

# Convergence of a L2 regularized Policy Gradient Algorithm for the Multi Armed Bandit

**Ștefana-Lucia Anița, Gabriel Turinici**



CEREMADE Université Paris Dauphine - PSL, Paris, France



27<sup>th</sup> International Conference on  
Pattern Recognition  
December 01-05, 2024, Kolkata, India



27<sup>th</sup> ICPR conference, Kolkata, India, Dec. 1-5, 2024

# Introduction to Multi-Armed Bandits

**Definition:** The Multi-Armed Bandit (MAB) problem is a classic framework in **reinforcement learning**, where:

- An agent repeatedly chooses among multiple actions (e.g., pulling different slot machine arms).
- The goal is to maximize the total reward over time.
- A balance must be struck between **exploration** (trying unknown actions) and **exploitation** (choosing the best-known action).

## Multiple Applications:

- Optimizing online ad placements or A/B testing.
- Choosing the **best (happy hour) drink** in a bar without spending too much money.
- Allocating research funds to promising areas for maximum impact.

# Softmax Policy Gradient MAB with $L_2$ Regularization

**Notations:** The MAB has  $k$  arms (choices), each choice  $a \leq k$  outcome is stochastic with a reward distribution of mean  $q_*(a)$ . At each time  $t$ , the agent selects an arm  $A_t$  and observes a reward  $R_t$  samples from the distribution of the arm  $A_t$ . The objective is to maximize the cumulative reward.

**Softmax Policy:** The agent maintains a preference vector  $H \in \mathbb{R}^k$ , defining the probability of selecting arm  $A$  through:  $\Pi_H(A) = \frac{e^{H(A)}}{\sum_{a=1}^k e^{H(a)}}$ .

**Our version : regularized loss:** The goal is to maximize:  $\mathcal{L}_\gamma(H) = \mathbb{E}_{A \sim \Pi_H} [R(A) - \frac{\gamma}{2} \|H\|^2]$ , where  $\gamma > 0$  is the  $L_2$  regularization coefficient.

**Optimization algorithm : Policy gradient, a specific variant of Stochastic Gradient Ascent:**

$$H_{t+1}(a) = H_t(a) + \rho_t [(R_t - \bar{R}_t)(\mathbb{1}_{a=A_t} - \Pi_{H_t}(a)) - \gamma H_t(a)],$$

with  $\rho_t =$  learning rate.

# Theoretical results : fixed or variable $\rho_t$ , large $\gamma$

## Proposition (Convergence conditions, fixed or variable $\rho_t$ , large $\gamma$ )

Assume  $\mu := \gamma - (\max_a q_*(a) - \min_a q_*(a)) > 0$ . Under appropriate hypotheses on distributions of arms  $A$ , there exists a unique optimum  $H_*$ . Moreover if

$$\rho_t \rightarrow 0 \text{ and } \sum_{t \geq 1} \rho_t = \infty. \quad (1)$$

then

$$\lim_{t \rightarrow \infty} H_t \stackrel{L^2}{=} H_*. \quad (2)$$

# Theoretical results : “linear decay schedule” and small $\gamma$

## Proposition (Convergence rate for “linear decay schedule”)

Let  $\beta_1, \beta_2 > 0$  two positive constants and take  $\rho_t = \frac{\beta_1}{1+\beta_2 t}$ . Under the same hypotheses as before for  $\gamma$  large enough : there is a unique solution  $H_*$  and the L2 regularized policy gradient MAB algorithm converges with the rate :

$$\mathbb{E}[\|H_t - H_*\|^2] = O\left(\frac{1}{t}\right) \text{ as } t \rightarrow \infty. \quad (3)$$

**Question** : optimum  $H^*$  will depend on  $\gamma$ , may pose problems, what about small  $\gamma$  ?

## Lemma

Let  $V(\gamma) := \max_{H \in \mathbb{R}^k} \mathcal{L}_\gamma(H)$ . Then  $\lim_{\gamma \rightarrow 0} V(\gamma) = V(0)$ .

# Numerical Simulations: Setup

## Experiment Design:

- $M = 1000$  tests of 2000 steps each.
- $k = 10$  arms with rewards  $R(A)$  normally distributed:  $R(A) \sim \mathcal{N}(q_*(A), 1)$ , where  $q_*(A)$  has mean 4 and unit variance.
- Initial distributions tested:
  - **Uniform:**  $\Pi_{H_0}$  with  $H_0 = (0, \dots, 0)$ .
  - **Biased:**  $\Pi_{H_0}$  with  $H_0 = (5, \dots, 0)$ .

## Plot Reward Normalization:

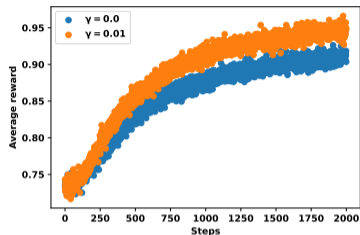
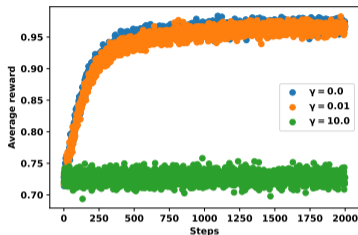
- Reward is scaled relative to the maximum possible reward:

$$R_{\text{scaled}} = \frac{R}{\max_a q_*(a)} \quad (\text{best reward} = 1).$$

## Parameter Testing:

- **Regularization coefficient:**  $\gamma \in \{0, 0.01, 10\}$ .
- **Learning rate:**  $\rho_t = 0.05$  or  $\rho_t = \frac{1}{1+0.05t}$ .

# Results with constant learning step



**Figure:** The average reward for  $\rho_t = 0.05$  (constant),  $\gamma$  is 0, 0.01 or 10 (see the legend). **Left :** start from a uniform distribution  $\Pi_{H_0}$  with  $H_0 = (0, \dots, 0)$ . **Right :** start from a biased distribution  $\Pi_{H_0}$  with  $H_0 = (5, \dots, 0)$ .

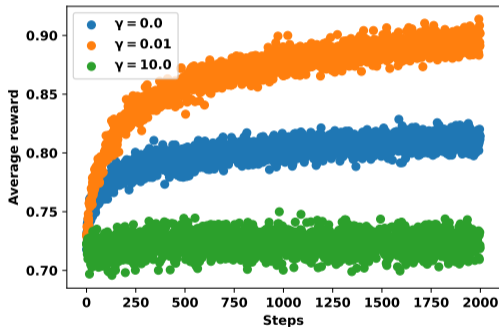
## Uniform Start ( $H_0 = (0, \dots, 0)$ ):

- $\gamma = 0.01$  achieves comparable performance to  $\gamma = 0$ .
- $\gamma = 10$  biases the solution and reduces performance.

## Biased Start ( $H_0 = (5, \dots, 0)$ ):

- $\gamma = 0.01$  improves performance over  $\gamma = 0$ .
- $\gamma = 10$  remains suboptimal.

# Results with linear decaying $\rho_t$



- **Biased Start** ( $H_0 = (5, \dots, 0)$ ):
- Learning rate:  $\rho_t = \frac{1}{1+0.05t}$  (linear decay schedule).
- Comparison of  $\gamma = 0$  or 0.01 or 10.

## Observations:

- $\gamma = 0.01$  achieves good performance, improving initial convergence.
- $\gamma = 10$  is too large, leading to suboptimal results.
- Decay of  $\rho_t$  helps transition from initial exploration to final convergence.



# Summary of numerical results, further tests

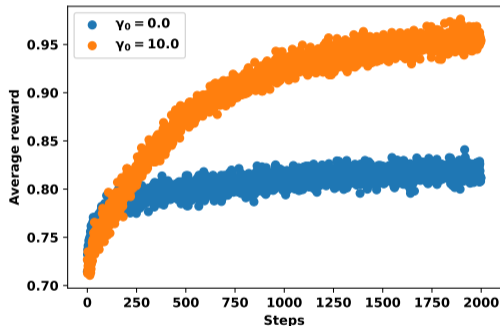
## Empirical results so far:

- Convergence occurs as predicted by the theory.
- For non-uniform initial guesses  $H_0$ : regularization i.e.,  $\gamma > 0$  improves convergence significantly.

- large  $\gamma$  may lead to suboptimal limit points  $H^*$ .

**Further idea:** best of the two worlds :  
variable  $\rho_t = \frac{1}{1+0.05 \cdot t}$  and  $\gamma_t = \frac{\gamma_0}{1+0.2 \cdot t}$ :

- Combines exploration and final convergence.
- $\gamma_0 > 0$  (orange in figure) outperforms non-regularized case  $\gamma_0 = 0$  (blue).



# Summary of theoretical and numerical results


**Objective:** we investigate a  $L_2$ -regularized policy gradient algorithm for Multi-Armed Bandit (MAB).

## Theoretical Results [1]:

- **Proposition 1:** Convergence established for both constant and variable step sizes ( $\rho_t$ ).
- **Proposition 2:** Convergence rate proven to be  $\mathcal{O}(1/t)$  for linear decay of  $\rho_t$ .
- **Lemma 1:** Regularization may shift the optimum but optimality is restored as  $\gamma \rightarrow 0$ .

**Key Takeaway (theoretical and numerical):** regularization helps improve convergence especially when starting far from optimality; it can be adjusted dynamically if needed.

**Future Work:** The optimal decay schedule for  $\gamma_t$  remains an open question: type, theoretical convergence.

-  Ștefana-Lucia Anița and Gabriel Turinici.  
*Convergence of a L2 Regularized Policy Gradient Algorithm for the Multi Armed Bandit*,  
page 407–422.  
Springer Nature Switzerland, December 2024.  
arXiv:2402.06388.