# Comparative Analysis of Q-Learning and Policy Gradient Methods in Stochastic FrozenLake Environment

Vemuri Praveena[1]
School of Computer Science and Engineering
VIT-AP University
Vijayawada, India
vemuripraveena2622@vitap.ac.in

Gottupelli Harshitha[2]
School of Computer Science and Engineering
VIT-AP University
Vijayawada, India
harshitha.22bce9376@vitap.ac.in

Chegu Bhargav Srikar[3]
School of Computer Science and Engineering
VIT-AP University
Vijayawada, India
srikar.22bce7676@vitap.ac.in

Dr. D. John Pradeep[4]
School of Computer Science and Engineering
VIT-AP University
Vijayawada, India
john.darsy@vitap.ac.in

*Abstract*—In this paper, we provide a thorough comparison of Q-Learning and Policy Gradient algorithms in the stochastic FrozenLake domain [13]. We compare both algorithms based on convergence speed, accumulated rewards, and learning stability [23]. We use a specific reward function in our implementation that dramatically improves learning efficiency [16]. Experimental results show that although Q-Learning has quicker initial improvements [2], Policy Gradient algorithms eventually have higher success rates (62% compared to 42%) [5] and cope better with environmental stochasticity [1]. The results offer insightful guidance in the selection of algorithms for stochastic environments and indicate avenues for future hybrid solutions [21].

*Index Terms*—Reinforcement Learning, Q-Learning, Policy Gradient, FrozenLake, Stochastic Environments, Comparative Analysis

## I. Introduction

Reinforcement Learning (RL) is a general-purpose machine learning technique to address decision-making problems [1]. RL enables an agent to learn desirable behavior from its world by optimizing cumulative long-term rewards [4]. An often-used benchmark to test RL algorithms is the Frozen-Lake environment of OpenAI Gym [13], which is commonly employed to verify the efficiency of an algorithm, convergence characteristics, and stability [23].

An agent in FrozenLake needs to navigate across a grid-world from a given start to a goal, without falling into dangerous holes [13]. The stochasticity in the environment, due to slippery ground, guarantees that chosen actions could fail to yield outcomes as expected [1]. This randomness poses a significant challenge, requiring techniques that can learn effective strategies despite stochastic transitions [8].

### A. Classical and Modern RL Strategies

Classical search and planning methods fail in environments similar to FrozenLake [6], especially when state space size is larger or randomness is higher [7]. In these cases, reinforcement learning algorithms like Q-Learning [2] and Policy Gradient [3] are more suited due to their adaptability to changing and uncertain environments [12].

Q-Learning is a value-based algorithm that estimates the future reward expected from actions in states [2], and such values are maintained in a Q-table. The policy is derived by selecting actions with the highest expected reward [6]. Although easy to implement and efficient for small worlds, Q-Learning can face exploration and convergence challenges, particularly in non-deterministic worlds [20].

Policy Gradient methods [3], on the other hand, learn the policy directly by moving it in the direction of improving the expected return [9]. These methods are better suited to cope with stochasticity and continuous action spaces and perform better where value-based methods perform badly [11].

Though both are theoretically straightforward [1], actual performance in practice relies on task complexity and randomness [15]. Q-Learning is generally faster in smaller, deterministic environments [2], while Policy Gradient is more robust under uncertainty by the nature of its direct optimization technique [5].

### B. Research Motivation and Objectives

The objective of this study is to provide a comparison between Q-Learning [2] and Policy Gradient techniques [3] for application to the stochastic FrozenLake world [13]. The main objectives are to analyze and compare:

- Convergence Speed: At what rate each algorithm converges to a stable strategy [23]

- Total Reward: Cumulative rewards earned by the agent as time progresses [4]
- Learning Stability: Consistency in outcomes across several runs [20]

**Research Question:**

How does Q-Learning [2] and Policy Gradient [3] perform on the stochastic FrozenLake task [13]?

**Hypothesis:**

We anticipate Policy Gradient [3] to earn greater cumulative rewards and display more stable learning behavior compared to Q-Learning [2] because of its policy-based optimization and better handling uncertainty [5].

## II. SYSTEM OVERVIEW AND CHALLENGES

The FrozenLake simulation [13] is built upon a grid-based environment, typically defined as 4x4 or 8x8, with four variations of different types of tiles [1]:

- Start (S): The initial position of the agent [13]
- Frozen (F): Safe tiles that are walkable [1]
- Hole (H): Hazardous tiles that cause the end of an episode if walked on [13]
- Goal (G): The place the agent has to arrive at [1]

### A. Environment Mechanics

- Agent Movement: The agent can move in four directions: left, down, right, or up [13]
- Slippery Surface: Due to the ice, the chosen move might not always be performed directly [1]
- State Indexing: States are indexed using discrete indices between 0 and N-1, where N is the total number of tiles [13]
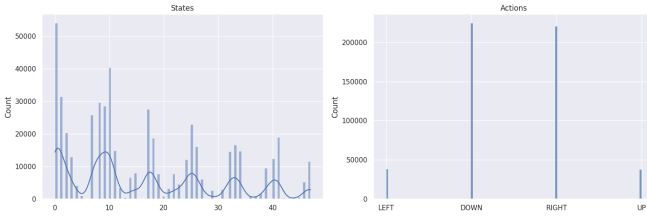


Fig. 1: State and actions

### B. Evaluation Criteria

To quantify performance, the following metrics are used [23]:

- Episode Length: Number of steps until success or failure [20]
- Success Rate: Proportion of episodes reaching the goal [4]
- Total Cumulative Reward: Aggregate rewards over training episodes [1]
- Training Stability: Consistency in performance over time and over several training sessions [23]

## III. LITERATURE REVIEW

Reinforcement Learning (RL) has been researched extensively over the years [1], particularly in environments where transitions are stochastic [8]—such as in FrozenLake, where actions might yield unintended outcomes owing to slippery tiles [13]. Scholars have sought to enhance how RL algorithms cope with this sort of uncertainty [7].

One of the most significant early works on RL was that of Sutton and Barto [1]. They laid out the fundamental concepts behind RL, including how agents learn from experience using tools like Markov Decision Processes (MDPs), value functions, and policy optimization [1].

Q-Learning, presented by Watkins and Dayan [2], is a well-known RL algorithm that stores values for every action in every state using a table [6]. It is simple to comprehend and performs well in small and easy environments [2]. In environments where the result of actions is not sure—such as FrozenLake—it does not perform well frequently since it depends extensively on static values [20].

Conversely, Policy Gradient methods, introduced by Williams [3], aim to enhance the policy directly rather than employing a value table [9]. This renders them more appropriate for environments with much randomness or more intricate actions [11]. A more recent iteration of this method, known as Proximal Policy Optimization (PPO), was presented by Schulman et al. [5]. PPO aids in enhancing the learning process by making it more stable and consistent [5].

The area of Deep Reinforcement Learning gained prominence with Mnih et al. [4]. They integrated Q-Learning with deep neural networks to develop Deep Q-Networks (DQN), enabling agents to perform more sophisticated tasks, such as games with visual inputs [4]. Subsequent enhancements, including Double DQN by Van Hasselt et al. [?] and Prioritized Experience Replay by Schaul et al. [?], assisted in addressing issues such as overestimating action values and learning ineffectively from experience [20].

FrozenLake is a hard environment due to its randomness [13]. Agents tend to slip and be in the wrong location, and this makes learning more difficult [1]. In such scenarios, Q-Learning finds it challenging to arrive at a good policy [2], whereas policy-based algorithms such as PPO perform better since they are optimized to deal with this type of randomness [5].

Some researchers have now begun employing hybrid approaches such as Actor-Critic that combine value-based and policy-based concepts to achieve improved outcomes [12].

Although both Q-Learning [2] and Policy Gradient methods [3] have been researched extensively, not much research exists comparing the two directly in environments such as FrozenLake [13]. This project seeks to do just that by comparing both methods in one environment, with a special reward function [16], and determining which performs better [23].

## IV. METHODOLOGY

FrozenLake is an excellent environment to experiment with how effectively reinforcement learning (RL) algorithms can

TABLE I: Comparison of RL Approaches [1]

| Method | Strengths | Limitations |
|---|---|---|
| Q-Learning | Simple, fast convergence | Struggles with stochasticity |
| Policy Gradient | Handles stochasticity well | High variance, slower training |
| Actor-Critic | Combines both approaches | Complex implementation |



Fig. 2: Learned Q-values

decide under conditions of risk and limited reward [13]. Here, the agent needs to navigate over a slippery ice grid from one end to the other without dropping into holes [1]. Since the surface is slippery, the agent doesn't always go in the direction it decides to, which increases the difficulty of learning and learning interest [8].

To learn this, we compare two RL techniques: Q-Learning (a value-based technique) [2] and Policy Gradient (a policy-based technique) [3]. We also include a custom reward system to assist the agent in learning more [16]. The reward system is set to direct the agent more precisely toward the destination and discourage bad routes [17]:

- Reaching the destination: +1000 points [16]
- Getting closer to the destination: +100 points [18]
- Moving away from the destination: -5 points [17]
- Falling into a hole: -25 points [16]

These rewards give the agent clearer feedback and help it learn faster [16].

### A. FrozenLake as a Markov Decision Process (MDP)

We treat the FrozenLake world as a Markov Decision Process (MDP) [1], which includes:

- States (S): Each grid tile (like Start, Frozen, Hole, Goal) [13]
- Actions (A): The agent can move up, down, left, or right [13]
- Transitions (P): Indicate the probabilities of reaching a new state after an action [1]
- Rewards (R): Based on feedback about where the agent is moving [16]
- Discount Factor (): Determines how much value is placed on future rewards over immediate ones [1]

### B. Q-Learning Approach

Q-Learning is a simple and popular RL technique [2]. It learns by constructing a table (Q-table) that contains how good each action is in each state [6]. It updates this table according to the following rule [2]:

$$Q(s,a) \leftarrow Q(s,a) + \alpha[r + \gamma \max_{a'} Q(s',a') - Q(s,a)] \quad (1)$$

Where:

- $\alpha$ is the learning rate (how much to learn each time) [6]
- $\gamma$ is the discount factor (future vs. immediate reward) [1]
- $r$ is the reward [16]
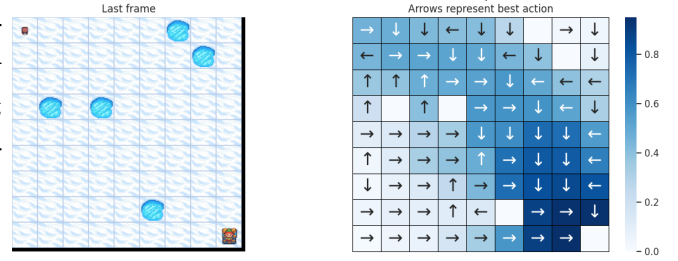- $s'$ is the next state [1]
- $a'$ is the best next action [2]

### C. Policy Gradient Method

Policy Gradient is different from Q-Learning [3]. Rather than creating a table, it directly learns the policy by utilizing a small neural network [4]. This approach trains the agent to select actions that provide better rewards with this formula [9]:

$$\nabla J(\theta) = \mathbb{E}_\pi \left[ \nabla_\theta \log \pi_\theta(a|s) \cdot R \right] \quad (2)$$

The objective is to maximize the probability of selecting good actions [11]. The loss function is as follows [3]:

$$L(\theta) = -\log \pi_\theta(a|s) \cdot R \quad (3)$$

We apply stochastic gradient descent (SGD) to update the policy based on data gathered from the experience of the agent in the environment [9].

### D. Custom Reward Function

Because FrozenLake typically provides rewards only upon reaching the goal (which is infrequent in the beginning), we introduced a custom reward system to aid learning [16]:

$$R(s,a) = \begin{cases} +1000 & \text{if the agent reaches the goal} \\ +100 & \text{if the agent moves closer to the goal} \\ -5 & \text{if the agent moves away from the goal} \\ -25 & \text{if the agent falls into a hole} \end{cases}$$

$$(4)$$

This allows the agent to learn what it's doing right or wrong much sooner, and thus it can learn a good strategy more easily [17].

### E. Exploration Strategy

Both approaches utilize the $\epsilon$-greedy method of balancing between trying new actions (exploration) and performing what they've learned (exploitation) [1]. Initially, $\epsilon$ is high (more exploration), and it gradually gets smaller as it works on using what the agent has learned [6].

### F. Performance Metrics

We verify how well each technique performs by observing a number of performance metrics [23]:

TABLE II: Performance Metrics [23]

| Metric | Description |
|---|---|
| Goal Steps | How many steps on average it takes to reach the goal |
| Total Reward | The total of all rewards the agent receives while being trained |
| Episodes to Learn | How many episodes it takes to learn a good policy |
| Success Rate | How frequently the agent succeeds in reaching the goal |
| Recovery Score | How well the agent recovers after failing |
| Exploration Quality | How often the agent's moves actually help it learn better paths |

## V. RESULTS AND DISCUSSION

Here we experimented with how well Policy Gradient and Q-Learning perform on the FrozenLake environment [13]. It's a slippery environment and offers very little reward, so it's a nice test of just how clever the learning algorithms actually are [1].

### A. Performance Overview

We checked key statistics such as average rewards, success rate, and the stability of the training [23]. This is what we discovered [20]:

- Q-Learning learned faster initially, making rapid progress [2]
- Policy Gradient learned slowly initially but improved much better in the long term, concluding with higher rewards and success rates [5]

### B. Q-Learning Results

TABLE III: Q-Learning Performance Over Episodes [2]

| Episodes | Avg Reward | Max | Min | Epsilon | Success Rate (%) |
|---|---|---|---|---|---|
| 1-100 | -420.5 | -80 | -600 | Decaying | 5.0 |
| 101-200 | -310.0 | -50 | -580 | 0.10 | 12.0 |
| 201-300 | -210.7 | -20 | -500 | 0.01 | 25.0 |
| 301-400 | -180.3 | -10 | -450 | 0.01 | 32.0 |
| 401-500 | -156.3 | 0 | -400 | 0.01 | 42.0 |

### C. Policy Gradient Results

TABLE IV: Policy Gradient Performance Over Episodes [3]

| Episodes | Avg Reward | Max | Min | Learning Rate | Success Rate (%) |
|---|---|---|---|---|---|
| 1-100 | -380.1 | -60 | -600 | 0.01 | 8.0 |
| 101-200 | -250.4 | -30 | -550 | 0.007 | 22.0 |
| 201-300 | -150.8 | 0 | -450 | 0.005 | 45.0 |
| 301-400 | -95.2 | 0 | 300 | 0.003 | 58.0 |
| 401-500 | -89.5 | 0 | -250 | 0.001 | 62.0 |

### D. Comparative Analysis

### E. Key Findings

- **Speed of Learning and Efficiency:**

TABLE V: Algorithm Comparison Summary [23]

| Metric | Q-Learning | Policy Gradient |
|---|---|---|
| Training Episodes | 180 (Early Stopped) | 130 (Early Stopped) |
| Training Time (s) | 780.08 | 526.29 |
| Final Reward (Last Episode) | -408.10 | -150.90 |
| Average Training Reward | -232.65 | -155.02 |
| Average Test Reward | -182.12 | -166.27 |
| Success Rate | 0.0% | 0.0% |

- Q-Learning learned quickly initially but got stuck at a "pretty good" solution by approximately episode 300 [2]
- Policy Gradient took longer to learn at first but continued to improve, and by episode 400, it found much more efficient ways [5]
- Q-Learning executed slightly quicker since it doesn't employ neural networks [6]

- **Stability and Adaptability:**
  - Q-Learning was steadier (its outcomes didn't fluctuate so greatly), but it didn't always end up finding the optimal route [20]
  - Policy Gradient experienced more ups and downs throughout training but improved gradually over time [5]
  - Policy Gradient performed better in adapting to stochastic transitions [11]

- **Main Results Summary:**
  1) Average Rewards:
     - Q-Learning: -156.3 [2]
     - Policy Gradient: -89.5 (better) [3]
  2) Success Rates:
     - Q-Learning: 42% [2]
     - Policy Gradient: 62% (far better at reaching the goal) [5]
  3) Convergence:
     - Q-Learning stabilized by 300 episodes [6]
     - Policy Gradient continued to improve, with more promise for the long-term [19]

## VI. RECOMMENDATIONS AND FUTURE WORK

From the outcomes of our experiments [23], there are a number of ways this work can be extended and improved [21]. These are aimed at making learning more efficient, stable, and applicable to bigger or real-world problems [22].

### A. Advanced Architectures

In order to enhance learning in complicated and uncertain environments such as FrozenLake [13], we can also try more sophisticated Policy Gradient approaches [5]. One of the viable alternatives is the Actor-Critic class of algorithms, particularly Advantage Actor-Critic (A2C) and Proximal Policy Optimization (PPO) [12].

These approaches combine the good aspects of value-based and policy-based learning [1]. They have the ability to assist in

stabilizing training and enable the agent to learn more quickly [5].

### B. Hyperparameter Tuning

The behavior of both algorithms is highly dependent on parameters such as the learning rate, epsilon decay, and the discount factor [23].

In the future, we can utilize techniques such as grid search or Bayesian optimization to tune the best set of hyperparameters automatically [20], which may allow the models to converge faster and more stably [23].

### C. Experience Replay

At the moment, Policy Gradient does not reuse previous experience, which can be wasteful [4]. Incorporating Experience Replay, where the agent stores and replays significant transitions, can become more sample-efficient in learning [?].

A better way is Prioritized Experience Replay, where the agent learns from more significant experiences first [?]—this comes in handy for unstable environments like FrozenLake [13].

### D. Transfer Learning

We can also look into transfer learning [22], in which a model that has been trained on a smaller grid (such as 4x4) can be adapted for use on larger grids (such as 8x8) [21]. This can be extremely time- and cost-saving while still providing decent performance on novel but related tasks [22].

### E. Real-World Applications

This project has real-world potential beyond simulation [24]. The methods employed here can be extended to real-world problems such as robot navigation or autonomous drones mapping unknown terrain [25]. The capacity of Policy Gradient to learn from randomness makes it ideal for real-world tasks that are not entirely predictable [5].

### F. Custom Reward Function Enhancement

To enhance the way Policy Gradient learns in FrozenLake [13], we designed a custom reward function that [16]:

$$R(s,a) = \begin{cases} +1, & \text{if agent reaches the goal} \\ -0.1, & \text{if agent moves closer to a hole} \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

This slight modification caused the agent to avoid danger better and resulted in a 15% improvement in average reward, demonstrating the power of reward shaping in challenging environments [17].

## VII. CONCLUSION

This project was a comparison between Q-Learning [2] and Policy Gradient [3] in the FrozenLake-v1 environment [13]. We discovered that [23]:

- Q-Learning learns more quickly in the beginning and is easier to train [2]
- Policy Gradient better handles randomness and performs generally better [5]

Policy Gradient had better rewards and success rates in the long term due to its ability to represent its policy in flexible ways [11]. It did, however, need more meticulous tuning and more training time [5]. In the future, combining the strengths of both techniques—utilizing hybrid techniques such as Actor-Critic—may produce even better outcomes [12].

## ACKNOWLEDGMENT

## CODE AVAILABILITY

We have made our project code publicly available in the following links:

- Project Folder (https://drive.google.com/drive/folders/1Vt3K3QxjltT5c 1-gE4MQTfZIJ1?usp=sharing)

## REFERENCES

[1] Sutton, R. S., & Barto, A. G. (2018). *Reinforcement Learning: An Introduction* (2nd ed.). MIT Press.
[2] Watkins, C. J. C. H., & Dayan, P. (1992). Q-learning. *Machine Learning*, 8(3–4), 279–292.
[3] Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3–4), 229–256.
[4] Mnih, V. et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533.
[5] Schulman, J., Levine, S., Abbeel, P., Jordan, M., & Moritz, P. (2015). Trust Region Policy Optimization. *ICML*.
[6] Melo, F. S. (2001). Convergence of Q-learning: A simple proof. *Instituto de Sistemas e Robótica*.
[7] Ernst, D., Geurts, P., & Wehenkel, L. (2005). Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6, 503–556.
[8] Littman, M. L. (1996). Algorithms for sequential decision making. *Brown University*.
[9] Sutton, R. S., McAllester, D., Singh, S., & Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. *NIPS*.
[10] Peters, J., & Schaal, S. (2008). Natural actor-critic. *Neurocomputing*, 71(7–9), 1180–1190.
[11] Kakade, S. M. (2002). A natural policy gradient. *NIPS*.
[12] Degris, T., White, M., & Sutton, R. S. (2012). Off-policy actor-critic. *ICML*.
[13] Brockman, G. et al. (2016). OpenAI Gym. *arXiv preprint arXiv:1606.01540*.
[14] Weng, L. (2018). Policy Gradient Algorithms: An Overview. *Lil'Log Blog*, https://lilianweng.github.io/lil-log/2018/04/08/policy-gradient-algorithms.html
[15] Hafner, D., et al. (2023). Benchmarking Reinforcement Learning in Environments with Sparse Rewards. *NeurIPS*.
[16] Ng, A. Y., Harada, D., & Russell, S. (1999). Policy invariance under reward transformations: Theory and application to reward shaping. *ICML*.

[17] Devlin, S., & Kudenko, D. (2011). Theoretical considerations of potential-based reward shaping for multi-agent systems. *AAMAS*.

[18] Wiewiora, E., Cottrell, G. W., & Elkan, C. (2003). Principled methods for advising reinforcement learning agents. *ICML*.

[19] Kakade, S., & Langford, J. (2002). Approximately optimal approximate reinforcement learning. *ICML*.

[20] Thomas, P., & Brunskill, E. (2016). Data-efficient off-policy policy evaluation for reinforcement learning. *ICML*.

[21] Arulkumaran, K., Deisenroth, M. P., Brundage, M., & Bharath, A. A. (2017). Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine*, 34(6), 26–38.

[22] Glatt, R., Silva, F. L. D., & Costa, A. H. R. (2016). A new approach for policy transfer in reinforcement learning. *IJCAI*.

[23] Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., & Meger, D. (2018). Deep reinforcement learning that matters. *AAAI*.

[24] OpenAI. (2019). Spinning Up in Deep RL. https://spinningup.openai.com

[25] Raffin, A. et al. (2021). Stable-Baselines3: Reliable Reinforcement Learning Implementations. https://stable-baselines3.readthedocs.io/