
EXP4-DFDC: A Non-Stochastic Multi-Armed Bandit for Cache Replacement

Anonymous Author(s)

Abstract

1 In this work we study a variant of the well-known multi-armed bandit (MAB)
2 problem, which has the properties of a delay in feedback, and a loss that declines
3 over time. We introduce an algorithm, EXP4-DFDC, to solve this MAB variant,
4 and demonstrate that the regret vanishes as the time increases. We also show
5 that LeCaR, a previously published machine learning-based cache replacement
6 algorithm, is an instance of EXP4-DFDC. Our results can be used to provide insight
7 on the choice of hyperparameters, and optimize future LeCaR instances.

8 1 Introduction

9 The *Multi-Armed Bandit* (MAB) problem [6] can be thought of as a sequence of “pull” requests for
10 which each request must elicit a response. MAB can be defined as a game where a player plays a
11 T -round prediction game: the player has access to a panel of N experts, each of whom is following a
12 specific strategy. In any round t , the player consults the N experts. Expert i recommends an action
13 $\xi_i(t) \in \{0, \dots, K\}$ from K available options. Each action is associated with a loss value, which is
14 revealed as feedback to the player either immediately or after a delay.

15 We consider a variant of the MAB problem, which takes into account cases for which as the delay in
16 feedback gets larger, the loss associated with the chosen action decreases. This variant differs from
17 any existing formulations of MAB with a delayed feedback [8, 1, 4, 5, 2, 3] where greater feedback
18 delays incur higher loss.

19 The cache replacement problem in storage systems can be considered as an instance of the MAB
20 problem: feedback is available on cache misses due to recent page evictions, and misses not found in
21 the page eviction history correspond to evictions from a distant past and are ignored.

22 2 Problem Statement

23 For our version of MAB, each cache replacement algorithm is thought of as an “expert” that
24 recommends which of the K items in the cache to evict on every miss. We assume that each
25 algorithm is allowed a limited size “history” to keep track of recent evictions, and to help in its
26 decision making.

27 The history also provides a natural aging mechanism for recent evictions. If the item is still in history,
28 feedback on its eviction is considered valuable and is used; otherwise, it is assumed that a significant
29 amount of time has passed, and the feedback is ignored. For an eviction that occurred at time (or
30 step or round) t' , the delayed feedback offered at time $t \geq t'$ is denoted by the indicator function,
31 $f_{(t',t)} \in \{0, 1\}$. If feedback arrives at time t , then $f_{(t',t)} = 1$, and is 0 otherwise. The delay function
32 is given by $d(t) = t - t' \in [1, m]$, where m is the delay threshold (as well as the history size).
33 Feedback provided at time t , for an action $\xi_i(t')$ taken at time t' , is based on the recommendation
34 of expert i , and observes a *loss* that is denoted by $l_{\xi_i(t')}(t, t')$. Note that feedback provided after a
35 delay threshold m incurs a loss of 0, and we consider the oblivious scenario where the loss function
36 is fixed, not adaptive.

37 3 Technical Contributions

38 Expert i is associated with a weight w_i , and all
 39 experts start off with equal weights. For each
 40 request, each expert recommends an action; a
 41 final action is chosen to be expert i 's recom-
 42 mendations with probability proportional to w_i .
 43 The weights of the experts get updated when-
 44 ever $f_{i(t',t)} = 1$, i.e., when feedback is provided.
 45 Shorter delays incur higher loss, and the loss
 46 decays with delay. Thus, regret is highest when
 47 feedback is immediate, i.e., when $t = t' + 1$ for
 48 the chosen action. If the maximum loss for a de-
 49 cision is $L = l_{\xi_i(t')}(t' + 1)$, then the estimated
 50 loss after delay $d(t)$ is given by $L/d(t)$, until the
 51 delay reaches the predefined threshold m , after
 52 which it becomes negligible. (See Figure 1.) As
 53 $t - t'$ increases, the loss for the chosen action
 54 decays. Thus, the Estimated loss calculation in
 55 EXP4-DFDC the major difference from EXP4.

56 The regret, $R_A(T)$, of any adaptive learning al-
 57 gorithm A after round T is used to measure its
 58 performance. The main objective of an adap-
 59 tive learning strategy is to minimize this regret,
 60 achieved by ensuring it vanishes over time. An
 61 explicit bound for the regret function can show
 62 that the algorithm EXP4-DFDC can be com-
 63 puted and that it has a vanishing regret. The
 64 regret of EXP4-DFDC can be shown to be of
 65 the form $O(\sqrt{KT \ln N})$, and as a result, the average regret will vanish with increasing time as
 66 $\bar{R}_A(t) = R_A(t)/T = O(T^{-1/2})$. Thus, as $T \rightarrow \infty$, $\bar{R}_A(t) \rightarrow 0$. We focus on applications where
 67 the loss of an action decays with a delay in feedback, but only up to a limit, after which it becomes
 68 zero. Details of the proof are omitted from this short abstract.

69 4 Discussion

70 We consider the cache management problem, which is a classic resource allocation problem from
 71 storage systems. The LeCaR algorithm is a recent cache replacement algorithm that uses online
 72 reinforcement learning to help decide which entry to evict when a new item is to be stored in the
 73 cache following a “cache miss” [7]. LeCaR relies on only two very fundamental cache replacement
 74 policies, namely the *Least Recently Used* (LRU) policy and the *Least Frequently Used* (LFU) policy.
 75 As a result, $N = 2$ represents the number of experts in LeCaR. LeCaR assumes that the best strategy
 76 at any given time is a probabilistic mix of the two policies and attempts to “learn” the optimal mix
 77 using a regret minimization strategy. With a slight modification, the LeCaR [7] cache replacement
 78 algorithm is an instance of EXP-DFDC: the position of an item in history is a reflection of the delay
 79 in feedback for its eviction decision, and is bounded by the history size. The regret is greatest if an
 80 item is requested immediately after its eviction. Else, as time wears on, until a threshold is reached,
 81 the regret continues to decay. If it is not present in the history, then the feedback is ignored, thus
 82 making its regret equal to zero. Cache management can be thought of as a special case of the MAB
 83 problem, and LeCaR can be considered as an application-specific version of EXP4-DFDC with two
 84 experts, delayed feedback and decaying losses, and has been shown to have vanishing regret over
 85 time. We note that LeCaR outperforms other state-of-the-art cache replacement algorithms for small
 86 cache sizes [7] with a fixed learning rate of 0.45. From our analysis, we demonstrate that it is possible
 87 to infer the learning rate in this class of algorithm without it being fixed a priori. Details of the
 88 mathematical proof of vanishing regret can be found in a forthcoming manuscript.

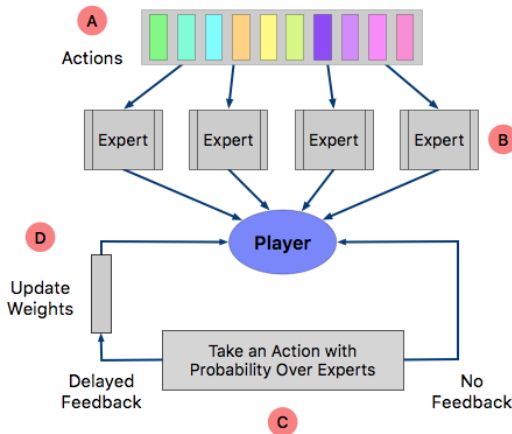


Figure 1: MAB problem variant. A) The set of K actions that are available to the panel of N experts. B) Each expert recommends an action to the player, which selects one. C) The feedback to the player can be immediate, or delayed, in which case the weights of the experts are updated based on the loss the player observed (D).

89 **References**

- 90 [1] Alekh Agarwal and John C Duchi. Distributed delayed stochastic optimization. In *Advances in*
91 *Neural Information Processing Systems*, pages 873–881, 2011.
- 92 [2] Thomas Desautels, Andreas Krause, and Joel W Burdick. Parallelizing exploration-exploitation
93 tradeoffs in gaussian process bandit optimization. *The Journal of Machine Learning Research*,
94 15(1):3873–3923, 2014.
- 95 [3] Miroslav Dudik, Daniel Hsu, Satyen Kale, Nikos Karampatziakis, John Langford, Lev
96 Reyzin, and Tong Zhang. Efficient optimal learning for contextual bandits. *arXiv preprint*
97 *arXiv:1106.2369*, 2011.
- 98 [4] John Langford, Alexander J Smola, and Martin Zinkevich. Slow learners are fast. *Advances in*
99 *Neural Information Processing Systems*, 22:2331–2339, 2009.
- 100 [5] Gergely Neu, Andras Antos, András György, and Csaba Szepesvári. Online markov decision
101 processes under bandit feedback. In *Advances in Neural Information Processing Systems*, pages
102 1804–1812, 2010.
- 103 [6] Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins*
104 *Selected Papers*, pages 169–177. Springer, 1985.
- 105 [7] Giuseppe Vietri, Liana V. Rodriguez, Wendy A. Martinez, Steven Lyons, Jason Liu, Raju
106 Rangaswami, Ming Zhao, and Giri Narasimhan. Driving cache replacement with ML-based
107 LeCaR. In *10th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 18)*,
108 Boston, MA, 2018. USENIX Association.
- 109 [8] Marcelo J Weinberger and Erik Ordentlich. On delayed prediction of individual sequences. *IEEE*
110 *Transactions on Information Theory*, 48(7):1959–1976, 2002.