Sequential Learning of the Pareto Front in Multi-objective Bandits



Wouter M. Koolen





ELLIS ILIR Workhop Oberwolfach Tuesday 27th February, 2024



UNPA ENS DE LYON

Élise Crepon



Aurélien Garivier



1. Motivation

2. Setting

- 3. Our Results
- 4. Those Computations
- 5. Conclusion



Almost all optimisation is multi-objective when you think about it.

- Vacation : sunny and tasty
- Drug trial : efficacy and toxicity
- Product dev: cost and sustainability
- . . .



Almost all optimisation is multi-objective when you think about it.

- Vacation : sunny and tasty
- Drug trial : efficacy and toxicity
- Product dev: cost and sustainability
- ...

Today: not in the mood to scalarise

Pareto Front



metric 1

Pareto Front



Pareto front is $\{4, 3, 6, 2\}$.



1. Motivation

2. Setting

- 3. Our Results
- 4. Those Computations
- 5. Conclusion

K-armed multi-objective bandit $\vec{\mu} = (\mu_1, \dots, \mu_K)$.

K-armed multi-objective bandit $ec{\mu}=(\mu_1,\ldots,\mu_{\mathcal{K}}).$

Each arm k represented by a mean vector μ_k in \mathbb{R}^d .

K-armed multi-objective bandit $\vec{\mu} = (\mu_1, \dots, \mu_K)$.

Each arm k represented by a mean vector μ_k in \mathbb{R}^d .

Observations from arm k are i.i.d. multivariate Gaussian $\mathcal{N}(\mu_k, l)$.

K-armed multi-objective bandit $\vec{\mu} = (\mu_1, \dots, \mu_K)$. Each arm k represented by a mean vector μ_k in \mathbb{R}^d . Observations from arm k are i.i.d. multivariate Gaussian $\mathcal{N}(\mu_k, I)$. We assume all μ_k are different. *K*-armed multi-objective bandit $\vec{\mu} = (\mu_1, \dots, \mu_K)$.

Each arm k represented by a mean vector μ_k in \mathbb{R}^d .

Observations from arm k are i.i.d. multivariate Gaussian $\mathcal{N}(\mu_k, I)$. We assume all μ_k are different.

We say arm k dominates arm i, denoted $\mu_k \succeq \mu_i$, if $\mu_k^j \ge \mu_i^j$ in every dimension j = 1, ..., d.

K-armed multi-objective bandit $\vec{\mu} = (\mu_1, \dots, \mu_K)$.

Each arm k represented by a mean vector μ_k in \mathbb{R}^d .

Observations from arm k are i.i.d. multivariate Gaussian $\mathcal{N}(\mu_k, l)$. We assume all μ_k are different.

We say arm k dominates arm i, denoted $\mu_k \succeq \mu_i$, if $\mu_k^j \ge \mu_i^j$ in every dimension j = 1, ..., d.

The **Pareto front** is the set of non-dominated arms:

$$S^*(ec{\mu}) \coloneqq \{k \mid orall i
eq k : \mu_i
eq \mu_k\}$$

We work in the setting of **fixed-confidence** $\delta \in (0, 1)$.

Protocol

For $t = 1, 2, ..., \tau$:

- Learner picks an arm $I_t \in [K]$.
- Learner sees $X_t \sim \mathcal{N}(\boldsymbol{\mu}_{I_t}, I)$

Learner recommends Pareto front $\hat{S} \subseteq [K]$



Learner is $\delta\text{-correct}$ if for any bandit instance $\vec{\mu}$

$$\mathbb{P}_{ec{\mu}}\left\{ au < \infty \land \hat{S}
eq S^*(ec{\mu})
ight\} \leq \delta$$

Goal: minimise sample complexity $\mathbb{E}_{\vec{\mu}}[\tau]$ over all δ -correct strategies.

Define the *alternatives* to $\vec{\mu}$ by

$$\mathsf{Alt}(ec{\mu}) \ \coloneqq \ \left\{ec{\lambda} \in \mathbb{R}^{K imes d} \mid S^*(ec{\lambda})
eq S^*(ec{\mu})
ight\}.$$

NB recall S^* is Pareto front

Define the *alternatives* to $\vec{\mu}$ by

$$\mathsf{Alt}(\vec{\mu}) \ \coloneqq \ \big\{ \vec{\lambda} \in \mathbb{R}^{K \times d} \ \big| \ S^*(\vec{\lambda}) \neq S^*(\vec{\mu}) \big\}.$$

NB recall S^* is Pareto front

Theorem (Garivier and Kaufmann 2016)

Fix a δ -correct strategy. Then for every bandit model $\vec{\mu}$

$$\mathbb{E}_{ec{m \mu}}[au] \ \geq \ {\mathcal T}^*(ec{m \mu}) \ln rac{1}{\delta}$$

where the characteristic time $T^*(\vec{\mu})$ is given by

$$rac{1}{T^*(ec{\mu})} \;=\; \max_{oldsymbol{w}\in riangle_K} \min_{oldsymbol{\lambda}\in riangle riangle$$

Idea is consider the oracle weight map

$$m{w}^*(ec{\mu}) \ \coloneqq \ rgmax_{w\in riangle \kappa} \min_{ec{\lambda}\in \mathsf{Alt}(ec{\mu})} \ rac{1}{2}\sum_{k=1}^K w_k \left\| \mu_k - oldsymbol{\lambda}_k
ight\|^2$$

and track the plug-in estimate: sample arm $l_t \sim w^* \left(\hat{ec{\mu}}(t-1)
ight).$



Idea is consider the oracle weight map

$$oldsymbol{w}^*(ec{\mu}) \ \coloneqq \ rgmax_{w\in riangle \kappa} \min_{oldsymbol{\lambda} \in \mathsf{Alt}(ec{\mu})} \ rac{1}{2} \sum_{k=1}^K w_k \left\| oldsymbol{\mu}_k - oldsymbol{\lambda}_k
ight\|^2$$

and track the plug-in estimate: sample arm $I_t \sim w^* \left(\hat{ec{\mu}}(t-1)
ight).$

Theorem (Degenne and Koolen, 2019)

Take set-valued interpretation of arg max defining w^* . Then $\vec{\mu} \mapsto w^*(\vec{\mu})$ is upper-hemicontinuous and convex-valued. Suitable tracking ensures that as $\hat{\vec{\mu}}(t) \to \vec{\mu}$, any choice $w_t \in w^*(\hat{\vec{\mu}}(t-1))$ have

$$\min_{oldsymbol{w}\inoldsymbol{w}^*(ec{\mu})} \left\|oldsymbol{w}_t - oldsymbol{w}
ight\|_{\infty} o 0$$

Track-and-Stop is asymptotically optimal: $\limsup_{\delta \to 0} \frac{\mathbb{E}_{\vec{\mu}}[\tau]}{\ln \frac{1}{\delta}} = T^*(\vec{\mu}).$





- 1. Motivation
- 2. Setting
- 3. Our Results
- 4. Those Computations
- 5. Conclusion

Kone, Kaufmann, and Richert (2023) consider identifying the Pareto Front among K arms in d dimensions.

- Asymptotically optimal algorithm for Pareto Front Identification.
- Computations in exponential $O(d^K)$ time per round.

Our Contribution

• Computations in polynomial $O(K^d)$ time per round.



- 1. Motivation
- 2. Setting
- 3. Our Results
- 4. Those Computations
- 5. Conclusion

Degenne, Koolen, and Ménard (2019): sufficient to implement best-response oracle (= gradient)

$$ec{\mu}, oldsymbol{w} \mapsto \min_{oldsymbol{\lambda} \in \mathsf{Alt}(ec{\mu})} rac{1}{2} \sum_{k=1}^{K} w_k \left\| oldsymbol{\mu}_k - oldsymbol{\lambda}_k
ight\|^2$$

Degenne, Koolen, and Ménard (2019): sufficient to implement best-response oracle (= gradient)

$$ec{\mu}, w \mapsto \min_{ec{\lambda} \in \mathsf{Alt}(ec{\mu})} \frac{1}{2} \sum_{k=1}^{K} w_k \left\| \mu_k - \lambda_k \right\|^2$$

Objective is convex, but domain $Alt(\vec{\mu})$ is not.

Degenne, Koolen, and Ménard (2019): sufficient to implement best-response oracle (= gradient)

$$ec{\mu}, w \mapsto \min_{ec{\lambda} \in \mathsf{Alt}(ec{\mu})} \frac{1}{2} \sum_{k=1}^{K} w_k \left\| \mu_k - \lambda_k \right\|^2$$

Objective is convex, but domain $Alt(\vec{\mu})$ is not.

Optimal transport problem

Recall

$$ec{\lambda}\in \mathsf{Alt}(ec{\mu})$$
 i.e. $S^*(ec{\lambda})
eq S^*(ec{\mu})$

Having a different Pareto front means either

- An arm on the front in $\vec{\mu}$ is off the front in $\vec{\lambda}$, or
- An arm off the front in $\vec{\mu}$ is on the front in $\vec{\lambda}$.

Taking arm 4 off the Pareto Front



Taking arm 4 off the Pareto Front



Example: we dominate arm 4 using arm 6 by moving each to the **weighted mid-point** in non-dominated coordinates.

Putting arm 1 on the Pareto Front



Putting arm 1 on the Pareto Front



Example: we make point 1 dominant by moving it north-east, and then moving all dominators **out of the way**.

The cost for moving point 1 onto the front is:

$$\min_{\boldsymbol{\lambda}_1} \ \frac{w_1}{2} || \mu_1 - \boldsymbol{\lambda}_1 ||^2 + \sum_{k \in S^*(\vec{\mu})} \frac{w_k}{2} \min_{j \in [d]} (\mu_k^j - \boldsymbol{\lambda}_1^j)_+^2$$

The cost for moving point 1 onto the front is:

$$\min_{\boldsymbol{\lambda}_1} \ \frac{w_1}{2} || \boldsymbol{\mu}_1 - \boldsymbol{\lambda}_1 ||^2 + \sum_{k \in S^*(\vec{\mu})} \frac{w_k}{2} \min_{j \in [d]} (\boldsymbol{\mu}_k^j - \boldsymbol{\lambda}_1^j)_+^2$$

and that is

$$\min_{\phi:S^*(\vec{\mu})\to[d]} \underbrace{\min_{\lambda_1} \frac{w_1}{2} ||\mu_1 - \lambda_1||^2}_{\text{control} loc} + \sum_{k\in S^*(\vec{\mu})} \frac{w_k}{2} (\mu_k^{\phi(k)} - \lambda_1^{\phi(k)})_+^2}_{\text{control} loc}$$

separable convex problem

The cost for moving point 1 onto the front is:

$$\min_{\boldsymbol{\lambda}_1} \ \frac{w_1}{2} || \boldsymbol{\mu}_1 - \boldsymbol{\lambda}_1 ||^2 + \sum_{k \in S^*(\vec{\mu})} \frac{w_k}{2} \min_{j \in [d]} (\boldsymbol{\mu}_k^j - \boldsymbol{\lambda}_1^j)_+^2$$

and that is

$$\min_{\phi:S^*(\vec{\mu})\to[d]} \underbrace{\min_{\lambda_1} \frac{w_1}{2} ||\mu_1 - \lambda_1||^2 + \sum_{k\in S^*(\vec{\mu})} \frac{w_k}{2} (\mu_k^{\phi(k)} - \lambda_1^{\phi(k)})_+^2}_{\text{separable convex problem}}$$

Not all $\phi: S^*(\vec{\mu}) \to [d]$ need to be attempted. Only $\binom{K+d-1}{d-1}$ due to geometry of \mathbb{R}^d .



- 1. Motivation
- 2. Setting
- 3. Our Results
- 4. Those Computations
- 5. Conclusion

With that, everything slots in place and we obtain an algorithm for Pareto Front Identification with

- asymptotically optimal sample complexity
- polynomial time cost per round

Now interested in going beyond

- Gaussian
- $\epsilon = 0$
- independence

Thanks!

С Degenne, R. and W. M. Koolen (Dec. 2019). "Pure Exploration with Multiple Correct Answers". In: Advances in Neural Information Processing Systems (NeurIPS) 32. C C Degenne, R., W. M. Koolen, and P. Ménard (Dec. 2019). "Non-Asymptotic Pure Exploration by Solving Games". In: Advances in Neural Information Processing Systems (NeurIPS) 32. C Garivier, A. and E. Kaufmann (June 2016). "Optimal Best Arm Identification with Fixed Confidence". In: 29th Annual Conference on Learning Theory. Vol. 49. Proceedings of Machine Learning Research. C Kone, C., E. Kaufmann, and L. Richert (2023). "Adaptive Algorithms for Relaxed Pareto Set Identification". In: arXiv preprint arXiv:2307.00424.