

Baby Tetris First Part Report

ERGIN Seçkin Yağmur and GROUSSON Dylan

November 28, 2025

1 Question 1: Discounted MDP For Player

The objective of this part is to choose the optimal policy that maximizes the expected cumulative discounted reward over time for the "player".

States

The state is represented by the current configuration of the grid, and the next piece to be placed.

$S :=$ (current grid, next piece to place)

Actions

The actions correspond to a position and orientation to place the piece on the grid.

$A :=$ (position, orientation)

Transition Function

The transition function is defined as $P(s'|s, a) := 1/2$

There are two possible pieces (I or L) as the next piece to place, and considering our state representation, and our assumption of uniform distribution between the pieces, each transition has an equal probability of 1/2.

Reward Function and Discount Factor

The reward function is computed based on the number of lines cleared after placing the piece, and a discount factor λ is used to weigh future rewards.

$R(s, a, s') :=$ number of lines cleared and $\lambda \in [0, 1]$

1.1 Finding Optimal Policy

We use the value iteration method to find the optimal policy that maximizes the expected rewards, iteratively updating the value function for each state based on the Bellman equation until it converges to the optimal value (margin ϵ).

$$V(s) = \max_{a \in A} \sum_{s' \in S} P(s'|s, a) \cdot (R(s, a, s') + \lambda V(s'))$$

Value iteration is better suited for maximizing the expected discounted reward since it computes the optimal value then extracts the policy and easier to implement

1.2 Implementation

1.2.1 Main Architecture

Based on the our design choices from Section 1, we decided to use an object oriented architecture (C++) for our implementation.

For our structure we have declared 7 classes:

1. *Point*: classical point (x,y) class used for grid coordinates
2. *Field*: the class that contain a boolean grid s.t. $\text{grid}[i][j] = 1$ if the corresponding cell is fill in the field
3. *Tromino*: the class that correspond to the baby Tetris pieces (I Piece or L Piece)
4. *State*: the class that contain both a field grid and an incoming Tromino
5. *Action*: that represent an action choosable on a given state (combination for placement and rotation of an incoming Tromino)
6. *Game*: the class that contain the configuration (completed lines scores), the current state and score of a baby Tetris game (the configuration is loaded from a file : `config.txt`)
7. *MDP*: the class that all the necessary methods in order to compute the value iteration over a game

Now that we overviewed the global implementation architecture, we will describe our implementation choices.

1.2.2 Implementation Choices

- **Piece teleported to their final position:** in order to make the number of intermediate states when choosing an action, we compute on a state all available actions on this state, then according to the policy we make a choice over these actions and we teleport the incoming piece in the current state directly to its final position (with the right rotation)
- **Reward Function:** The reward function as defined in Section 1 take both the current state s and an available action a_s . Here since we can calculate the available actions, we can know in which state s' we're going to arrive and we can compute directly in s' the number of completed lines before remove them right to the game. Then we abstract the reward function by counting the number of complete lines in the resulting state
- **Game Ending:** According to the subject, the game should end whenever a piece exceed the height limit of the grid. Again there, we can compute all available actions on a given state, since we cannot place a Tromino outside the grid, such an action will not appear in the set of available actions. Then the game end when there is no available action for the current game state
- **Max Iteration:** Since the value iteration isn't an exact method, if we fix an ϵ which is too small, we may iterate for a very long moment. For this reason we fix a maximum number of iteration for the expected value improvement, then we always a value iteration which compute in a reasonable time even if the expected value vector is less qualitative than the ϵ margin we fixed

1.2.3 Incoming Improvement

For now we compute naively the value iteration algorithm on all state combinations, which represents

$$= 2^{\text{gridHeight} \times \text{gridWidth}} * \text{nbOfTrominoTypes} = 2^{4 \times 4} * 2 = 2^{17} = 131072$$

states. In the near future we want to considerably reduce this number by computing the value iteration only on states that are reachable from a given starting state s_0 . For this purpose, we will make use of a Breadth First Search algorithm which will explore iteratively all reachable states from the current one.