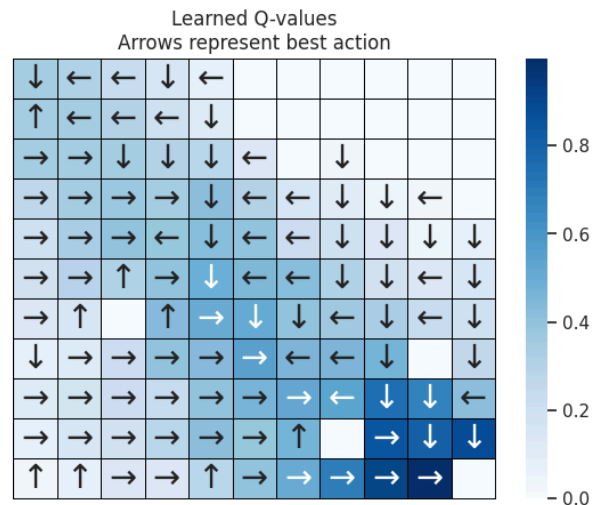
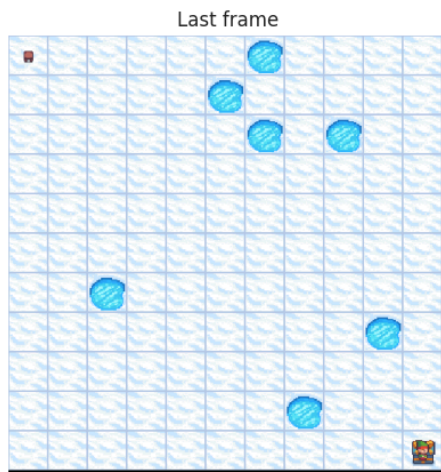


LAB 6, V5: Frozen Lake

GR 5

Introduction

This report presents the results of applying Q-learning to solve the FrozenLake-v1 environment with different hyperparameter settings. The FrozenLake environment represents a grid-world navigation problem where an agent must traverse a frozen lake from a starting point to a goal while avoiding holes. The environment uses a map size of 8x8 and the non-slippery version was selected. The Q-learning algorithm is a model-free reinforcement learning technique that learns the value of an action in a particular state. This algorithm uses an exploration-exploitation strategy to balance between trying new actions and using known good actions.



Experimental Setup

We conducted a series of experiments to investigate how different hyperparameters affect the performance of the Q-learning algorithm in the FrozenLake environment. Each experiment was run for 20,000 episodes with different configurations of the following parameters:

- **Learning rate (α):** Controls how much new information overrides old information (0.7, 0.8, 0.9)
- **Discount factor (γ):** Determines the importance of future rewards (0.98, 0.99, 0.995)

- **Epsilon decay rate:** Controls how quickly the exploration rate decreases (0.9995, 0.9999, 0.99995)

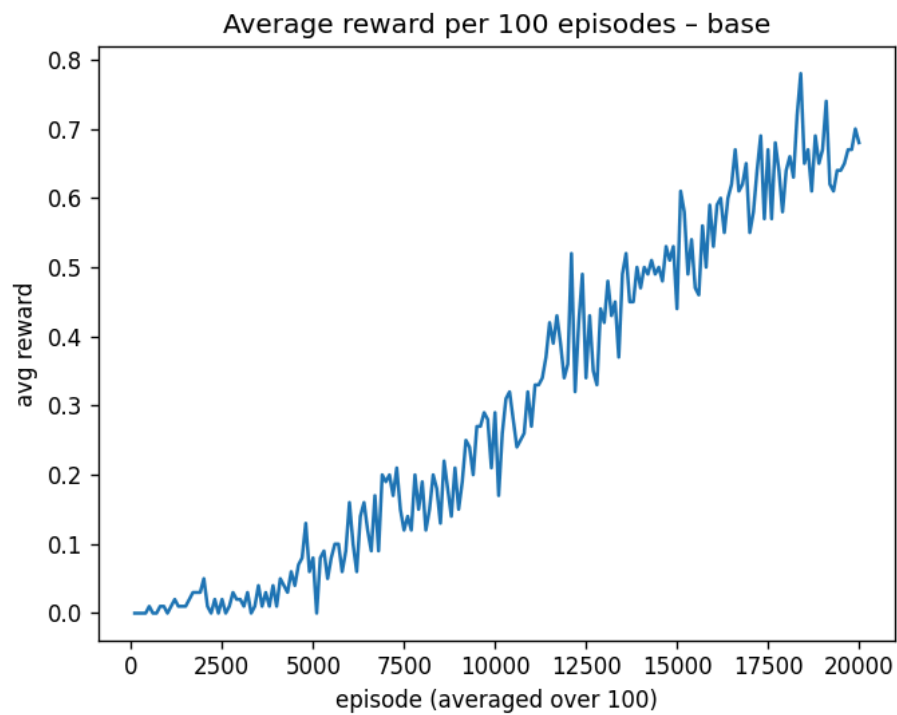
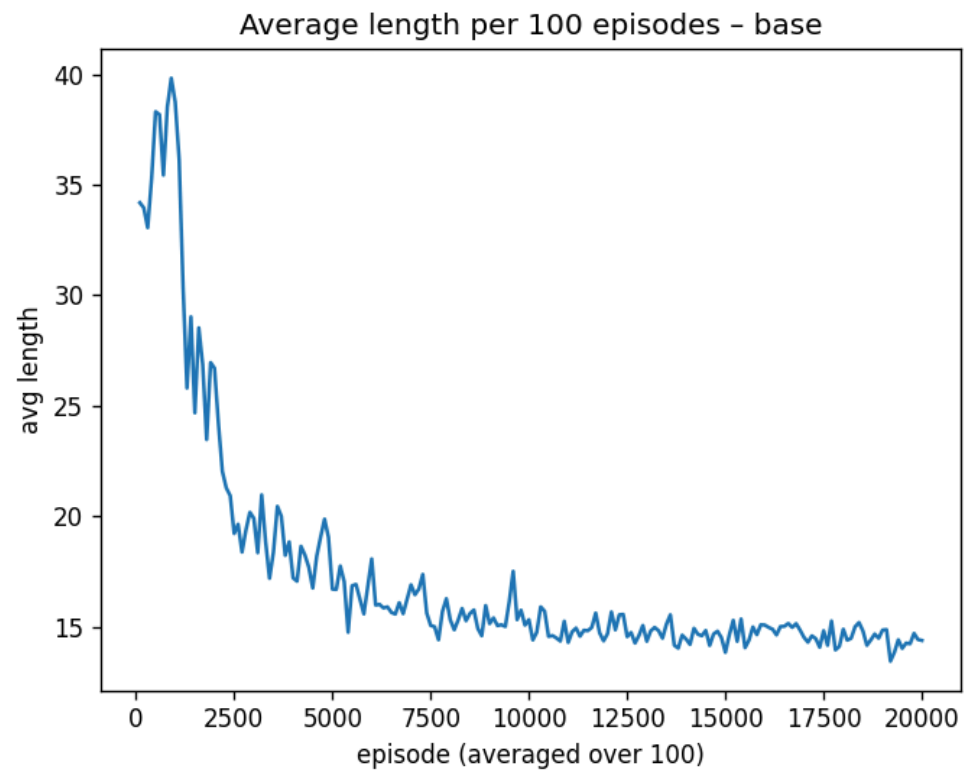
All experiments used the same starting epsilon ($\epsilon_0 = 1.0$) and random seed (42) for reproducibility. Performance was evaluated based on:

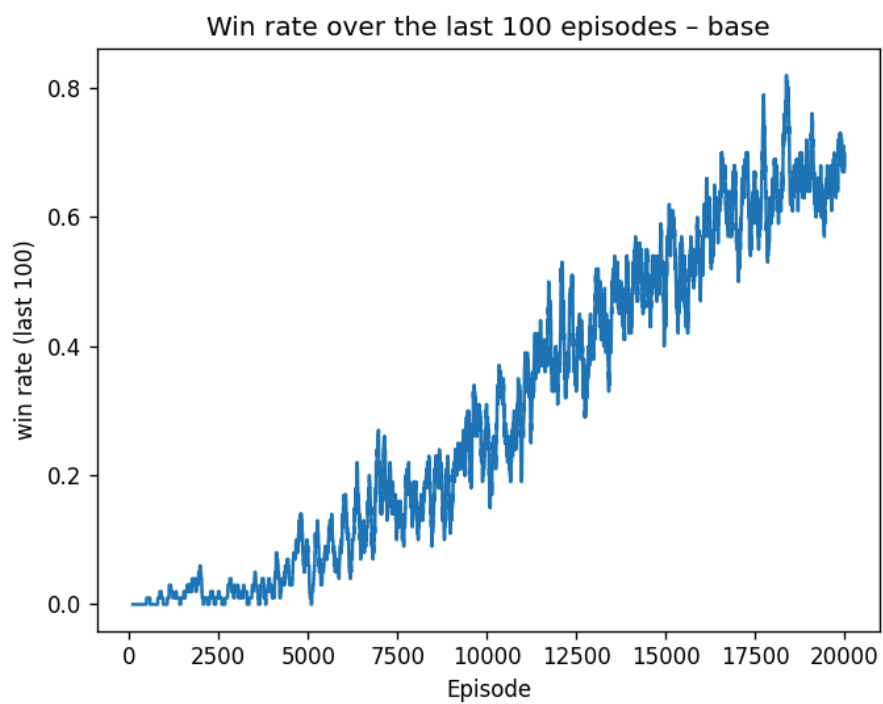
1. Reward graph
2. Average win rate per batch
3. Average length per batch
4. Final win rate during evaluation (1000 trials)

Results

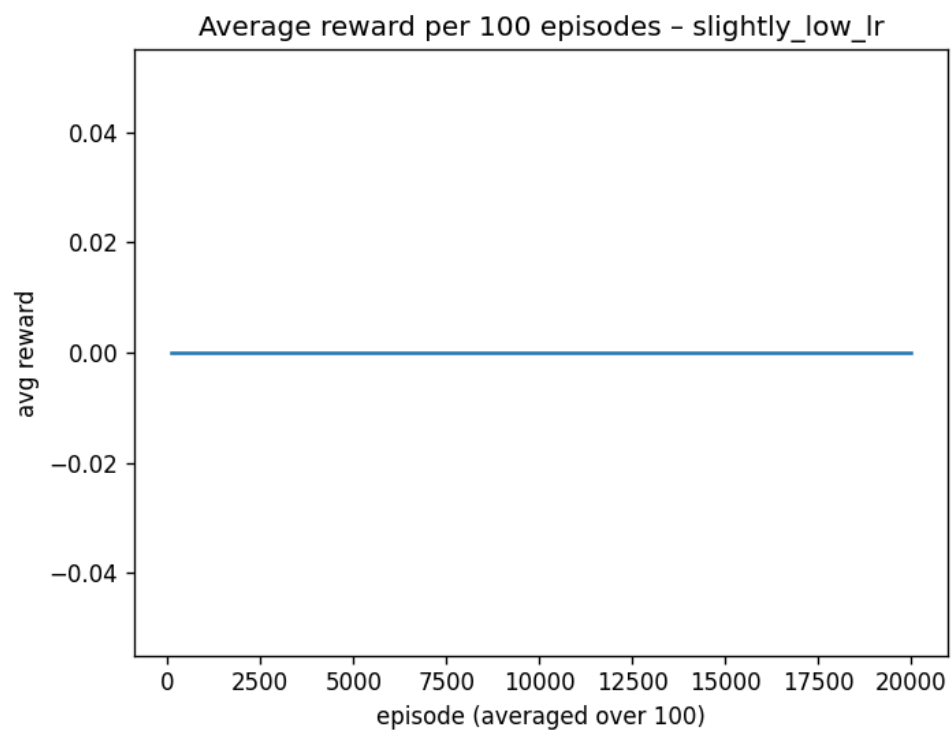
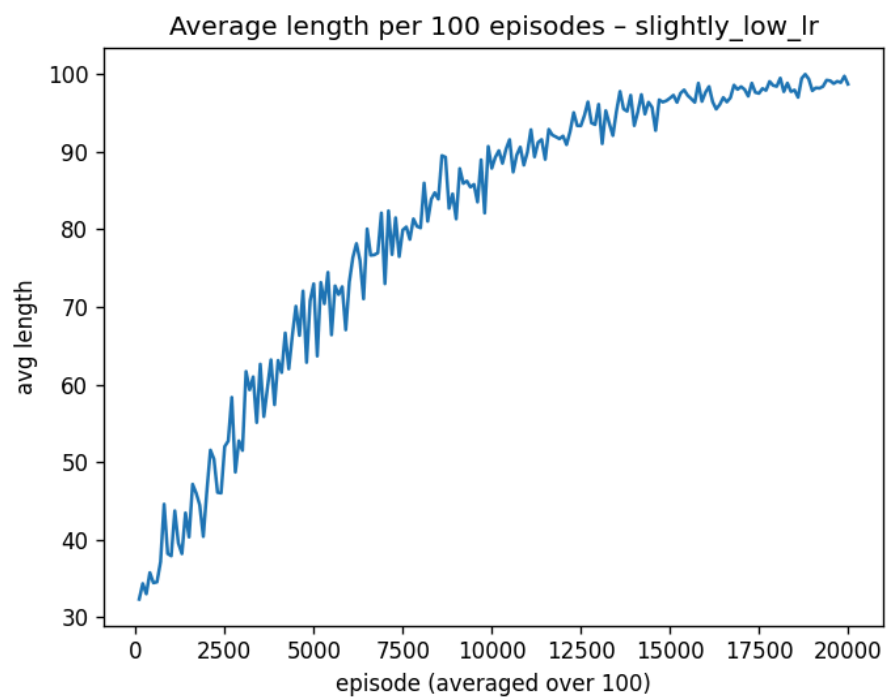
Learning Rate Variation

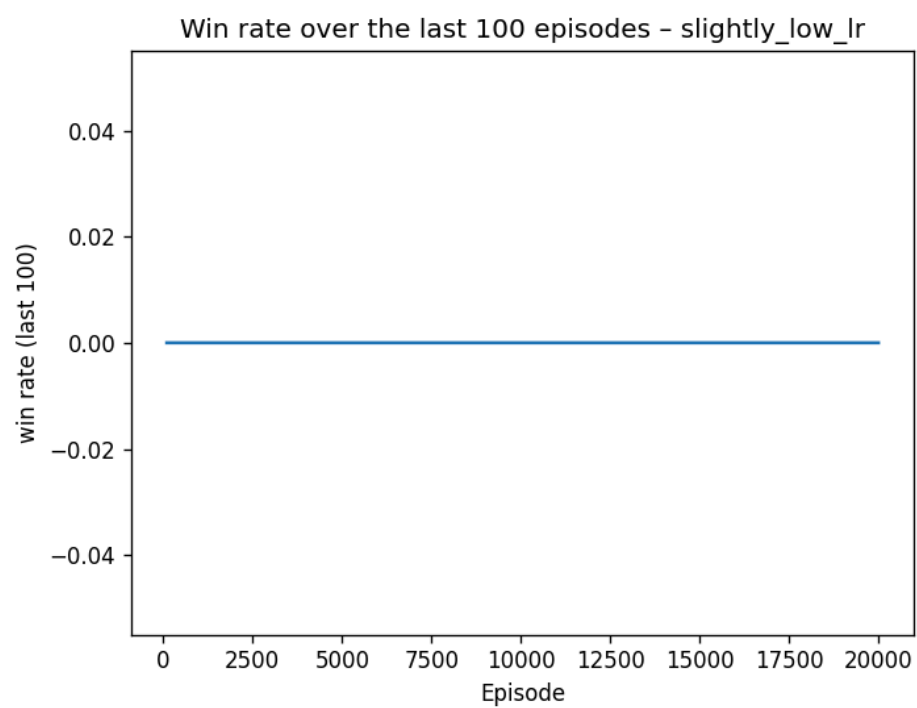
Base Case ($\alpha = 0.8$)



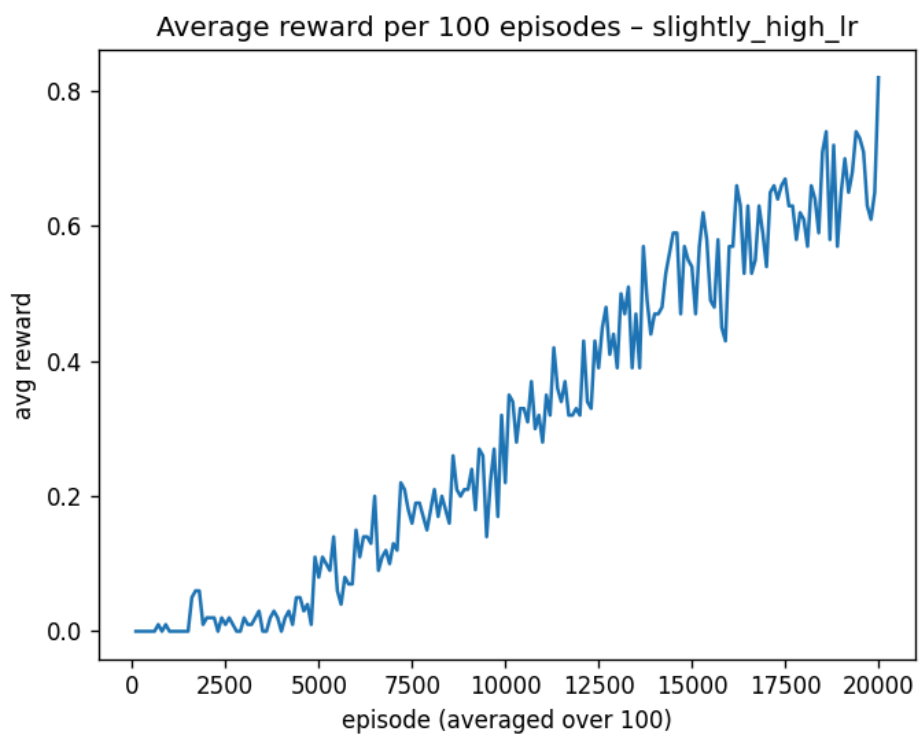
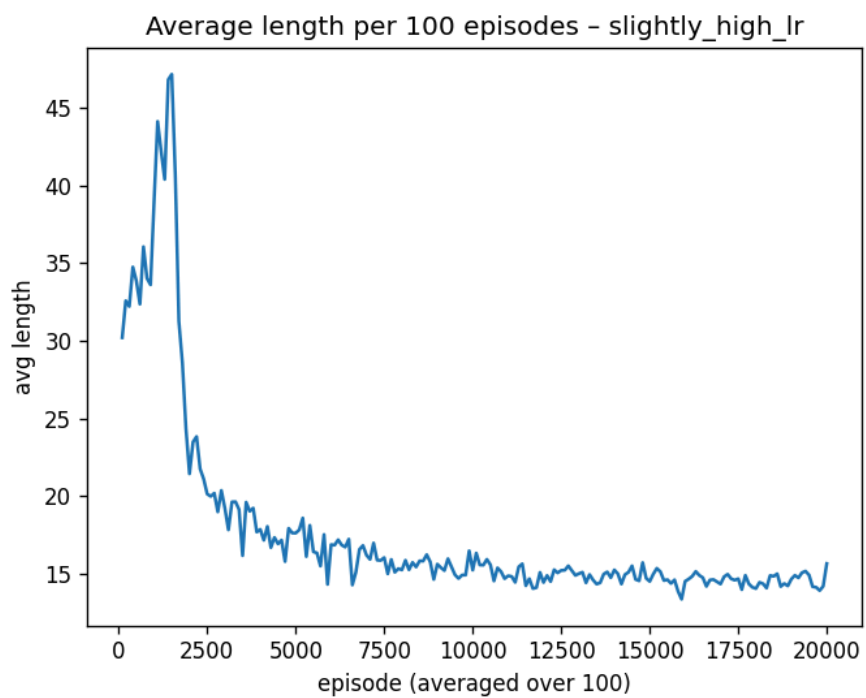


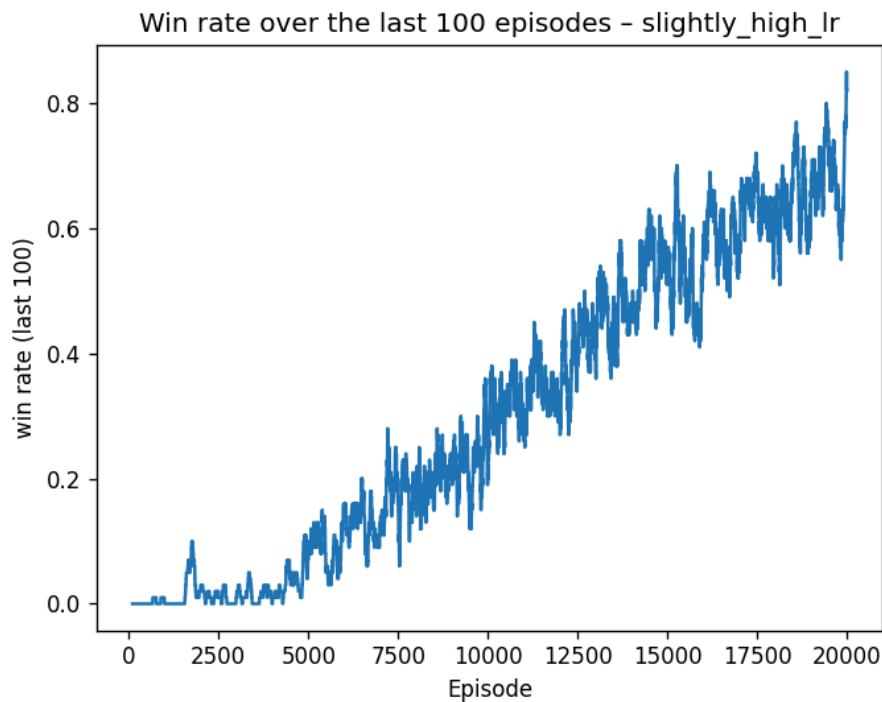
Slightly Lower Learning Rate ($\alpha = 0.7$)





Slightly Higher Learning Rate ($\alpha = 0.9$)



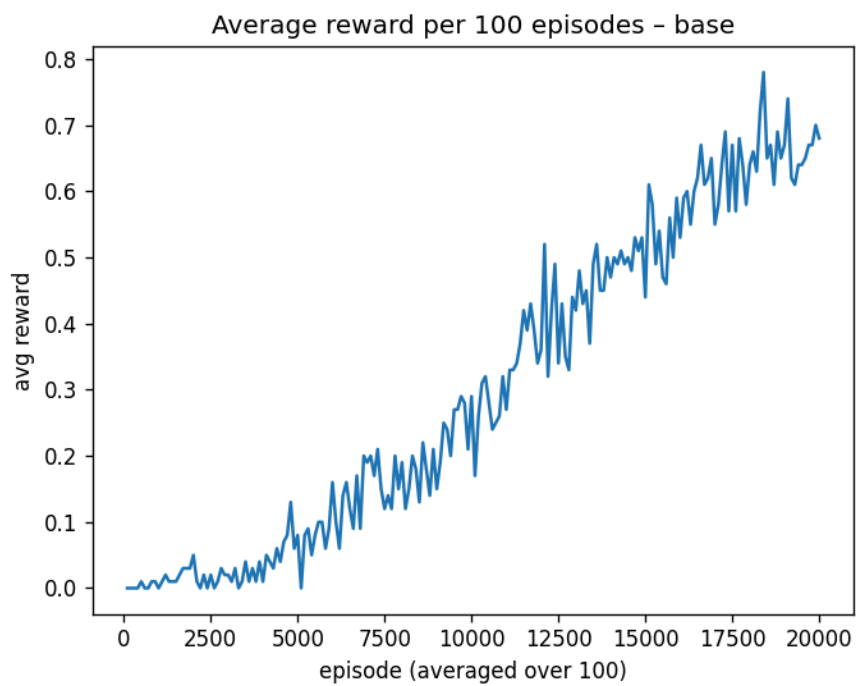
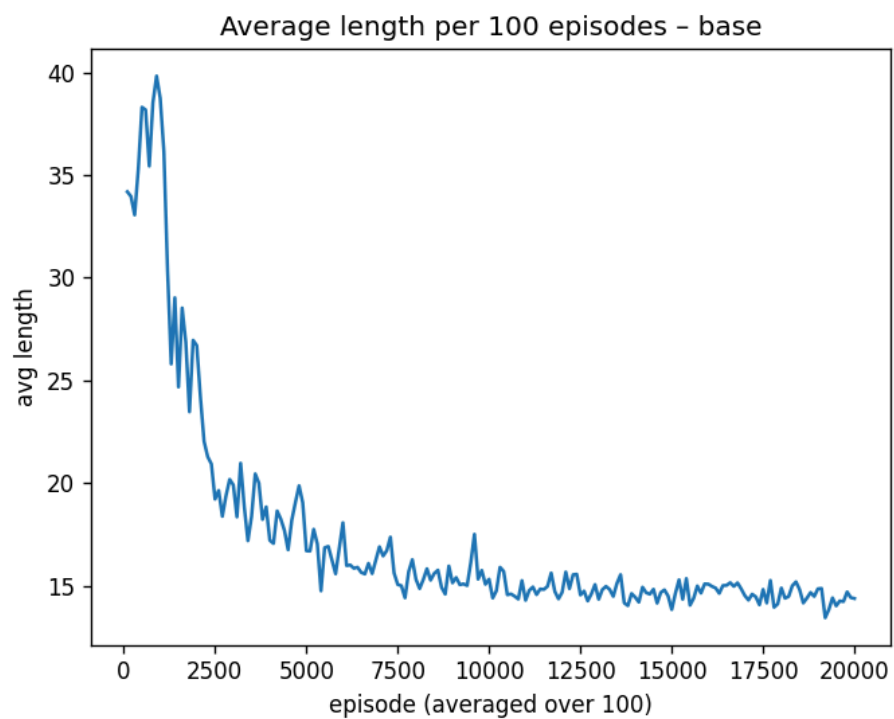


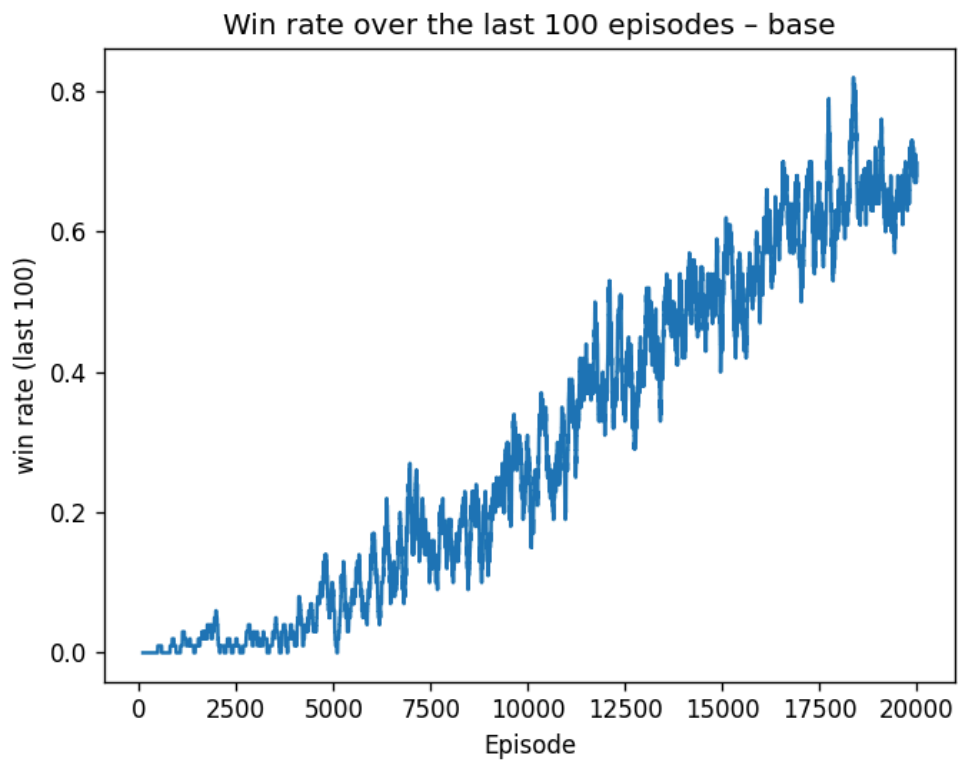
The learning rate had a dramatic effect on algorithm performance. The base case ($\alpha = 0.8$) achieved a 100% win rate during evaluation with an average episode length of 14 steps. The slightly higher learning rate ($\alpha = 0.9$) performed similarly, also reaching a 100% win rate during evaluation.

However, the slightly lower learning rate ($\alpha = 0.7$) failed completely, with a 0% win rate and consistently reached the maximum episode length, indicating that the agent was unable to learn an effective policy. This suggests that in this environment, having a sufficiently high learning rate is critical for successful learning.

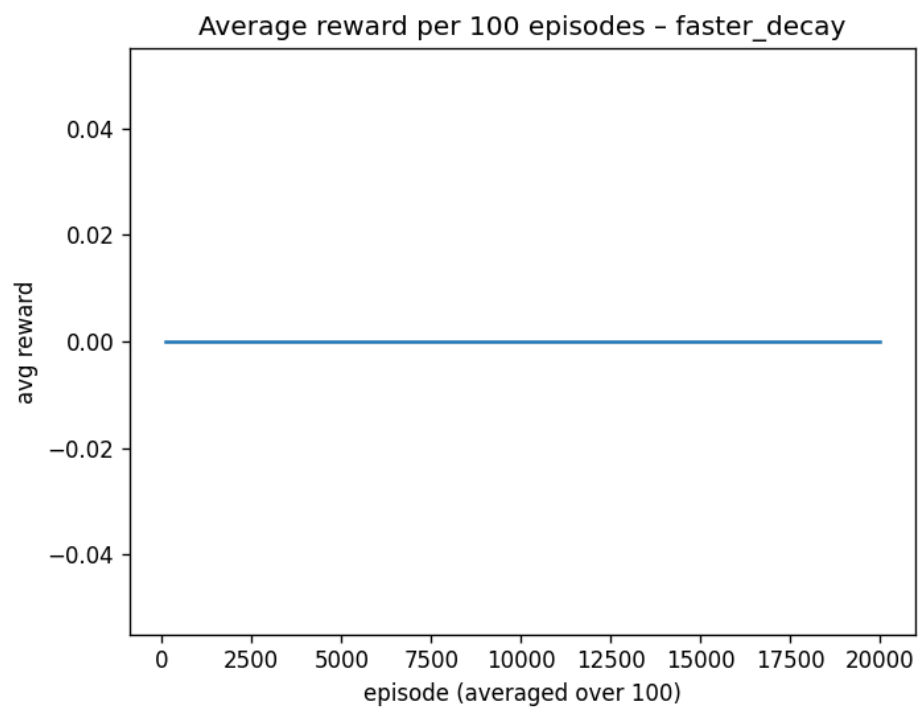
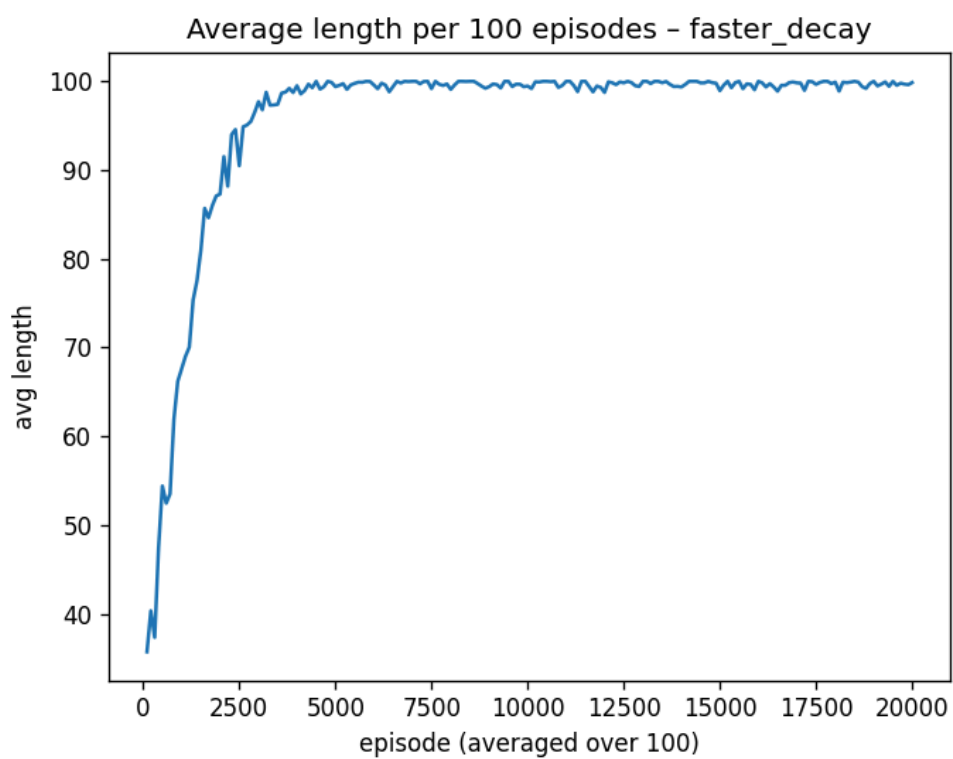
Epsilon Decay Rate Variation

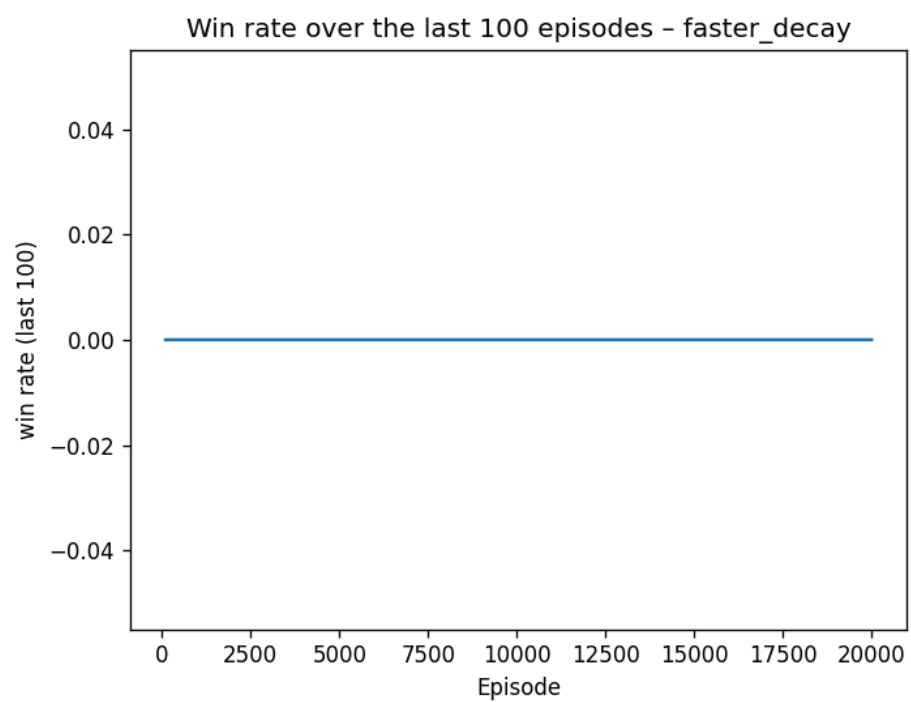
Base Case (decay = 0.9999)



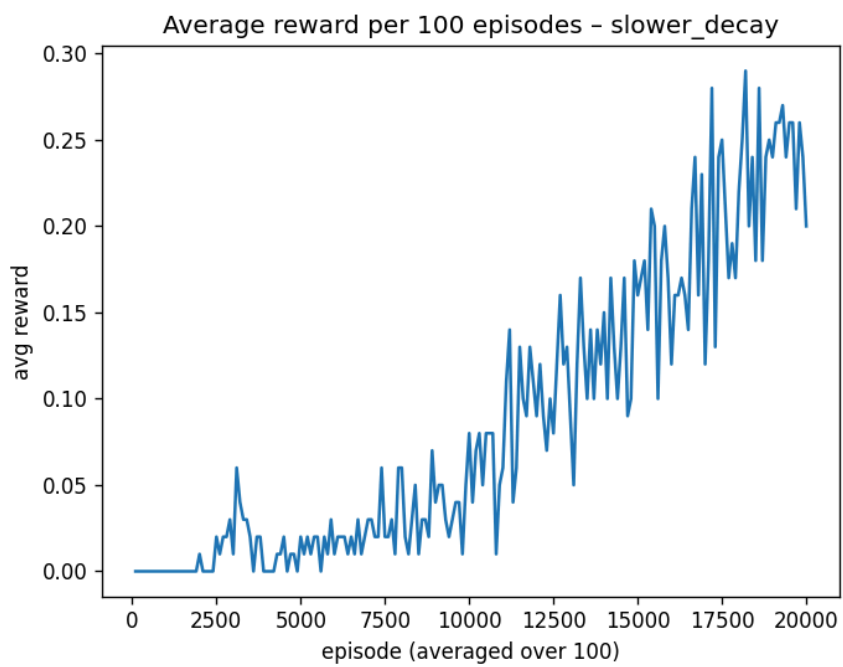
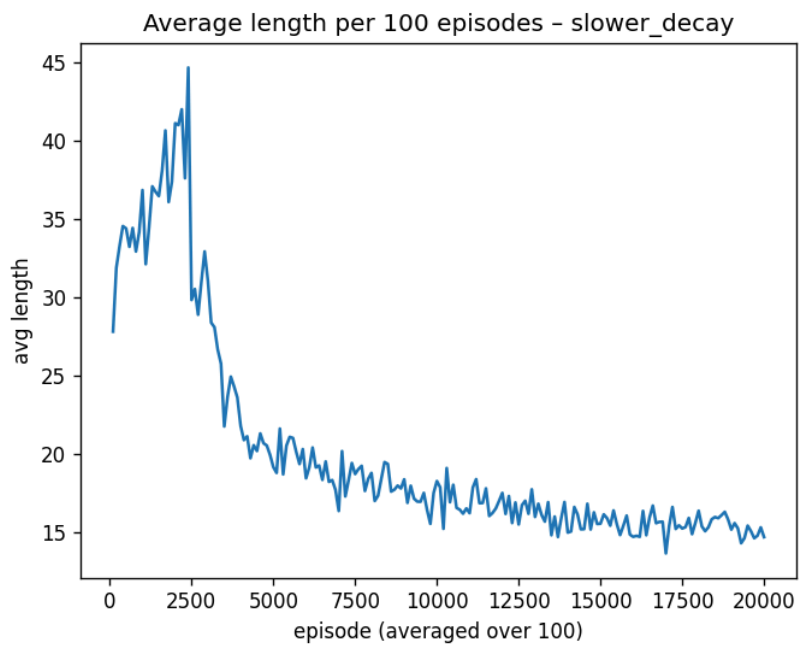


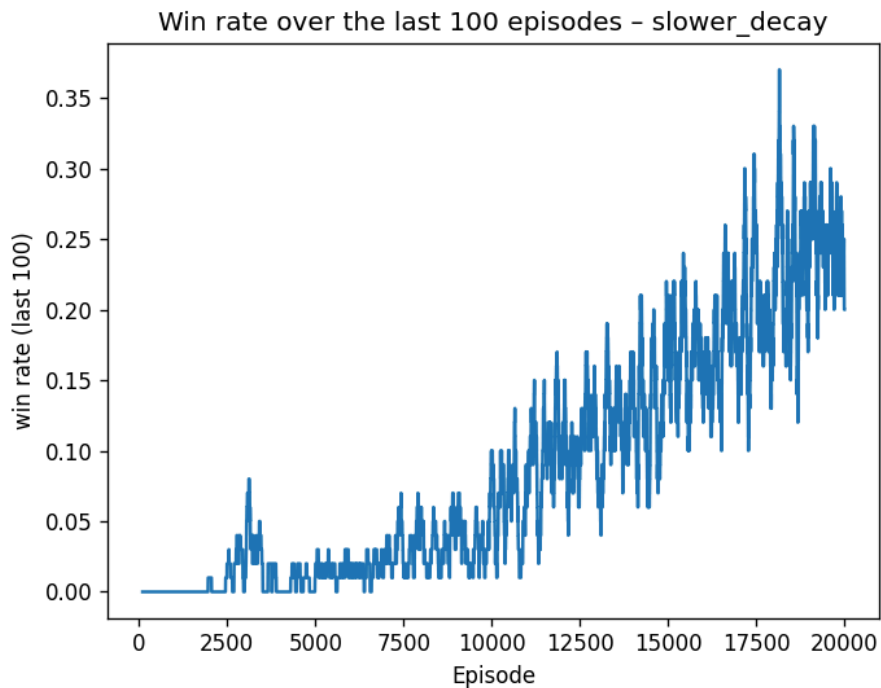
Faster Decay (decay = 0.9995)





Slower Decay (decay = 0.99995)





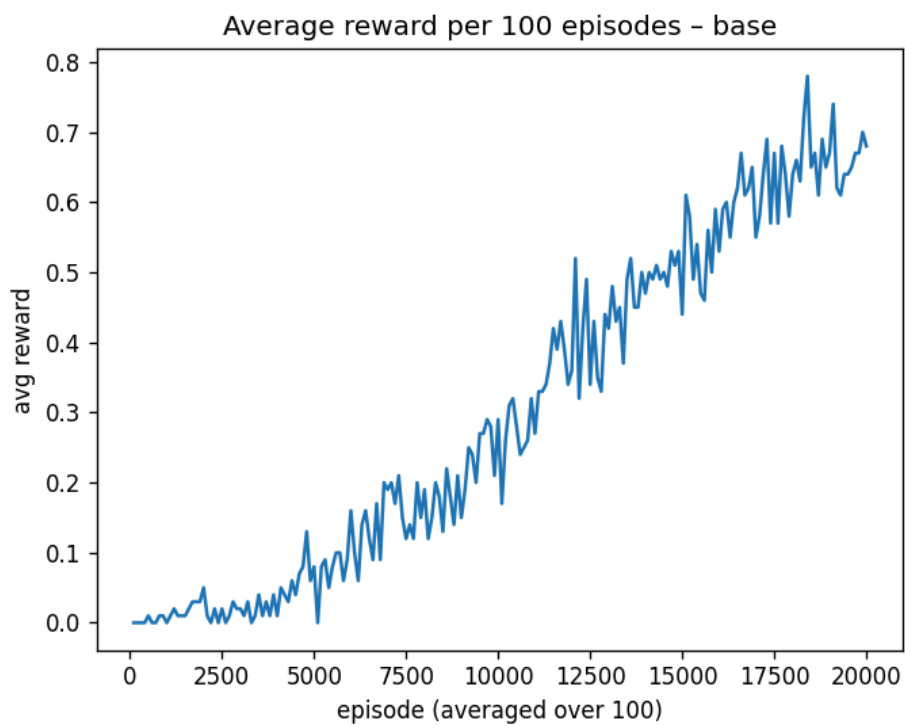
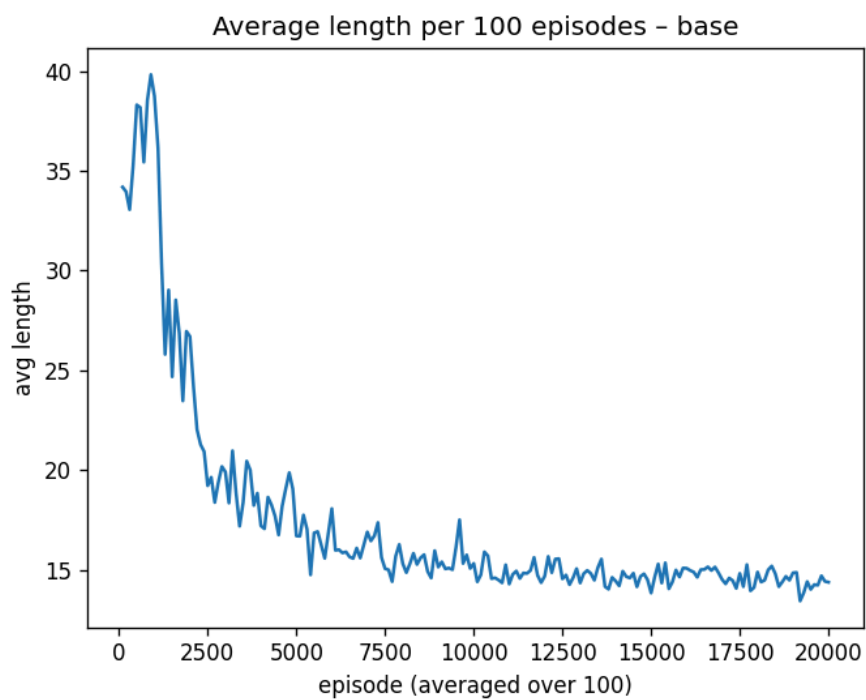
The epsilon decay rate significantly influenced exploration behavior and learning outcomes. The base decay rate (0.9999) allowed sufficient exploration while gradually transitioning to exploitation, resulting in successful learning.

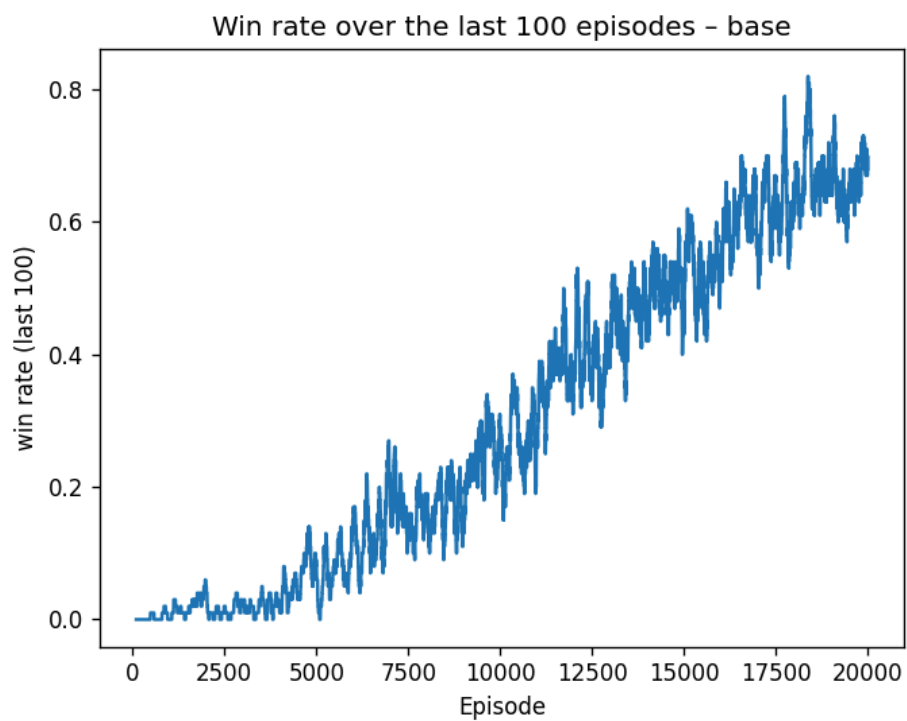
Interestingly, the faster decay rate (0.9995) led to complete failure, likely because exploration was curtailed too quickly before the agent could discover effective paths to the goal. The agent reached the minimum epsilon value (0.1) around episode 4000, much earlier than in other experiments.

The slower decay rate (0.99995) still achieved a 100% win rate during evaluation but showed much slower convergence during training, with average rewards reaching only 0.27 by episode 20000, compared to 0.708 in the base case. This demonstrates that maintaining higher exploration rates for longer can delay convergence but still lead to optimal policies eventually.

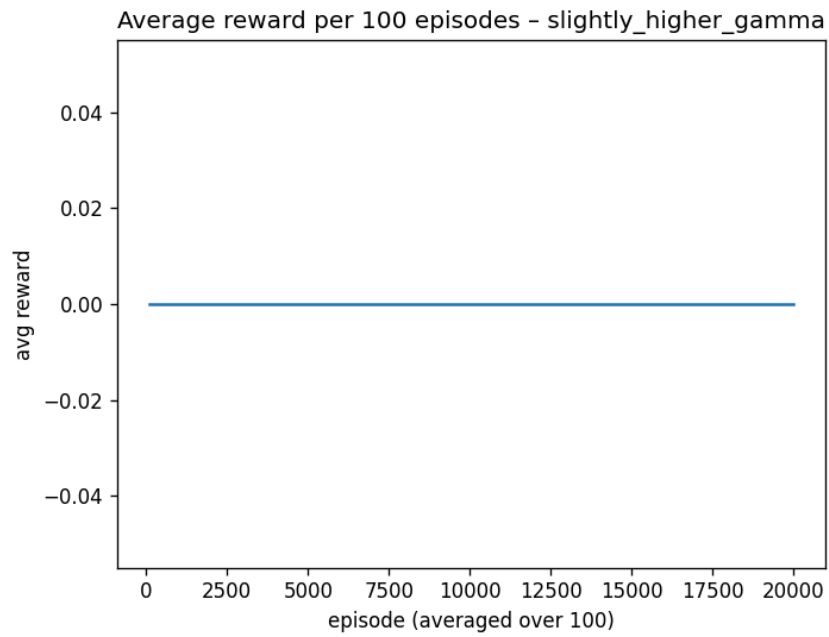
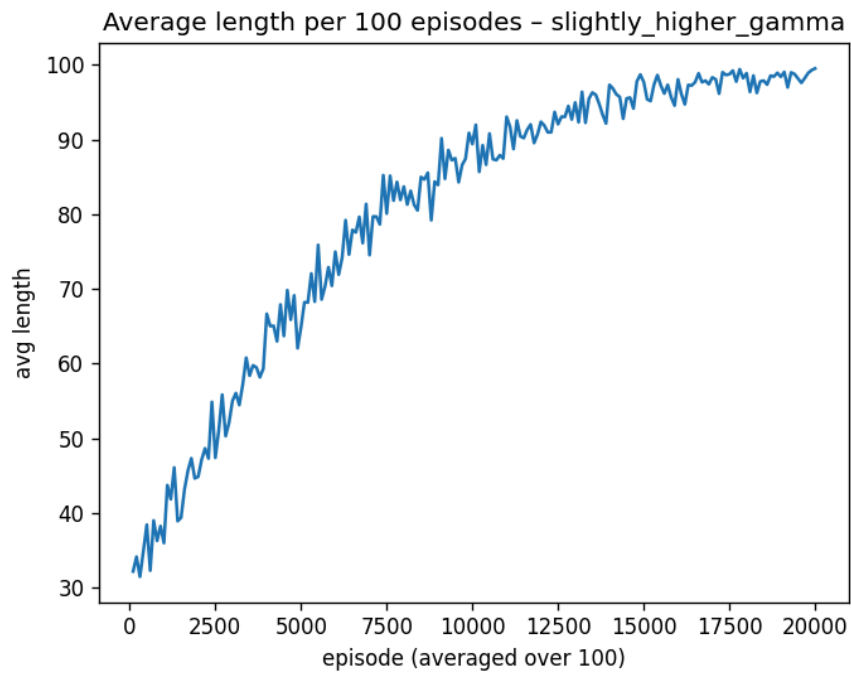
Discount Factor Variation

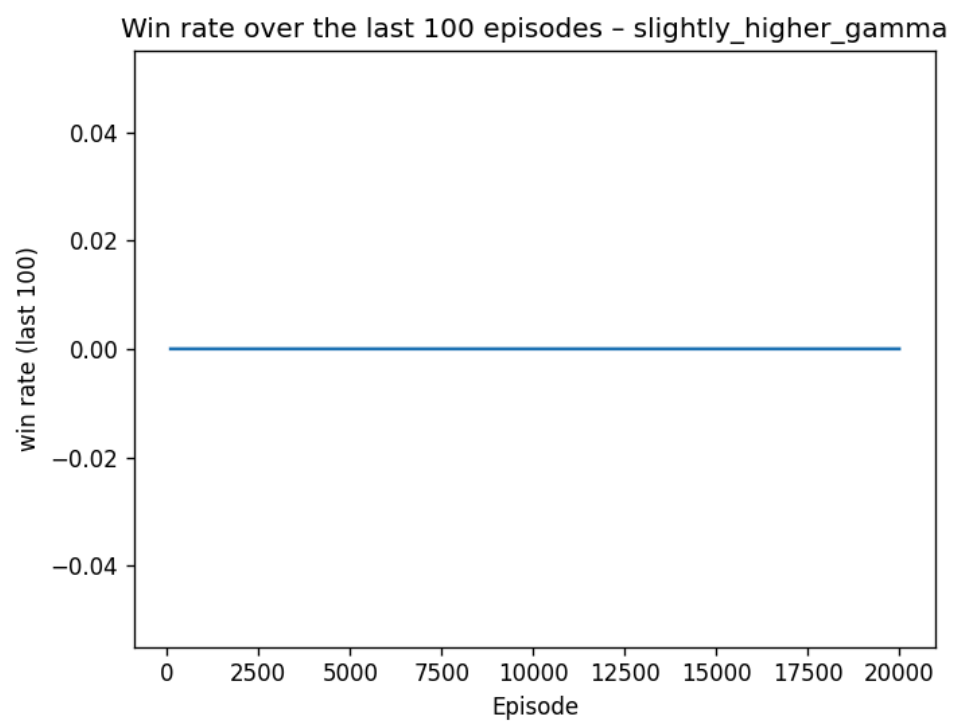
Base Case ($\gamma = 0.99$)



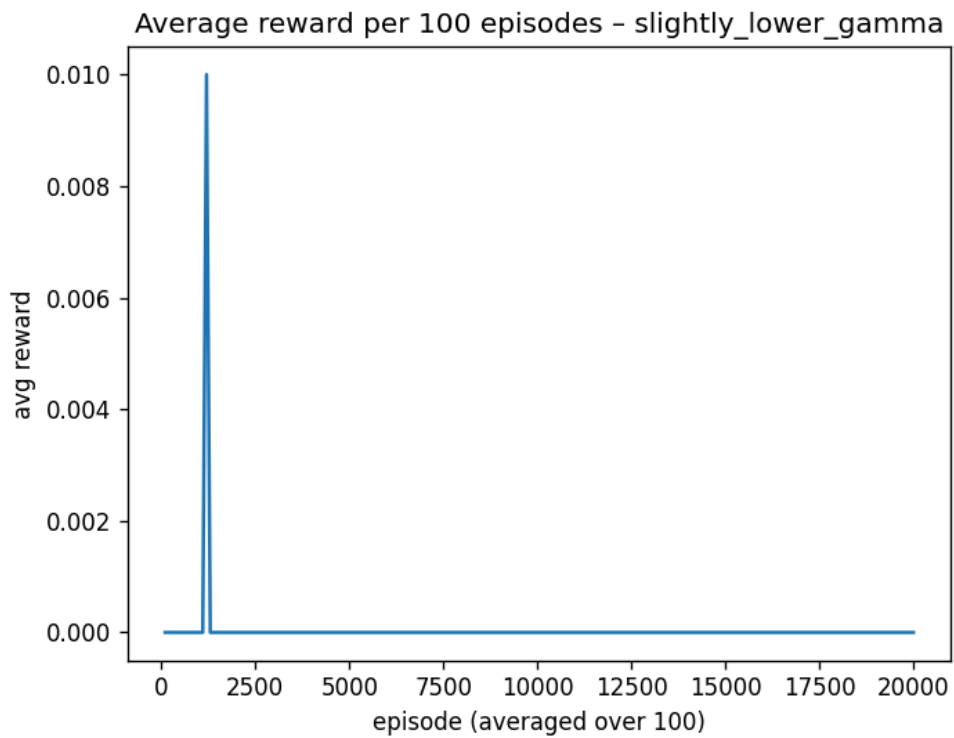
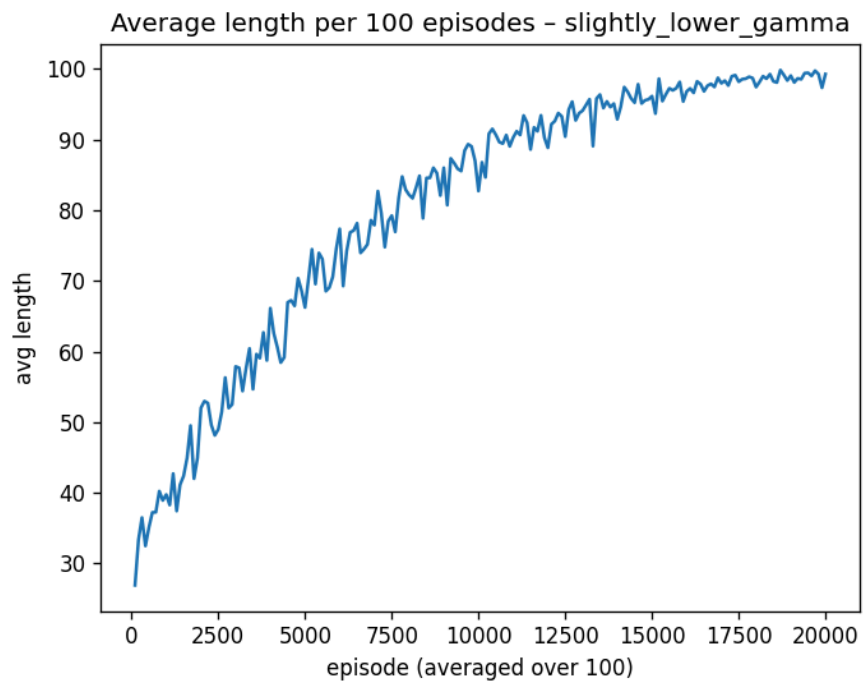


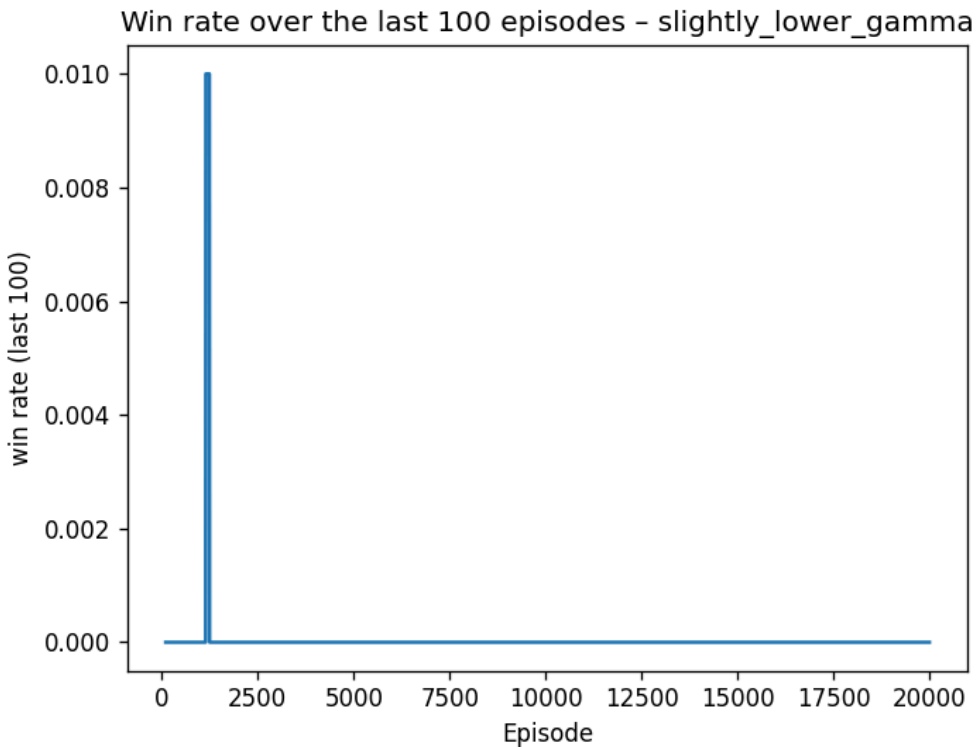
Slightly Higher Discount Factor ($\gamma = 0.995$)





Slightly Lower Discount Factor ($\gamma = 0.98$)





The discount factor determines how much the agent values future rewards. Interestingly, both the slightly higher (0.995) and slightly lower (0.98) discount factors resulted in complete learning failure (0% win rate) compared to the base case (0.99) which achieved perfect performance.

The results suggest that the optimal policy in this specific environment requires a precise valuation of future rewards, with $\gamma = 0.99$ striking the right balance.

Key Findings

1. **Hyperparameter Sensitivity:** The Q-learning algorithm exhibits extreme sensitivity to hyperparameters in the FrozenLake-8x8 environment. Only specific combinations led to successful learning, with even small deviations causing complete failure.
2. **Critical Learning Rate Threshold:** There appears to be a critical threshold for the learning rate between 0.7 and 0.8, below which learning fails to occur effectively.
3. **Balanced Exploration-Exploitation Trade-off:** The results emphasize the importance of a properly balanced exploration-exploitation schedule. Decaying exploration too quickly or maintaining too much persistent exploration both led to suboptimal performance.
4. **Precise Discount Factor Required:** The discount factor showed surprising sensitivity, with the apparently "ideal" value of 0.99 outperforming both slightly higher and lower values significantly.

5. **Convergence Patterns:** Successful configurations showed a characteristic pattern of increasing rewards and decreasing episode lengths.

Conclusion

This study demonstrates the critical importance of careful hyperparameter tuning for Q-learning in discrete environments. The FrozenLake-8x8 environment, despite its apparent simplicity, requires precise parameter settings to achieve optimal performance.

The base configuration ($\alpha=0.8$, $\gamma=0.99$, $\text{decay}=0.9999$, $\epsilon_{\text{min}}=0.1$) proved most effective, achieving a perfect 100% win rate during evaluation. This suggests that these parameters provide an optimal balance of learning rate, future reward discounting, and exploration-exploitation trade-off for this particular environment.