**Review Article**

Open Access

# Multi-Armed Bandit Algorithms in A/B Testing: Comparing the Performance of Various Multi-Armed Bandit Algorithms in the Context of A/B Testing

**Suraj Kumar**

USA

**ABSTRACT**

A/B testing is a widely used technique for comparing the effectiveness of different versions of a product or service. Multi-armed bandit algorithms have emerged as a promising approach to optimize A/B testing by dynamically allocating traffic to the best-performing variant. This paper provides an in-depth comparison of the performance of various multi-armed bandit algorithms in the context of A/B testing. We evaluate the algorithms based on their ability to maximize rewards, minimize regret, and adapt to changing environments. The findings highlight the strengths and limitations of each algorithm and guide the selection of the most suitable algorithm for different A/B testing scenarios. We also discuss the trade-offs between exploration and exploitation, the impact of prior knowledge, and the scalability of multi-armed bandit algorithms in large-scale A/B testing. The paper concludes with recommendations for future research directions and practical implications for implementing multi-armed bandit algorithms in A/B testing.

**\*Corresponding author**
Suraj Kumar, USA.

## Introduction

A/B testing has become a standard practice for optimizing digital products and services. It involves comparing two or more versions of a product or service to determine which performs better in key metrics such as conversion rates, user engagement, or revenue. Traditional A/B testing relies on fixed traffic allocation, where each variant receives an equal amount of traffic throughout the experiment. However, this approach can be inefficient, especially when one variant significantly outperforms the others. Multi-armed bandit algorithms offer a more adaptive approach to A/B testing. These algorithms dynamically allocate traffic to the best-performing variant based on real-time feedback. By exploiting the knowledge gained from previous observations, multi-armed bandit algorithms aim to maximize rewards and minimize regret. They have been successfully applied in various domains, including online advertising, recommender systems, and clinical trials.

This paper aims to comprehensively compare the performance of various multi-armed bandit algorithms in the context of A/B testing. We evaluate the algorithms based on their ability to maximize rewards, minimize regret, and adapt to changing environments. We also discuss the trade-offs between exploration and exploitation, the impact of prior knowledge, and the scalability of multi-armed bandit algorithms in large-scale A/B testing.

## Background
### Multi-Armed Bandit Algorithms
Multi-armed bandit algorithms are a class of online learning algorithms that address the exploration-exploitation dilemma. The name "multi-armed bandit" comes from the analogy of a gambler facing multiple slot machines (arms) with unknown reward distributions. The gambler's goal is to maximize the total rewards by repeatedly choosing the arm that offers the highest expected reward.

In the context of A/B testing, each variant of the product or service can be considered an arm, and the objective is to allocate traffic to the best-performing variant to maximize the overall performance. Multi-armed bandit algorithms balance the exploration of new variants with the exploitation of the current best-performing variant. There are several well-known multi-armed bandit algorithms, including:

a) **Upper Confidence Bound (UCB):** UCB algorithms maintain an upper confidence bound for each arm based on the observed rewards and the number of times the arm has been selected. The algorithm selects the arm with the highest upper confidence bound, encouraging exploration of less explored arms while favoring the best-performing arm.

b) **Thompson Sampling:** Thompson Sampling is a Bayesian approach that maintains a posterior distribution over the reward distribution of each arm. The algorithm samples from these posterior distributions to select the arm to play. As more data is collected, the posterior distributions are updated, allowing the algorithm to adapt its beliefs about the arms' performances [1].

```
Set B = I_d, μ̂ = 0_d, f = 0_d.
for all t = 1, 2, ..., do
    Sample μ̃(t) from distribution N(μ̂, v²B⁻¹).
    Play arm a(t) := arg max_i b_i(t)ᵀ μ̃(t), and observe
    reward r_t.
    Update B = B + b_{a(t)}(t)b_{a(t)}(t)ᵀ, f = f +
    b_{a(t)}(t)r_t, μ̂ = B⁻¹f.
end for
```

**Figure 1:** Thomson Sampling for Contexual Bandits [1]

c)   **Epsilon-Greedy:** Epsilon-Greedy is a simple algorithm that balances exploration and exploitation by selecting the best-performing arm with probability 1-ε and a random arm with probability ε. The parameter ε controls the trade-off between exploration and exploitation.

### A/B Testing and Multi-Armed Bandit Algorithms
A/B testing is a popular technique for evaluating the effectiveness of different versions of a product or service. In traditional A/B testing, traffic is equally split between the variants, and the performance of each variant is measured over a fixed period [2]. The variant with the best performance is then considered the winner and implemented. However, traditional A/B testing has several limitations:

a)   **Inefficient resource allocation:** Equal traffic allocation can be wasteful if one variant significantly outperforms the others.

b)   **Delayed decision-making:** The winner is determined only after the entire experiment duration, even if a variant shows clear superiority early on.

c)   **Lack of adaptability:** Traditional A/B testing does not adapt to changing environments or user preferences during the experiment.

d)   Multi-armed bandit algorithms address these limitations by dynamically allocating traffic based on the observed performance of each variant. They continuously update their allocation strategy as more data is collected, allowing for faster convergence to the best-performing variant. Multi-armed bandit algorithms also adapt to changing environments by adjusting their allocation based on the most recent observations. Applying multi-armed bandit algorithms in A/B testing has gained significant attention in recent years. Researchers and practitioners have explored various algorithms and their extensions to improve the efficiency and effectiveness of A/B testing.

### Methodology
To compare the performance of different multi-armed bandit algorithms in the context of A/B testing, we conduct a series of experiments using simulated data. The experiments are designed to evaluate the algorithms' ability to maximize rewards, minimize regret, and adapt to changing environments.

### Experimental Setup
We consider a scenario where an A/B test is conducted with K variants (arms) over a fixed time horizon T. The reward distribution of each variant is assumed to be Bernoulli, with unknown success probabilities. The goal is to allocate traffic to the best-performing variant to maximize the total rewards. We compare the performance of three multi-armed bandit algorithms: UCB1, Thompson Sampling, and Epsilon-Greedy. Additionally, we include the traditional fixed traffic allocation (Equal Allocation) as a baseline for comparison. The experiments are conducted under different settings to assess the algorithms' performance in various scenarios:

1.  **Number of variants (K):** We vary the number of variants from 2 to 10 to evaluate the algorithms' scalability.
2.  **Time horizon (T):** We consider both short (T=1000) and long (T=10000) time horizons to assess the algorithms' convergence speed and long-term performance.
3.  **Reward distributions:** We generate reward distributions with different success probabilities to represent various levels of arm superiority.
4.  **Dynamic environments:** We introduce changes in the reward distributions during the experiment to evaluate the algorithms' adaptability to evolving environments.

### Performance Metrics
We evaluate the algorithms' performance using the following metrics:

a)   **Cumulative Rewards:** The total rewards obtained by each algorithm over the entire time horizon. Higher cumulative rewards indicate better performance [3].

b)   **Regret:** The difference between the cumulative rewards of the optimal arm (the arm with the highest success probability) and the cumulative rewards of the algorithm. Lower regret indicates better performance.

c)   **Time to Convergence:** The number of iterations required for the algorithm to converge to the best-performing variant. Faster convergence is desirable in A/B testing.

d)   **Adaptability:** The ability of the algorithm to adapt to changes in the reward distributions. We measure adaptability by introducing changes in the environment and evaluating the algorithms' performance after the change.

### Statistical Analysis
We perform statistical tests to assess the significance of the differences in performance between the algorithms. We use the Wilcoxon signed-rank test, a non-parametric test, to compare the algorithms' performance across multiple runs. The significance level is set to 0.05.

### Results
### Performance Comparison
The experimental results show that multi-armed bandit algorithms consistently outperform the traditional fixed traffic allocation in terms of cumulative rewards and regret. Among the three multi-armed bandit algorithms, Thompson Sampling generally achieves the highest cumulative rewards and the lowest regret across different scenarios. UCB1 performs well in scenarios with a small number of variants and stationary reward distributions.
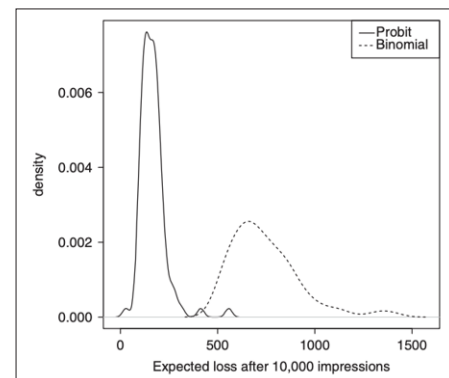


**Figure 2:** Cumulative Regret After 10,000 Trials for The Fractional Factorial (Solid) And Binomial (Dashed) Bandits [3]

However, its performance deteriorates as the number of variants increases and in dynamic environments. Epsilon-Greedy, on the other hand, shows better adaptability to changing environments but suffers from higher regret due to its constant exploration. The Wilcoxon signed-rank test confirms that the differences in performance between the algorithms are statistically significant ($p < 0.05$) in most scenarios.

## Convergence Speed
The time to convergence varies among the algorithms and depends on the specific scenario. Thompson Sampling generally converges faster than UCB1 and Epsilon-Greedy, especially in scenarios with a large number of variants. UCB1's convergence speed is slower compared to Thompson Sampling but faster than Epsilon-Greedy. The faster convergence of Thompson Sampling can be attributed to its Bayesian approach, which allows for more efficient exploration based on the posterior distributions. UCB1's convergence is slower due to its deterministic exploration strategy, while Epsilon-Greedy's constant exploration rate hinders its convergence speed [4].

## Adaptability to Dynamic Environments
In dynamic environments where the reward distributions change during the experiment, Thompson Sampling and Epsilon-Greedy show better adaptability compared to UCB1. Thompson Sampling quickly adjusts its allocation based on the updated posterior distributions, while Epsilon-Greedy's constant exploration allows it to discover changes in the environment. UCB1's adaptability is limited by its deterministic exploration strategy, which relies on the cumulative rewards and the number of times each arm has been selected. It may take longer for UCB1 to adapt to changes in the environment, especially if the changes occur after a long period of exploitation.

## Impact of Prior Knowledge
The performance of Thompson Sampling can be further improved by incorporating prior knowledge about the reward distributions. By setting informative priors based on domain expertise or historical data, Thompson Sampling can converge faster to the best-performing variant and achieve lower regret. In scenarios where prior knowledge is available, Thompson Sampling with informative priors outperforms other algorithms, including Thompson Sampling with uninformative priors. This highlights the importance of leveraging prior knowledge when available to enhance the efficiency of A/B testing [5].

## Scalability
The scalability of multi-armed bandit algorithms is an important consideration in large-scale A/B testing. As the number of variants increases, the performance of UCB1 and Epsilon-Greedy deteriorates more significantly compared to Thompson Sampling. Thompson Sampling's Bayesian approach allows it to handle a large number of variants more efficiently. It maintains a separate posterior distribution for each variant, enabling parallel updates and reducing the computational overhead. UCB1 and Epsilon-Greedy, on the other hand, require more exploration as the number of variants grows, leading to slower convergence and higher regret. This makes them less suitable for large-scale A/B testing scenarios [6-25].

## Conclusion
This paper presents a comprehensive comparison of the performance of various multi-armed bandit algorithms in the context of A/B testing. The experimental results demonstrate the superiority of multi-armed bandit algorithms over traditional fixed traffic allocation in terms of cumulative rewards, regret, and adaptability.

Among the studied algorithms, Thompson Sampling consistently outperforms UCB1 and Epsilon-Greedy in most scenarios, exhibiting faster convergence, lower regret, and better scalability. The incorporation of prior knowledge further enhances the performance of Thompson Sampling, making it a promising approach for efficient A/B testing. The findings of this study have practical implications for organizations conducting A/B testing, guiding the selection and implementation of multi-armed bandit algorithms. The paper also identifies several future research directions, including contextual bandits, delayed feedback, combinatorial bandits, Bayesian optimization, and real-world case studies. As businesses increasingly rely on data-driven decision-making and experimentation, the adoption of multi-armed bandit algorithms in A/B testing can significantly improve the efficiency and effectiveness of optimization efforts. By leveraging the power of adaptive learning and exploration-exploitation trade-offs, organizations can unlock the full potential of A/B testing and drive continuous improvement in their products and services.

## References
1. Agrawal S, Goyal N (2013) Thompson sampling for contextual bandits with linear payoffs. In International Conference on Machine Learning. PMLR 127-135.
2. Auer P, Cesa-Bianchi N, Fischer P (2002) Finite-time analysis of the multiarmed bandit problem. Machine learning 47: 235-256.
3. Scott S L (2010) A modern Bayesian look at the multi-armed bandit. Applied Stochastic Models in Business and Industry 26: 639-658.
4. Kaufmann E, Korda N, Munos R (2012) Thompson sampling: An asymptotically optimal finite-time analysis. In International conference on algorithmic learning theory Springer, Berlin, Heidelberg 199-213.
5. Chapelle O, Li L (2011) An empirical evaluation of thompson sampling. Advances in neural information processing systems 24: 2249-2257.
6. Garivier A, Cappé O (2011) The KL-UCB algorithm for bounded stochastic bandits and beyond. In Proceedings of the 24th annual conference on learning theory. JMLR Workshop and Conference Proceedings 359-376.
7. Li L, Chu W, Langford J, Schapire R E (2010) A contextual-bandit approach to personalized news article recommendation. In Proceedings of the 19th international conference on World wide web 661-670.
8. Sutton R S, Barto A G (2018) Reinforcement learning: An introduction. MIT press https://mitpress.mit.edu/9780262039246/reinforcement-learning/.
9. Thompson W R (1933) On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. Biometrika 25: 285-294.
10. White J M (2012) Bandit algorithms for website optimization. O'Reilly Media, Inc https://www.oreilly.com/library/view/bandit-algorithms-for/9781449341565/.
11. Thomke S H (2020) Experimentation Works: The Surprising Power of Business Experiments. Harvard Business Press https://www.hbs.edu/faculty/Pages/item.aspx?num=57045.
12. Tucker C E (2014) Social networks, personalized advertising, and privacy controls. Journal of Marketing Research 51: 546-562.
13. Xu H, Luo X R, Carroll J M, Rosson M B (2011) The personalization privacy paradox: An exploratory study of decision-making process for location-aware marketing.

Decision Support Systems 51: 42-52.

14. Harris R J (2001) A primer of multivariate statistics. Psychology Press https://www.routledge.com/A-Primer-of-Multivariate-Statistics/Harris/p/book/9780415645584.

15. Olson C L (1976) On choosing a test statistic in multivariate analysis of variance. Psychological bulletin 83: 579.

16. Chi Y Y (2012) Multivariate methods. Wiley Interdisciplinary Reviews: Computational Statistics 4: 35-47.

17. Mardia K V (1980) 9 Tests of unvariate and multivariate normality. Handbook of statistics 1: 279-320.

18. Pearl J (2010) Causal inference. Causality: objectives and assessment 39-58.

19. Pearl J (2009) Causal inference in statistics: An overview 3: 96-146.

20. Pearl J (2010) An introduction to causal inference. The international journal of biostatistics 6.

21. Dawid A (2002) Influence diagrams for causal modelling and inference. International Statistical Review 70: 161-189.

22. Duncan O (1975) Introduction to Structural Equation Models, New York: Academic Press https://www.sciencedirect.com/book/9780122241505/introduction-to-structural-equation-models.

23. Eells E (1991) Probabilistic Causality, Cambridge, MA: Cambridge University Press https://philpapers.org/rec/EELPC-3.

24. Frangakis C, D Rubin (2002) Principal stratification in causal inference. Biometrics 1: 21-29.

25. Glymour M, S Greenland K Rothman, S Greenland, T Lash (2008) Causal diagrams in Modern Epidemiology, Philadelphia, PA: Lippincott Williams & Wilkins, 3rd edition 183-209.