Reward**



	-1	-1
-1	-1	-1
-1	-10	+10



Q	a1	a2	a3	a4
S1	Q(S1,a1)	Q(S1, a2)	Q(S1, a3)	Q(S1, a4)
S2	Q(S2,a1)	Q(S2,a2)	Q(S2, a3)	Q(S2, a4)
S3	Q(S3,a1)	Q(S3,a2)	Q(S3, a3)	Q(S3, a4)
.

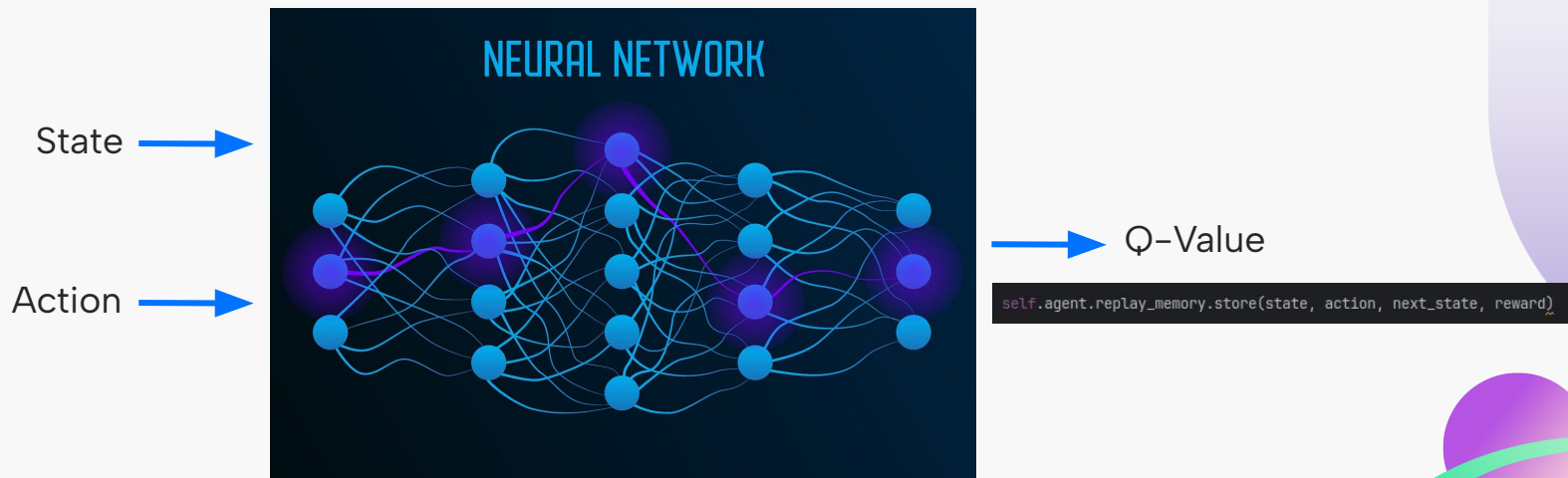
The DQN replace the Q-table with a neural network instead

Architecture

DQN:

The DQN Replace Q-Table with NN prediction.

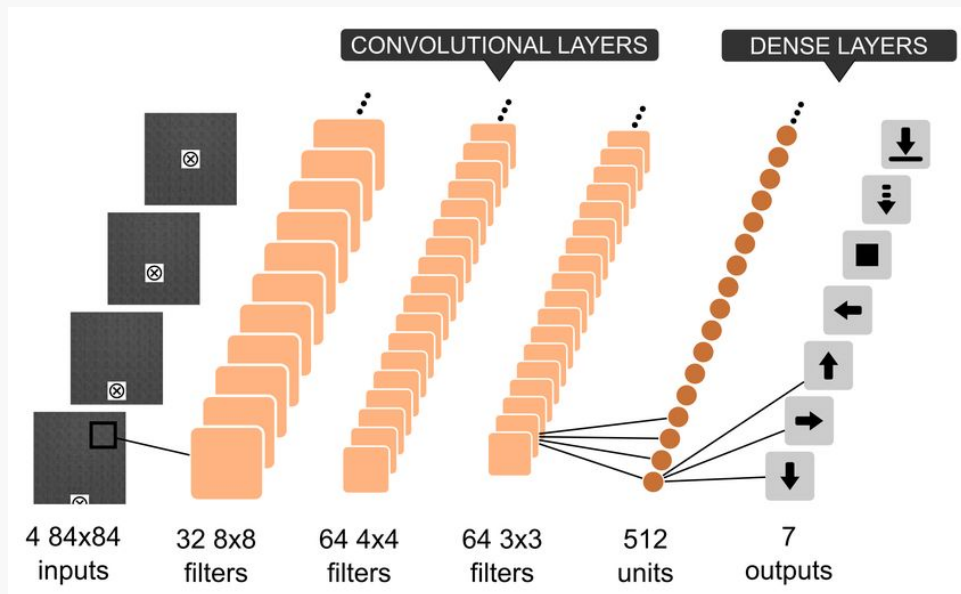
Performing with more robust and scalable solution rather than using a big action-state table



Architecture

Agent & NN

1. **Input Layer:**
 - State representation (game screen pixels)
2. **Convolutional Neural Networks (CNNs):**
 - Extract spatial features from inputs
3. **Fully Connected Layers:**
 - Combine features to predict Q-values
 - ReLU Activation function
4. **Output Layer:**
 - Q-values for each possible action



Architecture

Optimal Q action-value function

s – state (observation)

a – action

r – reward

t – index for step

γ – discount

π – policy $P(\text{action}|\text{state})$

$$Q^*(s,a) = \max_{\pi} \mathbb{E} [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi]$$

The data for the NN

The data is the states and actions that were happened in the environment

e – experience

t – index for step

D_t – collection of e, data set

$$e_t = (s_t, a_t, r_t, s_{t+1})$$

$$D_t = \{e_1, \dots, e_t\}$$

Architecture

Loss function

i – index for iteration

θ – weights

$U(D)$ – unified the data set

s' – state, s – previous state

a' – action, a – previous action

r – reward

γ – discount

Q – action-value function

$$L_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim U(D)} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a; \theta_i) \right)^2 \right]$$

Optimization algorithm

SGD – stochastic gradient descent

The background features several large, overlapping organic shapes in various colors: a large orange shape in the top left, a blue shape in the top right, a pink shape on the left, a purple shape in the bottom left, and a blue shape in the bottom right. A small pink sphere and a green curved line are also visible near the bottom left.

03

Articles

Articles



Human-level control through deep reinforcement learning

LETTER

doi:10.1038/nature14236

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. Fiedjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹



Playing Atari with 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

The background features several large, soft-edged, organic shapes in various colors: a large orange shape in the top left, a blue shape in the top right, a pink shape on the left, a purple shape in the bottom left, and a blue shape in the bottom right. A small pink sphere is positioned above a green curved line that arches over the purple shape.

**0
4**

Environment

Froze Lake

1. Game Description:

- Frozen Lake is a grid-based game where an agent must navigate from a start point to a goal point on a slippery surface, avoiding holes.

2. Environment Setup:

- Grid Layout: Typically an $n \times n$ grid (4x4 or 8x8)
 - S: Start point
 - G: Goal point
 - H: Holes (fall into one, and the game is over)
 - F: Frozen surface (safe to walk on)

3. Borders:

- The grid is bounded, meaning the agent cannot move outside the grid.
- Actions attempting to move out of bounds result in no movement.



Froze Lake

4. States and Actions:

- States (S): Each cell in the grid represents a unique state.
- Actions (A): Four possible actions – left, right, up, down.
- Rewards (R):
 - +1 for reaching the goal.
 - 0 for all other transitions (including falling into a hole).

5. Slippery Surface Dynamics:

- **Slippery Nature:** Due to the slippery surface, actions do not always result in moving to the intended adjacent cell. Instead, there is a chance the agent will slide to an unintended cell, adding randomness to the game and increasing the complexity of navigation.

The background features several large, overlapping organic shapes in various colors: a large orange shape in the top left, a blue shape in the top right, a pink shape on the left, a purple shape in the bottom left, and a blue shape in the bottom left. A green curved line and a small pink circle are also present near the center.

05

METHODOLOGY

Project Goal

Train an AI agent using Deep Q-Networks (DQN) to navigate the Frozen Lake environment and successfully reach the goal state while avoiding falling into holes.

Topic	OBJECTIVE
Problem Definition	<ul style="list-style-type: none">• Train DQN agent to navigate Frozen Lake, reaching goal while avoiding holes.
Environment Setup	<ul style="list-style-type: none">• Define Frozen Lake grid with start, goal, holes, and frozen surfaces.
State and Action Representation	<ul style="list-style-type: none">• States: Agent's grid position; Actions: Left, right, up, down movements.
DQN Implementation	<ul style="list-style-type: none">• Use neural network for Q-function approximation.
Training Process	<ul style="list-style-type: none">• Implement DQN model
Evaluation and Testing	<ul style="list-style-type: none">• Assess agent's performance: average reward, success rate.
Hyperparameter Tuning	<ul style="list-style-type: none">• Optimize parameters like learning rate, discount factor.

The background features several large, soft-edged, organic shapes in various colors: a large orange shape in the top left, a blue shape in the top right, a pink shape on the left, a purple shape in the bottom left, and a blue shape in the bottom right. A small red shape is also visible near the top center.

06

Project Plans

Project Stages

Implement RL:
DQN model and Environment

Stage 1



Train and Test RL Model

Stage 2



Collect and Analyze
Results

Stage 3



Thanks



Aviv Salomon
avivsalo@gmail.com

Ben Gornizky
bengornizky@gmail.com

