

MOTIVATION

- Product recommendation (e.g., movies)
- Adaptive hypothesis testing under linear assumption
- Optimization of a stochastic linear function

SETTING

The linear stochastic bandit model

- Set of arms $\mathcal{X} \subseteq \mathbb{R}^d$, $|\mathcal{X}| = K$, $\|x\|_2 \leq L$, $\forall x \in \mathcal{X}$.
- Linear reward model

$$r(x) = x^\top \theta^* + \varepsilon$$

- with $\theta^* \in \mathbb{R}^d$ **unknown** parameter and noise $\varepsilon \sim \mathcal{N}(0, \sigma^2)$.
- The (unique) best arm in \mathcal{X} :

$$x^* = \arg \max_{x \in \mathcal{X}} x^\top \theta^*$$

The best-arm identification problem ((0, δ)-PAC setting)

- $\hat{x}(n)$ – recommended best arm after n steps.
- Given a *fixed confidence* δ , design an *allocation strategy* and a *stopping criterion* such that:

$$\mathbb{P}(\hat{x}(n) = x^*) \geq 1 - \delta \text{ and } n \text{ as small as possible.}$$

TOOLS

Ordinary Least-Squares estimate

- Sequence of arms $\mathbf{x}_n = (x_1, \dots, x_n) \in \mathcal{X}^n$
- Sequence of rewards (r_1, \dots, r_n)
- OLS estimate, $A_{\mathbf{x}_n} = \sum_{t=1}^n x_t x_t^\top$, $b_{\mathbf{x}_n} = \sum_{t=1}^n x_t r_t$

$$\hat{\theta}_n = A_{\mathbf{x}_n}^{-1} b_{\mathbf{x}_n}$$

Prediction errors

- Fixed sequence and OLS estimate (w.p. $1 - \delta$)

$$\|x^\top \theta^* - x^\top \hat{\theta}_n\| \leq c \|x\|_{A_{\mathbf{x}_n}^{-1}} \sqrt{\log_n(K/\delta)}$$

- Adaptive sequence (Thm.2 in [1]) for η -regularized OLS (w.p. $1 - \delta$)

$$\|x^\top (\theta^* - \hat{\theta}_n^*)\| \leq \|x\|_{A_{\mathbf{x}_n}^{-1}} \left(\sigma \sqrt{d \log \left(\frac{1+nL^2/\eta}{\delta} \right)} + \eta^{1/2} \|\theta^*\| \right)$$

REFERENCES

- [1] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances Neural Information Processing Systems 24*, 2011.
- [2] E. Even-Dar, S. Mannor, and Y. Mansour. PAC bounds for multi-armed bandit and markov decision processed. In *Computational Learning Theory*, 2002.
- [3] V. Gabillon, M. Ghavamzadeh, and A. Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In *Advances Neural Information Processing Systems 25*, 2012.

ACKNOWLEDGEMENTS

GAPS AND DIRECTIONS

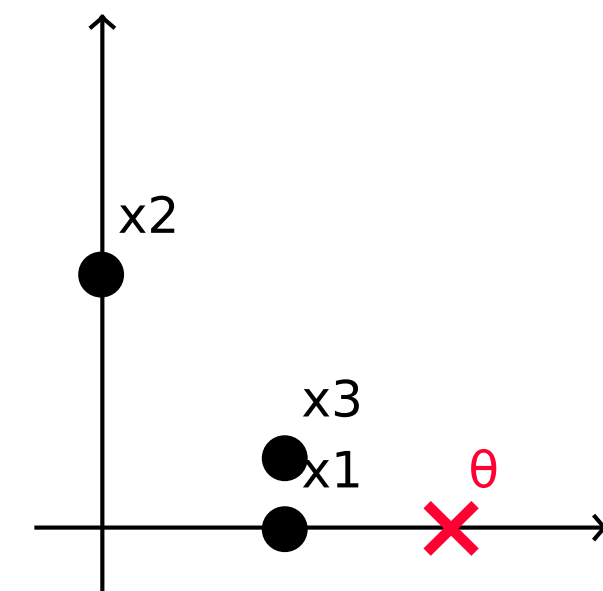
Value gaps

- For any pair $x, x' \in \mathcal{X}$, $\Delta(x, x') = (x - x')^\top \theta^*$
- Smallest gap $\Delta_{\min} = \min_{x \in \mathcal{X}} \Delta(x^*, x)$

The smaller the gaps the more difficult the problem.

Example

- x_1 and x_3 are very close, while x_2 is clearly suboptimal
- Only direction $y = x_1 - x_3$ (i.e., θ_2^*) must be estimated accurately
- x_2 provides *much* information about direction y (θ_2^*)

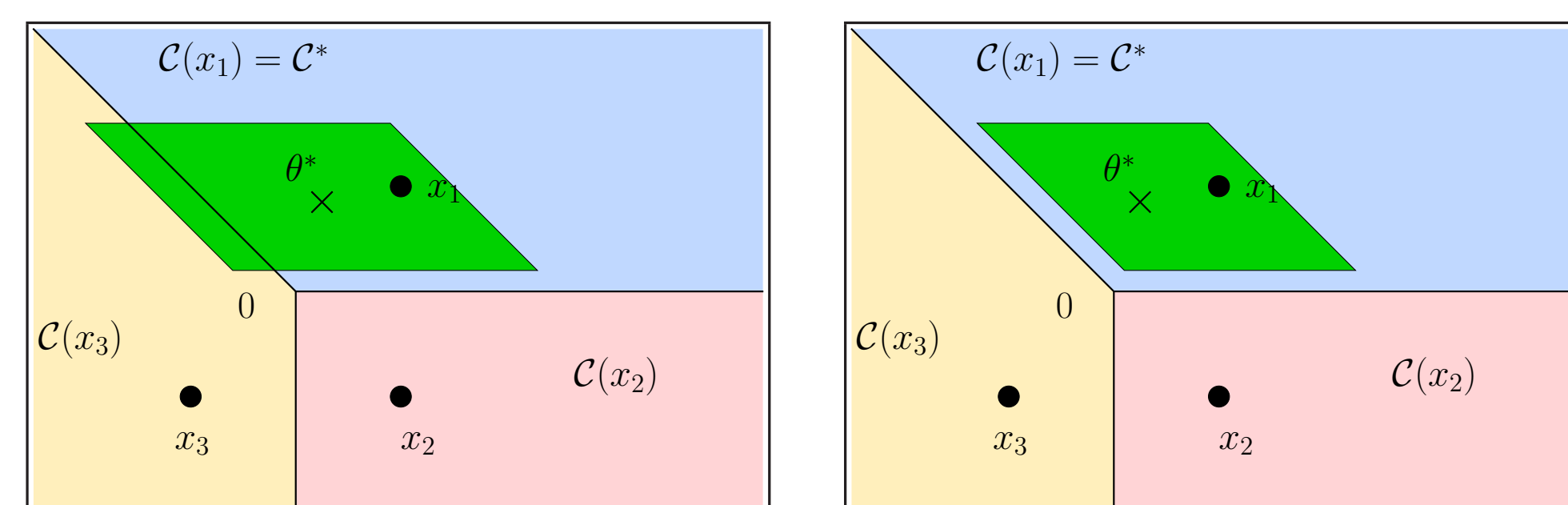


Use one arm to learn about the others (exploit the linear structure)!

Sets of directions

$$\mathcal{Y} = \{y = x - x'\}, \quad \mathcal{Y}^* = \{y = x^* - x\}$$

ILLUSTRATION



Optimality cones

$$\mathcal{C}(x) = \cap_{x' \in \mathcal{X}} \{\theta \in \mathbb{R}^d, (x - x')^\top \theta \geq 0\}$$

Confidence set

$$\mathcal{S}(\mathbf{x}_n) = \left\{ \theta, \forall y, y^\top (\theta^* - \theta) \leq c \|y\|_{A_{\mathbf{x}_n}^{-1}} \sqrt{\log(K^2/\delta)} \right\}$$

\mathcal{XY} -ORACLE

Intuition: select arms so that the confidence set shrinks into one optimality cone as soon as possible.

Stopping rule

$$\exists x \in \mathcal{X} \text{ s.t. } \mathcal{S}(\mathbf{x}_n) \subset \mathcal{C}(x)$$

Allocation rule

$$x_n^* = \arg \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}^*} \frac{\|y\|_{A_{x_n}^{-1}}}{\Delta(y)}$$

Oracle sample complexity

$$N^* = c^2 H_{\text{LB}} \log_n(K^2/\delta)$$

Complexity of Linear Best-Arm Identification

$$H_{\text{LB}} = \min_{\lambda \in \mathcal{D}^k} \max_{y \in \mathcal{Y}^*} \frac{\|y\|_{\Lambda_\lambda^{-1}}^2}{\Delta^2(y)}$$

N^ is the lower-bound on the sample complexity of any fixed allocation strategy*

Remarks

$$\max_{y \in \mathcal{Y}^*} \frac{\|y\|_{\Lambda_{\min}^{-1}}^2}{L \Delta_{\min}^2} \leq H_{\text{LB}} \leq \frac{4d}{\Delta_{\min}^2} \quad H_{\text{MAB}} \leq H_{\text{LB}} \leq 2H_{\text{MAB}}$$

STATIC ALLOCATIONS

G-Optimal Design: estimate θ^* uniformly well over all arms

$$x_n^G = \arg \min_{x_n} \max_{x \in \mathcal{X}} \|x\|_{A_{x_n}^{-1}}$$

\mathcal{XY} -Design: estimate the value of the gaps uniformly well over all the directions in \mathcal{Y}

$$x_n^{\mathcal{XY}} = \arg \min_{x_n} \max_{y \in \mathcal{Y}} \|y\|_{A_{x_n}^{-1}}$$

Empirical stopping criterion:

$$\exists x \in \mathcal{X}, \forall x' \in \mathcal{X}, \forall \theta \in \hat{S}(\mathbf{x}_n) \\ (x - x')^\top (\hat{\theta}_n - \theta) \leq \hat{\Delta}_n(x, x')$$

Sample complexity $O(d/\Delta_{\min}^2)$.

FROM \mathcal{Y} TO \mathcal{Y}^*

The minimum number of steps needed by the static \mathcal{XY} -allocation to discard all *suboptimal* directions.

$$M^* = \min\{n \in \mathbb{N}, \forall x \neq x^*, \forall x' \neq x^*,$$

$$S^*(x_n^{\mathcal{XY}}) \cap (\mathcal{C}(x) \cap \mathcal{C}(x')) = \emptyset\}.$$

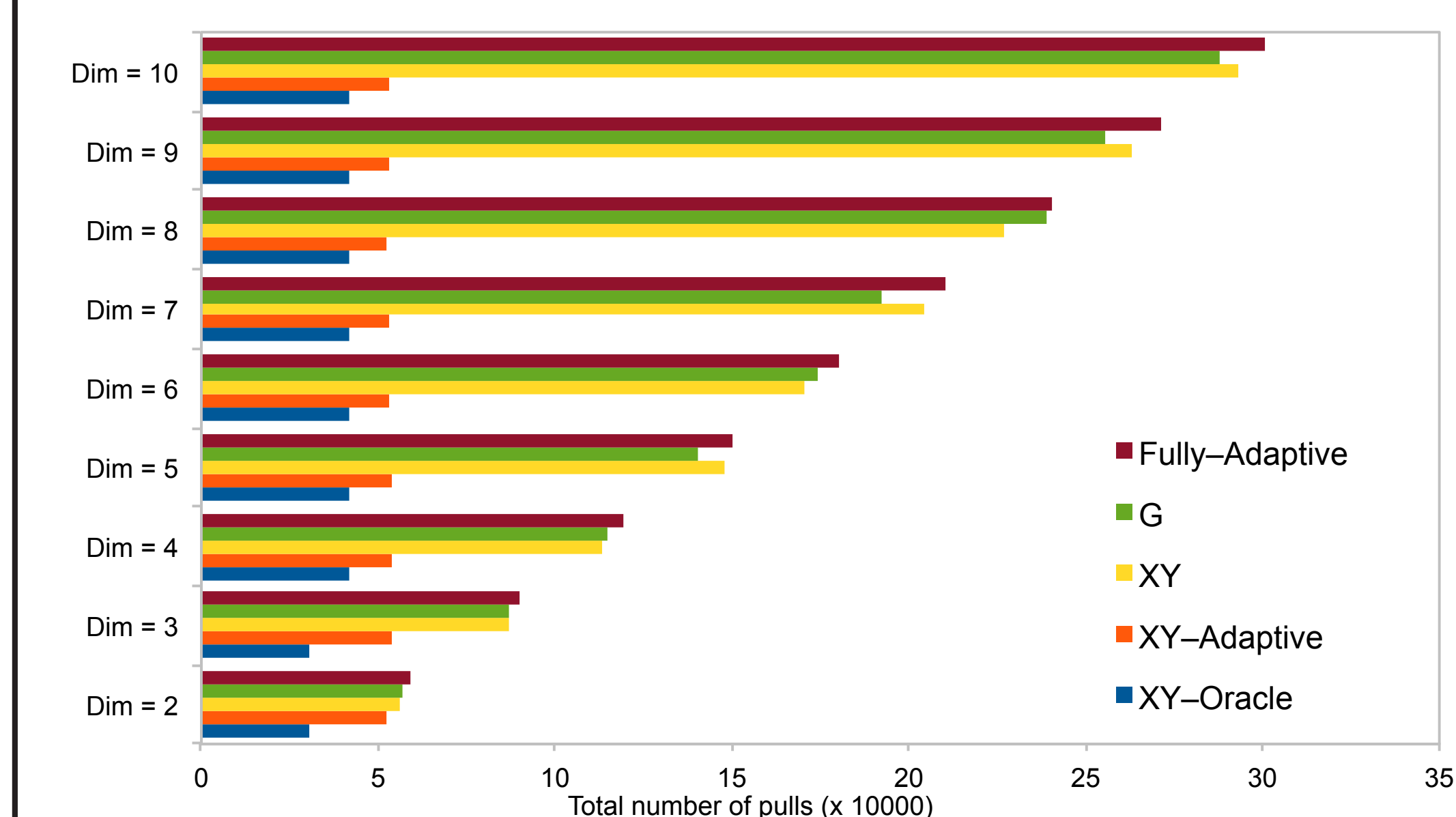
MAIN RESULT

Theorem 1. If the \mathcal{XY} -adaptive allocation strategy is implemented with a β -approximate method then

$$\mathbb{P} \left[N \leq \frac{(1 + \beta) \max\{M^*, \frac{16}{\alpha} N^*\}}{\log(1/\alpha)} \log \left(\frac{c \sqrt{\log_n(K^2/\delta)}}{\Delta_{\min}} \right) \wedge (\hat{x}_N = x^*) \right] \geq 1 - \delta.$$

The bound holds for any $(1 + \beta)$ -approximate allocation strategy: e.g., continuous relaxation, greedy incremental allocation.

EXPERIMENTS – SAMPLE COMPLEXITY AND ARM PULLS



Arm	\mathcal{XY} -oracle	\mathcal{XY} -adapt.	\mathcal{XY}	G	Fully-adapt.
x_1	207	263	29523	28014	740
x_2	41440	52713	29524	28015	149220
x_3	2	3	29524	28015	1
x_4	2	5	29524	28015	1
x_5	1	2	29524	28015	1
x_6	0	2	1	1	1
Total	41652	52988	147620	140075	149964

\mathcal{XY} -ADAPTIVE ALGORITHM

Input: $\mathcal{X} \in \mathbb{R}^d$; confidence δ ;
Phase length given by an α improvement;
Set $j = 1$; $\hat{\mathcal{X}}_j = \mathcal{X}$; $\hat{\mathcal{Y}}_1 = \mathcal{Y}$; $n = 0$;

STOPPING RULE

while $|\hat{\mathcal{X}}_j| > 1$ **do**
Start a new phase: $j = j + 1$, $t = 1$; $A_0 = I$
while $\rho^j/t \geq \alpha \rho^{j-1} (\mathbf{x}_{n_{j-1}}^j) / n_{j-1}$ **do**

ALLOCATION RULE

$$x_t = \arg \min_{x \in \mathcal{X}} \max_{y \in \hat{\mathcal{Y}}_j} y^\top (A + x x^\top)^{-1} y$$

Update $A_t = A_t + x_t x_t^\top$; $t = t + 1$; $n = n + 1$
 $\rho^j = \max_{y \in \hat{\mathcal{Y}}_j} y^\top A_t^{-1} y$

end while

$$b = \sum_{s=1}^t x_s r_s; \hat{\theta}_j = A_t^{-1} b$$

Recompute the set of potential optimal arms:

$$\hat{\mathcal{X}}_j = \{x, \exists x' : \|x - x'\|_{A_t^{-1}} \sqrt{\log_n(K^2/\delta)} \leq \hat{\Delta}_j(x', x)\}$$

Recompute the set of directions of interest:

$$\hat{\mathcal{Y}}_j = \{y = (x - x'); x, x' \in \hat{\mathcal{X}}_j\}$$

end while

RECOMMENDATION RULE

Return $\hat{x}(n)$ – the only arm remaining in $\hat{\mathcal{X}}_j$.

Setting:

- Fixed confidence $\delta = 0.05$.
- Set of arms: $\mathcal{X} \in \mathbb{R}^d$, $|\mathcal{X}| = d + 1$ and $d = 2, \dots, 10$.
- Canonical basis (x_1, \dots, x_d) and additional arm x_{d+1} very close to x_1 .
- $\theta^* = [2 \ 0 \ 0 \ \dots \ 0]^\top \rightarrow \Delta_{\min} = (x_1 - x_{d+1})^\top \theta^*$ much smaller than the other gaps.
- Identifying the best arm \rightarrow reducing uncertainty in the direction $\tilde{y} = (x_1 - x_{d+1})$.
- x_2 is almost aligned with $\tilde{y} \rightarrow$ **the most informative arm**.

The sample complexity **grows linearly with the dimension**:

- Fully-adaptive – despite pulling only the informative arms, the additional d term in the bound prevents a good performance.
- G and \mathcal{XY} – always consider the complete set \mathcal{Y} .

The sample complexity remains **constant**:

- \mathcal{XY} -Adaptive and \mathcal{XY} -Oracle – exclusively pull the two most informative arms, independently of the number of dimensions.