

CS 221 Project Report

Applying Varied Artificial Intelligence Techniques to Play 2048

Zhiyang He, Charles Lu, Stephen Ou
{hzyjerry,clu8,sdou}@stanford.edu

December 11th, 2015

1 Introduction

Released in 2014, the single-player puzzle game 2048 [1] quickly became a worldwide sensation: in addition to hijacking classrooms and quickly occupying people's phone screens, it also spawned countless spin-offs, strategy guides, and self-professed gurus. Meanwhile, the game offers a good platform to experiment and compare both traditional and cutting edge artificial intelligence techniques.

2048 is played on a 4 by 4 grid. Every turn, a random tile with a value 2 or 4 will appear on an empty cell. Then the player can choose to slide left, right, up, or down. All the tiles will slide all the way towards that direction. If any two tiles collide and have the same value, they will merge into one tile, with a new value being the sum of the two old tiles. Additionally, when two tiles collide, the score will be incremented by the value of the newly merged tile. There are two ways to end the game. If the value of one tile reaches 2048, the game is considered a win (though the player can continue playing). Otherwise, if all the cells on the grid are occupied with mismatched tiles and no move is possible, the game is lost.

2 Task Definition

Our goal for this project is to build and investigate a number of models which play to maximize the game score at a reasonable runtime. The input-output behavior is as follows: on each turn, given the state of the 2048 game board (as well as move history when necessary), our models will return one of four moves (left, right, up, down) to attempt to maximize the score.

Next, we would like to discuss how we model the

game of 2048 into specific game states. For each game state, there are two things to keep track of: the board and the score. The board is a four by four matrix, and each cell contains an integer (that is a power of 2) that indicates the current value of that cell. 0 is used to indicate unoccupied cells. The score variable is used to keep track of all the points received so far. The rule of the game of 2048 states that score is incremented by the sum of the two tiles when two tiles with the same value get merged.

There are two main methods that are available through the 2048 game state: `getLegalActions()` and `generateSuccessor()`. We will describe each of them in details below.

`getLegalActions()`: If the current agent is the human, there are four possible moves. The human can swipe left, right, up, or down. There is one special case. A move is considered invalid if the board in the successor state is the same as the current state. For example, if the left three columns are all filled with tiles and they have different values, swiping left is not a valid action because the successor state will not change. Next, if the current agent is the computer, there are maximum of sixteen possible moves. The computer can add a new tile with a value of 2 into an unoccupied cell.

`generateSuccessor()`: If the current agent is the human, all the tiles will slide in the direction specified. If the two neighboring tiles in that direction have the same value, they will collide and form a new tile with a new value that is the sum of the two old values. For example, if the board currently consists of only two tiles, both with a value of 2, at the bottom row. Swiping left will result in a merged tile with a value of 4, sitting in the bottom left corner. Next, if the current agent is the computer, a new tile with the value of 2 will be added

at the specified cell given by the row and column number.

In order to maximize the game score of 2048, we initially built six agents, three of which are baseline algorithms and three of which are more advanced search algorithms that we learned in class.

- 1) Random agent. This baseline agent simply picks a random move.
- 2) Up down agent. This baseline agent alternates between moving up and down.
- 3) Up left agent. This baseline agent alternates between moving up and left.
- 4) Expectimax agent. This applies the expectimax agent that we learned in class.
- 5) Minimax agent. This applies the minimax agent that we learned in class.
- 6) Minimax agent with alpha-beta pruning. This applies the alpha-beta pruning agent that we learned in class.

For the expectimax, minimax, and alpha-beta pruning agents, which require a way to score the current game state, we experimented with a number of evaluation functions. Though the purpose of this project was not to find the ideal evaluation function or tune perfect hyperparameters for the evaluation functions we wrote, the best evaluation function we used with the expectimax agent (with a depth of 3) had a win rate of 100 percent (getting to 2048 on every game).

Though the expectimax agent yielded very good results, the main disadvantage of using it was its computation time. With a branching factor of up to 4 for the players move and up to 16 for the computers move (placing an empty tile), it is quite infeasible to search the game tree to a large depth without pruning. Specifically, on a modern laptop, it took the expectimax agent with depth 3 almost a second per move.

Therefore, we next attempted to use the optimal moves suggested by the expectimax agent (depth 3) to train various reflex models using different machine learning techniques. Specifically, to generate the training, validation, and test datasets, we simply ran multiple playthroughs of 2048 using the search agent, and saved each game board state and the corresponding recommended move as a single data point.

In addition, we want to apply Q-learning algorithm to implement a reinforcement learning agent. This algorithm has not been applied to solving 2048 game puzzle. Inspired by the Google Deep Learning algorithm [2] solving the Atari game, we want to apply similar approach to tackle 2048 puzzle, using Q learning algorithm based on neural network. The learning algorithm has the advantage of needing no prior knowledge of the game (e.g. how 4-value random grids and 2-value random grids are distributed).

3 Infrastructure

We built a frontend interface for 2048 setup. We forked the original 2048 repository by Gabriele Cirulli and added our custom Javascript functions that talk to a lightweight Python server that computes the optimal move. The setup used Flask [3], an open source Python web framework. We wrote a JavaScript function that serializes the current board as a string and passes it to the Flask server via an AJAX request. The server computes the optimal move using a specified approach, and returns a response back to the frontend. Then, a callback JavaScript function updates the board using the optimal move.

One problem we ran into while doing simulation is that the HTTP request and overhead of displaying the front-end visualization was a bottleneck in terms of speed. While it is interesting to see the game being solved in the real user interface, the speed becomes problematic when we wanted to do a lot of simulation to get abundant results.

Therefore, we built a more robust backend-only simulator so we can run the game quickly. It uses the logic in gameState.py to generate a successor based on a move picked by the agent specified. It starts from a board with only 1 tile and outputs the final score and number of moves once all tiles have been filled. Without involving the frontend which requires a lot of back and forth HTTP request, the backend simulator was able to finish one full iteration of the game in few seconds.

To further increase the speed of the simulation, we decided to take advantages of the myth machines on Stanford campus. We parallelized the game simulation across 30 machines and ran them concurrently to get results. We were able to run 100 full iterations of the 2048 game across all 30 servers and gather 100 data points in less than 1

minute.

4 Approaches

We applied several different agents to the 2048 game, each using a different technique either a strategy specific to 2048, a search strategy using the game tree, a reflex strategy, or a reinforcement learning strategy.

For each strategy, we simulate the game for 300 times and record their scores. Then we compute detailed statistics so we can compare which strategy performs the best.

4.1 Basic Strategies

The first agent that we wrote was a random agent that simply picks a random move without considering the game state. This is meant to serve as a baseline.

Minimum Score	140
Maximum Score	1456
Mean Score	571.255
Median Score	546
Standard Deviation	280.527

Table 1: Random Agent

The up down agent alternates between moving up and down. This method simulates a player who does not play randomly and simply swipes up and down.

Minimum Score	64
Maximum Score	156
Mean Score	94.775
Median Score	88
Standard Deviation	18.791

Table 2: Up Down Agent

The up left agent alternates between moving left and right. This method is employed by many novice player as a shortcut to get to a high score. This is meant to trap the highest value tile towards the top left corner and make it easy to combine similar valued tiles.

Minimum Score	656
Maximum Score	1708
Mean Score	946.457
Median Score	796
Standard Deviation	261.353

Table 3: Up Left Agent

4.2 Search Strategies

In terms of search strategies, we tried out expectimax, minimax, and minimax with alpha-beta pruning. In the normal game of 2048, the computer move (where the new tile is inserted) is randomized, so technically only the expectimax strategy makes sense in that environment. However, we wanted to experiment with minimax strategy as well to see how the score differs if the computer aims to minimize the players score when inserting a new tile.

Because of the number of path to explore in the game state is huge, our initial experimentation uses a depth of 2. Once the depth reaches 2, we utilized the several evaluation functions that we developed to provide an accurate estimate of how good the current game board is.

Regarding the evaluation functions once the specified depth has been reached, we have tried four different ones:

1. Snake evaluation function. This evaluation function values tiles on the top left corner and devalues tiles on the bottom right corner. In effect, this moves the higher value tiles towards the top left and lower value tiles towards the bottom right.
2. Game score evaluation function. This evaluation function looks at the current score of the board and uses that to determine the value of the board.
3. Monotonicity evaluation function. This evaluation function measures how much the rows and cols are sorted in increasing or descending order. Specifically we count the total number of times the rows and cols switch from increasing to decreasing or vice versa. Smaller is better so we take its reciprocal.
4. Smoothness evaluation function. This evaluation function sums the difference between each pair of adjacent tiles. Again, smaller is better so we take its reciprocal.

We use a similar version of the expectimax algorithm as the one presented in lecture. If the current agent is the human, we take the action that produces the maximum score from the recursion tree. If the current agent is the computer, we randomly pick one action (to be more specific, an unoccupied cell location).

Overall, the first two turned out to be a success and the last two did not work out so well.

Minimum Score	5580
Maximum Score	33900
Mean Score	20293.433
Median Score	20712
Standard Deviation	2956.264

Table 4: Expectimax Agent with depth 2 and snake evaluation function

Minimum Score	1496
Maximum Score	16364
Mean Score	7397.181
Median Score	7168
Standard Deviation	3250.412

Table 5: Expectimax Agent with depth 2 and game score evaluation function

The monotonicity and smoothness evaluation functions unfortunately did not give an accurate estimate of the board, as the mean score hovers around 390 for both of them.

However, we received astonishing results once we increased the depth to 3. This is the **best** results out of all the agent, depth, and evaluation function combination.

Minimum Score	81900
Maximum Score	172426
Mean Score	137770
Median Score	158984
Standard Deviation	48849.40651

Table 6: Expectimax Agent with depth 3 and snake evaluation function

This result makes sense because by considering an extra level, the agent has $4 \cdot 16 = 64$ times more game states to pick an optimal move from. This enables a more intelligent choice which was able to increase our mean score significantly.

As mentioned above, even though the minimax algorithm does not fit the game perfectly because in the real 2048 game the computer picks a move random, we still want to validate the fact that the score generated by the minimax agent is lower than that of the expectimax agent with the same depth and evaluation function because the computer will attempt to minimize the score in minimax. Our assumption is correct.

We use a similar version of the minimax algorithm as the one presented in lecture. If the current agent is the human, we take the action that produces the maximum score from the recursion tree. If the current agent is the computer, we take the action (considering all possible insertions of new tiles) with the minimum score in the recursion tree.

Minimum Score	2600
Maximum Score	21036
Mean Score	13943.141
Median Score	14804
Standard Deviation	5279.458

Table 7: Minimax Agent with depth 2 and snake evaluation function

We use a similar version of the alpha beta pruning algorithm as the one presented in lecture. If the current agent is the human, we take the action that produces the maximum score from the recursion tree. However, if at any point as we loop through all the actions, the beta value is smaller or equal to the alpha value, we can terminate and return the maximum action so far. If the current agent is the computer, we take the action (considering all possible insertions of new tiles) with the minimum score in the recursion tree. Similarly, if at any point as we loop through all the actions the beta value is smaller or equal to the alpha value, we can terminate and return the minimum action so far.

Minimum Score	2368
Maximum Score	24676
Mean Score	15064.333
Median Score	15660
Standard Deviation	4782.109

Table 8: Minimax with Alpha-Beta Pruning Agent with depth 2 and snake evaluation function

4.3 Reflex Strategies

However, as mentioned above, the computation required for a depth of 3 is quite high, taking up to one second for each move. Therefore, as summarized above, we used training data from the expectimax agents moves to train various reflex models in Theano and Torch. We had a total of over 20,000 data points (consisting of a input in R^{16} and a move label corresponding to left, up, right, down).

We first implemented and trained a multi-class logistic regression model in Theano. In essence, we trained four logistic regressions (one for each move) with the entire dataset, corresponding to labels of 0 or 1 depending on whether the move was recommended or not. During prediction, we simply fed the board state into each model and went with the move with the highest expected probability, or the next highest scoring move if the top move was invalid.

Minimum Score	320
Maximum Score	400
Mean Score	375.2
Median Score	382
Standard Deviation	23.61

Table 9: Logistic regression classifier with no feature transform (raw tile value features)

As demonstrated by the poor scores (even worse than random!), even with validation during training, the practical test error rate was very high.

The same technique was tried, but using the softmax function $\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$ for $j = 1, \dots, K$ [4], yielding similar results.

Noticing that the inputs (tiles on the game board) were all powers of two, we considered applying a simple feature transformation: taking the base-2 logarithms of all features. However, this actually worsened the prediction results:

Minimum Score	268
Maximum Score	272
Mean Score	270.4
Median Score	270.4
Standard Deviation	23.61

Table 10: Logistic regression classifier with feature transform (raw tile value features)

Specifically, with a single-layer multi-class clas-

sifier using logistic regression, using the feature transform resulted in the predicted moves to always be upagain, likely a result of dataset bias as well as nonlinearity.

The above results are obviously very poor, but make sense given the hypothesized nonlinearity of the problem. Therefore, we next constructed various neural network models using Torch with the same 16 inputs and 4 outputs, but an additional layer with a variable number of neurons. As with before, sigmoids were used in the first transformation and a softmax for the outputs. However, for all models, regardless of the size of the hidden layer, the model again failed to generalize and outputted similar prediction results as above.

4.4 Reinforcement Learning Strategies

We want to obtain an agent that maximizes the accumulative future reward. This feature is described by optimal Q value:

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

, which describes the maximum sum of rewards discounted by γ at future steps t , achieved by a behaviour policy $\pi(s|a)$ after making an observation (s) and making an action (a). Q value (action-value) contains all pertinent information about the playing field along with a possible action. Q function (utility function) is used to describe how well our Q values approximates the actual game utilities: $\hat{Q}_{\text{opt}(s,a)} \leftarrow (1 - \eta) \underbrace{\hat{Q}_{\text{opt}(s,a)}}_{\text{prediction}} + \eta \underbrace{(r + \gamma \hat{V}_{\text{opt}(s')})}_{\text{target}}$.

This value can be computed by using the Q Learning algorithm we learned in class. In Q-learning algorithm, action utilities are evaluated based on immediate reward gained from taking actions, with the possibility of a delayed reward led to by the action. Due to the large state space of 2048 puzzle ($O(12^n)$), our data cannot possibly exploit all the possibilities. Thus we use a Q value based on neural network that covers the unknown states:

By linking our trainer with reinforcement learning code using `popen()`, `write()` and `read()` in python, we apply real-time reinforcement learning to our game playing. Every time a game state is generated, we feed it into our Q-Learning Brain object, the Brain outputs the optimal movement for

the game. Then a reward is calculated based on the move, and provided back to the brain as feedback.

We use a framework based on Torch 7 [5] called DeepQLearning [6] to implement this process. The DeepQLearning framework provides the following Q-Learning interface:

- `Brain.init(num_inputs, num_outputs):` initialize the brain
- `action = Brain.forward(state):` generate movement based on game state
- `Brain.backward(reward):` train the neural network based on reward

In addition, the framework provides entry points to further configure the Brain Object. Here we provide our customization of 2048 puzzle game Brain:

- `Brain.temporal_window = 0.` Due to instantaneous randomization, we take no past state/action pairs
- `Brain.experience_size = 10000.` Maximum number of experiences that we will save for training, this leads to a storage of 20Mb. Further training experiences are not stored, but used to replace former ones.
- `Brain.gamma = 0.99.` Decay factor that controls how much plan-ahead the agent does. Our choice leads to a factor of 0.36 after every 100 move (number of moves for a game)
- `Brain.epsilon = 1.0, Brain.epsilon_min = 0.05.` Purely randomized to highly deterministic policy at end

Our reward function is

$$\text{reward} = \begin{cases} 100000000 & \text{if game is won} \\ 100000000 & \text{if new highest number found} \\ -\text{avg_score} & \text{if game is lost} \end{cases}$$

Based on former strategies, a randomized agent has average performance of 571 points, with maximum 1456. By applying the reinforcement learning, our agent quickly surpasses this benchmark and begins generating scores up to 800 900 frequently.

Around 5 minutes after training begins, the agent starts producing grid with 128 in most of the experiments. Around 12 20 minutes after training,

an interesting pattern can be observed: the agent starts gathering grids with large numbers and keep them close together.

The performance of reinforcement learning stabilizes at average 1200 points (highest 4548) after 2 hours of training.

5 Literature Review

One of the most widely known 2048 AIs was published by Matt Overlan, who used an algorithm with iterative deepening depth first alpha-beta search. The evaluation function tries to keep the rows and columns monotonic (either all decreasing or increasing) while aligning same-valued tiles and minimizing the number of tiles on the grid.

6 Error Analysis and Future Works

6.1 Evaluation functions

Though the goal of this project was not to find the ideal evaluation function for 2048, we saw how critical a good evaluation function was to yielding good performance. In this project, we experimented with several different evaluation functions for the expectimax agent, and the snake evaluation function was able to give us the best success.

However, we do believe that the monotonicity, smoothness, empty tiles, and the raw score evaluation functions can be beneficial. Though we tried using a weighted combination of all five evaluation functions, with the hyperparameters we tried, the variance in total scores was greatly increased and the mean score was only sometimes marginally increased. However, to better tune the evaluation function, in the future, we can potentially better tune these hyperparameters for a combined evaluation function using a parameter sweep methodology. Additionally, we can explore other evaluation functions based on deeper domain knowledge.

6.2 Reinforcement Learning

Even though the reinforcement learning agent achieved high scores that can never be obtained by randomized movements, the average performance is still low compared to our Expectimax Agent. We

think that this is partly due to the randomized nature of 2048 game. Because random grid are created at every move, the game variability is much bigger than Atari, or other similar puzzles. Our dataset (2k games) is hardly representative of the vast possibilities.

On the other hand, compared to the large game state space, there exists only one right way of solving the puzzle (lining up values in snake-like fashion). It is very hard for reinforcement learning algorithms to find this approach simply based on randomized exploration. Even if it happens to achieve a high score in one game, the effect is averaged out by large number of low-score-low-reward game experiments.

For possible future work, we think that better reward function should be invented. The most fitting reward function should be able to filter out the few correct movements out of large number of randomized trials.

6.3 Machine learning

The performance of the reflex agents was far below expected. We can see several areas for improvement:

- 1) The dataset used to train the reflex agents was likely very biased. Specifically, we used several playthroughs of 2048 using the expectimax agent, saving every single board state and its corresponding recommended move. However, later in the game, larger-valued tiles or tiles which are locked in tend to stay in the same position for hundreds of moves. For instance, considering that the evaluation function tends to put the largest tile in the top left square, there tend to be two times as many data points containing a largest tile of value n than data points containing a largest tile of value $\frac{n}{2}$. Therefore, we hypothesize that such tiles repeatedly used in the training set failed to allow the reflex models to generalize. This was especially true in practice, since the vast majority of data points were from late in the game, while the reflex agent always failed early on.
- 2) Our models also likely had high bias. As demonstrated by our results, the process of picking moves in 2048 is clearly nonlinear, which explains the failure of single-layer models to generalize the move-making process. Furthermore,

data from the expectimax agent with depth 3 likely had much deeper strategy than could be modelled by such small networks. (For instance, our team members noticed many signs of deeper thinking in the depth 3 agent as compared to even the depth 2 agent.) Therefore, to address the failure of our neural networks with one hidden layer to generalize the move-making process, we could look into larger and more advanced models, as well as explore different techniques.

- 3) Rather than simply unrolling the board into a vector in R^{16} and expecting our models to learn the underlying patterns to the moves, we could try using more descriptive features as inputs. For instance, we could take inspiration from our monotonicity, smoothness, and num-empty evaluation functions and use those as features to give our models more insight into the actual game state, rather than simply giving the raw tile values.
- 4) Instead of deterministically choosing the move with the highest output, we could easily implement an agent to randomly pick a move weighted by the outputs. This would be a simple way to avoid the problem of the models repeatedly picking the same move deterministically.

References

- [1] <https://github.com/ov3y/2048-AI>
- [2] <http://www.nature.com/nature/journal/v518/n7540/full/nature14236.html>
- [3] <http://flask.pocoo.org/>
- [4] https://en.wikipedia.org/wiki/Softmax_function
- [5] <http://torch.ch/>
- [6] <https://github.com/blakeMilner/DeepQLearning>