

# Linear Multi-Resource Allocation with Semi-Bandit Feedback

**Authored by:**

Koby Crammer  
Csaba Szepesvari  
Tor Lattimore

## **Abstract**

We study an idealised sequential resource allocation problem. In each time step the learner chooses an allocation of several resource types between a number of tasks. Assigning more resources to a task increases the probability that it is completed. The problem is challenging because the alignment of the tasks to the resource types is unknown and the feedback is noisy. Our main contribution is the new setting and an algorithm with nearly-optimal regret analysis. Along the way we draw connections to the problem of minimising regret for stochastic linear bandits with heteroscedastic noise. We also present some new results for stochastic linear bandits on the hypercube that significantly out-performs existing work, especially in the sparse case.

## **1 Paper Body**

Economist Thomas Sowell remarked that “The first lesson of economics is scarcity: There is never enough of anything to fully satisfy all those who want it.”<sup>1</sup> The optimal allocation of resources is an enduring problem in economics, operations research and daily life. The problem is challenging not only because you are compelled to make difficult trade-offs, but also because the (expected) outcome of a particular allocation may be unknown and the feedback noisy. We focus on an idealised resource allocation problem where the economist plays a repeated resource allocation game with multiple resource types and multiple tasks to which these resources can be assigned. Specifically, we consider a (nearly) linear model with  $D$  resources and  $K$  tasks. In each time step  $t$  the economist chooses an allocation of resources  $M_t \in \mathbb{R}^{D \times K}$  where  $M_{tk} \in \mathbb{R}^D$  is the  $k$ th column and represents the amount of each resource type assigned to the  $k$ th task. We assume that the  $k$ th task is completed successfully with probability  $\min\{1, \frac{1}{h} M_{tk}^T \theta_k\}$  and  $\theta_k \in \mathbb{R}^D$  is an unknown non-negative vector that determines how the success rate of a given task depends on the quantity and type

of resources assigned to it. Naturally we will limit the availability of resources PK by demanding that  $M_t$  satisfies  $k=1$   $M_t d_k \leq 1$  for all resource types  $d$ . At the end of each time step the economist observes which tasks were successful. The objective is to maximise the number of successful tasks up to some time horizon  $n$  that is known in advance. This model is a natural generalisation of the one used by Lattimore et al. [2014], where it was assumed that there was a single resource type only. <sup>1</sup> He went on to add that "The first lesson of politics is to disregard the first lesson of economics." Sowell [1993]

<sup>1</sup>

An example application might be the problem of allocating computing resources on a server between a number of Virtual Private Servers (VPS). In each time step (some fixed interval) the controller chooses how much memory/cpu/bandwidth to allocate to each VPS. A VPS is said to fail in a given round if it fails to respond to requests in a timely fashion. The requirements of each VPS are unknown in advance, but do not change greatly with time. The controller should learn which VPS benefit the most from which resource types and allocate accordingly. The main contribution of this paper besides the new setting is an algorithm designed for this problem along with theoretical guarantees on its performance in terms of the regret. Along the way we present some additional results for the related problem of minimising regret for stochastic linear bandits on the hypercube. We also prove new concentration results for weighted least squares estimation, which may be independently interesting. The generalisation of the work of Lattimore et al. [2014] to multiple resources turns out to be fairly non-trivial. Those with knowledge of the theory of stochastic linear bandits will recognise some similarity. In particular, once the nonlinearity of the objective is removed, the problem is equivalent to playing  $K$  linear bandits in parallel, but where the limited resources constrain the actions of the learner and correspondingly the returns for each task. Stochastic linear bandits have recently been generating a significant body of research (e.g., Auer [2003], Dani et al. [2008], Rusmevichientong and Tsitsiklis [2010], Abbasi-Yadkori et al. [2011, 2012], Agrawal and Goyal [2012] and many others). A related problem is that of online combinatorial optimisation. This has an extensive literature, but most results are only applicable for discrete action sets, are in the adversarial setting, and cannot exploit the additional structure of our problem. Nevertheless, we refer the interested reader to (say) the recent work by Kveton et al. [2014] and references there-in. Also worth mentioning is that the resource allocation problem at hand is quite different to the "linear semi-bandit" proposed and analysed by Krishnamurthy et al. [2015] where the action set is also finite (the setting is different in many other ways besides). Given its similarity, it is tempting to apply the techniques of linear bandits to our problem. When doing so, two main difficulties arise. The first is that our payoffs are non-linear: the expected reward is a linear function only up to a point after which it is clipped. In the resource allocation problem this has a natural interpretation, which is that over-allocating resources beyond a certain point is fruitless. Fortunately, one can avoid this difficulty rather easily by ensuring that with high probability resources are never over-allocated. The second problem concerns achieving good

regret regardless of the task specifics. In particular, when the number of tasks  $K$  is large and resources are at a premium the allocation problem behaves more like a  $K$ -armed bandit where the economist must choose the few tasks that can be completed successfully. For this kind of problem regret should scale in the worst case with  $K$  only [Auer et al., 2002, Bubeck and Cesa-Bianchi, 2012]. The standard linear bandits approach, on the other hand, would lead to a bound on the regret that depends linearly on  $K$ . To remedy this situation, we will exploit that if  $K$  is large and resources are scarce, then many tasks will necessarily be under-resourced and will fail with high probability. Since the noise model is Bernoulli, the variance of the noise for these tasks is extremely low. By using weighted least-squares estimators we are able to exploit this and thereby obtain an improved regret. An added benefit is that when resources are plentiful, then all tasks will succeed with high probability under the optimal allocation, and in this case the variance is also low. This leads to a poly-logarithmic regret for the resource-laden case where the optimal allocation fully allocates every task.

2

## Preliminaries

If  $F$  is some event, then  $\bar{F}$  is its complement (i.e., it is the event that  $F$  does not occur). If  $A$  is positive definite and  $x$  is a vector, then  $\|x\|_A^2 = x^T A x$  stands for the weighted 2-norm. We write  $\|x\|$  to be the vector of element-wise absolute values of  $x$ . We let  $R \in \mathbb{R}^{D \times K}$  be a matrix with columns  $r_1, \dots, r_K$ . All entries in  $R$  are non-negative, but otherwise we make no global assumptions on  $R$ . At each time step  $t$  the learner chooses an allocation matrix  $M_t \in \mathbb{R}^{K \times D}$  where  $M_t(d, k) \in [0, 1]$  for all  $d \in [D], k \in [K]$ .

The assumption that each resource type has a bound of 1 is non-restrictive, since the units of any resource can be changed to accommodate this assumption. We write  $M_t(k) \in [0, 1]^D$  for the  $k$ th

column of  $M_t$ . The reward at time step  $t$  is  $y_t = \sum_{k=1}^K M_t(k) Y_{tk}$  where  $Y_{tk} \in \{0, 1\}$  is sampled from a Bernoulli distribution with parameter  $\theta(h(M_t), \eta_k)$  where  $\theta(h(M_t), \eta_k) = \min\{1, h(M_t, \eta_k)\}$ . The economist observes all  $Y_{tk}$ , however, not just the sum. The optimal allocation is denoted by  $M^*$  and defined by  $M^* = \arg \max_{M \in \mathcal{M}} \sum_{k=1}^K \theta(h(M), \eta_k)$ .

$k=1$

We are primarily concerned with designing an allocation algorithm that minimises the expected (pseudo) regret of this problem, which is defined by  $R_n = \sum_{k=1}^K \sum_{t=1}^n \theta(h(M_t), \eta_k) - \sum_{t=1}^n \theta(h(M^*), \eta_k)$ .

$k=1$

where the expectation is taken over both the actions of the algorithm and the observed reward. Optimal Allocations If  $\theta$  is known, then the optimal allocation can be computed by constructing an appropriate linear program. Somewhat surprisingly it may also be computed exactly in  $O(K \log K + D \log D)$  time using Algorithm 1 below. The optimal allocation is not so straightforward as, e.g., simply allocating resources to the incomplete task for which the corresponding  $\theta$  is largest in some dimension. For example, for  $K = 2$  tasks and  $d = 2$  resource types:

$\theta(0.5, 0.5) = \theta(1, 0) = \theta(0, 1) = \theta(1, 1) = 0.5$ . We see that



```

// Compute optimistic allocation:  $M_t = \arg \max_{M \in \mathcal{M}} \sum_k C_{tk} \mathbb{E}[h_{Mtk} | \mathcal{H}_t]$ 
5:
9: 10:
// Observe success indicators  $Y_{tk}$  for all tasks  $k$ :  $Y_{tk} \sim \text{Bernoulli}(\mathbb{E}[h_{Mtk} | \mathcal{H}_t])$ 
11: // Compute weights for all tasks  $k$ :  $w_{tk} = \arg \max_{w \in \mathcal{W}} \sum_k C_{tk} w_k$ 
(1  $\leq$   $h_{Mtk} | \mathcal{H}_t$ ) 0  $h_{Mtk} | \mathcal{H}_t$  13: end for

```

**Computational Efficiency** We could not find an efficient implementation of Algorithm 2 because solving the bilinear optimisation problem in Line 8 is likely to be NP-hard (Bennett and Mangasarian [1993] and also Petrik and Zilberstein [2011]). In our experiments we used a simple algorithm based on optimising for  $M$  and  $w$  in alternative steps combined with random restarts, but for large  $D$  and  $K$  this would likely not be efficient. In the supplementary material we present an alternative algorithm that is efficient, but relies on the assumption that  $k \neq k_1 \neq 1$  for all  $k$ . In this regime it is impossible to over-allocate resources and this fact can be exploited to obtain an efficient and practical algorithm with strong guarantees. Along the way, we are able to construct an elegant algorithm for linear bandits on the hypercube that enjoys optimal regret and adapts to sparsity. Computing the weights  $w_{tk}$  (Line 12) is (somewhat surprisingly) straight-forward. Define  $p = \mathbb{E}[h_{Mtk} | \mathcal{H}_t]$  and  $p_{tk} = h_{Mtk} | \mathcal{H}_t$ .  $p_{tk} \in [0, 1]$ .  $p_{tk} = \mathbb{E}[h_{Mtk} | \mathcal{H}_t]$  if  $p_{tk} \in [0, 1]$  otherwise.

Then the weights can be computed by  $w_{tk} = \arg \max_{w \in \mathcal{W}} \sum_k C_{tk} w_k$ .  
 $w_{tk} = \arg \max_{w \in \mathcal{W}} \sum_k C_{tk} w_k$  if  $p_{tk} \in [0, 1]$  otherwise.

(2)

A curious reader might wonder why the weights are computed by optimising within confidence set  $\mathcal{C}_{tk}$ , which has double the radius of  $\mathcal{C}_{tk}$ . The reason is rather technical, but essentially if the true parameter  $\theta_k$  were to lie on the boundary of the confidence set, then the corresponding weight could become infinite. For the analysis to work we rely on controlling the size of the weights. It is not clear whether or not this trick is really necessary.

4

#### Worst-case Regret for Algorithm 2

We now analyse the regret of Algorithm 2. First we offer a worst-case bound on the regret that depends on the time-horizon like  $O(\sqrt{n})$ . We then turn our attention to the resource-laden case where the optimal allocation satisfies  $h_{Mk} | \mathcal{H}_t = 1$  for all  $k$ . In this instance we show that the dependence on the horizon is only poly-logarithmic, which would normally be unexpected when the

action-space is continuous. The improvement comes from the weighted estimation that exploits the fact that the variance of the noise under the optimal allocation vanishes.

**Theorem 2.** Suppose Algorithm 2 is run with bound  $B \geq \max_k k^2$ . Then

$$\begin{aligned}
& p R_n \leq 1 + 4D \sum_{k=1}^n \max_{k \in K} k^2 + 4 \sum_{k=1}^n \log(1 + 4n^2) \cdot k \\
& \text{Choosing } \eta = B \sum_{k=1}^n \log \\
& 6nN \sum_{k=1}^n \\
& 3nN \sum_{k=1}^n \\
& 2 \\
& \text{and assuming that } B \leq O(\max_{k \in K} k^2), \text{ then} \\
& q \sum_{k=1}^n \log \leq O(D \sum_{k=1}^n \max_{k \in K} k^2 \log n \cdot \log \\
& k
\end{aligned}$$

The proof of Theorem 2 will follow by carefully analysing the width of the confidence sets as the algorithm makes allocations. We start by proving the validity of the confidence sets, and then prove the theorem. **Weighted Least Squares Estimation** For this sub-section we focus on the problem of estimating a single unknown  $\theta = \theta_k$ . Let  $n$   $M_1, \dots, M_n$  be a sequence of allocations to task  $k$  with  $M_t \in \mathcal{R}^D$ . Let  $\{F_t\}_{t=0}$  be a filtration with  $F_t$  containing information available at the end of round  $t$ , which means that  $M_t$  is  $F_{t-1}$  measurable. Let  $\eta_1, \dots, \eta_n$  be the sequence of weights chosen by Algorithm 2. The sequence of outcomes is  $Y_1, \dots, Y_n \in \{0, 1\}$  for which  $E[Y_t | F_{t-1}] = \eta(h(M_t), \theta)$ . The weighted regularised gram matrix is  $G_t = \sum_{s=1}^t \eta_s M_s M_s^T$  and the corresponding weighted least squares estimator is  $\hat{\theta}_t = G_t^{-1} \sum_{s=1}^t \eta_s M_s^T Y_s$ .

**Theorem 3.** If  $\|k\|_2 \leq B$  and  $\eta$  is chosen as in Eq. (1), then  $\|k - \hat{\theta}_t\|_2 \leq \eta \|G_t\|_2^{-1}$  for all  $t \leq n$  with probability at least  $1 - \frac{1}{nK}$ . Similar results exist in the literature for unweighted least-squares estimators (for example, Dani et al. [2008], Rusmevichientong and Tsitsiklis [2010], Abbasi-Yadkori et al. [2011]). In our case, however,  $G_t$  is the weighted gram matrix, which may be significantly larger than an unweighted version when the weights become large. The proof of Theorem 3 is unfortunately too long to include in the main text, but it may be found in the supplementary material. **Analysing the Regret**

We start with some technical lemmas. Let  $F$  be the failure event that  $\|k - \hat{\theta}_t\|_2 \geq \eta \|G_t\|_2^{-1}$  for some  $t \leq n$  and  $1 \leq k \leq K$ . Lemma 4 (Abbasi-Yadkori et al. [2012]). Let  $x_1, \dots, x_n$  be an arbitrary sequence of vectors with  $\|x_t\|_2 \leq 1$  and let  $G_t = I + \sum_{s=1}^t x_s x_s^T$ . Then  $\sum_{t=1}^n \log \det G_t \leq D \log(1 + n) + \frac{1}{2} \sum_{t=1}^n \|x_t\|_2^2$ .

**Corollary 5.** If  $F$  does not hold, then

$$\begin{aligned}
& \sum_{t=1}^n \sum_{k=1}^K \eta_t \|M_t\|_2^2 \leq 8D \log(1 + 4n^2) \cdot \sum_{k=1}^K \eta_t \\
& \sum_{t=1}^n
\end{aligned}$$

The proof is omitted, but follows rather easily by showing that  $\eta_t \|M_t\|_2^2$  can be moved inside the minimum at a price of increasing the loss at most by a factor of four, and then applying Lemma 4. See the supplementary material for the formal proof.

**Lemma 6.** Suppose  $F$  does not hold, then  $\sum_{k=1}^K \sum_{t=1}^n \eta_t \|M_t\|_2^2 \leq D \sum_{k=1}^K \sum_{t=1}^n \eta_t \|M_t\|_2^2 + 4 \sum_{k=1}^K \sum_{t=1}^n \eta_t \log(1 + 4n^2) \cdot k$ .

$$\begin{aligned}
& k=1 \\
& 5
\end{aligned}$$









Figure 3:  $\Delta$ Gap? dependence

Figure 2: Weights

80,000

30,000

40

20,000

20

Weighted Estimator Unweighted Estimator

0 0

1,000,000 t

7

Regret

40,000

?

Regret

60,000

$\Delta t_1 \Delta t_2$

0 0

1,000,000 t

20,000 10,000 0 0.0

0.5 ?

1.0

Conclusions and Summary

We introduced the stochastic multi-resource allocation problem and developed a new algorithm that enjoys near-optimal worst-case regret. The main drawback of the new algorithm is that its computation time is exponential in the dimension parameters, which makes practical implementations challenging unless both  $K$  and  $D$  are relatively small. Despite this challenge we were able to implement that algorithm using a relatively brutish approach to solving the optimisation problem, and this was sufficient to present experimental results on synthetic data showing that the algorithm is behaving as the theory predicts, and that the use of the weighted least-squares estimation is leading to a real improvement. Despite the computational issues, we think this is a reasonable first step towards a more practical algorithm as well as a solid theoretical understanding of the structure of the problem. As a consolation (and on their own merits) we include some other results: ? An efficient (both in terms of regret and computation) algorithm for the case where overallocation is impossible. ? An algorithm for linear bandits on the hypercube that enjoys optimal regret bounds and adapts to sparsity. ? Theoretical analysis of weighted least-squares estimators, which may have other applications (e.g., linear bandits with heteroscedastic noise). There are many directions for future research. The most natural is to improve the practicality of the algorithm. We envisage such an algorithm might be obtained by following the program below: ? Generalise the Thompson sampling analysis for linear bandits by Agrawal and Goyal [2012]. This is a highly non-trivial step, since it is no longer straight-forward to show that such an algorithm is optimistic with high probability. Instead it will be

necessary to make do with some kind of local optimism for each task. ? The method of estimation depends heavily on the algorithm over-allocating its resources only with extremely low probability, but this significantly slows learning in the initial phases when the confidence sets are large and the algorithm is acting conservatively. Ideally we would use a method of estimation that depended on the real structure of the problem, but existing techniques that might lead to theoretical guarantees (e.g., empirical process theory) do not seem promising if small constants are expected. It is not hard to think up extensions or modifications to the setting. For example, it would be interesting to look at an adversarial setting (even defining it is not so easy), or move towards a non-parametric model for the likelihood of success given an allocation. 8

## 2 References

Yasin Abbasi-Yadkori, Csaba Szepesvári, and David Tax. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011. Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvári. Online-to-confidence-set conversions and application to sparse stochastic bandits. In *AISTATS*, volume 22, pages 1–9, 2012. Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. *arXiv preprint arXiv:1209.3352*, 2012. Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *The Journal of Machine Learning Research*, 3:397–422, 2003. Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47:235–256, 2002. Kristin P Bennett and Olvi L Mangasarian. Bilinear separation of two sets in space. *Computational Optimization and Applications*, 2(3):207–227, 1993. Sébastien Bubeck and Nicolò Cesa-Bianchi. *Regret Analysis of Stochastic and Nonstochastic Multiarmed Bandit Problems*. Foundations and Trends in Machine Learning. Now Publishers Incorporated, 2012. ISBN 9781601986269. Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback. In *COLT*, pages 355–366, 2008. Akshay Krishnamurthy, Alekh Agarwal, and Miroslav Dudík. Efficient contextual semi-bandit learning. *arXiv preprint arXiv:1502.05890*, 2015. Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvári. Tight regret bounds for stochastic combinatorial semi-bandits. *arXiv preprint arXiv:1410.0949*, 2014. Tor Lattimore, Koby Crammer, and Csaba Szepesvári. Optimal resource allocation with semi-bandit feedback. In *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence (UAI)*, 2014. Marek Petrik and Shlomo Zilberstein. Robust approximate bilinear programming for value function approximation. *The Journal of Machine Learning Research*, 12:3027–3063, 2011. Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010. Thomas Sowell. *Is Reality Optional?: And Other Essays*. Hoover Institution Press, 1993.