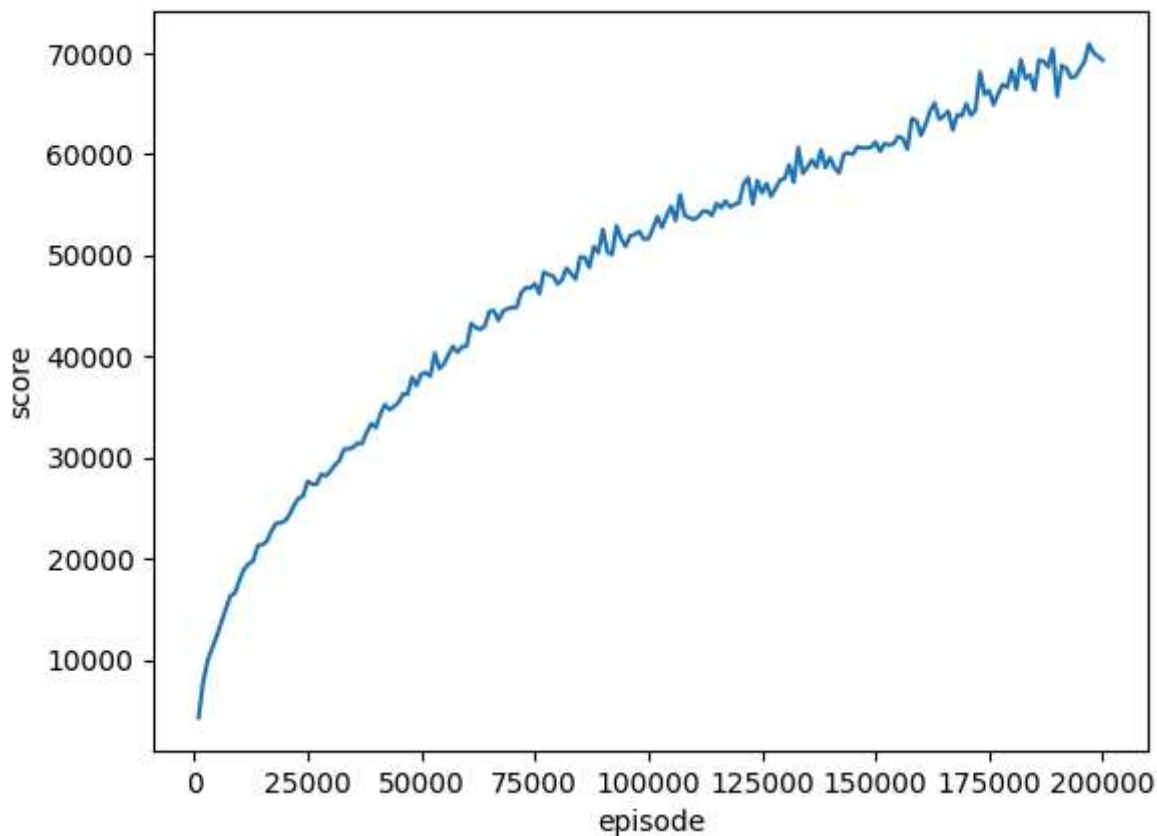


LAB1

A plot shows scores (mean) of at least 100k training episodes



Describe the implementation and the usage of n -tuple network.

因為 2048 的狀態總數非常巨大，所以需要 value function approximator 來估計盤面的 score，而 n -tuple network 在這次 LAB 就是作為這個角色。

n -tuple network 在 2048 的實作大致如下：

1. 取出幾個盤面位置形成一個 tuple
2. 以 tuple 為單位去建出一張表
3. 以盤面上對應位置的資訊去 encode 出 index，對表上該 index 進行查詢、更新

index of

- 分別取出盤面對應位置的值，並用左移將這些4-bit 的資訊 concat 成一個 mask (index)

```
1  size_t indexof(const std::vector<int>& patt, const board& b) const {
2      // TODO
3      size_t index = 0;
4      for (int t : patt) {
5          index = (index << 4) | b.at(t);
6      }
7      return index;
8  }
```

estimate

- 給定一個盤面，對所有該 pattern 的同構查表，並累加回傳

```
1  virtual float estimate(const board& b) const {
2      // TODO
3      float result = 0;
4      for (int i = 0; i < iso_last; i++) {
5          size_t index = indexof(isomorphic[i], b);
6          result += operator[](index);
7      }
8      return result;
9  }
```

update

- 給定一個盤面和更新值，先將更新值除以同構數，對每個同構進行更新，並累加更新後的值回傳



```
1  virtual float update(const board& b, float u) {
2      // TODO
3      float u_split = u / iso_last;
4      float value = 0;
5      for (int i = 0; i < iso_last; i++) {
6          size_t index = indexof(isomorphic[i], b);
7          operator[](index) += u_split;
8          value += operator[](index);
9      }
10     return value;
11 }
```

Explain the mechanism of TD(0).

TD(0) 是 Temporal Difference Learning 的一種方法，特色是不用等一個完整的 episode 結束，就可以利用 S_t 和 S_{t+1} 進行更新。

詳細的更新公式如下：

$$V(S_t) \leftarrow V(S_t) + \alpha[r_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Describe your implementation in detail including action selection and TD-backup diagram

next

- 我寫了一個 helper function，方便獲得一個當前 after-state popup 後可能的各個新盤面以及機率

```

1  std::vector<std::pair<board, float>> next(const board& b) const {
2      std::vector<std::pair<board, float>> next;
3      std::vector<int> spaces;
4      for (int i = 0; i < 16; i++) {
5          if (b.at(i) == 0) {
6              spaces.push_back(i);
7          }
8      }
9      static std::vector<std::pair<int, float>> tiles = { {1, 0.9}, {2, 0.1} };
10     for (auto &space : spaces) {
11         for (auto [tile, prob] : tiles) {
12             board next_board = b;
13             next_board.set(space, tile);
14             next.emplace_back(next_board, prob / spaces.size());
15         }
16     }
17     return next;
18 }

```

select best move

- 利用 helper function 獲得所有可能的 state 以及機率後，分別 estimate 這些 state 並以機率作為權重，即可得到 $V(S_{t+1})$ 的期望值，選一個最大的回傳。

```

1  state select_best_move(const board& b) const {
2      state after[4] = { 0, 1, 2, 3 }; // up, right, down, left
3      state* best = after;
4      for (state* move = after; move != after + 4; move++) {
5          if (move->assign(b)) {
6              // TODO
7              float value = 0;
8              std::vector<std::pair<board, float>> all_possible_state = next(move->after_state());
9              for (auto [next, prob] : all_possible_state) {
10                 float v = estimate(next);
11                 value += prob * v;
12             }
13             move->set_value(value);
14
15             if (move->value() + move->reward() > best->value() + best->reward())
16                 best = move;
17         } else {
18             move->set_value(-std::numeric_limits<float>::max());
19         }
20         // debug << "test " << *move;
21     }
22     return *best;
23 }

```

update episode

- TD(0) 可以用相鄰兩步更新，從頭到尾跟從尾到頭更新理論上都可行，只要分別計算 $estimate(S_t)$, $estimate(S_{t+1})$ 就能算出 error 並且 update。

- 如果從頭開始更新 $S_t \rightarrow S_{t+1} \rightarrow S_{t+2} \dots$ 。用 S_{t+1} 去更新 S_t ，update weight 後要更新 S_{t+1} 時，需要再重算一次權重更新後的 S_{t+1} ，造成計算的浪費。
- 因為 update 會順便回傳更新後的 $V(S_t)$ ，所以如果從尾巴來，更新 S_t 後可以存起來，順勢繼續往下更新 S_{t-1} ，可以少計算一次。
- 這邊就是把上面的公式實做出來

$$V(S_t) \leftarrow V(S_t) + \alpha[r_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

```
1 void update_episode(std::vector<state>& path, float alpha = 0.1) const {
2     // TODO
3     path.pop_back();
4     float target = 0;
5     for (auto it = path.rbegin(); it != path.rend(); it++) {
6         auto s = it;
7         float error = target + s->reward() - estimate(s->before_state());
8         target = update(s->before_state(), alpha * error);
9     }
10 }
```