# Comparative analysis of Sliding Window UCB and Discount Factor UCB in non-stationary environments: A Multi-Armed Bandit approach

**Haocheng Liu**

Boston University, Art & Sciences, Boston, 02215, United State

easonlhc@bu.edu

**Abstract.** The Multi-Armed Bandit (MAB) problem is a well-studied topic within stationary environments, where the reward distributions remain consistent over time. Nevertheless, many real-world applications often fall within non-stationary contexts, where the rewards from each arm can evolve. In light of this, our research focuses on examining and contrasting the effectiveness of two leading algorithms tailored for these shifting environments: the Sliding Window Upper Confidence Bound (SW-UCB) and the Discount Factor UCB (DF-UCB). By harnessing both simulated and real-world datasets, our evaluation encompasses adaptability, computational efficiency, and the potential for regret minimization. Our findings reveal that the SW-UCB is adept at swiftly adjusting to sudden shifts, whereas the DF-UCB emerges as the more resource-efficient option amidst gradual transitions. Notably, when pitted against conventional UCB algorithms within non-stationary contexts, both contenders exhibit substantial advancements. Such insights bear significant relevance to fields like online advertising, healthcare, and finance, where the capacity to nimbly adapt to dynamic environments is paramount.

**Keywords:** Multi-Armed Bandit Algorithms, Non-Stationary Environments, UCB.

## 1. Introduction

Multi-Armed Bandit problems are foundational frameworks for decision-making across diverse sectors such as healthcare, finance, and online advertising. These problems focus on optimizing choices among various options, termed "arms," to achieve the maximum expected reward.

Originally rooted in probability theory, the MAB problem has evolved into a fundamental concept in disciplines like machine learning and computational economics. Classical algorithms like Upper Confidence Bound and Thompson Sampling have been predominantly utilized to address MAB problems, especially in stationary environments where each arm's reward distribution remains unchanged. However, the central concern of this research centers on the inadequacies of traditional MAB algorithms in non-stationary environments. In these dynamic contexts, changing reward distributions over time can compromise the effectiveness of standard algorithms. Key challenges emerge in efficiently adapting to these variable reward distributions and in defining performance metrics, such as regret, within these mutable scenarios.

This paper delves into an analysis and comparison of two UCB algorithm variants: Sliding Window UCB and Discount Factor UCB, specifically in non-stationary settings. It is posited that these methods might outperform traditional UCB algorithms in adaptability and regret minimization. Optimizing decision-making in fluid environments has implications that span from automated trading systems to personalized healthcare. Consequently, this research could serve as a linchpin in tailoring MAB algorithms for real-world, non-stationary applications. After this introduction, the paper reviews pertinent literature, elaborates on the theoretical framework, outlines the methodology, presents findings, and concludes with an in-depth discussion. With an objective to bridge the existing gap between standard MAB algorithms and their effectiveness in dynamic environments, this research endeavors to furnish both theoretical knowledge and pragmatic directives for impending applications.

## 2. Literature Review

Evolution of MAB Algorithms: The Multi-Armed Bandit (MAB) paradigm holds its ground as an iconic model for decision-making amid uncertainty [1]. Historically, stalwarts like Upper Confidence Bound algorithms and Thompson Sampling catered predominantly to environments with static reward distributions. Yet, the capricious nature of real-world scenarios necessitated innovation within the MAB framework.

Adaptation Techniques for Non-Stationary Environments: It's become evident that conventional MAB strategies falter in addressing the inherent non-stationarity of many real-world applications [2]. In scenarios where reward dynamics are fluid, the prowess of an algorithm to identify and adjust to shifts takes center stage. Here, piecewise-stationary contexts arise, marking environments wherein reward distributions undergo sudden changes at specific intervals. The latest MAB research focuses intently on mechanisms that recognize and respond to these transitions.

Sliding Window UCB: A refined answer to the challenges of non-stationarity is the Sliding Window UCB. By emphasizing recent data within a predetermined 'window', this algorithm seeks to mirror current reward tendencies, making it agile in fluctuating landscapes [3]. Such an approach proves invaluable in contexts where historical insights might not reliably foretell upcoming trends.

Discount Factor UCB: Conversely, the Discount Factor UCB embraces a more measured tactic. Instead of completely overshadowing older data points as its counterpart does, it moderates their influence via a discounting technique [4]. This balanced consideration of past data ensures that the algorithm stays attuned to gradual shifts, making it proficient in situations marked by slowly adapting reward distributions [5].

The Concept of Regret in MAB: Regardless of the algorithm or its adaptive capabilities, the gold standard for success often hinges on 'regret'— the potential loss incurred from not consistently opting for the best arm throughout the decision-making timeline [6]. In the labyrinth of non-stationary environments, grasping the regret intricacies associated with diverse algorithms is paramount. Within this sphere, 'minimax regret' stands out, epitomizing the harshest relative performance of an algorithm when juxtaposed with the ideal strategy.

Gaps and Limitations in Existing Research: Notwithstanding the progress in adaptive MAB solutions, a comprehensive comparative study gauging the efficacy of tools like Sliding Window UCB and Discount Factor UCB across varied non-stationary backdrops is notably limited. Moreover, when faced with intricate challenges such as multi-objective decisions in unpredictable settings, the current suite of algorithms and assessment benchmarks requires further refinement [7].

## 3. Theoretical Framework

### 3.1. Definitions

Before diving into the complex algorithms and methodologies, it is essential to clarify some foundational terminology:

Arms: In the Multi-Armed Bandit context, arms represent the different options or choices available for selection during each round of decision-making.

Rewards: The numerical benefit obtained after pulling an arm, often modeled as random variables with underlying distributions.

Bandits: A colloquial term for the MAB problem setting, representing a metaphorical 'one-armed bandit' or slot machine with multiple arms to pull [8].

Regret: A performance metric that quantifies the cumulative difference between the rewards accrued by the algorithm and the rewards that could have been accrued by always choosing the optimal arm.

### 3.2. Mathematical Models

Sliding Window UCB: The core idea here is to consider only the most recent N rounds, where N is the size of the 'window'. The algorithm dynamically updates the upper confidence bounds using the following equation:

$$USB_{SW}(a) = \overline{x}a + \sqrt{\frac{2\log t}{n_a}} \qquad (1)$$

Here, $\overline{x}a$ is the average reward of arm a over the last N rounds, $n_a$ is the number of times arm a has been pulled in those N rounds, and t is the total number of rounds played [9]. Discount Factor UCB: In this model, past rewards are discounted by a factor $\gamma$ where $0 < \gamma \leq 1$. The discounted average reward for arm a is then:

$$\overline{x}DF, a = \sum_{t=1}^{T} \Upsilon(T-t) \, xa, t \qquad (2)$$

The UCB equation becomes:

$$UCBDF(a) = \overline{x}DF, a + \sqrt{\frac{2\log t}{n_a}} \qquad (3)$$

Where $n_a$ is the number of times arm a has been pulled so far.

### 3.3. Assumptions and Limitations

Assumptions: Recency of Data: The most salient assumption of Sliding Window UCB is that recent observations are the most informative about the current state of the system. Abrupt Changes: The algorithm inherently assumes that the environment can change suddenly and that such shifts are best captured by focusing on a small, recent set of observations [10]. Window Size: Assumes that an appropriate window size N is known or can be empirically determined, which represents the most relevant time frame for decision-making.

Limitations: Tuning Sensitivity: The performance is highly sensitive to the choice of the window size N. Too small a window might capture noise as signal, while too large a window may miss out on actual changes. Computational Burden: Maintaining a sliding window of observations for each arm can become computationally expensive, especially as the number of arms and the window size increase. Long-term Trends: The algorithm may struggle to capture long-term trends due to its focus on recent observations, potentially increasing regret in certain scenarios.

### 3.4. Discount Factor UCB

Assumptions: Decaying Relevance: Assumes that older observations still contain some level of relevant information, but their importance decays over time, captured by a discount factor $\gamma$. Gradual Changes: Best suited for environments where changes occur gradually rather than abruptly. Parameter Knowledge: Assumes that an optimal or near-optimal discount factor $\gamma$ can be determined, either theoretically or empirically.

Limitations: Hyperparameter Tuning: Like the Sliding Window UCB, the performance of Discount Factor UCB is highly sensitive to the choice of $\gamma$, and selecting an inappropriate value can severely impair its effectiveness. Delay in Adaptation: The method may lag in adapting to sudden, drastic changes in the environment due to its weighted inclusion of past data. Complexity in Nonlinear Changes: In

environments with nonlinear or cyclical trends, the constant discounting mechanism may not capture the complexities effectively, leading to suboptimal performance.

## 4. Methodology

### 4.1. Research Design
The core objective of this study is to critically assess the effectiveness of Sliding Window UCB and Discount Factor UCB algorithms in dealing with non-stationary environments. To achieve this, a multi-pronged approach will be employed: Comparative Analysis: Side-by-side evaluation of the two algorithms using both synthetic and real-world datasets to analyze various aspects such as regret, computational efficiency, and adaptability. Case Studies: Application of the algorithms to specific, well-documented non-stationary scenarios to observe their real-world efficacy and limitations. Simulations: Monte Carlo simulations will be conducted to understand how both algorithms perform across a range of different non-stationary environments, with varying levels of volatility and complexity.

### 4.2. Experiment Setup
Parameters and Settings Sliding Window UCB: Different window sizes (N) will be tested, ranging from small to large, to analyze sensitivity and adaptability. Discount Factor UCB: Various discount factors ($\gamma$) will be explored, from values close to zero to values approaching one. Environmental Complexity: The study will incorporate both smooth and abrupt changes in the reward distribution to test the algorithms' robustness. Number of Arms: Experiments will be conducted with varying numbers of arms to assess scalability and computational efficiency.

Data Collection and Sources: Simulated Data: Simulated datasets will be generated using statistical models to replicate different types of non-stationary environments, from gradual to abrupt changes in reward distributions. Real-world Data: Available time-series datasets from domains like finance, healthcare, and energy will be used to validate the algorithms' applicability in real-world scenarios.

### 4.3. Metrics and Evaluation Criteria
Regret: Cumulative regret will be measured to assess the opportunity cost of not selecting the optimal arm over time.

Computational Efficiency: CPU time and memory usage will be recorded to evaluate the algorithms' computational costs.

Adaptability: Metrics like 'time-to-adapt' after a change in the environment will be measured. This will assess how quickly each algorithm can adjust to new conditions.

Sensitivity Analysis: For each algorithm, how the performance metrics change with respect to their hyperparameters (N for Sliding Window UCB and $\gamma$ for Discount Factor UCB) will be scrutinized.

Confidence Intervals: For all metrics, 95% confidence intervals will be computed to provide a robust assessment of performance variability.

## 5. Results

### 5.1. Traditional UCB in Non-Stationary Settings
As a starting point, Traditional UCB was deployed in various non-stationary settings to establish a baseline for performance comparison. The results demonstrated significant limitations, chiefly reflected in the high cumulative regret and slow adaptation to changes, thus justifying the need for more adaptive techniques.

### 5.2. Sliding Window UCB
Performance:The Sliding Window UCB showcased a dramatic decrease in regret compared to Traditional UCB, particularly in environments characterized by abrupt changes. For smaller window sizes, the algorithm effectively adapted to rapid shifts but was sensitive to noise.

Adaptability: Sliding Window UCB's time-to-adapt was notably less than that of the Traditional UCB, especially when the window size was optimally tuned. However, performance decayed when the window size was misaligned with the rate of environmental changes.

Computational Efficiency: Although generally more efficient than Traditional UCB in terms of regret minimization, Sliding Window UCB incurred additional computational overhead due to the necessity to manage and update the sliding windows for each arm.

### 5.3. *Discount Factor UCB*

Performance: In environments where changes were more gradual, Discount Factor UCB outperformed both Traditional and Sliding Window UCB in terms of regret minimization, owing to its smoother adjustment to shifting reward distributions.

Adaptability: The adaptability of Discount Factor UCB was significant but varied depending on the selected discount factor. For slowly changing environments, higher values of $\gamma$ yielded better results.

Computational Efficiency: Discount Factor UCB demonstrated slightly better computational efficiency compared to Sliding Window UCB, as it avoids the overhead of window management by employing a mathematical decay function.

### 5.4. *Comparative Analysis*

Both tabular and graphical representations were employed to compare the performance, adaptability, and computational efficiency of Sliding Window UCB and Discount Factor UCB: Tables: Detailed tables were constructed to list the specific metrics, including regret, time-to-adapt, and computational costs.

Charts: Line charts for regret and bar charts for computational time were utilized for visual comparison.

Statistical Tests: T-tests were conducted to ascertain whether differences in performance metrics between the two algorithms were statistically significant.

## 6. Discussion

The analysis paints a nuanced picture where neither Sliding Window UCB nor Discount Factor UCB stands out as categorically superior. Instead, their efficacy is closely tied to the specific attributes of the non-stationary environments they navigate. While Sliding Window UCB thrives in fast-evolving contexts, it bears the weight of heightened computational demands and a heightened vulnerability to noise. Conversely, Discount Factor UCB shines in environments undergoing slower transitions and offers superior computational efficiency.

Sliding Window UCB's adaptability is notably impressive when the window size mirrors the pace of environmental shifts. However, this adaptability introduces computational complexities. Continually updating and maintaining individual sliding windows for each arm amplifies the algorithm's computational load. Thus, while adaptability remains Sliding Window UCB's hallmark, it might not be optimal for contexts with limited computational bandwidth. On the other end, Discount Factor UCB offers adaptability, albeit in a more understated, gradual fashion. This quality becomes particularly beneficial in slowly transitioning environments. With its lower computational demands compared to Sliding Window UCB, it utilizes a decay function, eliminating the need for data window management.

These insights enrich the discourse on the regret bounds associated with these algorithms. For instance, the analysis hints at the possibility that tweaking the discount factor or window size could spawn new or enhanced regret bounds, contingent on the intricacies of the non-stationary context. Such revelations pave the way for additional theoretical exploration aimed at crafting more precise or adaptive regret models. Practically, these findings carry significant weight. In the turbulent realm of financial markets, Sliding Window UCB might be ideal for high-frequency trading, where swift adaptability becomes paramount. In contrast, Discount Factor UCB's knack for navigating gradual shifts might make it more apt for long-term portfolio optimization. In the healthcare domain, the task of monitoring abrupt changes in patient data might be best entrusted to Sliding Window UCB. Yet, for chronicling long-term

patient metrics, a more measured approach via Discount Factor UCB could be more fitting. In summary, the choice between Sliding Window UCB and Discount Factor UCB should be dictated by the specific requirements of the application, including the rate of environmental change and computational resources available. Each has its merits and drawbacks, making them complementary tools in the toolbox for tackling multi-armed bandit problems in non-stationary environments.

## 7. Conclusion

The research provides a thorough evaluation of Sliding Window UCB and Discount Factor UCB algorithms in addressing the multi-armed bandit problem in non-stationary environments. Both algorithms exhibit marked enhancements over Traditional UCB in aspects of regret minimization and adaptability under fluctuating conditions. Sliding Window UCB stands out in rapidly changing environments but incurs higher computational overhead and increased sensitivity to noise. In contrast, Discount Factor UCB emerges as superior in scenarios with gradual environmental transitions and demonstrates greater computational efficiency. However, this research has certain limitations. Firstly, the datasets, both simulated and from real-world scenarios, might not encompass all intricacies inherent to non-stationary environments. Secondly, the primary focus on regret as a metric might overshadow other pertinent measures like fairness or robustness. Additionally, the selected parameters for each algorithm, though diversified, might not be all-encompassing, leaving potential performance facets uninvestigated.

Considering these limitations, it is recommended that subsequent studies engage more varied data sources for further validation. Broadening the metrics beyond regret could furnish a comprehensive perspective on algorithmic efficacy. Research dedicated to fine-tuning the parameters of Sliding Window and Discount Factor UCB for distinct non-stationary settings would be of notable worth. Furthermore, the potential of hybrid models that fluidly transition between the two algorithms based on real-time evaluations presents an intriguing prospect for future investigations.

## References

[1]   Cavenaghi, E., Sottocornola, G., Stella, F., & Zanker, M. (2021). Non stationary multi-armed bandit: Empirical evaluation of a new concept drift-aware algorithm. Entropy, 23(3), 380.
[2]   Cho, S. (2023). Use of Logarithmic Rates in Multi-Armed Bandit-Based Transmission Rate Control Embracing Frame Aggregations in Wireless Networks. Applied Sciences, 13(14), 8485.
[3]   Chaudhary, A., Rai, A., & Gupta, A. (2023). Maximizing Success Rate of Payment Routing using Non-stationary Bandits. arXiv preprint arXiv:2308.01028.
[4]   Nilsson, J. (2022). Multi-Armed Bandit to optimize the pricing strategy for consumer loans.
[5]   Russac, Y., Vernade, C., & Cappé, O. (2019). Weighted linear bandits for non-stationary environments. Advances in Neural Information Processing Systems, 32.
[6]   Kangas, A. (2021). Towards Embedded Implementations of Multi-Armed Bandit Algorithms for Optimised Channel Selection.
[7]   Cho, S. (2022, December). Multi-armed Bandit-Based Rate Control with Logarithmic Rates in CSMA/CA Wireless Networks. In International Conference on Computer Science and its Applications and the International Conference on Ubiquitous Information Technologies and Applications (pp. 631-637). Singapore: Springer Nature Singapore.
[8]   Qi, H., Wang, Y., & Zhu, L. (2023). Discounted Thompson Sampling for Non-Stationary Bandit Problems. arXiv preprint arXiv:2305.10718.
[9]   Balef, A. R., & Maghsudi, S. (2023). Piecewise-Stationary Multi-Objective Multi-Armed Bandit with Application to Joint Communications and Sensing. IEEE Wireless Communications Letters.
[10]  Jedor, M., Louëdec, J., & Perchet, V. (2020). Lifelong Learning in Multi-Armed Bandits. arXiv preprint arXiv:2012.14264.