



國立陽明交通大學

NATIONAL YANG MING CHIAO TUNG UNIVERSITY

# Deep Learning

## Lab2: Temporal Difference Learning

Name: 許子駿

Student ID: 311551166

Institute of Computer Science and Engineering

2023 年 4 月 9 日

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Details</b>	<b>3</b>
2.1	$n$ -tuple network . . . . .	3
2.2	TD(0) algorithm . . . . .	5
2.3	TD-backup diagram . . . . .	6
2.4	Action selection . . . . .	7
2.5	Implementation details . . . . .	7
<b>3</b>	<b>Results</b>	<b>11</b>
<b>4</b>	<b>Conclusion</b>	<b>13</b>

# Chapter 1

## Introduction

This task is to implement Temporal Difference (TD) learning algorithm to solve 2048 game using an  $n$ -tuple network.

# Chapter 2

## Details

Details are listed in this chapter including 5 parts,  $n$ -tuple network, TD(0) algorithm, TD-backup diagram, action selection, and implementation details.

### 2.1 $n$ -tuple network

1. Isomorphic patterns are created in constructor of class `pattern` (See Listing 2.1).

```
1 for (int i = 0; i < 8; i++) {
2     board idx = 0xfedcba9876543210ull;
3     if (i >= 4) idx.mirror();
4     idx.rotate(i);
5     // Get indices of isomorphic patterns on the board.
6     for (int t : p) {
7         isomorphic[i].push_back(idx.at(t));
8     }
9 }
```

Listing 2.1: C++ code of class `pattern` constructor.

2. Get index of an isomorphic pattern of the board.

```
1 size_t indexof(const std::vector<int>& patt, const board& b) const {
2     // TODO
3     size_t idx = 0;
4     for (size_t i = 0; i < patt.size(); i++)
5         idx |= b.at(patt[i]) << (4 * i); // Shift left 16 bits.
6
7     return idx;
8 }
```

Listing 2.2: C++ code of `indexof` function of class `pattern`.

3. Values of isomorphic patterns are sum up to represent the value of the board (See Listing 2.3).

```
1 virtual float estimate(const board& b) const {  
2     // TODO  
3     // Sum up values of all isomorphic patterns.  
4     float sum_iso_patterns = 0;  
5     for (int i = 0; i < iso_last; i++)  
6         sum_iso_patterns += weight[indexof(isomorphic[i], b)];  
7  
8     return sum_iso_patterns;  
9 }
```

Listing 2.3: C++ code of `estimate` function of class `pattern`.

4. Update isomorphic patterns with mean error, and then return the updated weight (See Listing 2.4).

```
1 virtual float update(const board& b, float u) {  
2     // TODO  
3     // Update all isomorphic patterns with mean.  
4     float mean = u / iso_last;  
5     float sum_iso_patterns = 0;  
6     for (int i = 0; i < iso_last; i++) {  
7         size_t idx = indexof(isomorphic[i], b);  
8         weight[idx] += mean;  
9         sum_iso_patterns += weight[idx];  
10    }  
11  
12    return sum_iso_patterns;  
13 }
```

Listing 2.4: C++ code of `update` function of class `pattern`.

## 2.2 TD(0) algorithm

This algorithm, starting from second last move of the episode, calculates error between before-state of the move  $s$  and next move  $s''$ , then updates before-state of  $s$  with the error, and finally calculates TD target with sum of reward of  $s$  and updated weight of before-state of  $s$ . Keep the TD target for next step, repeat steps above till the start of the episode (See Figure 2.1 and Listing 2.5).

For each episode,

```

Initialize (before-)state  $s$ 
While  $s$  is not terminal do
   $a \leftarrow \operatorname{argmax}_a \text{EVALUATE}(s, a')$ 
   $r, s', s'' \leftarrow \text{MAKE\_MOVE}(s, a)$ 
  STORE( $s, a, r, s', s''$ )
   $s \leftarrow s''$ 
End While
For ( $s, a, r, s', s''$ ) from terminal down to initial do
  LEARN_EVALUATION( $s, a, r, s', s''$ )
End For

```

← perform TD backup

Figure 2.1: Pseudo-code of TD(0) algorithm.

```

1 void update_episode(std::vector<state>& path, float alpha = 0.1) const {
2     // TODO
3     // Starts from second last move, and update state backwards iteratively.
4     float target = 0.0f;
5     for (path.pop_back(); path.size(); path.pop_back()) {
6         state& move = path.back();
7         target = move.reward() +
8             update(move.before_state(), alpha * (target - (move.value() - move.
9             reward())));
10    }

```

Listing 2.5: C++ code of `update_episode` function of class `learning`.

## 2.3 TD-backup diagram

1. Before-state: use before-state of move  $s$  as value function, and use reward  $r$  and before-state of next move  $s''$  as TD target to update before-state of the move  $s$  (See Figure 2.2).

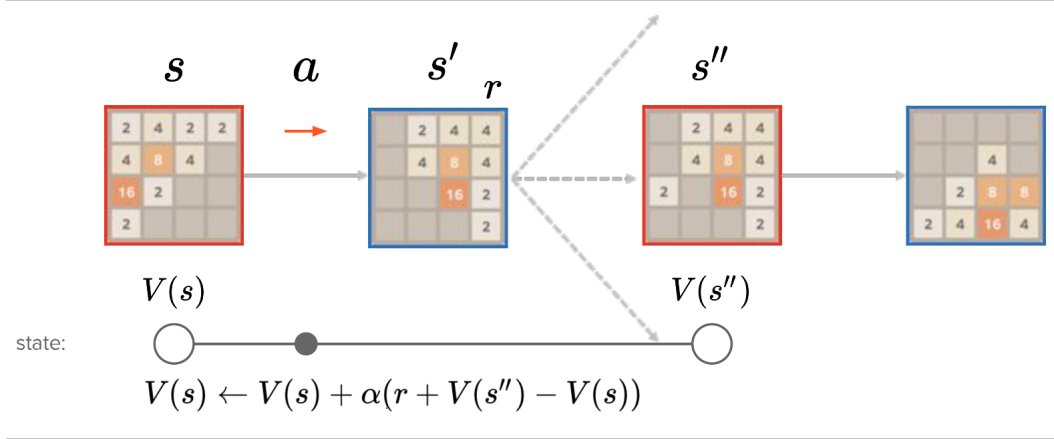


Figure 2.2: TD backup diagram of before-state.

2. After-state: use after-state of move  $s'$  as value function, and use reward  $r_{next}$  and after-state of next move  $s'_{next}$  as TD target to update after-state of the move  $s'$  (See Figure 2.3).

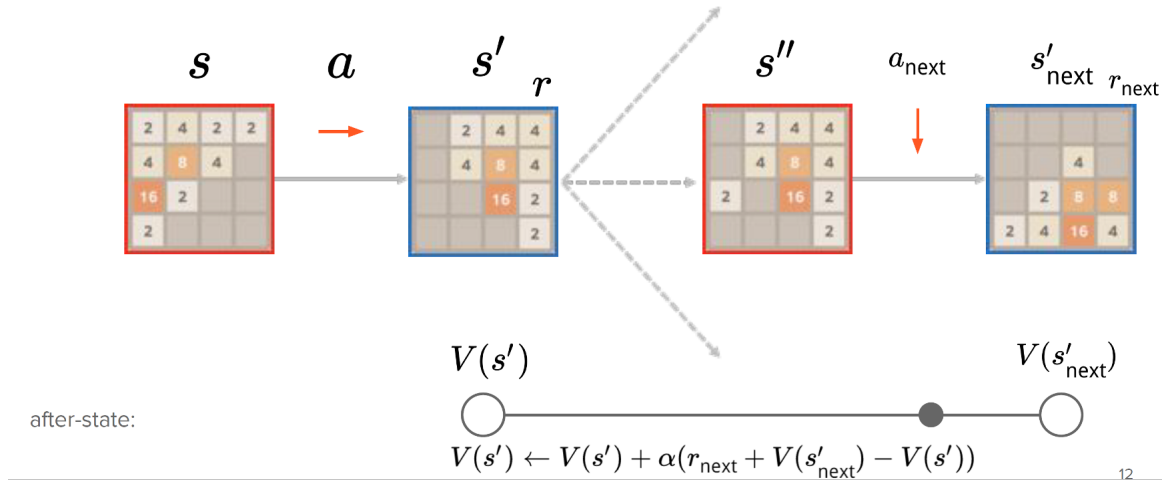


Figure 2.3: TD backup diagram of after-state.

## 2.4 Action selection

1. Before-state: after choosing an action  $A_t$ , randomly generate one 2-tile or 4-tile, then calculate and weighted sum up value functions of each new state  $S_{t+1}$  by the probability that 2-tile and 4-tile pop-up and number of empty tiles (See Figure 2.4).

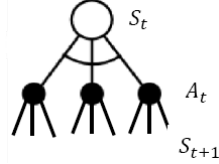


Figure 2.4: Action selection of before-state.

2. After-state: consider that before-state of next move  $S_{t+1}$  is fixed, and when choosing action of next move  $A_{t+1}$  is the same situation, so this approach chooses the action that results in highest value of after-state. (See Figure 2.5).

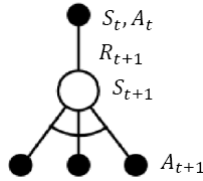


Figure 2.5: Action selection of after-state.

## 2.5 Implementation details

1. `indexof` (See Listing 2.2), `estimate` (See Listing 2.3), `update` (See Listing 2.4), `update_episode` (See Listing 2.5) functions have been listed in previous sections.
2. First, find all empty tiles during after-state, then randomly generate 2-tile or 4-tile, and calculate weighted sum to be the value of before-state of the move, and finally update the highest value for further comparisons, and return best move (See Listing 2.6).

```

1 state select_best_move(const board& b) const {
2     state after[4] = { 0, 1, 2, 3 }; // up, right, down, left
3     float weight_2_tile = 0.9f;
4     float weight_4_tile = 1.0f - weight_2_tile;
5

```



```

6   state* best = after;
7   float best_move_val = -std::numeric_limits<float>::max();
8   for (state* move = after; move != after + 4; move++) {
9       if (move->assign(b)) {
10          // TODO
11          board after_state = move->after_state();
12          int tiles[16];
13          int numEmptyTiles = 0;
14
15          // Find all empty tiles during after-state.
16          for (int i = 0; i < 16; i++)
17              if (after_state.at(i) == 0)
18                  tiles[numEmptyTiles++] = i;
19
20          // Find best move during after-state.
21          float val = move->reward();
22          for (int i = 0; i < numEmptyTiles; i++) {
23              board* cur_board = new board(uint64_t(after_state));
24
25              cur_board->set(tiles[i], 1); // 2-tile.
26              val += weight_2_tile * (estimate(*cur_board) /
27              numEmptyTiles);
28
29              cur_board->set(tiles[i], 2); // 4-tile.
30              val += weight_4_tile * (estimate(*cur_board) /
31              numEmptyTiles);
32          }
33
34          move->set_value(estimate(move->before_state()));
35
36          if (val > best_move_val) {
37              best = move;
38              best_move_val = val;
39          }
40          else {
41              move->set_value(-std::numeric_limits<float>::max());
42          }
43          debug << "test " << *move;
44      }
45      return *best;
46  }

```

Listing 2.6: C++ code of `select_best_move` function of class `learning`.

3. With original 4 patterns (See Listing 2.7 and Figure 2.6), the average reached 2048-tile percentage is about **90%** after 350K iterations. However, with additional 4 patterns, totally 8 patterns (See Listing 2.8 and Figure 2.7), it's about **95%** (Details are in Chapter 3).

```

1 tdl.add_feature(new pattern({ 0, 1, 2, 3, 4, 5 }));
2 tdl.add_feature(new pattern({ 4, 5, 6, 7, 8, 9 }));
3 tdl.add_feature(new pattern({ 0, 1, 2, 4, 5, 6 }));
4 tdl.add_feature(new pattern({ 4, 5, 6, 8, 9, 10 }));

```

Listing 2.7: C++ code of original 4 patterns.

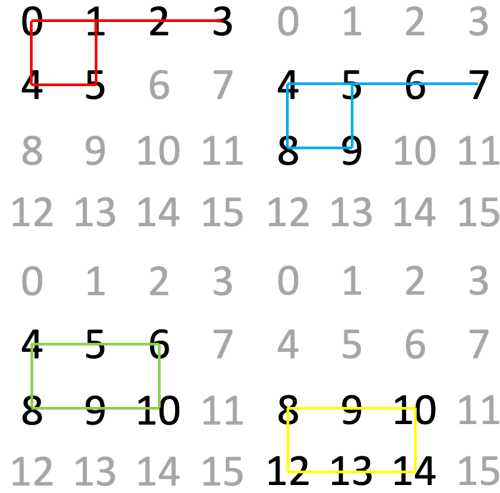


Figure 2.6: Illustration of original 4 patterns.

```

1 tdl.add_feature(new pattern({ 0, 2, 5, 10 }));
2 tdl.add_feature(new pattern({ 4, 6, 9, 14 }));
3 tdl.add_feature(new pattern({ 1, 4, 5, 6, 9, 13 }));
4 tdl.add_feature(new pattern({ 1, 4, 9, 14 }));

```

Listing 2.8: C++ code of new 4 patterns.

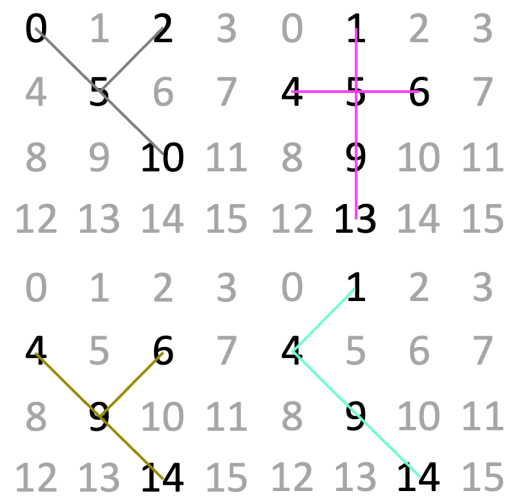


Figure 2.7: Illustration of new 4 patterns.

# Chapter 3

## Results

1. Mean score reaches about **77510** after 350K iterations training (See Figure 3.1).

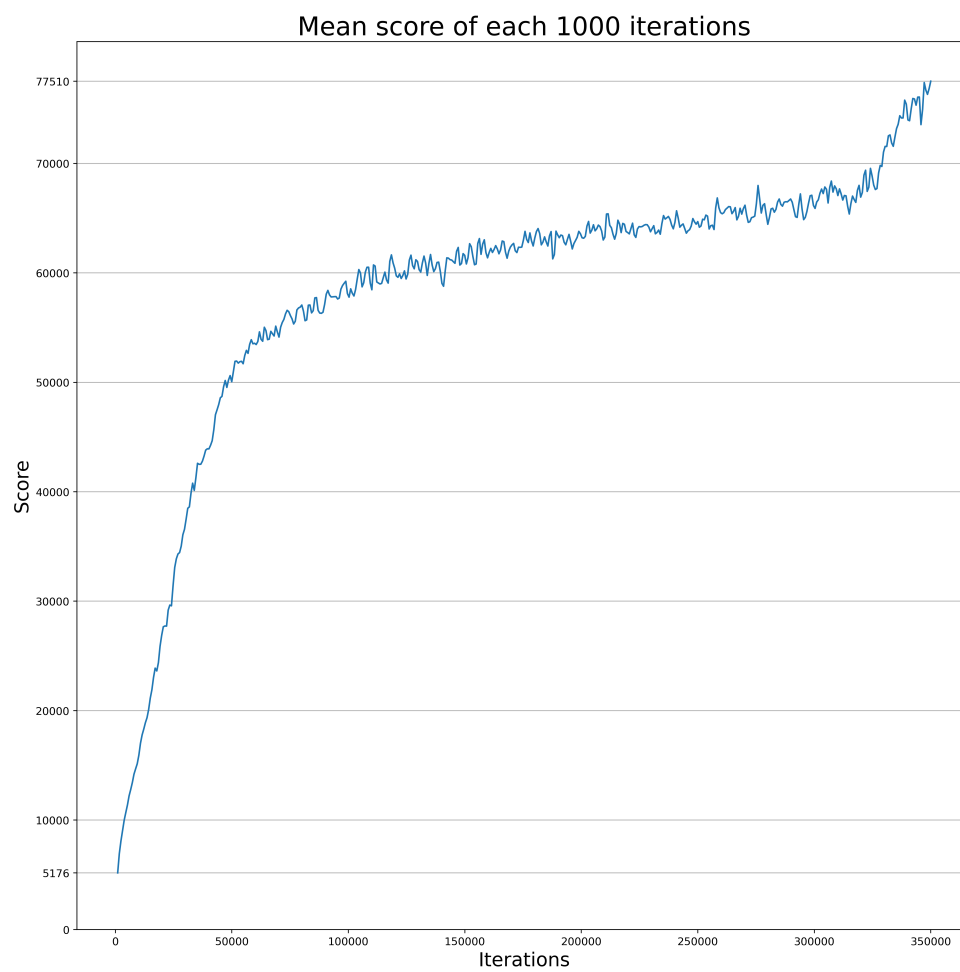


Figure 3.1: Mean score of each 1K iterations for total 350K iterations.

2. Average reached 2048-tile percentage is about **95%** after 350K iterations training (See Figure 3.2).

108000	mean = 76559.3	max = 155400
128	100%	(0.1%)
256	99.9%	(0.2%)
512	99.7%	(1.2%)
1024	98.5%	(4.9%)
2048	93.6%	(15.3%)
4096	78.3%	(54%)
8192	24.3%	(24.3%)
109000	mean = 76543	max = 159620
128	100%	(0.3%)
256	99.7%	(0.1%)
512	99.6%	(0.9%)
1024	98.7%	(3.2%)
2048	95.5%	(16.3%)
4096	79.2%	(55.8%)
8192	23.4%	(23.4%)
110000	mean = 77510.4	max = 154804
128	100%	(0.1%)
256	99.9%	(0.1%)
512	99.8%	(1.8%)
1024	98%	(3.7%)
2048	94.3%	(15.9%)
4096	78.4%	(53.4%)
8192	25%	(25%)

Figure 3.2: Score details after 350K training.<sup>1</sup>

---

<sup>1</sup>Note that the printed iterations here are wrong, since the model is trained 240K iterations at first and stop, then load the weight, and trained for 110K iterations for the second time. Totally 350K iterations training.

# Chapter 4

## Conclusion

After this experiment, I found that **patterns** play an important role in this task.

With original 4 patterns, the average reached 2048-tile percentage is only about **90%** after 350K iterations training. However, with another 4 patterns, totally 8 patterns are used, it reaches about **95%**.