# Online Resource Allocation with Limited Flexibility

Arash Asadpour, Xuan Wang, Jiawei Zhang

# Online Resource Allocation with Limited Flexibility

## 1 Introduction

In this paper, we consider a general class of resource allocation problems where requests (or items, tasks) arrive sequentially and need to be fulfilled by resources. There are $n$ different resources with the same initial capacity. There are $n$ types of requests, and a request is of type $j$ if it can be fulfilled by either resource $j$ or resource $j + 1$ (modulo $n$). This limited choices (substitution) brings flexibility to the system because the decision maker can choose how to fulfill the request at his discretion. We call such a pattern the *long chain substitution*, where the name comes from the graphical representation of the structure — if each resource and each request type are represented as a node and an edge between them means that the request can be fulfilled by the resource, then all the nodes are chained in a single cycle. A request reveals its type upon arrival, and the decision maker needs to choose one resource from the corresponding subset to fulfill this request and one unit of the chosen resource's capacity is consumed. In particular, we consider a *2-choice myopic online policy* — a type $j$ request is fulfilled by the resource with higher capacity among resources $j$ and $j + 1$. If both resources have zero capacity left, then this request fails to be fulfilled and we say a "lost sale" occurs. There are $m$ requests in total, and the goal is to fulfill as many requests as possible, i.e., to minimize the total expected lost sales.

The above resource allocation problem with limited flexibility has several interesting applications. One example is the inventory allocation problem when consumers' product substitution behavior is present. The increasing product variety and uncertain consumer demand are admittedly two prominent features in today's competitive market. It is common practice for manufacturers to change their products over attributes such as the color of shirts and the flavor of yoghurt to cater to consumers' heterogenous preferences while at the same time avoiding huge costs of new product development. Faced with a large number of similar variants to choose from, consumers' demand for a certain product can usually be satisfied by substitute products. The long chain substitution applies when consumers' individual taste favors only a small subset of similar products. For example, a customer who wishes to buy a black jacket may be willing to accept dark blue as a substitute, but not a yellow one. Moreover, when choosing among different products that involve price v.s. quality tradeoff, consumers usually would not sacrifice one to the other too much and often move up or down one step in the tradeoff. The long chain substitution may also be applicable to set-

tings in which the consumers' preferences are based on geographical conditions. For instance, an increasing number of retailers are now offering the "buy online and pick-up in store" service that aims to help reduce shipping and inventory costs. In such a service, customers place an order online and are offered a list of brick-and-mortar stores where they can pick up their order, and customers would usually prefer stores that are near their own location.

The resource allocation problem with limited flexibility also finds interesting applications in load balancing. In such a problem, tasks arrive sequentially and need to be assigned to a server. In the interest of response time of the tasks, the objective of an assignment procedure is to keep the maximum load of any server as small as possible. It is often the case that a task can be performed by a subset of servers with similar skills, and we can take advantage of the limited flexibility by assigning the task to the least loaded server among the two choices to achieve a balanced allocation.

## 2 An Abstraction: The Balls-into-Bins Model

The online resource allocation problem can be naturally modelled as a balls-into-bins problem. In a traditional *balls-into-bins* problem, $m$ balls are sequentially thrown into $n$ bins by some randomized placement algorithm. In general, the goal of such a placement procedure is to achieve a balanced allocation such that the load is distributed among the bins as evenly as possible. One metric to evaluate the evenness of the load distribution is the *maximum deviation*, which is the difference between the maximum load (i.e., the number of balls in the fullest bin) and the average load $m/n$ achieved by the optimal offline allocation. A classic balls-into-bins algorithm is the *d-choice scheme*. In such a process, each ball is placed in the least loaded of $d$ bins chosen independently and uniformly at random at the time of the placement. A well-known result in the CS literature is that even a small amount of choices (flexibility) can lead to significant improvement in load balancing. More specifically, it has been shown that the maximum deviation of the $d$-choice scheme is independent of the number of balls when $d \geq 2$, in contrast to the single choice scheme ($d = 1$) in which the maximum deviation diverges with the number of balls (see Raab and Steger 1998, Berenbrink et al. 2006).

The balls-into-bins model serves as a natural abstraction of our online resource allocation problem with limited flexibility. In our model, $m$ requests (balls) sequentially arrive, and each of them chooses $d = 2$ out of the $n$ resources (bins). Unlike the $d$-choice scheme where the choices are sampled independently and uniformly at random, each request's two choices follow the long chain

structure. The myopic online policy that uses the resource with higher capacity to fulfill the request is a reverse analog to the placement procedure in which the ball is thrown into the least loaded bin. Assume that each resource has an initial capacity of $m/n$, then it is clear that the total number of lost sales after all the $m$ requests have arrived is equivalent to the total excess load above the average $m/n$ when all the $m$ balls have been placed.

# 3   Main results

The results in Berenbrink et al. (2006) that the maximum deviation is independent of the number of balls immediately implies that the expected total number of lost sales is independent of the number of requests under random two choices. Our main result below shows that even under the long chain substitution, which is much less flexible than the random two choices, the expected total number of lost sales being independent of the number of requests continues to hold.

**Theorem 1.** *Assume there are n resources each with initial capacity m/n, and each of the m requests exhibits the long chain substitution. Then under the 2-choice myopic online allocation policy, the expected total number of lost sales is independent of m and bounded from above by $4n^2 \ln n$.*

We have also established the optimality of the long chain substitution in terms of minimizing the expected total lost sales among all structures with $2n$ arcs under the myopic online policy. In fact, if there exists some type of requests that can be fulfilled by only one resource, or if some resource can only serve one type of requests, then the expected total lost sales diverges with the number of requests. This implies that the limited choices (flexibility) is not only efficient, but also crucial.

# References

Berenbrink, P., Czumaj, A., Steger, A., and Vcking, B. (2006). Balanced allocations: The heavily loaded case. *SIAM Journal on Computing*, 35(6):1350–1385.

Raab, M. and Steger, A. (1998). "balls into bins" - a simple and tight analysis. In *Proceedings of the Second International Workshop on Randomization and Approximation Techniques in Computer Science*, RANDOM '98, pages 159–170, London, UK, UK. Springer-Verlag.