

Learning how to play 2048 using TD(λ) learning over n -tuple networks

Maadhav Gupta

Ashoka University

March 2025

Actions and States

- Actions = {up, down, left, right}
- States = $\{s_0, s_1, s_2, \dots\}$
#States $\approx 10^{16}$

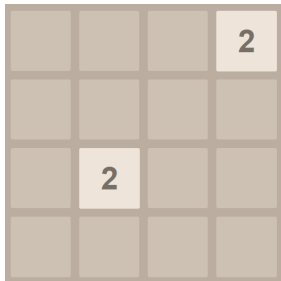


Figure: s_0

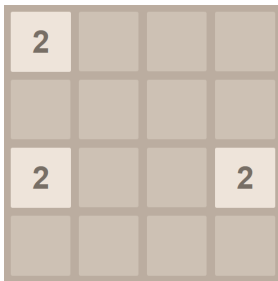


Figure: $s_1 = s_0$ after left

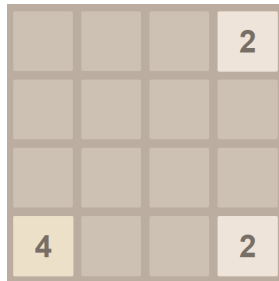


Figure: $s_2 = s_1$ after down

2048 Heuristics

- Monotonicity
- High Corner Value
- High number of blank tiles



A 4x4 grid representing a 2048 game board. The tiles are as follows:

			2
4			
8	4		
16	8	4	

n -Tuple Networks

- Uses fixed groups of tile positions (tuples).
- Each tuple maps observed tile patterns to a value.
- Multiple overlapping tuples are used.
- Captures local board structures and generalizes across similar configurations.
- Fast to evaluate and effective for learning non-linear patterns.

n -Tuple Networks

- Some states

2	4		
16		2	
128	32	4	
8	2	64	4

- Network of tuples

2
16
128
8

4
32
2

2
4
64

4

Value Function Approximation

- The value function $V(s)$ estimates expected future reward from state s .
- Approximated as:

$$V(s) = \sum_{i=1}^n w_i [f_i(s)]$$

- $f_i(s)$: index in the lookup table for tuple i based on observed pattern.
- w_i : weight table storing values for tuple i .
- Enables generalization by sharing weights across states with similar features.

TD(λ) with Eligibility Traces

- Temporal Difference learning with eligibility traces updates value estimates across sequences.
- TD error: $\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$
- Eligibility trace update:
 - $e_t(f) \leftarrow \gamma \lambda e_{t-1}(f) + 1$ (if feature f is active)
- Weight update:
 - $V(f) \leftarrow V(f) + \alpha \cdot \delta_t \cdot e_t(f)$
- Assigns credit to both recent and earlier features, improving long-term learning.
- Suits 2048 as reward is delayed and survival is crucial.

ϵ -Greedy Policy with Decay

- Balances exploration and exploitation during training.
- At each decision point:

$$a = \begin{cases} \text{random action} & \text{with probability } \epsilon \\ \arg \max_a Q(s, a) & \text{with probability } 1 - \epsilon \end{cases}$$

- ϵ decays from 1.0 to 0.01 over episodes.
- Encourages exploration early and stable policies later.

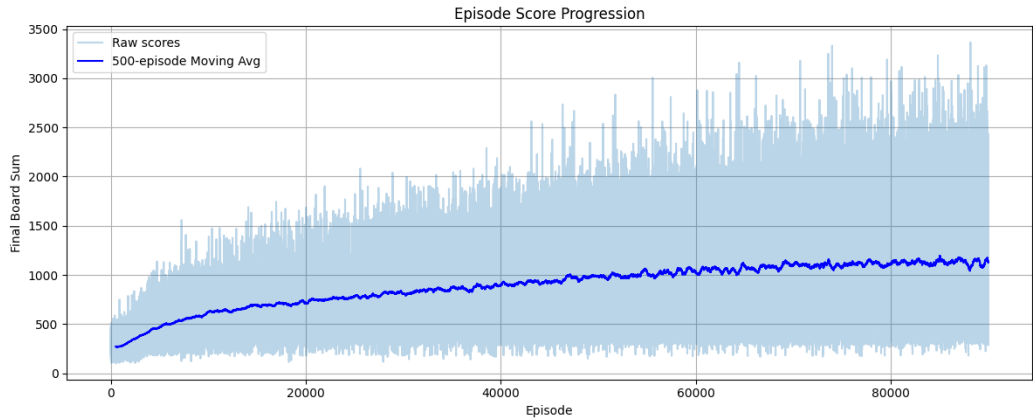
Experimentation

- TD(0) attempt: $\alpha = 0.01$, $\epsilon = 0.01$ fixed, $\gamma = 0.99$, *rows* only in n -tuple network
 - Episodes observed: 50000
 - Average score: below 600
 - Runtime: 6 hours
- TD(λ) attempt 1: $\alpha = 0.1$, $\epsilon = 0.1$ fixed, $\gamma = 1.0$, $\lambda = 0.9$, set of 30 n -tuples considered
 - Episodes observed: 50000
 - Average score: below 400, (as much as the random agent would score (or worse))
 - Runtime: 18 hours

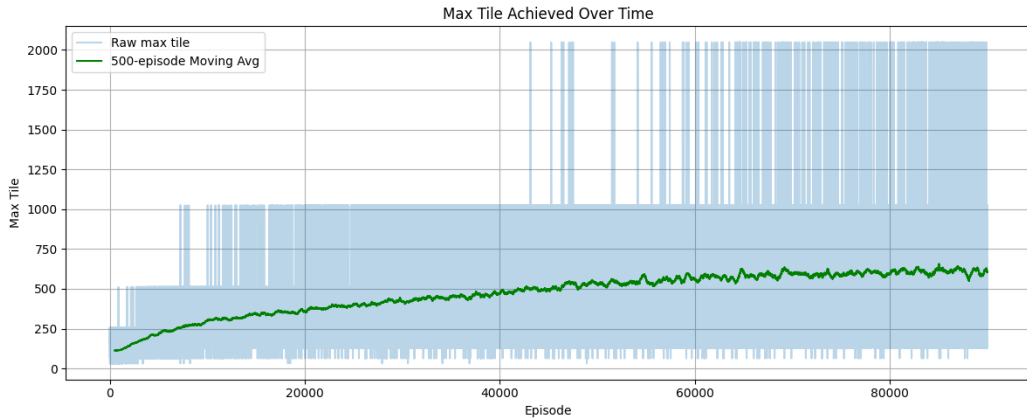
Experimentation

- TD(λ) attempt 2: $\alpha = 0.01$, $\epsilon = 0.01$ decaying to 0.1, $\gamma = 0.99$, $\lambda = 0.5$, n -tuple network = $\{\text{rows}, \text{columns}, \text{diagonals}\}$
 - Episodes observed: 100000
 - Average scores: reached 1200
 - 2048 achieved 5% of the times
 - Runtime 10 hours
 - ϵ decays to 0.01 in 11,500 episodes

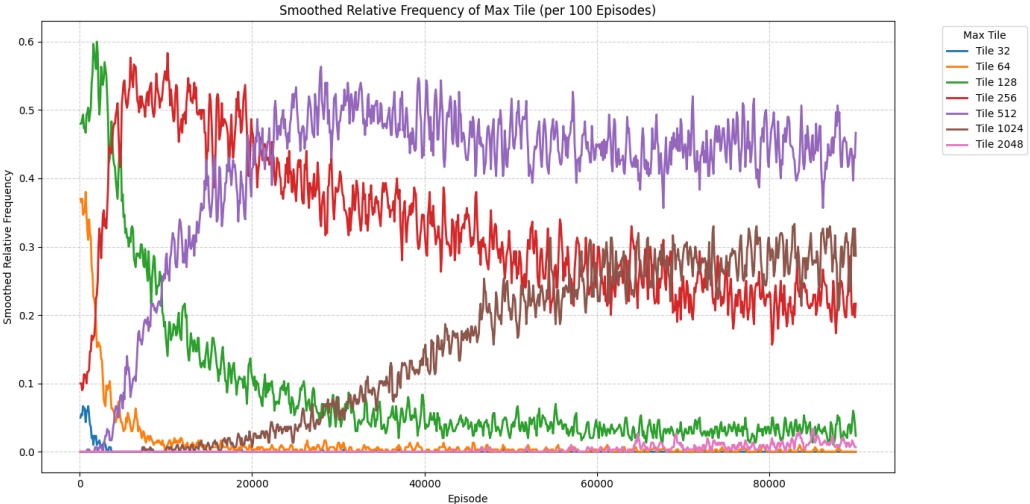
Results



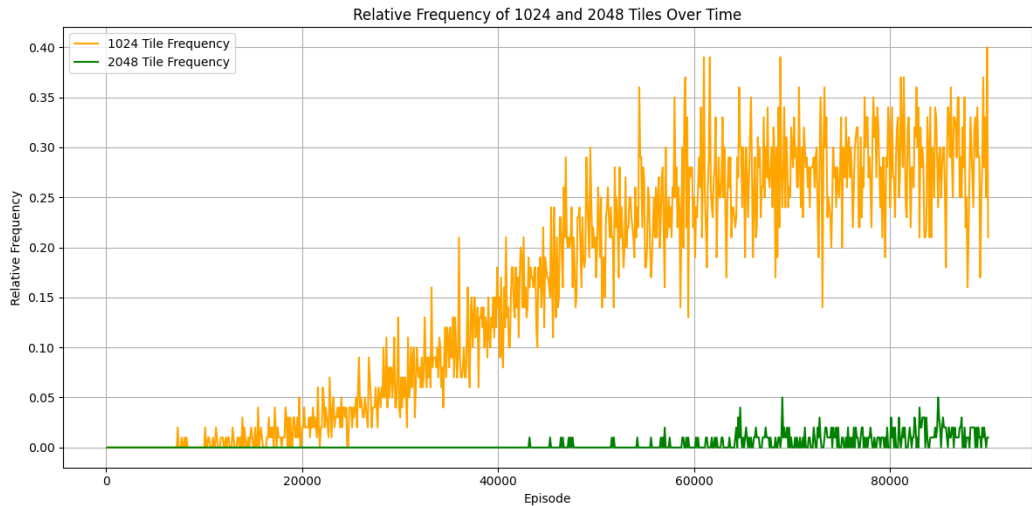
Results



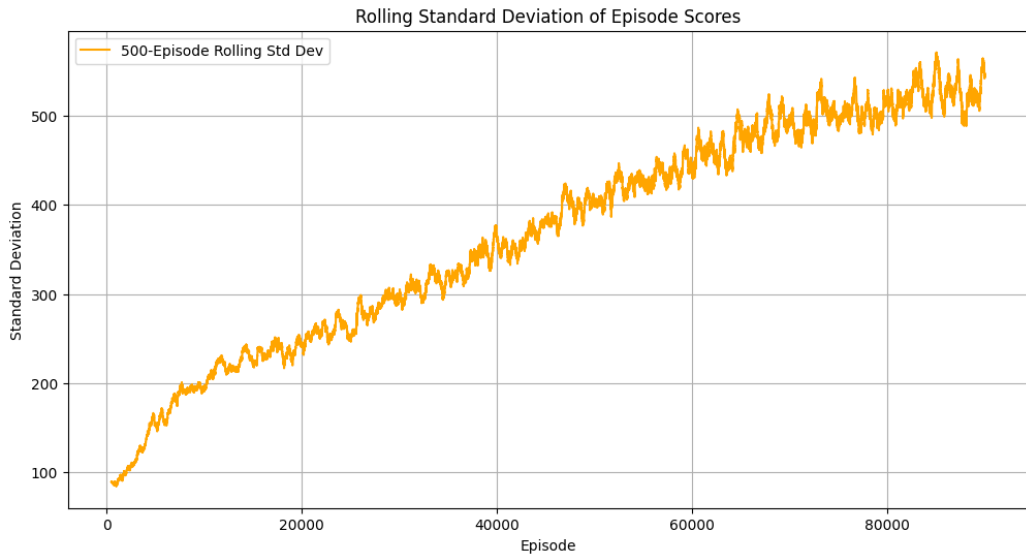
Results



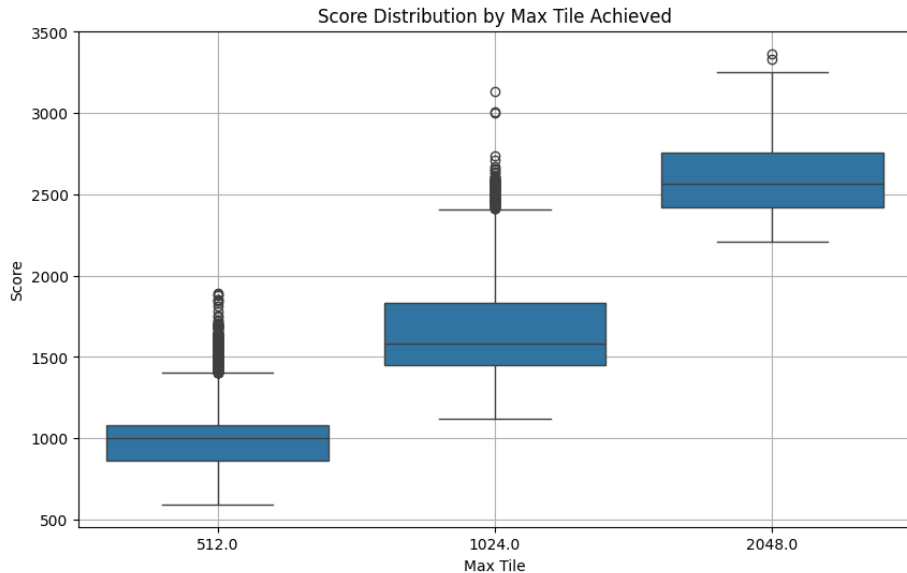
Results



Results



Results



Improvements

- Can use more n -tuples and try different parameters or for more episodes
- Introduce exploration again after a few, say 10000, episodes
- Need to encourage corner points, present value function focusing on sum of grid