

Design of Pac-Man Strategies with Embedded Markov Decision Process in a Dynamic, Non-Deterministic, Fully Observable Environment

Jiachang (Ernest) Xu

CONTENTS

I	Introduction	1
II	Properties of Environment	1
II-A	Dynamic Environment	1
II-B	Non-Deterministic Environment . . .	2
II-C	Fully Observable Environment	2
III	Markov Decision Process	2
III-A	Bellman's Equation	2
III-A.1	Generalized Formula . . .	2
III-A.2	Abstract Adaption	2
III-A.3	Detailed Adaptation	2
III-B	Initialize Data Structures	2
III-B.1	Inactive Ghostbuster Mode	3
III-B.2	Defensive Ghostbuster Mode	3
III-B.3	Offensive Ghostbuster Mode	3
III-C	Update Utilities	3
III-D	Definition of Convergence	3
III-D.1	Convergence Tolerance . .	4
III-D.2	Early Stopping Point . . .	4
IV	Workflow of Strategies	4
IV-A	Mapping Operation	4
IV-B	Value Iteration	4
IV-C	Maximum Expected Utility	4
V	Methodology and Evaluation	5
V-A	Early Warning System	5
V-A.1	Surprise Attack	5
V-A.2	Decaying Threat	5
V-A.3	Preliminary Finding	5
V-B	Object-Oriented Design	6
V-B.1	Debug Mode Master Switch	6
V-B.2	Logging for Fine Tuning .	6
V-C	Parameter Tuning	6
V-C.1	Discount Factor	6
V-C.2	Ghostbuster Mode	7
V-C.3	Safety Distance	7
V-C.4	Optimal Parameter Setting	7
VI	Conclusion	7
	References	8

Abstract—This paper iterates the design of Pac-Man strategies, whose decision-making protocol is solely based on Markov Decision Process, without the support of pathfinding algorithms nor heuristic functions, in a dynamic, non-deterministic, fully observable environment. This project provides the rare opportunity to refine the understanding of, and practice the application of Markov Decision Process in a classic arcade game setting. After significant number of hours of parameter tuning, my design achieved a win rate ranging between 50% and 60%. With the proven effectiveness of embedded Markov Decision Process, a spin-off Pac-Man AI project that incorporates the advantages of pathfinding algorithms, heuristic functions, and Markov Decision Process shall be on the agenda.

I. INTRODUCTION

Imagine that you can see into the future. Would this be valuable for your decision-making process at the moment? My answer is "yes". With the embedded Markov Decision Process, the Pac-Man agent of my design, named `MDPAgent`, is able to predict the utility of taking an action in a dynamic, non-deterministic, fully observable environment. In other word, the `MDPAgent` quantifies the usefulness of moving to a neighbor location by adapting Bellman's Equation to this project. This project refines my understanding and application of the Markov Decision Process, and value iteration of Bellman's Equation, and in turn proves the effectiveness of Markov Decision Process in decision-making tasks in a non-deterministic environment. This paper will provide a detailed discussion with visual aids about (1) properties of Pac-Man environment in Section II, (2) adaptation of Markov Decision Process to this project in Section III, (3) my design of strategies and workflow of `MDPAgent` in Section IV, and (4) the methodology to optimize and evaluate the effectiveness of my design in Section V.

II. PROPERTIES OF ENVIRONMENT

This section will discuss the properties of the environment, dynamic, non-deterministic, and fully observable, that this project works within in the three following subsections.

A. Dynamic Environment

The environment of this project is dynamic, because the ghosts are always on the move. One major difference between the game state at timestamp t and $t + 1$ is the location of the ghost. For example (Fig. 1), the ghost can moves either east or west.

Wall	Wall	Wall	Wall	Wall	Wall	Wall
Wall		Agent				Wall
Wall		Wall	Wall	Wall		Wall
Wall		Wall	Food			Wall
Wall		Wall	Wall	Wall		Wall
Wall	Food	← Ghost →				Wall
Wall	Wall	Wall	Wall	Wall	Wall	Wall

Fig. 1. Example of possible movements of the ghost

B. Non-Deterministic Environment

The environment of this project is non-deterministic, because the agent is not guaranteed to move in the direction in which it intends to move. More precisely, the probability for the agent to actually move in the direction in which it intends to move is 80%. There is a 10% probability each for the agent to move in the direction that is either left or right of its intended direction. Additionally, if the direction in which the agent actually moves leads to a wall, the agent stops at the original location.

C. Fully Observable Environment

The environment of this project is fully observable, because the version of `api.py` that this project works with provides information of every component in the maze.

III. MARKOV DECISION PROCESS

Markov Decision Process is the core technique that this project uses to design Pac-Man strategies in a dynamic, non-deterministic, fully observable environment. According to the book *Artificial Intelligence: A Modern Approach* (Third Edition) by Russel and Norvig, generalized Markov Decision Process consists of four key features:

- S : a set of states
- $A(s)$: a set of actions available at state s
- $\Pr(s'|s, a)$: a transition model that dictates that probability that action a from state s at timestamp t leads to state s' at timestamp $t + 1$
- $R(s)$: a reward function

This section will discuss how this project adapts Bellman's Equation to the Pac-Man problem in Subsection III-A, and how to perform value iteration on Bellman's Equation through Subsection III-B to III-D. The discussion of value iteration will include (1) how to initialize data structures of rewards and utilities for different components in the maze under different circumstances, (2) how to update utility values for all the non-wall locations in the maze, and (3) definition of convergence.

A. Bellman's Equation

1) *Generalized Formula*: Bellman's Equation is the center piece of the embedded Markov Decision Process in my design of Pac-Man strategies. The generalized Bellman's Equation (Formula 1) updates the utility value of state s based on the action a , which belongs to the set of available actions $A(s)$ at state s and results in maximum expected utility. The discount factor γ dictates how much Bellman's Equation sees into the future, which is a parameter that can be fine-tuned to maximize win rate through experimentation (tested in Paragraph V-C.1).

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} \Pr(s'|s, a) U(s') \quad (1)$$

2) *Abstract Adaption*: To adapt the generalized Bellman's Equation to the Pac-Man problem, Formula 2 updates the utility value for each non-wall location l at iteration i based on the maximum expected utility of its eastern, western, northern, and southern neighbors at iteration $i - 1$. The eastern, western, northern, and southern neighbor locations of a location l are denoted as e_l , w_l , n_l , and s_l .

$$U_i(l) = R(l) + \gamma \max[EU_{i-1}(e_l), EU_{i-1}(w_l), EU_{i-1}(n_l), EU_{i-1}(s_l)] \quad (2)$$

3) *Detailed Adaption*: Because of the non-deterministic nature of the environment (defined in Subsection II-B), by expanding the expression of expected utility of each neighbor location of a location l , Formula 3 gives the detailed adaption of Bellman's Equation to the Pac-Man problem. In case that a neighbor location of a location l is a wall, the utility of the neighbor location of the location l at iteration $i - 1$ (i.e. $U_{i-1}(e_l)$, $U_{i-1}(w_l)$, $U_{i-1}(n_l)$, or $U_{i-1}(s_l)$) shall be substituted by the utility of the location l itself (i.e. $U_{i-1}(l)$).

$$U_i(l) = R(l) + \gamma \max[0.8U_{i-1}(e_l) + 0.1U_{i-1}(n_l) + 0.1U_{i-1}(s_l), 0.8U_{i-1}(w_l) + 0.1U_{i-1}(s_l) + 0.1U_{i-1}(n_l), 0.8U_{i-1}(n_l) + 0.1U_{i-1}(w_l) + 0.1U_{i-1}(e_l), 0.8U_{i-1}(s_l) + 0.1U_{i-1}(e_l) + 0.1U_{i-1}(n_l)] \quad (3)$$

B. Initialize Data Structures

The first step of value iteration on Bellman's Equation is to initialize data structures of rewards and utilities for different components in the maze. The utility of each non-wall location is initialized as 0. To counter the problem of surprise attack, I invented an early warning mechanism, based on a technique called decaying threat, to preemptively evade from the ghosts that are hiding behind the foods (defined later in Subsection V-A). In summary, the early warning system keeps a safety distance between the agent and each ghost. Fig. 2 provides an example of reward initialization.

Wall	Wall	Wall	Wall	Wall	Wall	Wall
Wall	Free 0.0	Agent 0.0	Free 0.0	Free 0.0	Free 0.0	Wall
Wall	Free 0.0	Wall	Wall	Wall	Free 0.0	Wall
Wall	Free 0.0	Wall	Food +10.0	Free 0.0	Free 0.0	Wall
Wall	Free 0.0	Wall	Wall	Wall	Free 0.0	Wall
Wall	Food +10.0	Free -450.0	Ghost -500.0	Free -450.0	Free 0.0	Wall
Wall	Wall	Wall	Wall	Wall	Wall	Wall

Fig. 2. Initialized rewards of smallGrid layout at timestamp 0

Because the version of `api.py` that this project works with provides the edible-or-hostile state of each ghost, my design incorporates 3 different ghostbuster modes to experiment with to maximize win rate in the `mediumClassic` layout. The three following paragraphs will iterate in details about how different choices of the parameter `ghostbuster` mode initialize their corresponding rewards differently. Table I compares the key behavioral features of different ghostbuster modes. Table II lists how each ghostbuster mode initializes rewards for different components in the maze. The parameter `ghostbuster` mode is one that can be fine-tuned to maximize win rate through experimentation (tested in Paragraph V-C.2).

TABLE I

BEHAVIORAL COMPARISON OF DIFFERENT GHOSTBUSTER MODE

Ghostbuster mode:	Inactive	Defensive	Offensive
Evade from hostile ghosts?	Yes	Yes	Yes
Aim to eat edible ghosts?	No	Yes	Yes
Aim to eat capsules?	No	No	Yes

TABLE II

REWARD INITIALIZATION OF DIFFERENT GHOSTBUSTER MODE

Ghostbuster mode:	Inactive	Defensive	Offensive
R(Capsule) =	0.0	0.0	+50.0
R(Food) =	+10.0	+10.0	+10.0
R(Free space) =	0.0	0.0	0.0
R(Edible ghost) =	0.0	+400.0	+400.0
R(Hostile ghost) =	-500.0	-500.0	-500.0

1) *Inactive Ghostbuster Mode*: The inactive ghostbuster mode always runs away from the ghosts, no matter they are edible or hostile, and doesn't intentionally aim to eat the capsules. Therefore, the inactive ghostbuster mode only initializes rewards for foods, free spaces, and ghosts: +10.0 for each food; 0.0 for each free space; -500.0 for each ghost and decaying threats for its surrounding spaces.

2) *Defensive Ghostbuster Mode*: The defensive ghostbuster mode always aims to eat edible ghosts as a priority, and runs away from hostile ghosts, but doesn't intentionally aim to eat the capsules. Therefore, the defensive ghostbuster

mode initializes rewards for free spaces, and hostile ghosts: 0.0 for each free space; -500.00 for each hostile ghost and decaying threats for its surrounding spaces. The defensive ghostbuster mode also initializes rewards for either edible ghosts or foods, but never for both at the same time: +400.0 for each edible ghost; +10.0 for each food.

3) *Offensive Ghostbuster Mode*: The offensive ghostbuster mode always aims to eat edible ghosts as a priority, runs away from hostile ghosts, and aggressively aims to eat the capsules. Therefore, the offensive ghostbuster mode initializes rewards for free spaces, and hostile ghosts: 0.0 for each free space; -500.00 for each hostile ghost and decaying threats for its surrounding spaces. The offensive ghostbuster mode also initializes rewards for either capsules, or edible ghosts, or foods, but never for more than one type of them at the same time: +50.0 for each capsule; +400.0 for each edible ghost; +10.0 for each food.

C. Update Utilities

The second step of value iteration on Bellman's equation is to update utilities for all the non-wall locations until convergence. Each iteration makes a deep copy of the `utilities` data structure, named `previous_utilities`. It uses the adapted Bellman's Equation (Formula 3) to update the `utilities` (iteration i) data structure, based on the `previous_utilities` (iteration $i - 1$) data structure.

The technique of value iteration repeats the above process of updating utilities until convergence. An example of utility convergence is provided (Fig. 3). The next subsection will provide a computationally efficient definition of convergence.

Wall	Wall	Wall	Wall	Wall	Wall	Wall
Wall	Free +1.535	Agent +0.836	Free +0.556	Free +1.024	Free +1.882	Wall
Wall	Free +3.097	Wall	Wall	Wall	Free +3.562	Wall
Wall	Free +6.069	Wall	Food +24.99	Free +13.63	Free +6.981	Wall
Wall	Free +11.89	Wall	Wall	Wall	Free +3.805	Wall
Wall	Food +23.29	Free -498.7	Ghost -840.2	Free -511.1	Free +0.492	Wall
Wall	Wall	Wall	Wall	Wall	Wall	Wall

Fig. 3. Convergent utilities of smallGrid layout at timestamp 0

D. Definition of Convergence

To make the value iteration work in a computational efficient way, we need a definition of convergence. Because all the reward and utility values assigned are decimal, to reach exact convergence, meaning the difference of the utility value of a location l at iteration i and that at iteration $i - 1$ for each location equals exactly 0, is very computationally costly. Therefore, I came up with two techniques, (1) convergence tolerance and (2) early stopping point, to reduce computational cost.

1) *Convergence Tolerance*: The convergence tolerance is an arbitrary parameter, whose value is of the user's choice, to balance the trade-off between the exactness of convergence and the computational cost of value iteration on Bellman's Equation. If the sum of the differences of the utility value for location l at iteration i and that at iteration $i - 1$ is below the convergence tolerance, we mark that the utility values are fully convergent across the maze. The smaller the convergence tolerance is, the closer the final convergent utility values are to the theoretical convergent utility values. The convergence tolerance of my choice is 0.1.

2) *Early Stopping Point*: The early stopping point is another arbitrary parameter, whose value is of the user's choice, to balance the trade-off between the exactness of convergence and the computational cost of value iteration on Bellman's Equation. For an entire cycle of value iteration, even though the technique of convergence tolerance is in place, it might still possess the possibility of non-convergence within a reasonably small number of iterations. Therefore, by putting a ceiling on the number of iterations, value iteration can stop early. The ceiling is called early stopping point. The early stopping point of my choice is 100 for non-sparse targets, such as foods, and 200 for sparse targets, such as capsules and edible ghosts.

IV. WORKFLOW OF STRATEGIES

The Pac-Man agent I designed is named `MDPAgent`. With the Markov Decision Process embedded, the overall workflow of `MDPAgent`'s strategies is divided into three major functional components in sequence (Fig. 4): (1) mapping operation, (2) value iteration, and (3) decision-making based on maximum expected utility. The three following subsections will discuss these three major functional components with visual aid of UML-style diagrams.

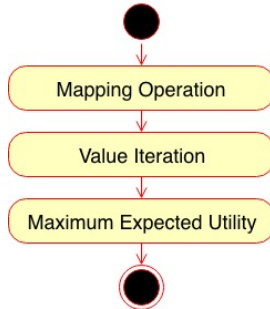


Fig. 4. Activity diagram of `MDPAgent.getAction()`

A. Mapping Operation

The first major functional component of `MDPAgent`'s workflow is the mapping operation. Although this project returns to the setting of a fully observable environment, maintaining a collection of internal memories about some aspects of information about the maze helps reduce computational redundancy. More specifically, that is to prevent from repeatedly calling `api.py` for information about the static or predictable components of the environment. The

static components of the environment, such as locations of walls, corners, and navigable spaces, are those that are fixed once the maze is instantiated. the predictable components of the environment, such as locations of capsules and foods, are those that can be tracked as long as the initial record of these components and travel record of the `MDPAgent` are provided constantly.

The mapping operation (Fig. 5) works as follows. At the initial state of each round, the `MDPAgent` initializes the internal memories to store and track information about the static and predictable components of the environment by calling the relevant functions from `api.py`. Then, every time the `getAction()` method is invoked, the mapping operation logs the game state history, and finds the location of the agent. If the location of the agent belongs to the internal memories of available capsules or foods, the mapping operation removes this location from the corresponding internal memories.

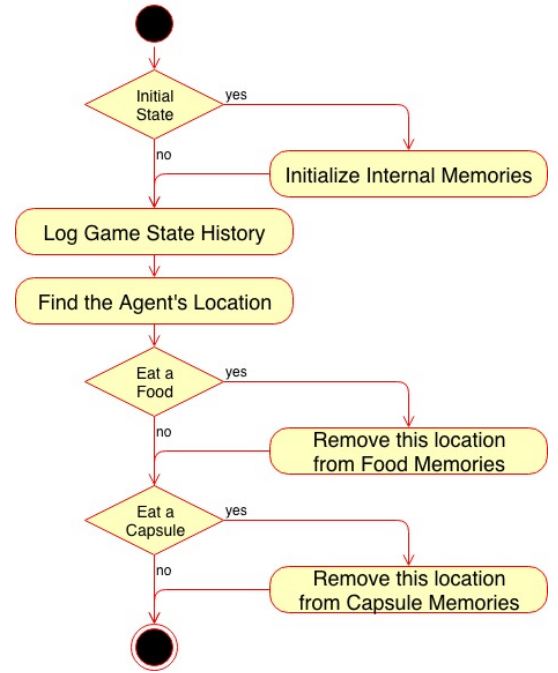


Fig. 5. Activity diagram of Mapping Operation

B. Value Iteration

The second major functional component of `MDPAgent`'s workflow is to apply value iteration on Bellman's Equation. To begin with, the `MDPAgent` initializes data structures of rewards and utilities for different components in the maze (defined in Subsection III-B). Then, the `MDPAgent` updates utility values for all the non-wall locations in the maze (defined in Subsection III-C), until reaching either full convergence or early stopping point (defined in Subsection III-D).

C. Maximum Expected Utility

The third major functional component of `MDPAgent`'s workflow is the decision-making process based on maximum

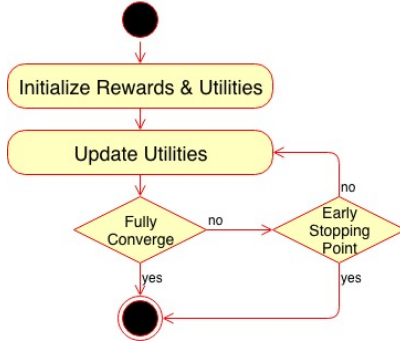


Fig. 6. Activity diagram of Value Iteration

expected utility. After the cycle of value iteration on the current game state, which means reaching either full convergence on all the free spaces or early stopping point, the MDPAgent uses the rule of maximum expected utility to decide which neighbor location to move towards (Fig. 7). More specifically, the MDPAgent looks at all of its neighbor locations, to which all the legal actions (the stop option removed) lead to, and chooses the neighbor location with the highest value of expected utility calculated by value iteration. Then, the MDPAgent returns the legal action that leads to that neighbor location.

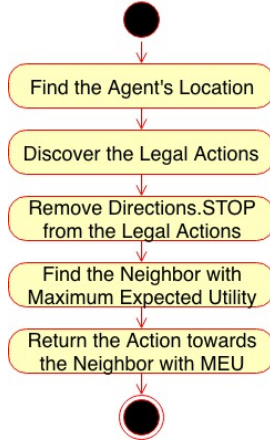


Fig. 7. Activity diagram of Maximum Expected Utility

V. METHODOLOGY AND EVALUATION

Due to non-deterministic nature of this project, the first version of my design had a catastrophic win rate of less than 20%. Therefore, I used my creativity to optimize my design, and analyzed the effectiveness of these new features based on the data generated by running thousands of rounds of games. The following subsections will discuss (1) the early warning system to evade from ghosts, (2) modular programming of object-oriented design, and (3) how I fine-tuned the key parameter setting of my MDPAgent strategies.

A. Early Warning System

In the very early stage of this project, I simply initialized reward values to different components of the environment

as follows: +10.0 for each food; 0.0 for each free space; -500.0 for each ghost. However, this naive way of initializing rewards only resulted in a win rate around 20%. After analyzing the printed utility values of the entire maze for those failed rounds, I discovered that the main cause was the surprise attack by the ghosts, and I provided a countermeasure to surprise attack, called decaying threat. The two following paragraphs will give the definitions of surprise attack and decaying threat.

1) *Surprise Attack*: Because the technique of value iteration on Bellman's Equation only updates utilities based on the neighbor with maximum expected utility, if the ghosts are hiding inside a cluster of foods (Fig. 8), it is impossible for the information of the ghosts to penetrate the barrier of foods and reach the agent. Therefore, the agent is not aware of the ghosts until it is too close to the ghosts to evade from them. To solve this problem, the agent needs an early warning system to preemptively evade from the ghosts.

Wall	Free 0.0	Free 0.0	Free 0.0	Free 0.0	Wall
Wall	Free 0.0	Wall	Wall	Wall	Wall
Wall	Free 0.0	Food +10.0	Ghost -500.0	Food +10.0	Free 0.0
Wall	Free 0.0	Wall	Food +10.0	Wall	Free 0.0
Wall	Free 0.0	Wall	Agent 0.0	Wall	Free 0.0

Fig. 8. Surprise Attack: ghosts hiding behind foods

2) *Decaying Threat*: I created an early warning system by initializing decaying threats for each ghost's surrounding spaces. Starting with -500.0 for each ghost, the early warning system recursively initializes negative reward values at an arbitrary decay rate as the radius to the ghost grows. The early warning system stops the decaying threats at a safety distance (a parameter that needs to be fine-tuned through experimentation, tested in Paragraph V-C.3). For example (Fig. 9), the decay rate is 100.0; the safety distance is 5 steps. The agent knows from which direction the ghost will attack it at that moment, at the maximum distance of 5. Therefore, in this example, the agent definitely would not move to east, because it is more dangerous (i.e. more negative utility).

3) *Preliminary Finding*: Table III finds that the early warning system of my design, which uses the solution of decaying threats, increases the stabilized win rate by 30 percentage point. This is merely a preliminary finding. The effectiveness of the early warning system will be formally tested in Paragraph V-C.3.

TABLE III
EVALUATE THE EFFECTIVENESS OF EARLY WARNING SYSTEM

	No decaying threat	Decaying threat
Stabilized win rate	20%	50%

Wall	Free 0.0	Free 0.0	Free 0.0	Free 0.0	Wall
Wall	Free 0.0	Wall	Wall	Wall	Wall
Wall	Agent 0.0	Food -100.0	Food -200.0	Food -300.0	Food -400.0
Wall	Free 0.0	Wall	Free -100.0	Wall	Ghost -500.0
Wall	Free 0.0	Wall	Free 0.0	Wall	Free -400.0

Fig. 9. Decaying Threat: a solution to counter surprise attack

B. Object-Oriented Design

In this project, I further put modular programming into effect. This good practice of object-oriented design improves modularity and readability of my code structure (Fig. 10). I breaks down the `getAction()` method, a level-1 function, into three level-2 functions, which correspond to the three major functional components of the overall workflow of `MDPAgent`'s strategies (defined in Section IV). Furthermore, the `_value_iteration()` method, a level-2 function, is broken down into three level-3 functions, including two functional sub-components, which represent the core functionalities inside value iteration, and a helper function that print the rewards or utilities data structure on a user-friendly display. This design of function call hierarchy provides convenience for debug mode and log mode, which I will discuss in the two following paragraphs.

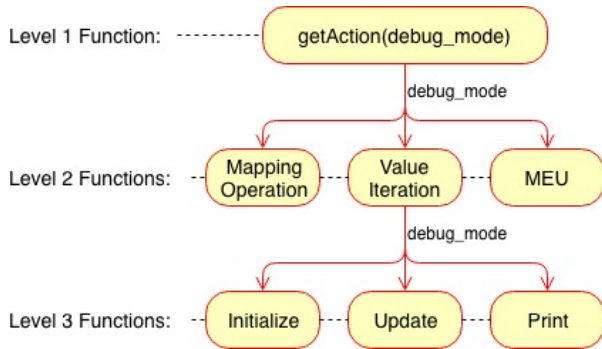


Fig. 10. Function call hierarchy of `MDPAgent.getAction()`

1) *Debug Mode Master Switch*: My design of the 3-level function call hierarchy (Fig. 10) allows me to embed a master switch for debug mode. I added an optional input parameter `debug_mode` for the `getAction()` method, and a required input parameter `debug_mode` for every level-2 function. The level-1 setting of `debug_mode` therefore automatically flows to all the level-2 functions. By changing the default value of the level-1 input parameter `debug_mode`, the `MDPAgent` switches between active and inactive debug mode across all its member methods, like flipping a master switch. If the debug mode is active, the `getAction()` method automatically print every well-formatted debug message, including the user-friendly display of the rewards

or utilities data structure. This design is more efficient and more elegant than spending time on toggling comments on all the `print()` lines.

2) *Logging for Fine Tuning*: Because the game calls the `final()` method on `MDPAgent` at the end of each round, I added an optional input parameter `log_mode` for the `final()`. By changing the default value of the input parameter `log_mode`, the `MDPAgent` switches between active and inactive log mode. If the log mode is active, the `MDPAgent` logs the parameter (i.e. discount factor, convergence tolerance, ghostbuster mode, etc.) setting of this round, along with the result (i.e. win or loss) and score to a `log.csv` file on the append mode. The log file allows me to aggregate all the game result and analyze the performance of different parameter settings.

C. Parameter Tuning

This subsection analyzes the performance of various parameter settings. There are three key parameters in my design that needs fine-tuning to find their optimal values to maximize win rate in both the `smallGrid` layout and the `mediumClassic` layout through experimentation: (1) the discount factor inside Bellman's Equation, (2) the ghostbuster mode that dictates the agent's behavioral pattern, and (3) the safety distance to initialize reward values for a ghost's surrounding spaces. I wrote a separate `test-case-generator.py` code file to generate all the combinations of selected values of the three key parameters (Table IV). Because there are no capsules present in the `smallGrid` layout, the tuning of ghostbuster mode is not applicable for the `smallGrid` layout. In total, there are 40 and 120 different combinations of test cases for the `smallGrid` layout and the `mediumClassic` layout respectively. I ran 50 rounds of games on each test case, and logged the record of each round into a log file for parameter analysis. The total number of testing data entry is 2,000 and 6,000 for the `smallGrid` layout and the `mediumClassic` layout respectively. The three following paragraphs will analyze each one of these key parameters above separately, including the reasons for their selected values for testing, and discuss the findings based on data gathered through experimentation.

TABLE IV
SELECTED VALUES OF THE THREE KEY PARAMETERS

Key parameter	Selected values
Discount factor	0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, 0.0
Ghostbuster mode	inactive, defensive, offensive
Safety distance	4 steps, 3 steps, 2 steps, 1 step

1) *Discount Factor*: The parameter discount factor (defined in Subsection III-A) dictates how much the agent sees into the future. The bigger the value of discount factor is, the more the agent sees into the future, and the longer it takes to converge. The parameter tuning of discount factor includes 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, and 0.0, which is a reasonable sampling from the domain of discount

factor between 0 (inclusive) and 1 (exclusive). According to data gathered through experimentation, discount factor being 0.6 results in the highest average win rate of 60.50% on the `smallGrid` layout; discount factor being 0.9, 0.8 and 0.7 results in the top 3 highest win rates between 34% and 35%. Because less time is required to converges for a lower discount factor, discount factor being 0.7 is preferred.

TABLE V
PARAMETER ANALYSIS OF DISCOUNT FACTOR ON `smallGrid`

Discount factor	Average win rate
0.9	46.00%
0.8	53.50%
0.7	56.50%
0.6	60.50%
0.5	58.00%
0.4	59.00%
0.3	49.50%
0.2	16.00%
0.1	18.50%
0.0	9.50%

TABLE VI
PARAMETER ANALYSIS OF DISCOUNT FACTOR ON `mediumClassic`

Discount factor	Average win rate
0.9	34.33%
0.8	34.83%
0.7	34.17%
0.6	32.33%
0.5	30.17%
0.4	23.67%
0.3	15.50%
0.2	11.50%
0.1	6.17%
0.0	0.00%

2) *Ghostbuster Mode*: The parameter ghostbuster mode (defined in Subsection III-B) dictates how aggressively the agent goes after the ghosts. The inactive ghostbuster mode is the least aggressive; the defensive ghostbuster mode is in the middle; the offensive ghostbuster mode is the most aggressive. The parameter tuning of ghostbuster mode includes all of its possible values. According to data gathered through experimentation, ghostbuster mode being inactive results in the highest average win rate of 32.15%. My first impression of ghostbuster mode is that the offensive ghostbuster mode should be the optimal choice, because it proactively destroys the ghosts. However, the experimentation tells the exactly opposite. One possible explanation for the optimal setting of ghostbuster mode being inactive is that, since eating all the foods wins the game, the inactive ghostbuster mode, which does not bother to hunt down capsules or ghosts, spends less time to win a game on average, decreasing the chance being eaten by the ghosts.

3) *Safety Distance*: The parameter safety distance (defined in Subsection V-A) dictates the range of the early warning system. The larger the value of safety distance is, the farther away the agent anticipates the incoming ghosts. The parameter tuning of safety distance only includes 4 steps, 3 steps, 2 steps, and 1 step, because any safety distance

TABLE VII
PARAMETER ANALYSIS OF GHOSTBUSTER MODE ON `mediumClassic`

Ghostbuster mode	Average win rate
Inactive	32.15%
Defensive	17.95%
Offensive	16.70%

larger than or equal to 5 steps is over-caution that might drastically increase the running time of each round. If the safety distance is 1 step, it is equivalent to initializing utility values without the technique of decaying threats. According to data gathered through experimentation, safety distance being 2 steps results in the highest average win rate of 50.40% on the `smallGrid` layout; safety distance being 4 steps results in the highest average win rate of 31.80% on the `mediumClassic` layout. Because for both the `smallGrid` layout and the `mediumClassic` layout, safety distance being 1 step results in a significantly lower win rate than any other options, it formally proves the effectiveness of the early warning system.

TABLE VIII
PARAMETER ANALYSIS OF SAFETY DISTANCE ON `smallGrid`

Safety distance	Average win rate
4 steps	40.60%
3 steps	42.60%
2 steps	50.40%
1 step	37.20%

TABLE IX
PARAMETER ANALYSIS OF SAFETY DISTANCE ON `mediumClassic`

Safety distance	Average win rate
4 steps	31.80%
3 steps	25.73%
2 steps	21.07%
1 step	10.47%

4) *Optimal Parameter Setting*: According to the findings of separate analyses of the three key parameters, the optimal parameter setting is specified in Table X. To evaluate the effectiveness of these optimal parameter settings, I ran 10,000 rounds on `smallGrid` layout plus 1,000 rounds on `mediumClassic` layout, achieving 61.34% and 56.40% win rate respectively (Table XI).

TABLE X
SPECIFICATION OF OPTIMAL PARAMETER SETTING

Key parameter	Optimal value for <code>smallGrid</code>	Optimal value for <code>mediumClassic</code>
Discount factor	0.6	0.7
Ghostbuster mode	N/A	Inactive
Safety distance	2 steps	4 steps

VI. CONCLUSION

This project provides the opportunities to apply Markov Decision Process in a non-deterministic environment. My design of Pac-Man strategies gives a practical sense of

TABLE XI

EVALUATION OF OPTIMAL PARAMETER SETTING ON `MEDIUMCLASSIC`

Simulation configuration	Win rate
10,000 rounds on <code>smallGrid</code> layout	61.34%
1,000 rounds on <code>mediumClassic</code> layout	56.40%

effectiveness of Markov Decision Process in decision-making process. In addition, this project highlights the importance of modular programming of object-oriented design, which rewarded me with robust code structures, and the necessary infrastructure in place to help debugging and logging. Finally, data through statistically significant counts of experiments provide evidence to hypotheses and findings. Considering the advantages of both pathfinding algorithms from my last Pac-Man project, and Markov Decision Process from this project, I believe a spin-off project that incorporate the fast searching capability of pathfinding algorithms, and the look-into-future magic of Markov Decision Process should be on the agenda.

REFERENCES

- [1] Russell, Stuart J.; Norvig, Peter (2010), Artificial Intelligence: A Modern Approach (3rd ed.), Upper Saddle River, New Jersey: Pearson Education, Inc., ISBN 0-13-604259-7, chpt. 17.
- [2] Bellman, R. (1957). "A Markovian Decision Process". *Journal of Mathematics and Mechanics*. 6.
- [3] Moore, Edward F. (1959). "The shortest path through a maze". *Proceedings of the International Symposium on the Theory of Switching*. Harvard University Press. pp. 285–292.