



UNIVERSITÀ  
DEGLI STUDI DI TRIESTE



DATA SCIENCE &  
SCIENTIFIC COMPUTING

# Exploring Reinforcement Learning in Pokemon Battles

Reinforcement Learning Project

Matteo Boi

---

A.A. 2022/2023

# Index



- **Introduction:**
  - Pokemon Battles as a RL environment
- **First approach**
  - Q-learning
  - Results
- **Second approach**
  - Multi-Armed bandit
  - Results
- **Improvements & Conclusions**

# Pokemon Battles as RL environment



## Complex problem:

- Current pokemon
- Status conditions
- HP level
- Stat changes
- Opponent's pokemon
- Move informations
- Move usage
- State of the player's team
- Possibility to switch pokemon
- Battlefield condition
- Items



# Pokemon Battles as RL environment



## Simplification:

- Agent = Pokemon.
- Each pokemon has one type, no stats, same HP.
- Only generic moves (fire, water, ...).
- No items.

**Goal:** to make the agent learn type advantages.

- Fire is supereffective against grass (2x damage).
- Fire is not very effective against water (0.5 damage).

# Environment V0



- 3 possible types (fire, water, grass)
- **Action space:** fire move, water move, grass move.
- **Observation:** type of the opponent pokemon.
- **Episode:** entire battle (until one of the two pokemons faints)
- **Two reward functions:**
  - Full reward ("guided"): +1 if the move used is supereffective  
0 if the move is standard  
-1 if the move is not very effective  
+ (remaining HP) of the agent pokemon at the end of the episode
  - HP reward: = (remaining HP) no guidance during the episode

# Environment V1



- 4 possible types (normal, fire, water, grass)
- **Action space:** normal move, fire move, water move, grass move.
- **Observation:** type of the opponent pokemon.
- **Episode:** entire battle (until one of the two pokemons faints)
- **Two reward functions:**
  - Full reward ("guided"): +1 if the move used is superffective  
0 if the move is standard  
-1 if the move is not very effective  
+ (remaining HP) of the agent pokemon at the end of the episode
  - HP reward: = (remaining HP) no guidance during the episode

# Environment V1



- 4 possible types (normal, fire, water, grass)
- **Action space:** normal move, fire move, water move, grass move.
- **Observation:** (type of the agent pokemon, type of the opponent pokemon).
- **Episode:** entire battle (until one of the two pokemons faints)
- **Additional damage:** if the move used is of the same type of the pokemon
- **Two reward functions:**
  - Full reward ("guided"):
    - As before.
    - **+0.5** if the move used is of the same type of the pokemon.
  - HP reward:
    - As before.

# Environment V2



- **6 possible types** (normal, fire, water, grass, fighting, flying)
- **Random moveset:** (move of the same type of the pokemon,  
3 random moves between the remaining types).
- **Action space:** 4 moves in the moveset
- **Observation:** (type of the agent pokemon,  
type of the opponent pokemon,  
types of the 3 moves in the moveset)

6 possible values
6 possible values
6 possible values for each move



# Q-learning

Off-policy temporal difference control with  $\epsilon$ -greedy action selection.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

With  $Q$  that directly approximates  $q_*$  (optimal action-value function).

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

        Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

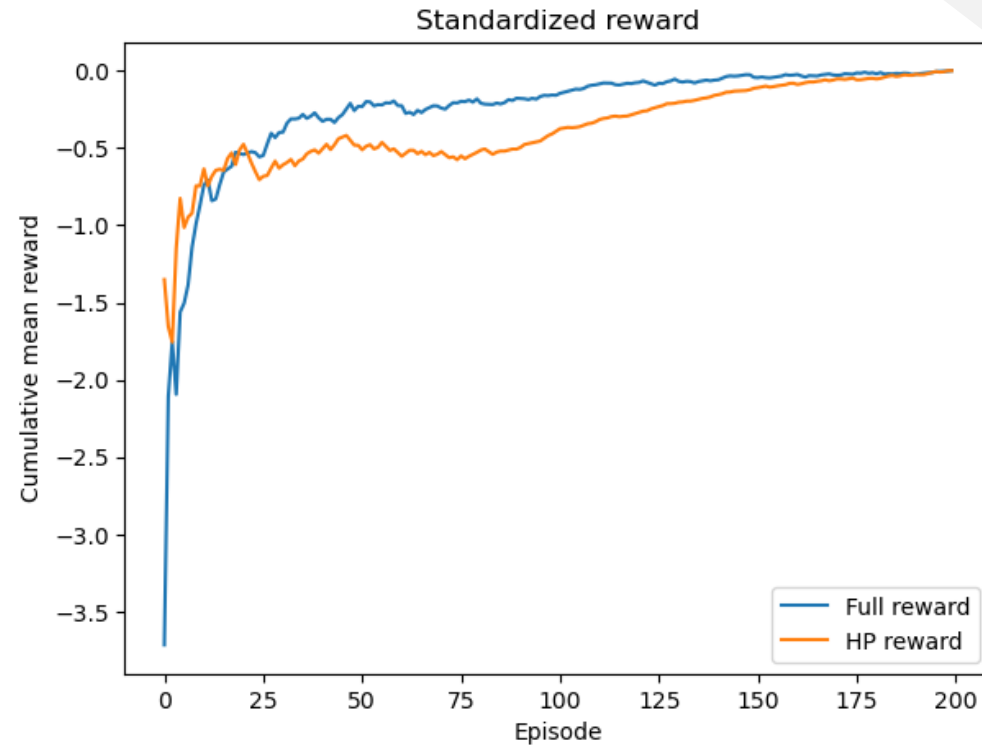
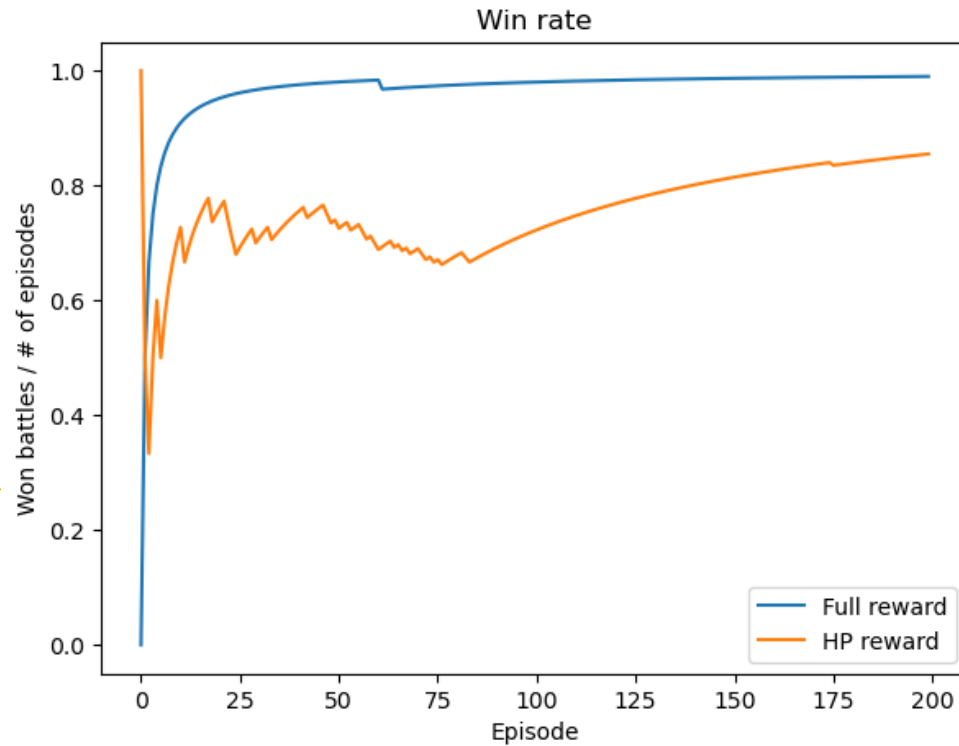
    until  $S$  is terminal



# Results: Environment V0



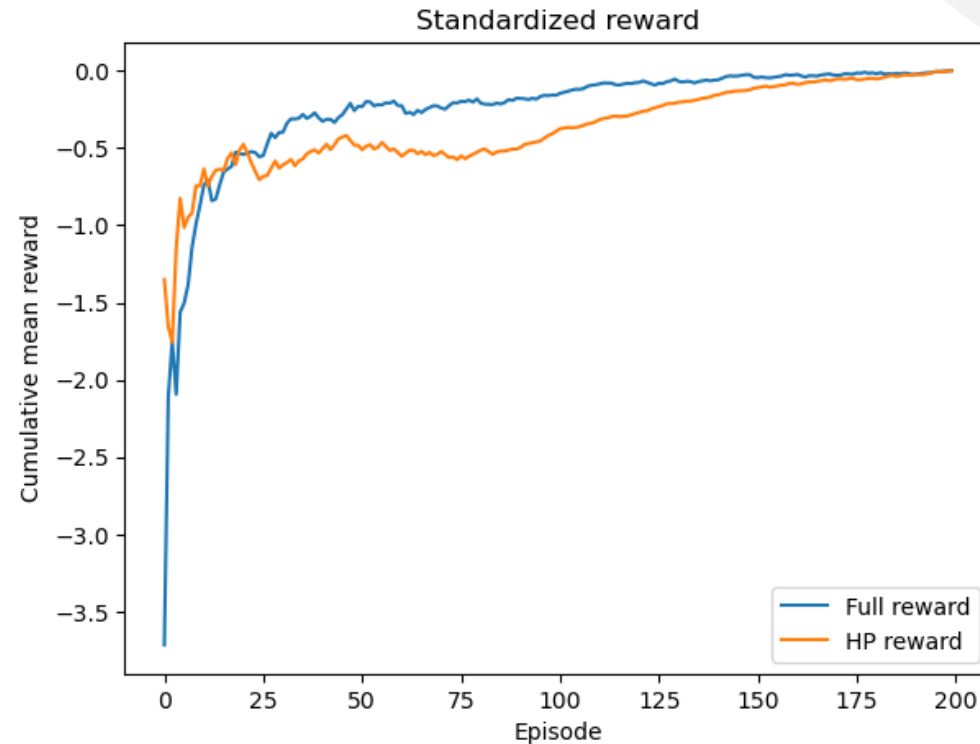
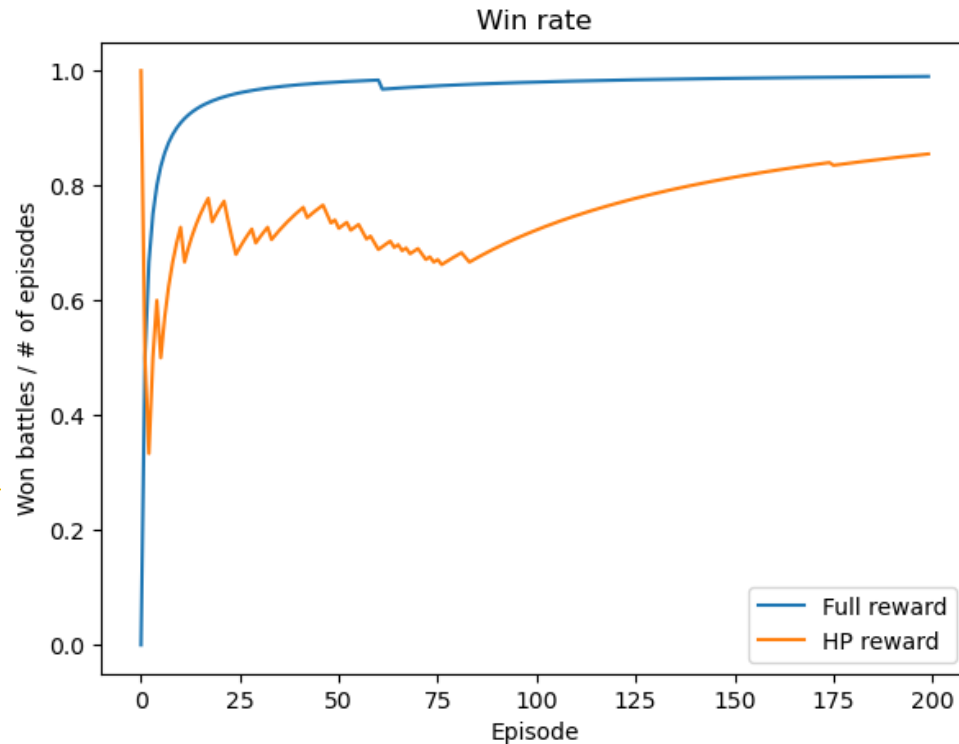
Comparison: Full reward vs. HP reward



# Results: Environment V0



Comparison: Full reward vs. HP reward



**Convergence:**     ~5 episodes for Full reward  
                             ~30 episodes for HP reward

# Results: Environment V0



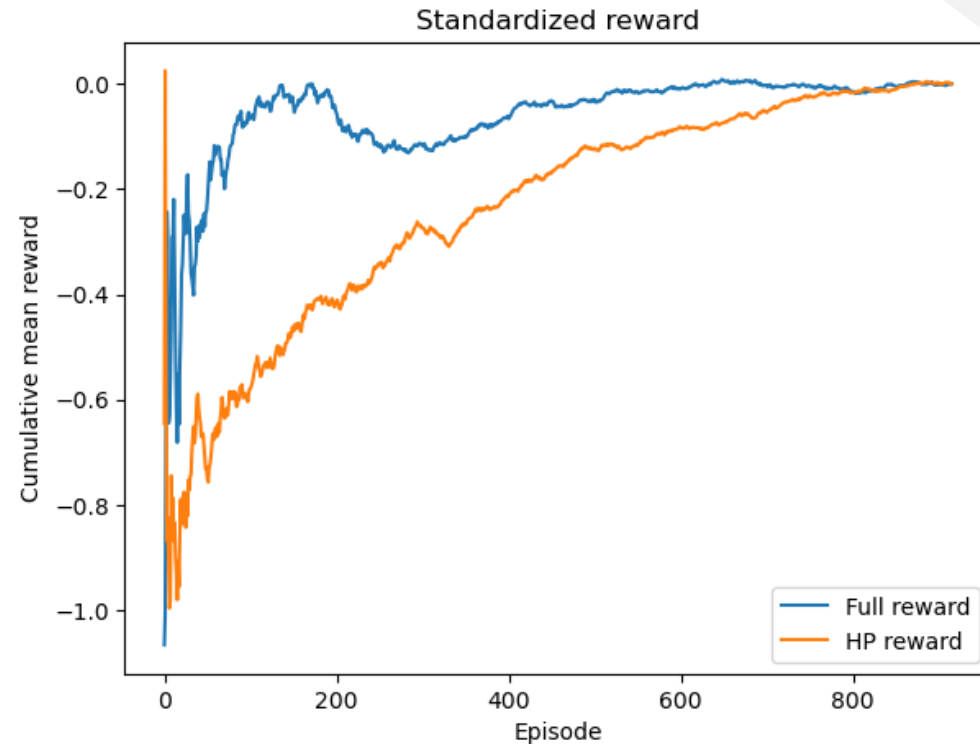
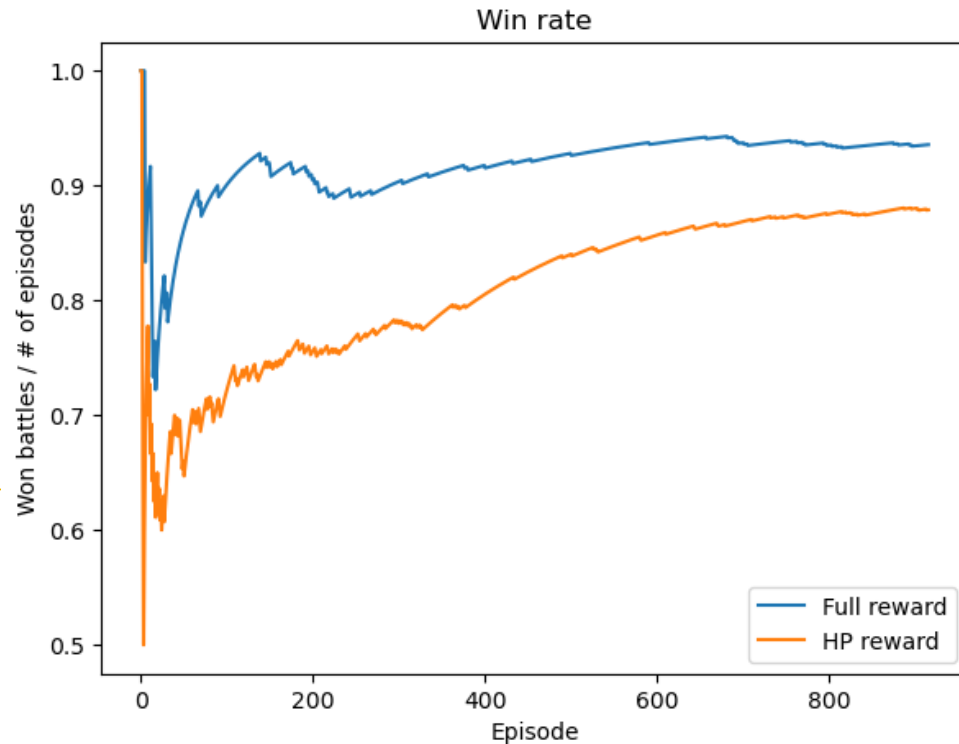
- Optimal Policy:

Agent's move type	Opponent type		
	Fire type	Water type	Grass type
	Water move	Grass move	Fire move

# Results: Environment V1



Comparison: Full reward vs. HP reward



**Convergence:** ~1100 episodes for Full reward  
~1300 episodes for HP reward



# Results: Environment V1

- Optimal Policy:

		Opponent type			
Agent type		Normal type	Fire type	Water type	Grass type
	Normal type	Normal move	Water move	Grass move	Fire move
	Fire type	Fire move	Water move	Grass move	Fire move
	Water type	Water move	Water move	Grass move	Fire move
	Grass type	Grass move	Water move	Grass move	Fire move



# Multi-Armed Bandit

**Observation:** the state never changes during the episode.

**Multi-Armed bandit algorithm:**  $\epsilon$ -greedy action selection.

- **Arms:** actions available to the pokemon

## A simple bandit algorithm

Initialize, for  $a = 1$  to  $k$ :

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Loop forever:

$$A \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with probability } 1 - \epsilon \\ \text{a random action} & \text{with probability } \epsilon \end{cases} \quad (\text{breaking ties randomly})$$

$$R \leftarrow \text{bandit}(A)$$

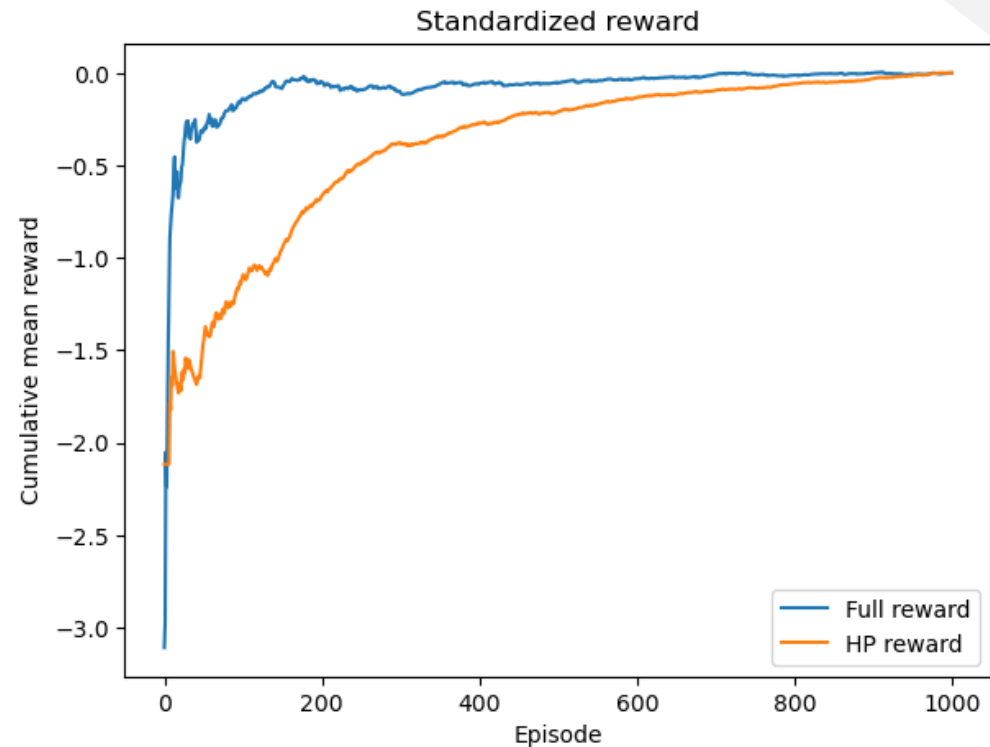
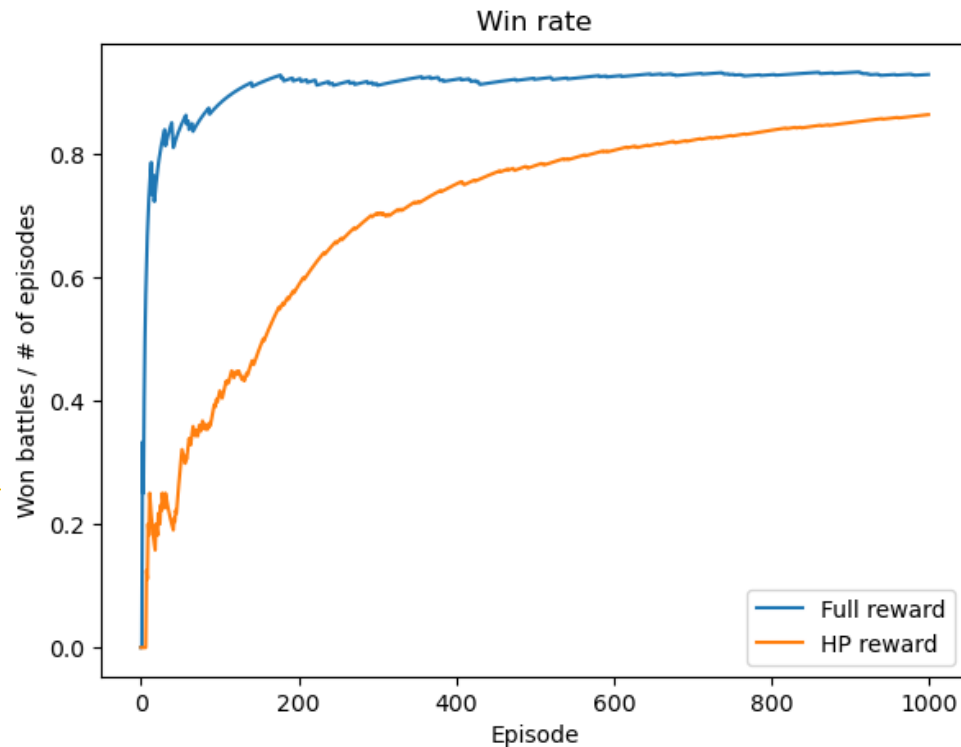
$$N(A) \leftarrow N(A) + 1$$

$$Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)]$$

# MAB Results: Environment V0



Comparison: Full reward vs. HP reward

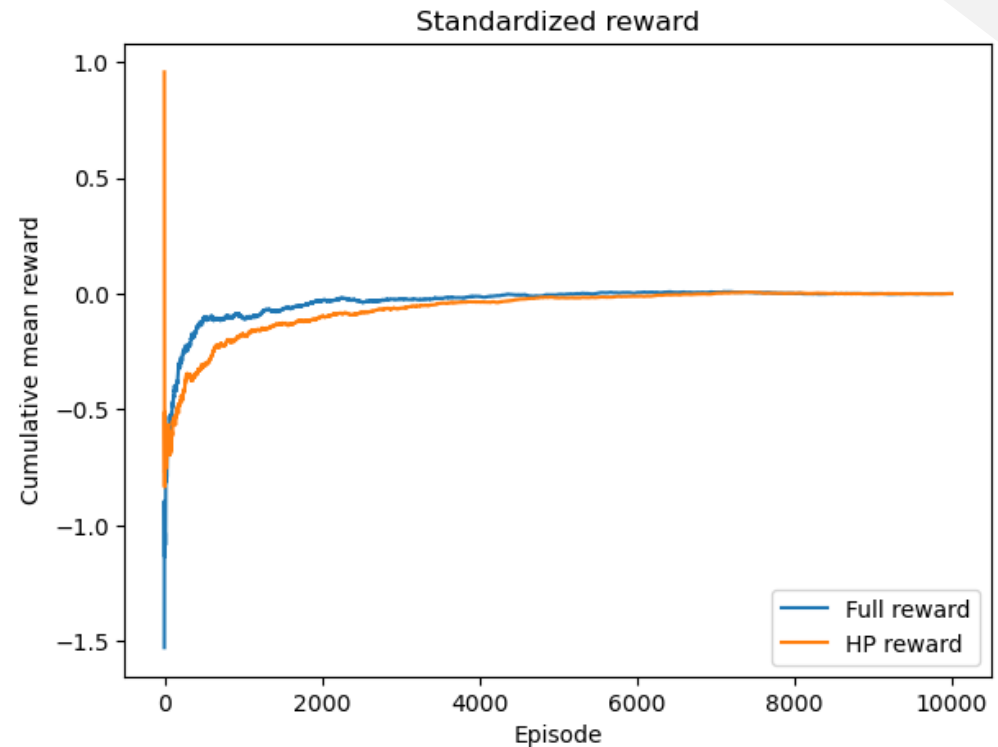
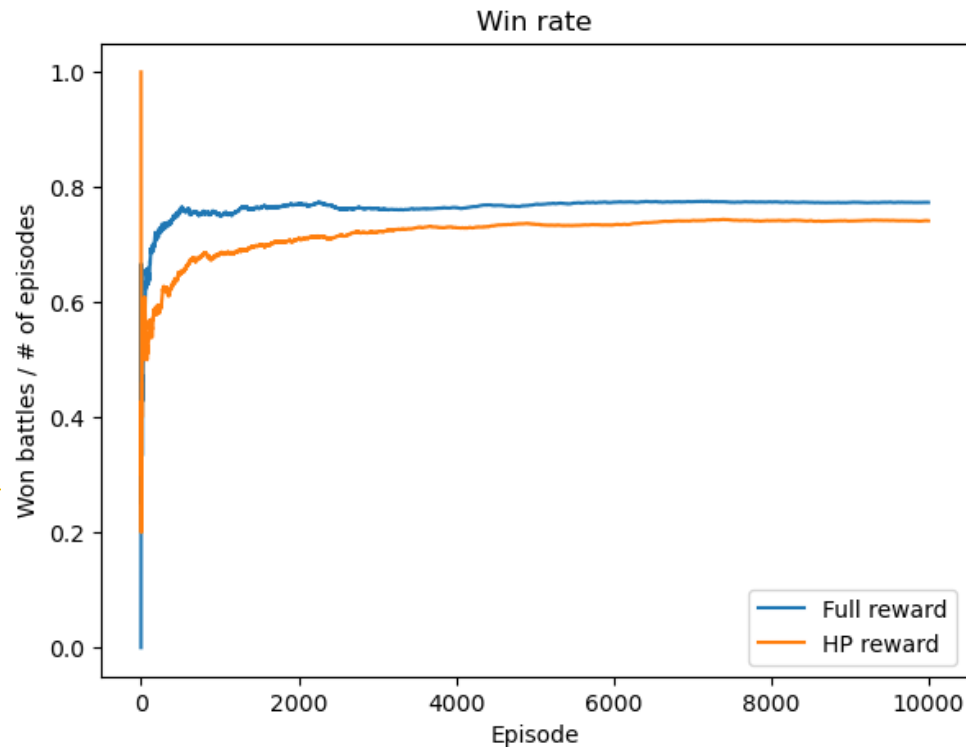




# MAB Results: Environment V1



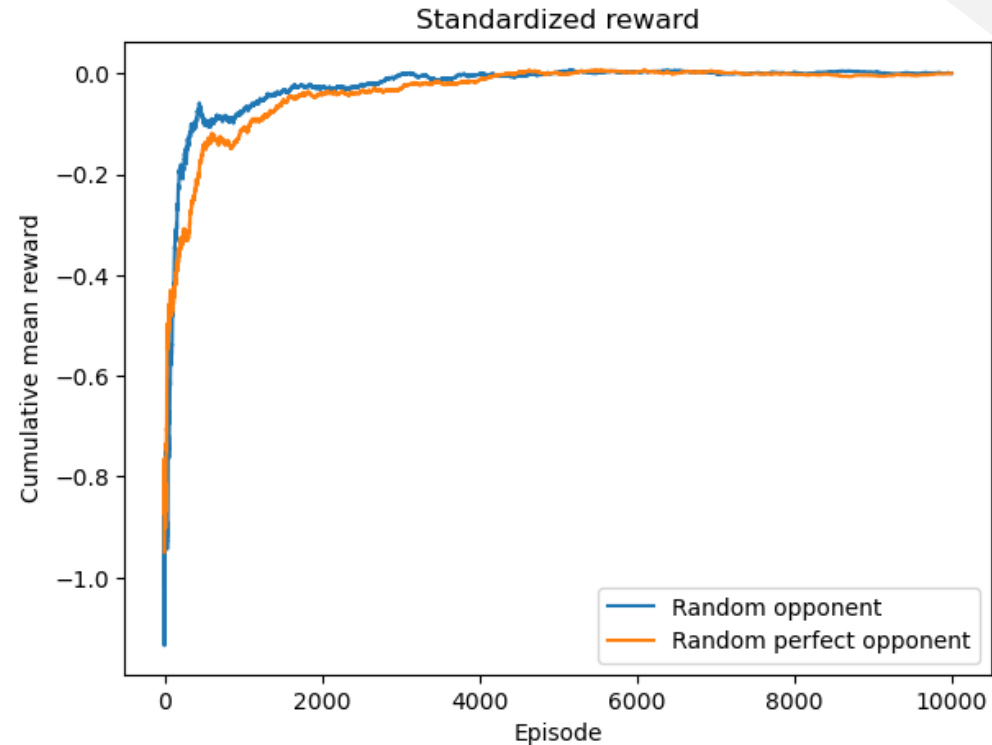
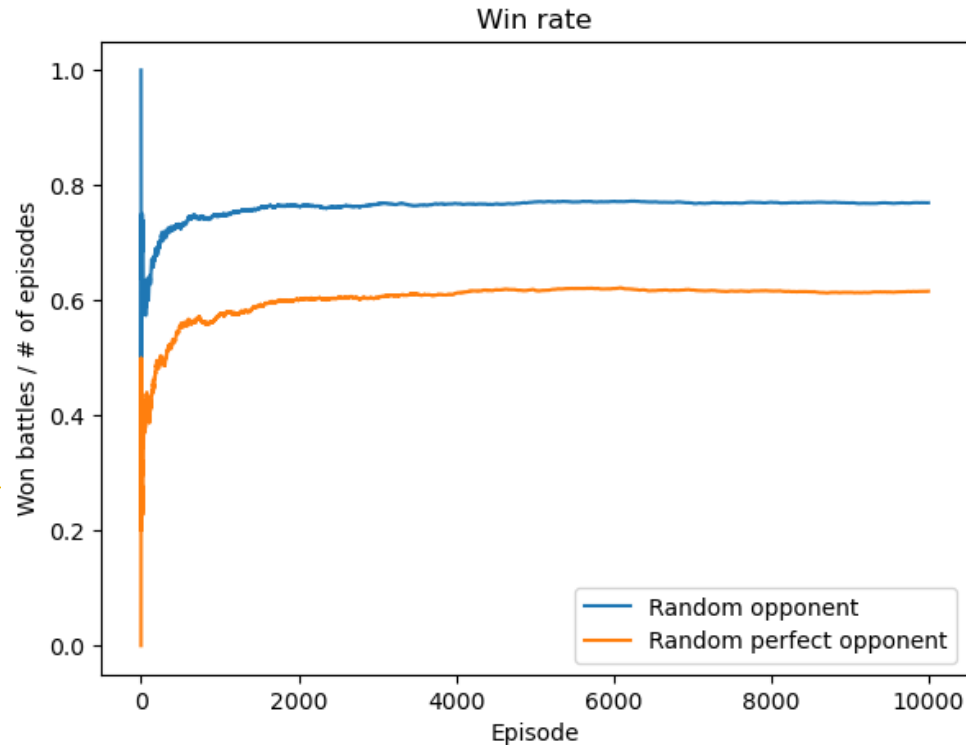
Comparison: Full reward vs. HP reward





# MAB Results: Environment V1

- **Random-perfect opponent:** chooses a random move with  $p = 0.5$   
chooses the supereffective move otherwise



**Observation:** a good portion of the lost battles is due to bad matchups

# Improvements



- **Learning opponent** instead of a dummy one.
- **More complex scenarios:**
  - HP as observation.
  - Integration of more pokemon types.
  - Strategic moves for more complex interactions.
- **Dynamic turn order system:** based on the pokemon speed.
- **Partial observability of the opponent:** hidden type and moves.



**Thank You.**