

STRUCTURED MULTI-ARMED BANDITS RLSS

July 02, Lille

Odalric-Ambrym Maillard

INRIA LILLE – NORD EUROPE

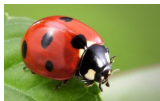
...SequeL...

Eco-sustainable decision making

▶ Plant health-care:



: $\mathcal{A} = \{$

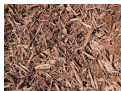


$\}$

▶ Ground health-care:



: $\mathcal{A} = \{$



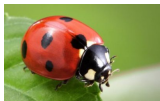
$\}$

Eco-sustainable decision making

▶ Plant health-care:



: $\mathcal{A} = \{$

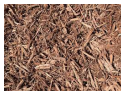


$\}$

▶ Ground health-care:



: $\mathcal{A} = \{$



$\}$

Medical decision companion

▶ Emergency admission filtering:



: $\mathcal{A} = \{$



,



,



,



$\}$



- Suggest medical consultation or treatment based on smart meters.



- ▶ Suggest medical consultation or treatment based on smart meters.
- ▶ Time series, hidden variables, risk-aversion.



- ▶ Recommend drug dosage w.r.t. genome of individuals.



- ▶ Recommend drug dosage w.r.t. genome of individuals.
- ▶ Huge dimension, Gene interactions.



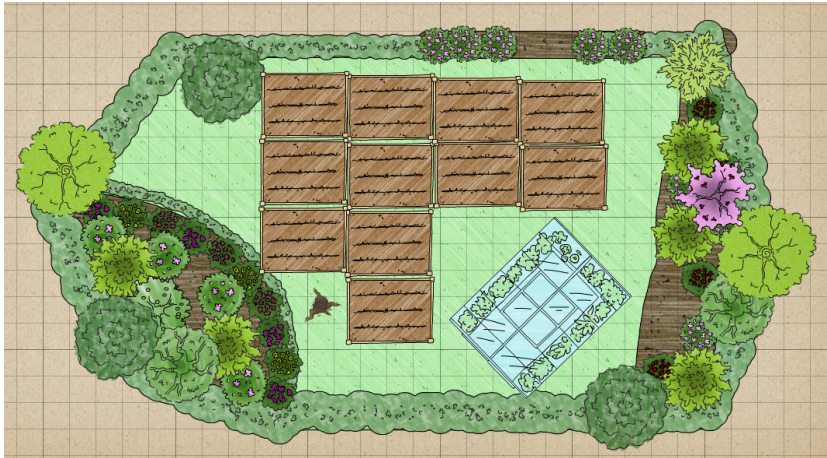
Massive Open Online Course

- ▶ Recommend exercises that maximize learning progression

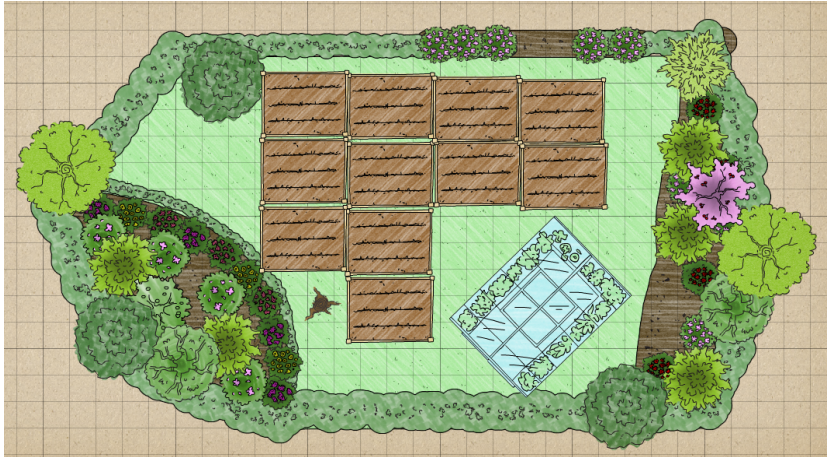


Massive Open Online Course

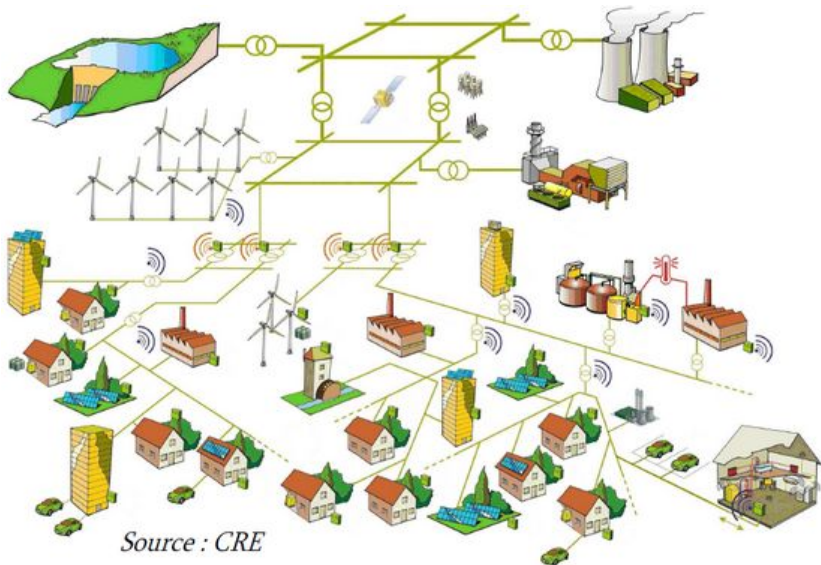
- ▶ Recommend exercises that maximize learning progression
- ▶ Non-stationary rewards, few interactions



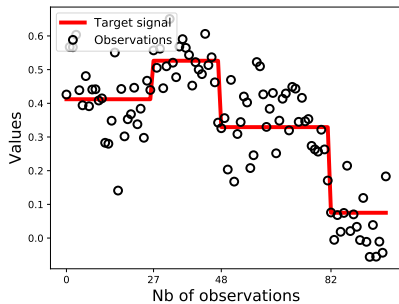
- ▶ Recommend good practice between farms/share knowledge.



- ▶ Recommend good practice between farms/share knowledge.
- ▶ Strong correlations, hidden variables, delayed feedback.



- ▶ Distributed Optimization, Cognitive Radio Networks, etc.



- Time Series, HMMs, Autoregressive models, etc.

STRUCTURES

LINEAR BANDITS

STRUCTURED LOWER BOUNDS

CONCLUSION, PERSPECTIVE



[Camera Calibration Toolbox for Matlab](#)

This is a release of a **Camera Calibration Toolbox** for Matlab® with a complete documentation. This document may also be used as a tutorial on **camera** ...
[www.vision.caltech.edu/bouguetj/calib_doc/](#) - 14k - [Cached](#)



[Omnivis 2003: Omnidirectional Vision and Camera Networks](#)

A complete paper, not longer than six (6) pages including figures and references, should be submitted in **camera-ready** IEEE 2-column format of single-spaced ...
[www.cs.wustl.edu/~pless/omnivis2003/](#) - 5k - [Cached](#)

[Camera Calibration Toolbox for Matlab](#)

A **Camera Calibration Toolbox** from the Institute of Robotics and Mechatronics, Germany - DLR CalDe and DLR CalLab is a very complete tool for **camera** ...
[www.vision.caltech.edu/bouguetj/calib_doc/htmls/links.html](#) - 16k - [Cached](#)

[The Page of Omnidirectional Vision](#)

ICCV 2005 Omnisixth Workshop on Omnidirectional Vision, **Camera** ... Automatic Surveillance Using Omnidirectional and Active **Cameras** at the PRIP Lab, ...
[www.cis.upenn.edu/~kostas/omni.html](#) - 39k - [Cached](#)

[Digital Camera Characteristics](#)

It is necessary to know your **camera** characteristics if you intend to make full use of the functions available on your **camera** ...
[www.ncsu.edu/science/junction/route/usetech/digitalcamera/](#) - 10k - [Cached](#)

[A Comparison of PMD-Cameras and Stereo-Vision for the Task of...](#)

File Format: PDF/Adobe Acrobat - [View as HTML](#)
systems and **PMD cameras** is discussed qualitatively and ... the stereo system as well as the **PMD camera** will be compared in section 4 based on those ...
[vision.middlebury.edu/conferences/bencos2007/pdf/beder.pdf](#)



Google Custom Search Search

UCSD Computer Vision Web Search

[Camera Calibration Toolbox for Matlab](#)
This is a release of a **Camera Calibration Toolbox** for Matlab® with a complete documentation. This document may also be used as a tutorial on **camera** ...
[www.vision.caltech.edu/bouguetj/calib_doc/](#) - 14k - [Cached](#)

[Omnivis 2003: Omnidirectional Vision and Camera Networks](#)
A complete paper, not longer than six (6) pages including figures and references, should be submitted in **camera-ready** IEEE 2-column format of single-spaced ...
[www.cs.wustl.edu/~pless/omnivis2003/](#) - 5k - [Cached](#)

[Camera Calibration Toolbox for Matlab](#)
A **Camera Calibration Toolbox** from the Institute of Robotics and Mechatronics, Germany - DLR CalDe and DLR CalLab is a very complete tool for **camera** ...
[www.vision.caltech.edu/bouguetj/calib_doc/htmls/links.html](#) - 16k - [Cached](#)

[The Page of Omnidirectional Vision](#)
ICCV 2005 Omnisixth Workshop on Omnidirectional Vision, **Camera** ... Automatic Surveillance Using Omnidirectional and Active **Cameras** at the PRIP Lab, ...
[www.cis.upenn.edu/~kostas/omni.html](#) - 39k - [Cached](#)

[Digital Camera Characteristics](#)
It is necessary to know your **camera** characteristics if you intend to make full use of all of the functions available on your **camera** ...
[www.ncsu.edu/science/junction/route/usetech/digitalcamera/](#) - 10k - [Cached](#)

[A Comparison of PMD-Cameras and Stereo-Vision for the Task of...](#)
File Format: PDF/Adobe Acrobat - [View as HTML](#)
systems and PMD **cameras** is discussed qualitatively and ... the stereo system as well as the PMD **camera** will be compared in section 4 based on those ...
[vision.middlebury.edu/conferences/bencos2007/pdf/beder.pdf](#)

► Actions: List of items.



[Camera Calibration Toolbox for Matlab](#)

This is a release of a **Camera Calibration Toolbox** for Matlab® with a complete documentation. This document may also be used as a tutorial on **camera** ...
[www.vision.caltech.edu/bouguetj/calib_doc/](#) - 14k - [Cached](#)



[Omnivis 2003: Omnidirectional Vision and Camera Networks](#)

A complete paper, not longer than six (6) pages including figures and references, should be submitted in **camera-ready** IEEE 2-column format of single-spaced ...
[www.cs.wustl.edu/~pless/omnivis2003/](#) - 5k - [Cached](#)

[Camera Calibration Toolbox for Matlab](#)

A **Camera Calibration Toolbox** from the Institute of Robotics and Mechatronics, Germany - DLR CalDe and DLR CalLab is a very complete tool for **camera** ...
[www.vision.caltech.edu/bouguetj/calib_doc/htmls/links.html](#) - 16k - [Cached](#)

[The Page of Omnidirectional Vision](#)

ICCV 2005 Omnisixth Workshop on Omnidirectional Vision, **Camera** ... Automatic Surveillance Using Omnidirectional and Active **Cameras** at the PRIP Lab, ...
[www.cis.upenn.edu/~kostas/omni.html](#) - 39k - [Cached](#)

[Digital Camera Characteristics](#)

It is necessary to know your **camera** characteristics if you intend to make full use of all of the functions available on your **camera** ...
[www.ncsu.edu/science/junction/route/usetech/digitalcamera/](#) - 10k - [Cached](#)

[A Comparison of PMD-Cameras and Stereo-Vision for the Task of...](#)

File Format: PDF/Adobe Acrobat - [View as HTML](#)
systems and PMD **cameras** is discussed qualitatively and ... the stereo system as well as the PMD **camera** will be compared in section 4 based on those ...
[vision.middlebury.edu/conferences/bencos2007/pdf/beder.pdf](#)

- ▶ Actions: List of items.
- ▶ Reward/loss: Ranking of preferred item.



[Camera Calibration Toolbox for Matlab](#)

This is a release of a **Camera Calibration Toolbox** for Matlab® with a complete documentation. This document may also be used as a tutorial on **camera** ...
[www.vision.caltech.edu/bouguetj/calib_doc/ - 14k - Cached](#)



[Omnivis 2003: Omnidirectional Vision and Camera Networks](#)

A complete paper, not longer than six (6) pages including figures and references, should be submitted in **camera-ready** IEEE 2-column format of single-spaced ...
[www.cs.wustl.edu/~pless/omnivis2003/ - 5k - Cached](#)

[Camera Calibration Toolbox for Matlab](#)

A **Camera Calibration Toolbox** from the Institute of Robotics and Mechatronics, Germany - DLR CalDe and DLR CalLab is a very complete tool for **camera** ...
[www.vision.caltech.edu/bouguetj/calib_doc/htmls/links.html - 16k - Cached](#)

[The Page of Omnidirectional Vision](#)

ICCV 2005 Omnisixth Workshop on Omnidirectional Vision, **Camera** ... Automatic Surveillance Using Omnidirectional and Active **Cameras** at the PRIP Lab, ...
[www.cis.upenn.edu/~kostas/omni.html - 39k - Cached](#)

[Digital Camera Characteristics](#)

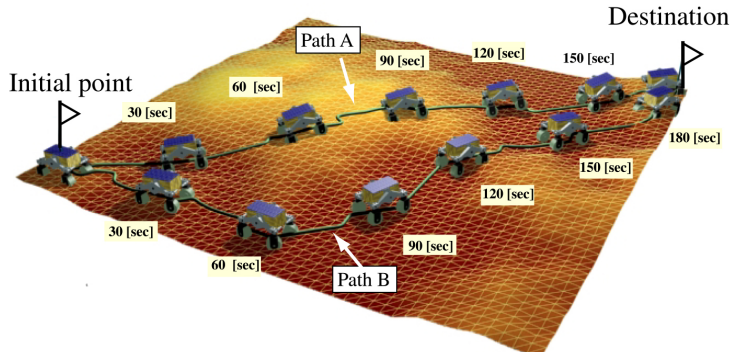
It is necessary to know your **camera** characteristics if you intend to make full use of all of the functions available on your **camera** ...
[www.ncsu.edu/science/junction/route/usetech/digitalcamera/ - 10k - Cached](#)

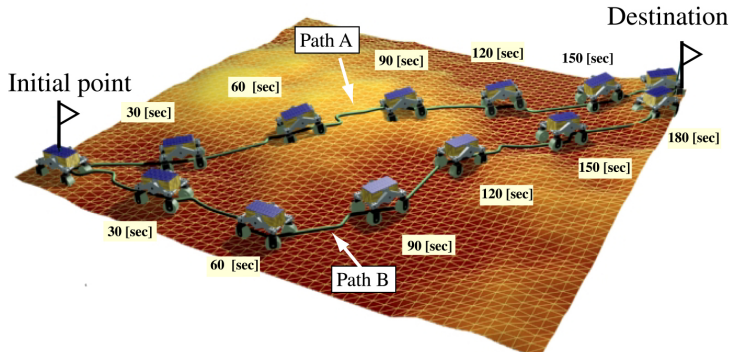
[A Comparison of PMD-Cameras and Stereo-Vision for the Task of...](#)

File Format: PDF/Adobe Acrobat - [View as HTML](#).
systems and PMD **cameras** is discussed qualitatively and ... the stereo system as well as the PMD **camera** will be compared in section 4 based on those ...
[vision.middlebury.edu/conferences/bencos2007/pdf/beder.pdf](#)

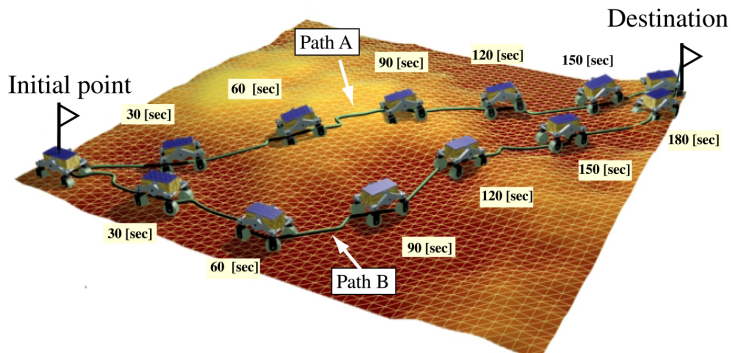
- ▶ Actions: List of items.
- ▶ Reward/loss: Ranking of preferred item.
- ▶ Ordering

STRUCTURE : PATHS

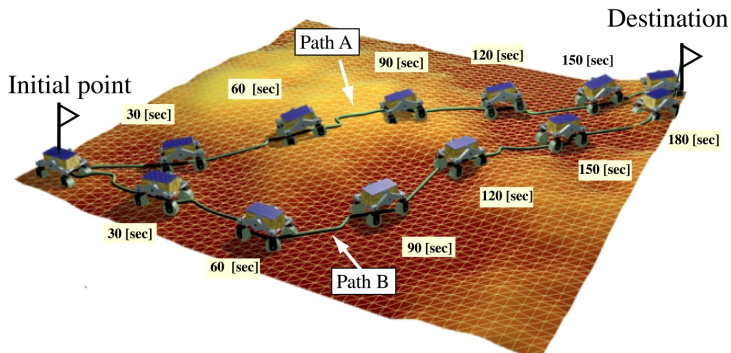




- ▶ Actions: (valued) Paths.

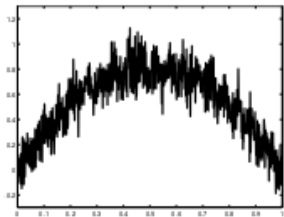
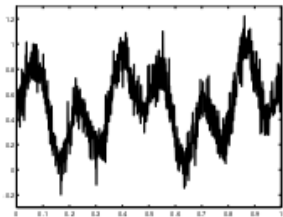
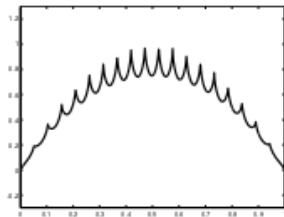
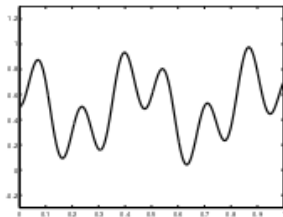


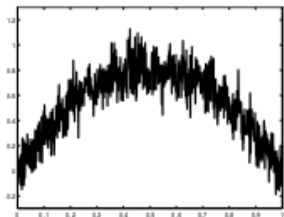
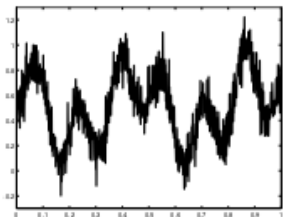
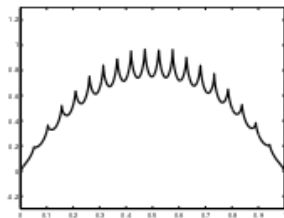
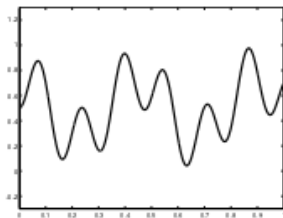
- ▶ Actions: (valued) Paths.
- ▶ Reward/loss: cumulative value on the path.



- ▶ Actions: (valued) Paths.
- ▶ Reward/loss: cumulative value on the path.
- ▶ Paths have edges in common.

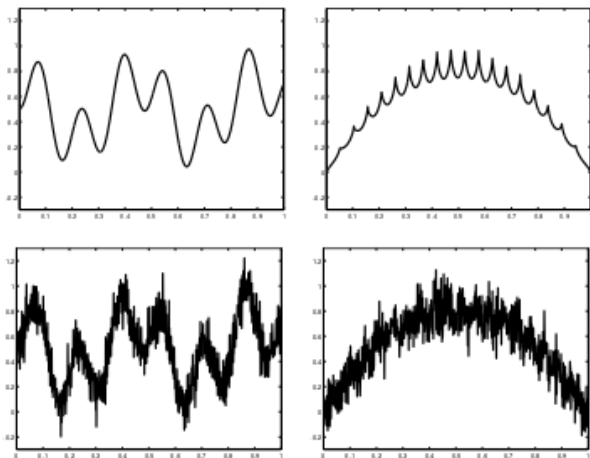
STRUCTURE: SMOOTH REWARDS





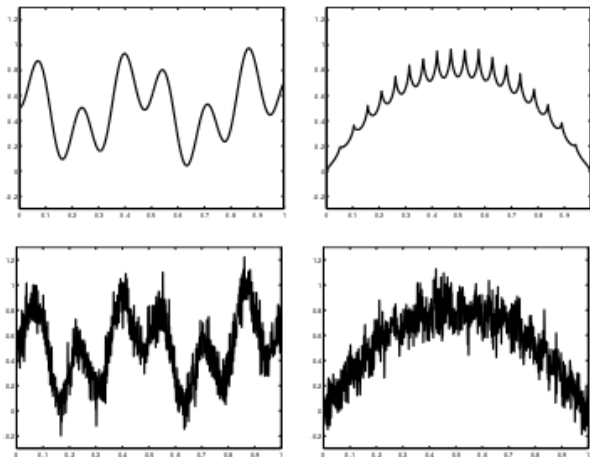
► Actions: $x \in \mathbb{R}$

STRUCTURE: SMOOTH REWARDS



- ▶ Actions: $x \in \mathbb{R}$
- ▶ Reward/loss: $f(x) + \xi$

STRUCTURE: SMOOTH REWARDS



- ▶ Actions: $x \in \mathbb{R}$
- ▶ Reward/loss: $f(x) + \xi$
- ▶ Regularity.

STRUCTURES

LINEAR BANDITS

STRUCTURED LOWER BOUNDS

CONCLUSION, PERSPECTIVE

STRUCTURES

LINEAR BANDITS

Regression

Linear UCB, Linear TS

Graph-linear Bandits

Extension to Kernels

STRUCTURED LOWER BOUNDS

CONCLUSION, PERSPECTIVE

Sequential optimization game

At each time $t \in \mathbb{N}$, sample at $x_t \in \mathcal{X}$, receive $y_t \in \mathbb{R}$, where

$$y_t = \underbrace{f_\star(x_t)}_{\text{target}} + \underbrace{\xi_t}_{\text{noise}}.$$

Goal: Minimize cumulative regret

$$\mathcal{R}_T \stackrel{\text{def}}{=} \sum_{t=1}^T f_\star(\star) - f_\star(x_t) \text{ where } \star \in \text{Argmax } f_\star(x).$$

Sequential optimization game

At each time $t \in \mathbb{N}$, sample at $x_t \in \mathcal{X}$, receive $y_t \in \mathbb{R}$, where

$$y_t = \underbrace{f_\star(x_t)}_{\text{target}} + \underbrace{\xi_t}_{\text{noise}}.$$

Goal: Minimize cumulative regret

$$\mathcal{R}_T \stackrel{\text{def}}{=} \sum_{t=1}^T f_\star(\star) - f_\star(x_t) \text{ where } \star \in \text{Argmax } f_\star(x).$$

► Actions : $x \in \mathcal{X}$.

Sequential optimization game

At each time $t \in \mathbb{N}$, sample at $x_t \in \mathcal{X}$, receive $y_t \in \mathbb{R}$, where

$$y_t = \underbrace{f_\star(x_t)}_{\text{target}} + \underbrace{\xi_t}_{\text{noise}}.$$

Goal: Minimize cumulative regret

$$\mathcal{R}_T \stackrel{\text{def}}{=} \sum_{t=1}^T f_\star(\star) - f_\star(x_t) \text{ where } \star \in \text{Argmax } f_\star(x).$$

- ▶ Actions : $x \in \mathcal{X}$.
- ▶ Means : $f_\star(x)$. Mean at x and x' not arbitrarily different !

- ▶ Set of arms \mathcal{X}

- ▶ Set of arms \mathcal{X}
- ▶ At time t , pick $X_t \in \mathcal{X}$, receive

$$Y_t = f_*(X_t) + \xi_t$$

where ξ_t is centered and further conditionally sub-Gaussian.

f_* belongs to a linear function space:

$$\mathcal{F}_\Theta = \left\{ f_\theta : x \mapsto \theta^\top \varphi(x), \theta \in \Theta \right\} \text{ where } \Theta \in \mathbb{R}^d, \varphi : \mathcal{X} \rightarrow \mathbb{R}^d.$$

θ : Parameter, φ : Feature function.

- ▶ Set of arms \mathcal{X}
- ▶ At time t , pick $X_t \in \mathcal{X}$, receive

$$Y_t = f_\star(X_t) + \xi_t$$

where ξ_t is centered and further conditionally sub-Gaussian.

f_\star belongs to a linear function space:

$$\mathcal{F}_\Theta = \left\{ f_\theta : x \mapsto \theta^\top \varphi(x), \theta \in \Theta \right\} \text{ where } \Theta \in \mathbb{R}^d, \varphi : \mathcal{X} \rightarrow \mathbb{R}^d.$$

θ : Parameter, φ : Feature function.

- ▶ Unknown parameter $\theta_\star \in \mathbb{R}^d$.

- ▶ Set of arms \mathcal{X}
- ▶ At time t , pick $X_t \in \mathcal{X}$, receive

$$Y_t = f_\star(X_t) + \xi_t$$

where ξ_t is centered and further conditionally sub-Gaussian.

f_\star belongs to a linear function space:

$$\mathcal{F}_\Theta = \left\{ f_\theta : x \mapsto \theta^\top \varphi(x), \theta \in \Theta \right\} \text{ where } \Theta \in \mathbb{R}^d, \varphi : \mathcal{X} \rightarrow \mathbb{R}^d.$$

θ : Parameter, φ : Feature function.

- ▶ Unknown parameter $\theta_\star \in \mathbb{R}^d$.
- ▶ Best arm $x_\star = \operatorname{argmax}_{x \in \mathcal{X}} \langle \theta_\star, \varphi(x) \rangle$

- ▶ *Polynomials*: $\mathcal{X} = \mathbb{R}$, $\varphi(x) = (1, x, x^2, \dots, x^{d-1})$, $\Theta = \mathcal{B}_{2,d}(0, 1)$ unit Euclidean ball of \mathbb{R}^d .
- ▶ *Bandits*: $\mathcal{X} = \mathcal{A} = \{1, \dots, A\}$, $\varphi(a) = e_a \in \mathbb{R}^A$, $\Theta = [0, 1]^A$.
- ▶ *Shortest path*: $\mathcal{X} \subset \mathcal{A}^L$ (paths of length L), $\varphi_{(a,\ell)}(x) = \mathbb{I}\{x_\ell = a\}$, $\Theta = [0, 1]^{|\mathcal{X}|}$.
 $\mathcal{X} \subset \{0, 1\}^d$, paths in graph with d edges, $\varphi(x) = x$, $\Theta \subset [0, 1]^d$ mean travel time for each edge (Combes et al. 2015).
- ▶ *Contextual bandits*: $\mathcal{X} = \mathcal{C} \times \mathcal{A}$, $\varphi((c, a)) = (1, c, a, ca, \dots)$
- ▶ *Smooth function on graph*: $\mathcal{X} =$ nodes of a graph with adjacency matrix G , $\varphi =$ eigenfunctions of the Graph-Laplacian.

► *Linear space*: $\mathcal{F} = \left\{ f_{\theta} : f_{\theta}(x) = \langle \theta, \varphi(x) \rangle, \theta \in \mathbb{R}^d, \theta \in \Theta \right\}$.

Ex: $\varphi(x) = (1, x, x^2)$, $f_{\theta}(x) = 2 + \frac{1}{2}x - 2x^2$, $\theta = (2, 1/2, -2)$.

► *Linear space*: $\mathcal{F} = \left\{ f_{\theta} : f_{\theta}(x) = \langle \theta, \varphi(x) \rangle, \theta \in \mathbb{R}^d, \theta \in \Theta \right\}$.

Ex: $\varphi(x) = (1, x, x^2)$, $f_{\theta}(x) = 2 + \frac{1}{2}x - 2x^2$, $\theta = (2, 1/2, -2)$.

► *Loss* : $\ell(y, y') = \frac{(y - y')^2}{2}$

- ▶ **Linear space**: $\mathcal{F} = \{f_\theta : f_\theta(x) = \langle \theta, \varphi(x) \rangle, \theta \in \mathbb{R}^d, \theta \in \Theta\}$.
Ex: $\varphi(x) = (1, x, x^2)$, $f_\theta(x) = 2 + \frac{1}{2}x - 2x^2$, $\theta = (2, 1/2, -2)$.
- ▶ **Loss**: $\ell(y, y') = \frac{(y - y')^2}{2}$
- ▶ **Objective**: from $(x_n, y_n)_{n \leq N}$ optimize

$$\min_{\theta \in \Theta} \sum_{n=1}^N \ell(y_n, f_\theta(x_n)).$$

$$\min_{\theta \in \Theta} \sum_{n=1}^N \left(y_n - \theta^\top \varphi(x_n) \right)^2. \quad (1)$$

- Any solution to (1) must satisfy

$$G_N \theta = \sum_{n=1}^N \varphi(x_n) y_n, \text{ where } G_N = \sum_{n=1}^N \varphi(x_n) \varphi(x_n)^\top \text{ (} d \times d \text{ matrix)}.$$

- ▶ Any solution to (1) must satisfy

$$G_N \theta = \sum_{n=1}^N \varphi(x_n) y_n, \text{ where } G_N = \sum_{n=1}^N \varphi(x_n) \varphi(x_n)^\top \text{ (} d \times d \text{ matrix).}$$

- ▶ *Matrix notations:*

$$Y_N = (y_1, \dots, y_N)^\top \in \mathbb{R}^N,$$

$$\Phi_N = (\varphi^\top(x_1), \dots, \varphi^\top(x_N))^\top \text{ (} N \times d \text{ matrix).}$$

$$G_N \theta = \Phi_N^\top Y_N, \text{ where } G_N = \Phi_N^\top \Phi_N.$$

- ▶ Specific solution: $\theta_N^\dagger = G_N^\dagger \Phi_N^\top Y_N$ where G_N^\dagger : pseudo-inverse of G_N .

- ▶ Specific solution: $\theta_N^\dagger = G_N^\dagger \Phi_N^\top Y_N$ where G_N^\dagger : pseudo-inverse of G_N .
- ▶ Solutions:

$$\begin{aligned}\Theta_N &= \{\theta \in \Theta : G_N(\theta_N^\dagger - \theta) = 0\} \\ &= \{\theta_N^\dagger + \ker(G_N)\} \cap \Theta.\end{aligned}$$

- ▶ Specific solution: $\theta_N^\dagger = G_N^\dagger \Phi_N^\top Y_N$ where G_N^\dagger : pseudo-inverse of G_N .
- ▶ Solutions:

$$\begin{aligned}\Theta_N &= \{\theta \in \Theta : G_N(\theta_N^\dagger - \theta) = 0\} \\ &= \{\theta_N^\dagger + \ker(G_N)\} \cap \Theta.\end{aligned}$$

- ▶ When $\Theta = \mathbb{R}^d$ and G_N is invertible, $G_N^\dagger = G_N^{-1}$,

$$(\textit{Ordinary Least-squares}) \quad \theta_N = G_N^{-1} \Phi_N^\top Y_N.$$

► *Error control:*

$$\forall x \in \mathcal{X}, \quad |f_{\star}(x) - f_{\theta_N}(x)| \leq \|\theta_{\star} - \theta_N\|_A \|\varphi(x)\|_{A^{-1}}. \quad (2)$$

for each invertible matrix A , where $\|x\|_A = \sqrt{x^T A x}$.

► *Error control:*

$$\forall x \in \mathcal{X}, \quad |f_{\star}(x) - f_{\theta_N}(x)| \leq \|\theta_{\star} - \theta_N\|_A \|\varphi(x)\|_{A^{-1}}. \quad (2)$$

for each invertible matrix A , where $\|x\|_A = \sqrt{x^T A x}$.

► Matrix $A = G_N$ has natural interpretation: for $\theta \in \Theta_N$ (solution),

$$\sum_{n=1}^N (f_{\star}(x_n) - f_{\theta}(x_n))^2 = \sum_{n=1}^N (\theta^{\star} - \theta)^T \varphi(x_n) \varphi(x_n)^T (\theta^{\star} - \theta) = \|\theta^{\star} - \theta\|_{G_N}^2.$$

(Over-fitting is $\forall \theta \in \Theta_N, \|\theta^{\star} - \theta\|_{G_N} = 0$).

Study $\|\theta_{\star} - \theta_N\|_{G_N}$

When G_N is not invertible, introduce regularization parameter $\lambda \in \mathbb{R}_*^+$.

When G_N is not invertible, introduce regularization parameter $\lambda \in \mathbb{R}_*^+$.

► *Regularized* solution

$$\theta_{N,\lambda} = G_{N,\lambda}^{-1} \Phi_N^\top Y_N \text{ where } G_{N,\lambda} = \Phi_N^\top \Phi_N + \lambda I_d.$$

When G_N is not invertible, introduce regularization parameter $\lambda \in \mathbb{R}_*^+$.

- ▶ *Regularized* solution

$$\theta_{N,\lambda} = G_{N,\lambda}^{-1} \Phi_N^\top Y_N \text{ where } G_{N,\lambda} = \Phi_N^\top \Phi_N + \lambda I_d.$$

- ▶ Bayesian interpretation:

For *Prior* $\theta \sim \mathcal{N}(0, \Sigma)$, i.i.d. setup, Gaussian noise ($\xi_n \sim \mathcal{N}(0, \sigma^2)$),

Posterior: $\hat{f}_N(x) | x, x_1, y_1, \dots, x_N, y_N \sim \mathcal{N}(\mu_N(x), \sigma_N^2(x))$ where

$$\begin{aligned} \mu_N(x) &= \varphi(x)^\top (\Phi_N^\top \Phi_N + \sigma^2 \Sigma^{-1})^{-1} \Phi_N^\top Y_N \\ \sigma_N^2(x) &= \sigma^2 \varphi(x)^\top (\Phi_N^\top \Phi_N + \sigma^2 \Sigma^{-1})^{-1} \varphi(x). \end{aligned}$$

When G_N is not invertible, introduce regularization parameter $\lambda \in \mathbb{R}_*^+$.

- ▶ *Regularized* solution

$$\theta_{N,\lambda} = G_{N,\lambda}^{-1} \Phi_N^\top Y_N \text{ where } G_{N,\lambda} = \Phi_N^\top \Phi_N + \lambda I_d.$$

- ▶ Bayesian interpretation:

For *Prior* $\theta \sim \mathcal{N}(0, \Sigma)$, i.i.d. setup, Gaussian noise ($\xi_n \sim \mathcal{N}(0, \sigma^2)$),

Posterior: $\hat{f}_N(x) | x, x_1, y_1, \dots, x_N, y_N \sim \mathcal{N}(\mu_N(x), \sigma_N^2(x))$ where

$$\begin{aligned} \mu_N(x) &= \varphi(x)^\top (\Phi_N^\top \Phi_N + \sigma^2 \Sigma^{-1})^{-1} \Phi_N^\top Y_N \\ \sigma_N^2(x) &= \sigma^2 \varphi(x)^\top (\Phi_N^\top \Phi_N + \sigma^2 \Sigma^{-1})^{-1} \varphi(x). \end{aligned}$$

- ▶ Prior $\Sigma = \frac{\sigma^2}{\lambda} I_d$ gives **regularized least-squares** $\mu_N(x) = \varphi(x)^\top \theta_{N,\lambda}$.

When G_N is not invertible, introduce regularization parameter $\lambda \in \mathbb{R}_*^+$.

- ▶ *Regularized* solution

$$\theta_{N,\lambda} = G_{N,\lambda}^{-1} \Phi_N^\top Y_N \text{ where } G_{N,\lambda} = \Phi_N^\top \Phi_N + \lambda I_d.$$

- ▶ Bayesian interpretation:

For *Prior* $\theta \sim \mathcal{N}(0, \Sigma)$, i.i.d. setup, Gaussian noise ($\xi_n \sim \mathcal{N}(0, \sigma^2)$),

Posterior: $\hat{f}_N(x) | x, x_1, y_1, \dots, x_N, y_N \sim \mathcal{N}(\mu_N(x), \sigma_N^2(x))$ where

$$\begin{aligned} \mu_N(x) &= \varphi(x)^\top (\Phi_N^\top \Phi_N + \sigma^2 \Sigma^{-1})^{-1} \Phi_N^\top Y_N \\ \sigma_N^2(x) &= \sigma^2 \varphi(x)^\top (\Phi_N^\top \Phi_N + \sigma^2 \Sigma^{-1})^{-1} \varphi(x). \end{aligned}$$

- ▶ Prior $\Sigma = \frac{\sigma^2}{\lambda} I_d$ gives **regularized least-squares** $\mu_N(x) = \varphi(x)^\top \theta_{N,\lambda}$.
- ▶ Interpret λ as prior value on variance.

Study $\|\theta_* - \theta_{N,\lambda}\|_{G_{N,\lambda}}$

Standard regression noise assumptions

- ▶ *iid samples* $(x_t)_t$ are i.i.d., $(\xi_t)_t$ are i.i.d., independent from $(x_t)_t$.

Standard regression noise assumptions

- ▶ *iid samples* $(x_t)_t$ are i.i.d., $(\xi_t)_t$ are i.i.d., independent from $(x_t)_t$.

Standard regression noise assumptions

- ▶ *iid samples* $(x_t)_t$ are i.i.d., $(\xi_t)_t$ are i.i.d., independent from $(x_t)_t$.
- ▶ *sub-Gaussian* noise: For some $\sigma^2 > 0$,

$$\forall t \in \mathbb{N}, \forall \gamma \in \mathbb{R}, \quad \ln \mathbb{E} \left[\exp(\gamma \xi_t) \right] \leq \frac{\gamma^2 \sigma^2}{2}.$$

Standard regression noise assumptions

- ▶ *iid samples* $(x_t)_t$ are i.i.d., $(\xi_t)_t$ are i.i.d., independent from $(x_t)_t$.
- ▶ *sub-Gaussian* noise: For some $\sigma^2 > 0$,

$$\forall t \in \mathbb{N}, \forall \gamma \in \mathbb{R}, \quad \ln \mathbb{E} \left[\exp(\gamma \xi_t) \right] \leq \frac{\gamma^2 \sigma^2}{2}.$$

- ▶ = for $\mathcal{N}(0, \sigma^2)$ [Exercice]

Sequential regression noise assumption

- ▶ *Predictable sequence* (not iid): x_t is \mathcal{H}_{t-1} -measurable and y_t is \mathcal{H}_t -measurable. \mathcal{H}_t : history.

Standard regression noise assumptions

- ▶ *iid samples* $(x_t)_t$ are i.i.d., $(\xi_t)_t$ are i.i.d., independent from $(x_t)_t$.
- ▶ *sub-Gaussian* noise: For some $\sigma^2 > 0$,

$$\forall t \in \mathbb{N}, \forall \gamma \in \mathbb{R}, \quad \ln \mathbb{E} \left[\exp(\gamma \xi_t) \right] \leq \frac{\gamma^2 \sigma^2}{2}.$$

- ▶ = for $\mathcal{N}(0, \sigma^2)$ [Exercice]

Sequential regression noise assumption

- ▶ *Predictable sequence* (not iid): x_t is \mathcal{H}_{t-1} -measurable and y_t is \mathcal{H}_t -measurable. \mathcal{H}_t : history.
- ▶ *Conditionally* sub-Gaussian noise: For some $\sigma^2 > 0$,

$$\forall t \in \mathbb{N}, \forall \gamma \in \mathbb{R}, \quad \ln \mathbb{E} \left[\exp(\gamma \xi_t) \middle| \mathcal{H}_{t-1} \right] \leq \frac{\gamma^2 \sigma^2}{2}.$$

STRUCTURES

LINEAR BANDITS

Regression

Linear UCB, Linear TS

Graph-linear Bandits

Extension to Kernels

STRUCTURED LOWER BOUNDS

CONCLUSION, PERSPECTIVE

- ▶ *Least-squares (regularized) estimate* of θ_* :

$$\theta_{t,\lambda} = \underbrace{[\Phi_t^\top \Phi_t + \lambda I_d]^{-1}}_{G_{t,\lambda}} \Phi_t^\top Y_t.$$

- ▶ Choose $X_{t+1} = \operatorname{argmax}_{x \in \mathcal{X}} \langle \theta_{t,\lambda}, \varphi(x) \rangle$.

- ▶ *Least-squares (regularized) estimate* of θ_* :

$$\theta_{t,\lambda} = \underbrace{[\Phi_t^\top \Phi_t + \lambda I_d]^{-1}}_{G_{t,\lambda}} \Phi_t^\top Y_t.$$

- ▶ Choose $X_{t+1} = \operatorname{argmax}_{x \in \mathcal{X}} \langle \theta_{t,\lambda}, \varphi(x) \rangle$.

⇒ Exploitation only !

Optimism in Face of Uncertainty - Linear

Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári "Improved Algorithms for Linear Stochastic Bandits"
NIPS, 2011.

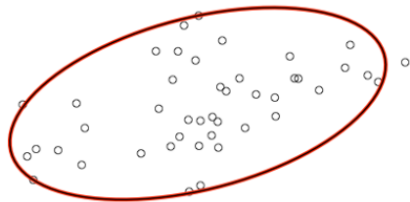
$$X_{t+1} = \operatorname{argmax}_{x \in \mathcal{X}} \max \left\{ f_{\theta}(x) : \theta \text{ is plausible} \right\}$$

$$X_{t+1} = \operatorname{argmax}_{x \in \mathcal{X}} \max \left\{ f_{\theta}(x) : \theta \text{ is plausible} \right\}$$

► Plausible: $C_t(\delta) = \left\{ \theta : \|\theta - \theta_{t,\lambda}\|_{G_{t,\lambda}} \leq B_t(\delta) \right\}$

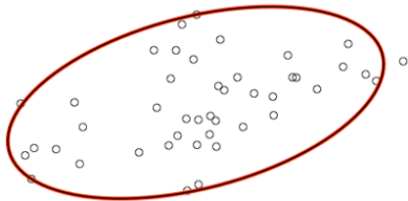
$$X_{t+1} = \operatorname{argmax}_{x \in \mathcal{X}} \max \left\{ f_{\theta}(x) : \theta \text{ is plausible} \right\}$$

- ▶ Plausible: $C_t(\delta) = \left\{ \theta : \|\theta - \theta_{t,\lambda}\|_{G_{t,\lambda}} \leq B_t(\delta) \right\}$
- ▶ Confidence ellipsoid such that $\mathbb{P}(\theta_{\star} \in C_t(\delta)) \geq 1 - \delta$.



$$X_{t+1} = \operatorname{argmax}_{x \in \mathcal{X}} \max \left\{ f_{\theta}(x) : \theta \text{ is plausible} \right\}$$

- ▶ Plausible: $C_t(\delta) = \left\{ \theta : \|\theta - \theta_{t,\lambda}\|_{G_{t,\lambda}} \leq B_t(\delta) \right\}$
- ▶ Confidence ellipsoid such that $\mathbb{P}(\theta_{\star} \in C_t(\delta)) \geq 1 - \delta$.



- ▶ Explicit solution

$$X_{t+1} = \operatorname{argmax}_{x \in \mathcal{X}} \langle \theta_{t,\lambda}, \varphi(x) \rangle + B_t(\delta) \|\varphi(x)\|_{G_{t,\lambda}^{-1}}.$$

\Rightarrow UCB-style exploitation and exploitation trade-off!

How to build $B_t(\delta)$?

How to build $B_t(\delta)$?

▶ (Dani, Kakade 2008) $B_t(\delta) = \sqrt{\max(128d \ln(t) \ln(t^2/\delta), 64/9 \ln^2(t^2/\delta))}$

How to build $B_t(\delta)$?

- ▶ (Dani, Kakade 2008) $B_t(\delta) = \sqrt{\max(128d \ln(t) \ln(t^2/\delta), 64/9 \ln^2(t^2/\delta))}$
- ▶ (Rusmevichientong, Tsitsiklis 2009)

$$B_t(\delta) = C \sqrt{\ln(t)} \sqrt{d \ln \left(\frac{36 \max_x \|\varphi(x)\|^2}{\lambda} t \right) + \ln(1/\delta)}$$

How to build $B_t(\delta)$?

▶ (Dani, Kakade 2008) $B_t(\delta) = \sqrt{\max(128d \ln(t) \ln(t^2/\delta), 64/9 \ln^2(t^2/\delta))}$

▶ (Rusmevichientong, Tsitsiklis 2009)

$$B_t(\delta) = C \sqrt{\ln(t)} \sqrt{d \ln \left(\frac{36 \max_x \|\varphi(x)\|^2}{\lambda} t \right) + \ln(1/\delta)}$$

▶ OFUL (Abbasi et al, 2011)

$$B_t(\delta) = \sqrt{\lambda} \|\theta^*\|_2 + \sqrt{2 \ln \left(\frac{\det(G_N + \lambda I)^{1/2}}{\delta \lambda^{d/2}} \right)}$$

$$|f_{\theta^*}(x) - f_{\theta_{N,\lambda}}(x)| \leq \|\theta_{\star} - \theta_{N,\lambda}\|_{G_{N,\lambda}} \|\varphi(x)\|_{G_{N,\lambda}^{-1}}$$

Decomposition lemma

$$\|\theta_{\star} - \theta_{N,\lambda}\|_{G_{N,\lambda}} \leq \sqrt{\lambda} \|\theta^*\|_2 + \|\Phi_N^{\top} E_N\|_{G_{N,\lambda}^{-1}}$$

where $E_N = (\xi_1, \dots, \xi_N)^{\top} \in \mathbb{R}^N$.

Key observation: sum of *conditionally centered* vector variables

$$\Phi_N^{\top} E_N = \sum_{n=1}^N \varphi(x_n) \xi_n \in \mathbb{R}^d.$$

\Rightarrow *Concentration inequality for vectors !*

Make use of *self-normalized* concentration inequalities.

$$\begin{aligned}
\theta^* - \theta_{N,\lambda} &= \theta^* - G_{N,\lambda}^{-1} \Phi_N^\top Y_N \\
&= \theta^* - G_{N,\lambda}^{-1} \Phi_N^\top (\Phi_N \theta^* + E_N) \\
&= (I - G_{N,\lambda}^{-1} G_N) \theta^* - G_{N,\lambda}^{-1} \Phi_N^\top E_N \\
&= G_{N,\lambda}^{-1} (G_{N,\lambda} - G_N) \theta^* - G_{N,\lambda}^{-1} \Phi_N^\top E_N. \\
&= \underbrace{\lambda G_{N,\lambda}^{-1} \theta^*}_{(1)} - \underbrace{G_{N,\lambda}^{-1} \Phi_N^\top E_N}_{(2)}.
\end{aligned}$$

$$\begin{aligned}
(1) \quad \|\lambda G_{N,\lambda}^{-1} \theta^*\|_{G_{N,\lambda}} &= \lambda \sqrt{\theta^{*\top} G_{N,\lambda}^{-1} G_{N,\lambda} G_{N,\lambda}^{-1} \theta^*} \\
&\leq \frac{\lambda}{\sqrt{\text{eig}_{\min}(G_{N,\lambda})}} \|\theta^*\|_2 \leq \sqrt{\lambda} \|\theta^*\|_2
\end{aligned}$$

$$(2) \quad \|G_{N,\lambda}^{-1} \Phi_N^\top E_N\|_{G_{N,\lambda}} = \|\Phi_N^\top E_N\|_{G_{N,\lambda}^{-1}}.$$

What it means to be *self-normalized* ?

In dimension $D = 1$, $\lambda = 0$, $G_N = \sum_{n=1}^N \varphi(x_n)^2$

$$\|\Phi_N^\top E_N\|_{G_{N,\lambda}^{-1}} = \frac{|\sum_{n=1}^N \varphi(x_n) \xi_n|}{\sqrt{\sum_{n=1}^N \varphi(x_n)^2}} = \frac{|\sum_{n=1}^N Z_n|}{\sqrt{\sum_{n=1}^N \sigma_n^2}}$$

Basic self-normalized (Gaussian) concentration inequality

For fixed t , Z_1, \dots, Z_t , independent, $Z_n \sim \mathcal{N}(0, \sigma_n^2)$, $\delta \in (0, 1]$

$$\mathbb{P}\left(\left|\frac{\sum_{n=1}^t Z_n}{\sqrt{\sum_{n=1}^t \sigma_n^2}}\right| \geq \sqrt{2 \ln(2/\delta)}\right) \leq \delta$$

Basic (Gaussian) concentration inequality For fixed t , Z_1, \dots, Z_t i.i.d. $\mathcal{N}(0, \sigma^2)$, $\delta \in (0, 1]$

$$\mathbb{P}\left(\frac{1}{t} \sum_{n=1}^t Z_n \geq \sqrt{\frac{2\sigma^2 \ln(1/\delta)}{t}}\right) \leq \delta$$

Likewise, using the Chernoff-method, we can show for fixed t , Z_1, \dots, Z_t , independent, $Z_n \sim \mathcal{N}(0, \sigma_n^2)$, $\delta \in (0, 1]$

$$\mathbb{P}\left(\sum_{n=1}^t Z_n \geq \sqrt{2 \sum_{n=1}^t \sigma_n^2 \ln(1/\delta)}\right) \leq \delta$$

Thus

$$\mathbb{P}\left(\frac{\sum_{n=1}^t Z_n}{\sqrt{\sum_{n=1}^t \sigma_n^2}} \geq \sqrt{2 \ln(1/\delta)}\right) \leq \delta$$

Extension to dimension d by the *Laplace method* (De la Peña et al., 2004).

Let $Z \in \mathbb{R}^d$ random *vector*, B a $d \times d$ random *matrix* such that

$$(\textit{Sub-Gaussian}) \quad \forall \gamma \in \mathbb{R}^d, \quad \ln \mathbb{E}[\exp(\gamma^\top Z - \frac{1}{2} \gamma^\top B \gamma)] \leq 0.$$

Then for any deterministic $d \times d$ matrix C , w.p. $\geq 1 - \delta$,

$$\|Z\|_{(B+C)^{-1}} \leq \sqrt{2 \ln \left(\frac{\det(B+C)^{1/2}}{\delta \det(C)^{1/2}} \right)}.$$

Extension to dimension d by the *Laplace method* (De la Peña et al., 2004).

Let $Z \in \mathbb{R}^d$ random *vector*, B a $d \times d$ random *matrix* such that

$$(\textit{Sub-Gaussian}) \quad \forall \gamma \in \mathbb{R}^d, \quad \ln \mathbb{E}[\exp(\gamma^\top Z - \frac{1}{2} \gamma^\top B \gamma)] \leq 0.$$

Then for any deterministic $d \times d$ matrix C , w.p. $\geq 1 - \delta$,

$$\|Z\|_{(B+C)^{-1}} \leq \sqrt{2 \ln \left(\frac{\det(B+C)^{1/2}}{\delta \det(C)^{1/2}} \right)}.$$

► Application: $Z = \sum_{n=1}^N \varphi(x_n) \xi_n$, $B = G_{N,0}$ $C = \lambda I_d$.

1) Quantity

$$M_t^\gamma = \exp \left(\langle \gamma, Z \rangle - \frac{1}{2} \|\lambda\|_B^2 \right)$$

is a super martingale such that for all t , $\mathbb{E}[M_t^\gamma] \leq 1$.

1) Quantity

$$M_t^\gamma = \exp \left(\langle \gamma, Z \rangle - \frac{1}{2} \|\lambda\|_B^2 \right)$$

is a super martingale such that for all t , $\mathbb{E}[M_t^\gamma] \leq 1$.

2) Choice of γ ? Replace optimization with **integration** (Laplace) !

Introduce distribution $\Lambda \sim \mathcal{N}(0, C^{-1})$, and M_t^Λ .

1) Quantity

$$M_t^\gamma = \exp \left(\langle \gamma, Z \rangle - \frac{1}{2} \|\lambda\|_B^2 \right)$$

is a super martingale such that for all t , $\mathbb{E}[M_t^\gamma] \leq 1$.

2) Choice of γ ? Replace optimization with **integration** (Laplace) !

Introduce distribution $\Lambda \sim \mathcal{N}(0, C^{-1})$, and M_t^Λ .

a) $\mathbb{E}[M_t^\Lambda] \leq 1$

b) $\mathbb{E}[M_t^\Lambda] = \mathbb{E}[\mathbb{E}[M_t^\Lambda | \mathcal{F}_\infty]]$ and

$$\mathbb{E}[M_t^\Lambda | \mathcal{F}_\infty] = \int_{\mathbb{R}^d} \exp \left(\langle \gamma, Z \rangle - \frac{1}{2} \|\lambda\|_B^2 \right) f(\lambda) d\lambda$$

where f denotes the pdf of $\Lambda \sim \mathcal{N}(0, C^{-1})$.

3) Direct calculations show that

$$\mathbb{E}[M_t^\wedge | \mathcal{F}_\infty] = \left(\frac{\det(C)}{\det(B+C)} \right)^{1/2} \exp \left(\frac{1}{2} \|Z\|_{(B+C)^{-1}}^2 \right)$$

$$\text{Then } \mathbb{E} \left[\left(\frac{\det(C)}{\det(B+C)} \right)^{1/2} \exp \left(\frac{1}{2} \|Z\|_{(B+C)^{-1}}^2 \right) \right] \leq 1$$

4) Markov inequality yields:

$$\begin{aligned} & \mathbb{P} \left(\|Z\|_{(B+C)^{-1}}^2 > 2 \ln \left(\frac{\det(B+C)^{1/2}}{\delta \det(B)^{1/2}} \right) \right) \\ &= \mathbb{P} \left(\exp \left(\frac{1}{2} \|Z\|_{(B+C)^{-1}}^2 \right) > \frac{\det(B+C)^{1/2}}{\delta \det(B)^{1/2}} \right) \leq \delta. \end{aligned}$$

► Application: $Z = \sum_{n=1}^N \varphi(x_n) \xi_n$, $B = G_{N,0}$ $C = \lambda I_d$.

$$\mathbb{P}\left(\|\Phi_N^\top E_N\|_{G_{N,\lambda}^{-1}} \geq 2 \ln \left(\frac{\det(G_{N,\lambda})^{1/2}}{\delta \lambda^{d/2}}\right)\right) \leq \delta.$$

- ▶ Application: $Z = \sum_{n=1}^N \varphi(x_n) \xi_n$, $B = G_{N,0}$, $C = \lambda I_d$.
- $$\mathbb{P}\left(\|\Phi_N^\top E_N\|_{G_{N,\lambda}^{-1}} \geq 2 \ln \left(\frac{\det(G_{N,\lambda})^{1/2}}{\delta \lambda^{d/2}}\right)\right) \leq \delta.$$
- ▶ Time-uniform bound ($\forall N$): handles random stopping time N .

- ▶ Application: $Z = \sum_{n=1}^N \varphi(x_n) \xi_n$, $B = G_{N,0}$, $C = \lambda I_d$.

$$\mathbb{P}\left(\|\Phi_N^\top E_N\|_{G_{N,\lambda}^{-1}} \geq 2 \ln \left(\frac{\det(G_{N,\lambda})^{1/2}}{\delta \lambda^{d/2}} \right)\right) \leq \delta.$$

- ▶ Time-uniform bound ($\forall N$): handles random stopping time N .
- ▶ Property:

$$\mathbb{E}[M_N^\wedge] = \mathbb{E}[\liminf_{m \rightarrow \infty} M_{\min(N,m)}^\wedge] \leq \liminf_{m \rightarrow \infty} \mathbb{E}[M_{\min(N,m)}^\wedge] \leq 1.$$

\Rightarrow *Confidence ellipsoid* on θ_* :

$$C_t(\delta) = \left\{ \theta : \|\theta - \theta_{t,\lambda}\|_{G_{t,\lambda}} \leq \sqrt{\lambda} \|\theta^*\|_2 + \sqrt{2 \ln \left(\frac{\det(G_t + \lambda I)^{1/2}}{\delta \lambda^{d/2}} \right)} \right\},$$

Information gain γ_T

Log-determinant Lemma

$$\gamma_T = \ln \left(\frac{\det(G_{T,\lambda})}{\det(\lambda I_d)} \right) = \sum_{t=1}^T \ln \left(1 + \|\varphi(x_t)\|_{G_{t-1,\lambda}^{-1}}^2 \right)$$

Information gain γ_T

Log-determinant Lemma

$$\gamma_T = \ln \left(\frac{\det(G_{T,\lambda})}{\det(\lambda I_d)} \right) = \sum_{t=1}^T \ln (1 + \|\varphi(x_t)\|_{G_{t-1,\lambda}^{-1}}^2)$$

- ▶ $\det(\lambda I_d)$: volume before observing data; $\det(G_{T,\lambda})$: volume after observing x_1, \dots, x_t .

Information gain γ_T

Log-determinant Lemma

$$\gamma_T = \ln \left(\frac{\det(G_{T,\lambda})}{\det(\lambda I_d)} \right) = \sum_{t=1}^T \ln \left(1 + \|\varphi(x_t)\|_{G_{t-1,\lambda}^{-1}}^2 \right)$$

- ▶ $\det(\lambda I_d)$: volume before observing data; $\det(G_{T,\lambda})$: volume after observing x_1, \dots, x_T .
- ▶ Captures how much the "*volume*" of information is modified by samples x_1, \dots, x_T .

Information gain γ_T

Log-determinant Lemma

$$\gamma_T = \ln \left(\frac{\det(G_{T,\lambda})}{\det(\lambda I_d)} \right) = \sum_{t=1}^T \ln \left(1 + \|\varphi(x_t)\|_{G_{t-1,\lambda}^{-1}}^2 \right)$$

- ▶ $\det(\lambda I_d)$: volume before observing data; $\det(G_{T,\lambda})$: volume after observing x_1, \dots, x_t .
- ▶ Captures how much the "*volume*" of information is modified by samples x_1, \dots, x_t .
- ▶ $\gamma_T = O(d \ln(T))$ for d -dimensional linear space.

$$\begin{aligned}
\det(G_{n,\lambda}) &= \det(G_{n-1,\lambda} + \varphi(x_n)\varphi(x_n)^\top) \\
&= \det(G_{n-1,\lambda}) \det\left(I + G_{n-1,\lambda}^{-1/2}\varphi(x_n)\left(G_{n-1,\lambda}^{-1/2}\varphi(x_n)\right)^\top\right) \\
&= \det(G_{n-1,\lambda})\left(1 + \|\varphi(x_n)\|_{G_{n-1,\lambda}^{-1}}^2\right) \\
&= \det(\lambda I) \prod_{t=1}^n \left(1 + \|\varphi(x_t)\|_{G_{t-1,\lambda}^{-1}}^2\right)
\end{aligned}$$

Thus,

$$\ln\left(\frac{\det(G_{n,\lambda})}{\lambda^d}\right) = \sum_{t=1}^n \ln\left(1 + \|\varphi(x_t)\|_{G_{t-1,\lambda}^{-1}}^2\right)$$

- ▷ We have good confidence bounds: let us exploit them!
- ▷ Simplest approach:

$$\begin{aligned} X_{t+1} &= \operatorname{argmax}_{x \in \mathcal{X}} \max \{ \langle \theta, \varphi(x) \rangle : \theta \in \mathcal{C}_t(\delta) \} . \\ &= \operatorname{argmax}_{x \in \mathcal{X}} f_t^+(x) \end{aligned}$$

Regret

If $f_\star(x) \in [-1, 1]$ for all x , then w.p. higher than $1 - \delta$,

$$\mathcal{R}_T = O\left(\sqrt{T\gamma_T} \left(\|\theta_\star\|_2 + \sigma \sqrt{2 \ln(1/\delta) + 2\gamma_T} \right)\right)$$

- ▷ Is this optimal way of exploiting linear structure?

Instantaneous regret r_t (note: $r_t \leq 2$)

$$\begin{aligned} r_t &= f_*(x_*) - f_*(x_t) \\ &\leq f_{t-1}^+(x_t) - f_*(x_t) \text{ with high probability} \\ &\leq |f_{t-1}^+(x_t) - f_{\lambda, t-1}(x_t)| + |f_{\lambda, t-1}(x_t) - f_*(x_t)| \\ &\leq 2\|\varphi(x_t)\|_{G_{t,\lambda}^{-1}} B_{t-1}(\delta). \end{aligned}$$

Thus, we deduce that with probability higher than $1 - \delta$:

$$\begin{aligned} \mathfrak{R}_T &= \sum_{t=1}^T r_t \leq \sum_{t=1}^T 2 \min\{\|\varphi(x_t)\|_{G_{t,\lambda}^{-1}} B_{t-1}(\delta), 1\} \\ &\leq 2B_T(\delta) \sum_{t=1}^T \min\{\|\varphi(x_t)\|_{G_{t,\lambda}^{-1}}, 1\} \\ &\leq 2B_T(\delta) \sqrt{T \sum_{t=1}^T \min\{\|\varphi(x_t)\|_{G_{t,\lambda}^{-1}}^2, 1\}}. \end{aligned}$$

We conclude remarking that $\min\{A, 1\} \leq \frac{\ln(1+A)}{\ln(2)}$ for all $A \geq 0$.

Thompson in Sampling for Linear - Bandits

Shipra Agrawal, Navin Goyal "Thompson Sampling for Contextual Bandits with Linear Payoffs"
arXiv:1209.3352, 2014.

► **Bayesian model:**

$$y_t = x_t^T \theta + \varepsilon_t, \quad \theta \sim \mathcal{N}(0, \kappa^2 I_d), \quad \varepsilon_t \sim \mathcal{N}(0, \sigma^2).$$

Explicit posterior: $p(\theta | x_1, y_1, \dots, x_t, y_t) = \mathcal{N}(\hat{\theta}(t), \Sigma_t)$.

► **Bayesian model:**

$$y_t = x_t^T \theta + \varepsilon_t, \quad \theta \sim \mathcal{N}(0, \kappa^2 I_d), \quad \varepsilon_t \sim \mathcal{N}(0, \sigma^2).$$

Explicit posterior: $p(\theta | x_1, y_1, \dots, x_t, y_t) = \mathcal{N}(\hat{\theta}(t), \Sigma_t)$.

► **Thompson Sampling**

$$\begin{aligned} \tilde{\theta}(t) &\sim \mathcal{N}(\hat{\theta}(t), \Sigma_t), \\ x_{t+1} &= \operatorname{argmax}_{x \in \mathcal{D}_{t+1}} x^T \tilde{\theta}(t). \end{aligned}$$

[Li et al. 12],[Agrawal & Goyal 13]

STRUCTURES

LINEAR BANDITS

Regression

Linear UCB, Linear TS

Graph-linear Bandits

Extension to Kernels

STRUCTURED LOWER BOUNDS

CONCLUSION, PERSPECTIVE

$\mathcal{G} = (\mathcal{V}, \mathcal{E})$ graph with set of nodes $\mathcal{V} = \{1, \dots, N\}$, and edges $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$.

$\mathcal{G} = (\mathcal{V}, \mathcal{E})$ graph with set of nodes $\mathcal{V} = \{1, \dots, N\}$, and edges $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$.

▶ $\mathbf{W} = (w_{i,j})_{i,j}$ Weight matrix (non-negative weights)

$\mathcal{G} = (\mathcal{V}, \mathcal{E})$ graph with set of nodes $\mathcal{V} = \{1, \dots, N\}$, and edges $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$.

- ▶ $\mathbf{W} = (w_{i,j})_{i,j}$ Weight matrix (non-negative weights)
- ▶ $\mathbf{D} = \text{Diag}((\sum_j w_{i,j})_i)$ Degree matrix

$\mathcal{G} = (\mathcal{V}, \mathcal{E})$ graph with set of nodes $\mathcal{V} = \{1, \dots, N\}$, and edges $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$.

- ▶ $\mathbf{W} = (w_{i,j})_{i,j}$ Weight matrix (non-negative weights)
- ▶ $\mathbf{D} = \text{Diag}((\sum_j w_{i,j})_i)$ Degree matrix
- ▶ $\mathbf{L} = \mathbf{D} - \mathbf{W}$ graph Laplacian matrix

A graph function is seen as a vector $f \in \mathbb{R}^N$ assigning values to nodes.

$$f^\top \mathbf{L} f = \frac{1}{2} \sum_{i,j \leq N} w_{i,j} (f_i - f_j)^2.$$

Properties:

A graph function is seen as a vector $f \in \mathbb{R}^N$ assigning values to nodes.

$$f^\top \mathbf{L} f = \frac{1}{2} \sum_{i,j \leq N} w_{i,j} (f_i - f_j)^2.$$

Properties:

- ▶ \mathbf{L} is symmetric, positive, semi-definite.

A graph function is seen as a vector $f \in \mathbb{R}^N$ assigning values to nodes.

$$f^\top \mathbf{L} f = \frac{1}{2} \sum_{i,j \leq N} w_{i,j} (f_i - f_j)^2.$$

Properties:

- ▶ \mathbf{L} is symmetric, positive, semi-definite.
- ▶ Smallest eigenvalue is 0, corresponding vector $\mathbf{1}_N$

A graph function is seen as a vector $f \in \mathbb{R}^N$ assigning values to nodes.

$$f^\top \mathbf{L} f = \frac{1}{2} \sum_{i,j \leq N} w_{i,j} (f_i - f_j)^2.$$

Properties:

- ▶ \mathbf{L} is symmetric, positive, semi-definite.
- ▶ Smallest eigenvalue is 0, corresponding vector $\mathbf{1}_N$
- ▶ Eigenvalues : $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$

Let $\mathbf{L} = \mathbf{Q}^\top \mathbf{\Lambda} \mathbf{Q}$ where

Let $\mathbf{L} = \mathbf{Q}^\top \mathbf{\Lambda} \mathbf{Q}$ where

- ▶ $\mathbf{\Lambda}$: $N \times N$ diagonal matrix with eigenvalues of \mathbf{L}

Let $\mathbf{L} = \mathbf{Q}^\top \mathbf{\Lambda} \mathbf{Q}$ where

- ▶ $\mathbf{\Lambda}$: $N \times N$ diagonal matrix with eigenvalues of \mathbf{L}
- ▶ \mathbf{Q} : $N \times N$ matrix whose columns are eigenvectors of \mathbf{L} .

Any graph-function f decomposes as $f = \mathbf{Q}\alpha$ for some α , that is

Let $\mathbf{L} = \mathbf{Q}^\top \mathbf{\Lambda} \mathbf{Q}$ where

- ▶ $\mathbf{\Lambda}$: $N \times N$ diagonal matrix with eigenvalues of \mathbf{L}
- ▶ \mathbf{Q} : $N \times N$ matrix whose columns are eigenvectors of \mathbf{L} .

Any graph-function f decomposes as $f = \mathbf{Q}\alpha$ for some α , that is

- ▶ $f(i) = \sum_{j \in \mathcal{V}} \alpha_j Q_{i,j} = \langle \alpha, q(i) \rangle$ where $q(i) = (Q_{i,j})_j$ is i^{th} eigenvector.

Let $\mathbf{L} = \mathbf{Q}^\top \mathbf{\Lambda} \mathbf{Q}$ where

- ▶ $\mathbf{\Lambda}$: $N \times N$ diagonal matrix with eigenvalues of \mathbf{L}
- ▶ \mathbf{Q} : $N \times N$ matrix whose columns are eigenvectors of \mathbf{L} .

Any graph-function f decomposes as $f = \mathbf{Q}\alpha$ for some α , that is

- ▶ $f(i) = \sum_{j \in \mathcal{V}} \alpha_j Q_{i,j} = \langle \alpha, q(i) \rangle$ where $q(i) = (Q_{i,j})_j$ is i^{th} eigenvector.
- ▶ Then, $f^\top \mathbf{L} f = \sum_{i \in \mathcal{V}} \lambda_i \alpha_i^2 = \|\alpha\|_{\mathbf{\Lambda}} \stackrel{\text{def}}{=} \|f\|_{\mathcal{G}}$

\Rightarrow *Linear space* induced by the Graph:

$$\mathcal{F}_{\mathcal{G}} = \{f : f(x) = \langle \alpha, q(x) \rangle, \|\alpha\|_{\mathbf{\Lambda}} \leq 1\}$$

Low-norm $\|f\|_{\mathcal{G}}$ means:

Let $\mathbf{L} = \mathbf{Q}^\top \mathbf{\Lambda} \mathbf{Q}$ where

- ▶ $\mathbf{\Lambda}$: $N \times N$ diagonal matrix with eigenvalues of \mathbf{L}
- ▶ \mathbf{Q} : $N \times N$ matrix whose columns are eigenvectors of \mathbf{L} .

Any graph-function f decomposes as $f = \mathbf{Q}\alpha$ for some α , that is

- ▶ $f(i) = \sum_{j \in \mathcal{V}} \alpha_j Q_{i,j} = \langle \alpha, q(i) \rangle$ where $q(i) = (Q_{i,j})_j$ is i^{th} eigenvector.
- ▶ Then, $f^\top \mathbf{L} f = \sum_{i \in \mathcal{V}} \lambda_i \alpha_i^2 = \|\alpha\|_{\mathbf{\Lambda}} \stackrel{\text{def}}{=} \|f\|_{\mathcal{G}}$

\Rightarrow *Linear space* induced by the Graph:

$$\mathcal{F}_{\mathcal{G}} = \{f : f(x) = \langle \alpha, q(x) \rangle, \|\alpha\|_{\mathbf{\Lambda}} \leq 1\}$$

Low-norm $\|f\|_{\mathcal{G}}$ means:

- ▶ $(f_i - f_j)^2$ is small if $w_{i,j}$ is large

Let $\mathbf{L} = \mathbf{Q}^\top \mathbf{\Lambda} \mathbf{Q}$ where

- ▶ $\mathbf{\Lambda}$: $N \times N$ diagonal matrix with eigenvalues of \mathbf{L}
- ▶ \mathbf{Q} : $N \times N$ matrix whose columns are eigenvectors of \mathbf{L} .

Any graph-function f decomposes as $f = \mathbf{Q}\alpha$ for some α , that is

- ▶ $f(i) = \sum_{j \in \mathcal{V}} \alpha_j Q_{i,j} = \langle \alpha, q(i) \rangle$ where $q(i) = (Q_{i,j})_j$ is i^{th} eigenvector.
- ▶ Then, $f^\top \mathbf{L} f = \sum_{i \in \mathcal{V}} \lambda_i \alpha_i^2 = \|\alpha\|_{\mathbf{\Lambda}} \stackrel{\text{def}}{=} \|f\|_{\mathcal{G}}$

\Rightarrow *Linear space* induced by the Graph:

$$\mathcal{F}_{\mathcal{G}} = \{f : f(x) = \langle \alpha, q(x) \rangle, \|\alpha\|_{\mathbf{\Lambda}} \leq 1\}$$

Low-norm $\|f\|_{\mathcal{G}}$ means:

- ▶ $(f_i - f_j)^2$ is small if $w_{i,j}$ is large
- ▶ *similar value* between *neighbor nodes*.

Further references for bandits on graphs:

- ▶ Michal Valko, Rémi Munos, Branislav Kveton, Tomáš Kocák: *Spectral Bandits for Smooth Graph Functions*, in International Conference on Machine Learning (ICML 2014).

Further references for bandits on graphs:

- ▶ Michal Valko, Rémi Munos, Branislav Kveton, Tomáš Kocák: *Spectral Bandits for Smooth Graph Functions*, in International Conference on Machine Learning (ICML 2014).
- ▶ Alexandra Carpentier, Michal Valko: *Revealing graph bandits for maximizing local influence*, in International Conference on Artificial Intelligence and Statistics (AISTATS 2016).

STRUCTURES

LINEAR BANDITS

Regression

Linear UCB, Linear TS

Graph-linear Bandits

Extension to Kernels

STRUCTURED LOWER BOUNDS

CONCLUSION, PERSPECTIVE

Let k be a kernel function (continuous, symmetric positive definite) on a compact \mathcal{X} with positive finite Borel measure μ .

There exists an at most *countable* sequence $(\sigma_i, \psi_i)_{i \in \mathbb{N}^*}$ where $\sigma_i \geq 0$, $\lim_{i \rightarrow \infty} \sigma_i = 0$ and $\{\psi_i\}$ form an orthonormal basis of $L_{2,\mu}(\mathcal{X})$, such that

$$k(x, y) = \sum_{j=1}^{\infty} \sigma_j \psi_j(x) \psi_j(y') \quad \text{and} \quad \|f\|_{\mathcal{K}}^2 = \sum_{j=1}^{\infty} \frac{\langle f, \psi_j \rangle_{L_{2,\mu}}^2}{\sigma_j}$$

Let $\varphi_i = \sqrt{\sigma_i} \psi_i$ (hence $\|\varphi_i\|_{L_2} = \sqrt{\sigma_i}$, $\|\varphi_i\|_{\mathcal{K}} = 1$.)

If $f = \sum_i \theta_i \varphi_i$, then $\|f\|_{\mathcal{K}}^2 = \sum_i \theta_i^2$.

Similar to parametric regression except with infinite parameter.

Let k be a kernel function.

In the parametric case, we built $\theta_{\lambda,t}$, then $f_{\lambda,t}(x) = \langle \theta_{\lambda,t}, \varphi(x) \rangle$.

After observing $Y_t = (y_1, \dots, y_t)^\top \in \mathbb{R}^t$, we now build directly:

$$\text{(Kernel estimate)} \quad f_{\lambda,t}(x) = k_t(x)^\top (\mathbf{K}_t + \lambda I_t)^{-1} Y_t,$$

where

Let k be a kernel function.

In the parametric case, we built $\theta_{\lambda,t}$, then $f_{\lambda,t}(x) = \langle \theta_{\lambda,t}, \varphi(x) \rangle$.

After observing $Y_t = (y_1, \dots, y_t)^\top \in \mathbb{R}^t$, we now build directly:

$$\text{(Kernel estimate)} \quad f_{\lambda,t}(x) = k_t(x)^\top (\mathbf{K}_t + \lambda I_t)^{-1} Y_t,$$

where

$$\blacktriangleright k_t(x) = (k(x, x_{t'}))_{t' \leq t} \in \mathbb{R}^t,$$

Let k be a kernel function.

In the parametric case, we built $\theta_{\lambda,t}$, then $f_{\lambda,t}(x) = \langle \theta_{\lambda,t}, \varphi(x) \rangle$.

After observing $Y_t = (y_1, \dots, y_t)^\top \in \mathbb{R}^t$, we now build directly:

$$\text{(Kernel estimate)} \quad f_{\lambda,t}(x) = k_t(x)^\top (\mathbf{K}_t + \lambda I_t)^{-1} Y_t,$$

where

$$\blacktriangleright k_t(x) = (k(x, x_{t'}))_{t' \leq t} \in \mathbb{R}^t,$$

Let k be a kernel function.

In the parametric case, we built $\theta_{\lambda,t}$, then $f_{\lambda,t}(x) = \langle \theta_{\lambda,t}, \varphi(x) \rangle$.

After observing $Y_t = (y_1, \dots, y_t)^\top \in \mathbb{R}^t$, we now build directly:

$$\text{(Kernel estimate)} \quad f_{\lambda,t}(x) = k_t(x)^\top (\mathbf{K}_t + \lambda I_t)^{-1} Y_t,$$

where

- ▶ $k_t(x) = (k(x, x_{t'}))_{t' \leq t} \in \mathbb{R}^t$,
- ▶ $\mathbf{K}_t = (k(x_s, x_{s'}))_{s, s' \leq t} \in \mathbb{R}^{t \times t}$,

for a parameter $\lambda \in \mathbb{R}$.

Theorem (Durand & M. 2017, Kernel estimation error)

$\forall \delta \in [0, 1]$, with probability higher than $1 - \delta$, it holds **simultaneously over all** $x \in \mathcal{X}$ and $\mathbf{t} \geq \mathbf{0}$,

$$|f_{\star}(x) - f_{\lambda, \mathbf{t}}(x)| \leq \sqrt{k_{\lambda, \mathbf{t}}(x, x)} \left[\|f_{\star}\|_k + \frac{\sigma}{\sqrt{\lambda}} \sqrt{2 \ln(1/\delta) + 2\gamma_{\mathbf{t}}(\lambda)} \right],$$

where

Theorem (Durand & M. 2017, Kernel estimation error)

$\forall \delta \in [0, 1]$, with probability higher than $1 - \delta$, it holds **simultaneously over all** $x \in \mathcal{X}$ and $\mathbf{t} \geq \mathbf{0}$,

$$|f_{\star}(x) - f_{\lambda, \mathbf{t}}(x)| \leq \sqrt{k_{\lambda, \mathbf{t}}(x, x)} \left[\|f_{\star}\|_k + \frac{\sigma}{\sqrt{\lambda}} \sqrt{2 \ln(1/\delta) + 2\gamma_{\mathbf{t}}(\lambda)} \right],$$

where

► $k_{\lambda, \mathbf{t}}(x, x) = k(x, x) - k_{\mathbf{t}}(x)^{\top} (\mathbf{K}_{\mathbf{t}} + \lambda I_{\mathbf{t}})^{-1} k_{\mathbf{t}}(x)$: *posterior variance*.

Theorem (Durand & M. 2017, Kernel estimation error)

$\forall \delta \in [0, 1]$, with probability higher than $1 - \delta$, it holds **simultaneously over all** $x \in \mathcal{X}$ and $\mathbf{t} \geq \mathbf{0}$,

$$|f_{\star}(x) - f_{\lambda, \mathbf{t}}(x)| \leq \sqrt{k_{\lambda, \mathbf{t}}(x, x)} \left[\|f_{\star}\|_k + \frac{\sigma}{\sqrt{\lambda}} \sqrt{2 \ln(1/\delta) + 2\gamma_{\mathbf{t}}(\lambda)} \right],$$

where

- ▶ $k_{\lambda, \mathbf{t}}(x, x) = k(x, x) - k_{\mathbf{t}}(x)^{\top} (\mathbf{K}_{\mathbf{t}} + \lambda I_{\mathbf{t}})^{-1} k_{\mathbf{t}}(x)$: *posterior variance*.
- ▶ $\gamma_{\mathbf{t}}(\lambda) = \frac{1}{2} \sum_{t'=1}^{\mathbf{t}} \ln \left(1 + \frac{1}{\lambda} k_{\lambda, t'-1}(x_{t'}, x_{t'}) \right)$: *information gain*.

Theorem (Durand & M. 2017, Kernel estimation error)

$\forall \delta \in [0, 1]$, with probability higher than $1 - \delta$, it holds **simultaneously over all** $x \in \mathcal{X}$ and $\mathbf{t} \geq \mathbf{0}$,

$$|f_{\star}(x) - f_{\lambda, \mathbf{t}}(x)| \leq \sqrt{k_{\lambda, \mathbf{t}}(x, x) B_{\lambda, \mathbf{t}-1}(\delta)},$$

where

- ▶ $k_{\lambda, \mathbf{t}}(x, x) = k(x, x) - k_{\mathbf{t}}(x)^{\top} (\mathbf{K}_{\mathbf{t}} + \lambda I_{\mathbf{t}})^{-1} k_{\mathbf{t}}(x)$: *posterior variance*.
- ▶ $\gamma_{\mathbf{t}}(\lambda) = \frac{1}{2} \sum_{t'=1}^{\mathbf{t}} \ln \left(1 + \frac{1}{\lambda} k_{\lambda, t'-1}(x_{t'}, x_{t'}) \right)$: *information gain*.
- ▶ $\|f_{\star}\|_k$: Reproducing Kernel Hilbert Space norm.

$k(x, x')$	Captures	γ_T
$\langle x, x' \rangle$	"Linear functions"	$O(d \ln(T))$
$\exp(-\frac{\ x-x'\ ^2}{2\ell^2})$	"Smooth functions"	$O(\ln(T)^{d+1})$
...

Many kernels, for different properties of the signal
(graph-smoothness, periodic, change points, etc.)

Minimize the regret: $\mathcal{R}_T = \sum_{t=1}^T f_*(\star) - f_*(x_t)$.

Kernel-UCB

$$x_t \in \operatorname{argmax}_{x \in \mathcal{X}} f_t^+(x) \quad \text{where} \quad f_t^+(x) = f_{\lambda, t-1}(x) + \sqrt{k_{\lambda, t-1}(x, x) B_{\lambda, t-1}(\delta)}.$$

Kernel-TS (on discrete set $\mathbb{X} \subset \mathcal{X}$)

$$x_t \in \operatorname{argmax}_{x \in \mathbb{X}} \tilde{f}_t(x) \quad \text{where} \quad \tilde{f}_t \sim \mathcal{N}(\hat{\mathbf{f}}_{t-1}, \hat{\Sigma}_{t-1}) \quad \text{with}$$

$$\hat{\mathbf{f}}_{t-1} = (f_{\lambda, t-1}(x))_{x \in \mathbb{X}}, \quad \hat{\Sigma}_{t-1} = (k_{\lambda, t-1}(x, x') B_{\lambda, t-1}(\delta)^2)_{x, x' \in \mathbb{X}}.$$

More info in (Durand et al., 2018, JMLR)

STRUCTURES

LINEAR BANDITS

STRUCTURED LOWER BOUNDS

CONCLUSION, PERSPECTIVE

STRUCTURES

LINEAR BANDITS

STRUCTURED LOWER BOUNDS

Lower bounds

Lipschitz bandits

Ranking bandits

Metric-graph of bandits

CONCLUSION, PERSPECTIVE

Set of optimal arms for $\nu = (\nu_a)_{a \in \mathcal{A}}$: $\mathcal{A}_*(\nu) = \text{Argmax}_{a \in \mathcal{A}} \mu_a(\nu)$.

Definition (Uniformly Good strategies)

A bandit strategy is *uniformly-good* on \mathcal{D} if

$$\forall \nu = (\nu_a)_{a \in \mathcal{A}} \in \mathcal{D}, \forall a \notin \mathcal{A}_*(\nu), \quad \mathbb{E}[N_T(a)] = o(T^\alpha) \quad \text{for all } \alpha \in (0, 1].$$

Theorem ((Lai, Robbins 85) "Price for being uniformly-good")

Any uniformly good strategy on $\mathcal{D} = \text{Bern}^{\mathcal{A}}$ must satisfy

$$\forall a \notin \mathcal{A}_*(\nu) \quad \liminf_{T \rightarrow \infty} \frac{\mathbb{E}_\nu[N_T(a)]}{\log(T)} \geq \frac{1}{\text{kl}(\mu_a(\nu), \mu_*(\nu))}.$$

REGRET LOWER BOUNDS

Set of optimal arms for $\nu = (\nu_a)_{a \in \mathcal{A}}$: $\mathcal{A}_*(\nu) = \text{Argmax}_{a \in \mathcal{A}} \mu_a(\nu)$.

Definition (Uniformly Good strategies)

A bandit strategy is *uniformly-good* on \mathcal{D} if

$$\forall \nu = (\nu_a)_{a \in \mathcal{A}} \in \mathcal{D}, \forall a \notin \mathcal{A}_*(\nu), \quad \mathbb{E}[N_T(a)] = o(T^\alpha) \quad \text{for all } \alpha \in (0, 1].$$

Theorem ((Lai, Robbins 85) "Price for being uniformly-good")

Any uniformly good strategy on $\mathcal{D} = \text{Bern}^{\mathcal{A}}$ must satisfy

$$\forall a \notin \mathcal{A}_*(\nu) \quad \liminf_{T \rightarrow \infty} \frac{\mathbb{E}_\nu[N_T(a)]}{\log(T)} \geq \frac{1}{\text{kl}(\mu_a(\nu), \mu_*(\nu))}.$$

Main tool: *Change of measure*

(Probability) $\forall \Omega, \forall c \in \mathbb{R}, \mathbb{P}_\nu \left(\Omega \cap \left\{ \log \left(\frac{d\nu}{d\tilde{\nu}}(X) \right) \leq c \right\} \right) \leq \exp(c) \mathbb{P}_{\tilde{\nu}}(\Omega).$

(Expectation) $\mathbb{E}_\nu \left[\log \left(\frac{d\nu}{d\tilde{\nu}}(X) \right) \right] \geq \sup_{g: \mathcal{X} \rightarrow [0,1]} \text{kl} \left(\mathbb{E}_\nu[g(X)], \mathbb{E}_{\tilde{\nu}}[g(X)] \right).$

Consider $\theta, \theta' \in \Theta$:

$$\hat{\mathcal{L}}_T = \sum_{s=1}^T \ln \left(\frac{\nu_{\theta'_{A_s}}(Y_s)}{\nu_{\theta_{A_s}}(Y_s)} \right) = \sum_{a \in \mathcal{A}} \sum_{s=1}^T \mathbb{I}\{A_s = a\} \ln \left(\frac{\nu_{\theta'_a}(Y_s)}{\nu_{\theta_a}(Y_s)} \right)$$

Consider $\theta, \theta' \in \Theta$:

$$\hat{\mathcal{L}}_T = \sum_{s=1}^T \ln \left(\frac{\nu_{\theta'_{A_s}}(Y_s)}{\nu_{\theta_{A_s}}(Y_s)} \right) = \sum_{a \in \mathcal{A}} \sum_{s=1}^T \mathbb{I}\{A_s = a\} \ln \left(\frac{\nu_{\theta'_a}(Y_s)}{\nu_{\theta_a}(Y_s)} \right)$$

For any event Ω it holds (*Change of measure*)

$$\begin{aligned} \mathbb{P}_{\theta'}[\Omega] &= \mathbb{E}_{\theta}[\exp(\hat{\mathcal{L}}_T)\mathbb{I}\{\Omega\}] = \mathbb{E}_{\theta}[\exp(\hat{\mathcal{L}}_T)|\Omega] \mathbb{P}_{\theta}[\Omega] \\ &\stackrel{\text{Jensen}}{\geq} \exp\left(\mathbb{E}_{\theta}[\hat{\mathcal{L}}_T|\Omega]\right) \mathbb{P}_{\theta}[\Omega] = \exp\left(\frac{\mathbb{E}_{\theta}[\hat{\mathcal{L}}_T\mathbb{I}\{\Omega\}]}{\mathbb{P}_{\theta}[\Omega]}\right) \mathbb{P}_{\theta}[\Omega], \end{aligned}$$

Reorganizing the terms, we get $-\mathbb{E}_{\theta}[\hat{\mathcal{L}}_T\mathbb{I}\{\Omega\}] \geq \mathbb{P}_{\theta}[\Omega] \ln \left(\frac{\mathbb{P}_{\theta}[\Omega]}{\mathbb{P}_{\theta'}[\Omega]} \right)$.

WHY KL? LOG-LIKELIHOOD (FROM WEYL 1940)

Consider $\theta, \theta' \in \Theta$:

$$\hat{\mathcal{L}}_T = \sum_{s=1}^T \ln \left(\frac{\nu_{\theta'_s}(Y_s)}{\nu_{\theta_s}(Y_s)} \right) = \sum_{a \in \mathcal{A}} \sum_{s=1}^T \mathbb{I}\{A_s = a\} \ln \left(\frac{\nu_{\theta'_a}(Y_s)}{\nu_{\theta_a}(Y_s)} \right)$$

For any event Ω it holds (*Change of measure*)

$$\begin{aligned} \mathbb{P}_{\theta'}[\Omega] &= \mathbb{E}_{\theta}[\exp(\hat{\mathcal{L}}_T)\mathbb{I}\{\Omega\}] = \mathbb{E}_{\theta}[\exp(\hat{\mathcal{L}}_T)|\Omega]\mathbb{P}_{\theta}[\Omega] \\ &\stackrel{\text{Jensen}}{\geq} \exp\left(\mathbb{E}_{\theta}[\hat{\mathcal{L}}_T|\Omega]\right)\mathbb{P}_{\theta}[\Omega] = \exp\left(\frac{\mathbb{E}_{\theta}[\hat{\mathcal{L}}_T\mathbb{I}\{\Omega\}]}{\mathbb{P}_{\theta}[\Omega]}\right)\mathbb{P}_{\theta}[\Omega], \end{aligned}$$

Reorganizing the terms, we get $-\mathbb{E}_{\theta}[\hat{\mathcal{L}}_T\mathbb{I}\{\Omega\}] \geq \mathbb{P}_{\theta}[\Omega] \ln \left(\frac{\mathbb{P}_{\theta}[\Omega]}{\mathbb{P}_{\theta'}[\Omega]} \right)$. Likewise for the complement Ω^c . Summing up the terms, we obtain

$$\begin{aligned} -\mathbb{E}_{\theta}[\hat{\mathcal{L}}_T] &= \sum_{a \in \mathcal{A}} \mathbb{E}_{\theta}[N_T(a)] \text{KL}(\theta_a, \theta'_a) \\ &\geq \mathbb{P}_{\theta}[\Omega] \ln \left(\frac{\mathbb{P}_{\theta}[\Omega]}{\mathbb{P}_{\theta'}[\Omega]} \right) + (1 - \mathbb{P}_{\theta}[\Omega]) \ln \left(\frac{1 - \mathbb{P}_{\theta}[\Omega]}{1 - \mathbb{P}_{\theta'}[\Omega]} \right). \end{aligned}$$

WHY KL? LOG-LIKELIHOOD (FROM WEYL 1940)

Consider $\theta, \theta' \in \Theta$:

$$\widehat{\mathcal{L}}_T = \sum_{s=1}^T \ln \left(\frac{\nu_{\theta'_a}(Y_s)}{\nu_{\theta_a}(Y_s)} \right) = \sum_{a \in \mathcal{A}} \sum_{s=1}^T \mathbb{I}\{A_s = a\} \ln \left(\frac{\nu_{\theta'_a}(Y_s)}{\nu_{\theta_a}(Y_s)} \right)$$

For any event Ω it holds (*Change of measure*)

$$\begin{aligned} \mathbb{P}_{\theta'}[\Omega] &= \mathbb{E}_{\theta}[\exp(\widehat{\mathcal{L}}_T)\mathbb{I}\{\Omega\}] = \mathbb{E}_{\theta}[\exp(\widehat{\mathcal{L}}_T)|\Omega]\mathbb{P}_{\theta}[\Omega] \\ &\stackrel{\text{Jensen}}{\geq} \exp\left(\mathbb{E}_{\theta}[\widehat{\mathcal{L}}_T|\Omega]\right)\mathbb{P}_{\theta}[\Omega] = \exp\left(\frac{\mathbb{E}_{\theta}[\widehat{\mathcal{L}}_T\mathbb{I}\{\Omega\}]}{\mathbb{P}_{\theta}[\Omega]}\right)\mathbb{P}_{\theta}[\Omega], \end{aligned}$$

Reorganizing the terms, we get $-\mathbb{E}_{\theta}[\widehat{\mathcal{L}}_T\mathbb{I}\{\Omega\}] \geq \mathbb{P}_{\theta}[\Omega] \ln \left(\frac{\mathbb{P}_{\theta}[\Omega]}{\mathbb{P}_{\theta'}[\Omega]} \right)$. Likewise for the complement Ω^c . Summing up the terms, we obtain

$$\sum_{a \in \mathcal{A}} \mathbb{E}_{\theta}[N_T(a)] \text{KL}(\theta_a, \theta'_a) \geq \mathbf{k}1(\mathbb{P}_{\theta}[\Omega], \mathbb{P}_{\theta'}[\Omega])$$

Hence for all suboptimal arm $a \neq \star_\theta$,

$$\mathbb{E}_\theta[N_T(a)] \geq \sup_{\Omega, \theta'} \frac{\text{kl}(\mathbb{P}_\theta[\Omega], \mathbb{P}_{\theta'}[\Omega]) - \sum_{a' \neq a} \text{KL}(\theta_{a'}, \theta'_{a'}) \mathbb{E}_\theta[N_T(a')]}{\text{KL}(\theta_a, \theta'_a)} .$$

Hence for all suboptimal arm $a \neq \star_\theta$,

$$\mathbb{E}_\theta[N_T(a)] \geq \sup_{\Omega, \theta'} \frac{\mathbf{k}1(\mathbb{P}_\theta[\Omega], \mathbb{P}_{\tilde{\theta}}[\Omega]) - \sum_{a' \neq a} \mathbf{KL}(\theta_{a'}, \theta'_{a'}) \mathbb{E}_\theta[N_T(a')]}{\mathbf{KL}(\theta_a, \theta'_a)}.$$

Choose θ' such that a is optimal. Let $\Omega = \{N_T(a) > T^\alpha\}$.

▶ $\mathbb{P}_\theta[\Omega] \leq \mathbb{E}_\theta[N_T(a)] T^{-\alpha} = o(1)$ (*Consistency*)

▶ $\sum_{a' \in \mathcal{A}} N_T(a') = T$ (*Construction*)

Thus $\mathbf{k}1(\mathbb{P}_\theta[\Omega], \mathbb{P}_{\tilde{\theta}}[\Omega]) \simeq \ln\left(\frac{1}{\mathbb{P}_{\tilde{\theta}}(N_T(a) \leq T^\alpha)}\right) \geq \ln\left(\frac{T - T^\alpha}{\sum_{a' \neq a} \mathbb{E}_{\tilde{\theta}}[N_T(a')]} \right) \simeq \ln(T)$.

Hence for all suboptimal arm $a \neq \star_\theta$,

$$\mathbb{E}_\theta[N_T(a)] \geq \sup_{\Omega, \theta'} \frac{\mathbf{k}1(\mathbb{P}_\theta[\Omega], \mathbb{P}_{\tilde{\theta}}[\Omega]) - \sum_{a' \neq a} \mathbf{KL}(\theta_{a'}, \theta'_{a'}) \mathbb{E}_\theta[N_T(a')]}{\mathbf{KL}(\theta_a, \theta'_a)}.$$

Choose θ' such that a is optimal. Let $\Omega = \{N_T(a) > T^\alpha\}$.

▶ $\mathbb{P}_\theta[\Omega] \leq \mathbb{E}_\theta[N_T(a)] T^{-\alpha} = o(1)$ (*Consistency*)

▶ $\sum_{a' \in \mathcal{A}} N_T(a') = T$ (*Construction*)

Thus $\mathbf{k}1(\mathbb{P}_\theta[\Omega], \mathbb{P}_{\tilde{\theta}}[\Omega]) \simeq \ln\left(\frac{1}{\mathbb{P}_{\tilde{\theta}}(N_T(a) \leq T^\alpha)}\right) \geq \ln\left(\frac{T - T^\alpha}{\sum_{a' \neq a} \mathbb{E}_{\tilde{\theta}}[N_T(a')]} \right) \simeq \ln(T)$.

▶ **No constraint** on $\theta'_{a'}$ for $a' \neq a$: $\theta'_{a'} = \theta_{a'}$ kills the blue terms.

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_\theta[N_T(a)]}{\ln(T)} \geq \frac{1 - 0}{\inf_{\tilde{\theta}_a} \{\mathbf{KL}(\theta_a, \theta'_a) : \mu'_{a'} > \mu_{\star_\theta}\}}$$

▷ Insight from lower bound: Any *uniformly-good* strategy on \mathcal{D} must satisfy:

$$\forall a \notin \mathcal{A}_*(\nu), \liminf_T \frac{\mathbb{E}[N_T(a)]}{\log(T)} \geq \sup \left\{ \frac{1}{\text{KL}(\nu_a, \tilde{\nu}_a)} : \underbrace{\tilde{\nu} = (\nu_1, \dots, \tilde{\nu}_a, \dots, \nu_A)}_{\text{most confusing (unstructured)}}, \mathcal{A}_*(\tilde{\nu}) = \{a\} \right\}$$

▷ Insight from lower bound: Any *uniformly-good* strategy on \mathcal{D} must satisfy:

$$\forall a \notin \mathcal{A}_*(\nu), \liminf_T \frac{\mathbb{E}[N_T(a)]}{\log(T)} \geq \sup \left\{ \frac{1}{\text{KL}(\nu_a, \tilde{\nu}_a)} : \underbrace{\tilde{\nu} = (\nu_1, \dots, \tilde{\nu}_a, \dots, \nu_A)}_{\text{most confusing (unstructured)}}, \mathcal{A}_*(\tilde{\nu}) = \{a\} \right\}$$

▷ KL-UCB plays arms *not pulled enough* for being *uniformly-good*:

$$a_{t+1} \in \operatorname{argmax}_{a \in \mathcal{A}} \max \left\{ \mathbb{E}_{\tilde{\nu}_a}[X] : N_T(a) \leq \frac{\log(T)}{\text{KL}(\hat{\nu}_{t,a}, \tilde{\nu}_a)}, \tilde{\nu} \text{ most confusing for } a \right\}$$

- ▷ Insight from lower bound: Any *uniformly-good* strategy on \mathcal{D} must satisfy:

$$\forall a \notin \mathcal{A}_*(\nu), \liminf_T \frac{\mathbb{E}[N_T(a)]}{\log(T)} \geq \sup \left\{ \frac{1}{\text{KL}(\nu_a, \tilde{\nu}_a)} : \underbrace{\tilde{\nu} = (\nu_1, \dots, \tilde{\nu}_a, \dots, \nu_A)}_{\text{most confusing (unstructured)}}, \mathcal{A}_*(\tilde{\nu}) = \{a\} \right\}$$

- ▷ KL-UCB plays arms *not pulled enough* for being *uniformly-good*:

$$a_{t+1} \in \operatorname{argmax}_{a \in \mathcal{A}} \max \left\{ \mathbb{E}_{\tilde{\nu}_a}[X] : N_T(a) \leq \frac{\log(T)}{\text{KL}(\hat{\nu}_{t,a}, \tilde{\nu}_a)}, \tilde{\nu} \text{ most confusing for } a \right\}$$

Play an arm in order to
rule-out a most confusing instance
 (Selects one causing maximal regret if not played.)

- ▷ Different from “expecting the best reward in the best world”: testing.

Following the same proof as for the *fundamental Lemma* one can obtain the following generalization:

Lemma (\mathcal{D} -constrained regret lower bound)

Let \mathcal{D} be any set of bandit configurations and $\nu \in \mathcal{D}$. Then any uniformly-good strategy on \mathcal{D} must incur a regret

$$\liminf_{T \rightarrow \infty} \frac{\mathfrak{R}_{T,\nu}}{\ln(T)} \geq \inf \left\{ \sum_{a \in \mathcal{A}} c_a (\mu_{\star}(\nu) - \mu_a(\nu)) : \right. \\ \left. \forall a \in \mathcal{A}, c_a \geq 0, \inf_{\nu' \in \tilde{\mathcal{D}}(\nu)} \sum_{a \in \mathcal{A}} c_a \text{KL}(\nu_a, \nu'_a) \geq 1 \right\}.$$

where we introduced the set of maximally confusing distributions

$$\tilde{\mathcal{D}}(\nu) = \left\{ \nu' \in \mathcal{D} : \mathcal{A}^*(\nu') \cap \mathcal{A}^*(\nu) = \emptyset, \forall a \in \mathcal{A}^*(\nu), \text{KL}(\nu_a, \nu'_a) = 0 \right\}.$$

- ▶ Solution to an *optimization* problem!
- ▶ Specialization to the multi-armed bandit setup of an even more general result from Graves&Lai, 97 (extending Agrawal 89).

Using similar steps as for unstructured lower bounds, we get

$\forall a \notin \mathcal{A}^*(\nu), \forall \nu' \in \mathcal{D}$ s.t. $\mathcal{A}^*(\nu') = \{a\}$

$$\liminf_T \frac{\sum_{a' \in \mathcal{A}} \mathbb{E}[N_T(a')] \text{KL}(\nu_{a'}, \nu'_{a'})}{\ln(T)} \geq \liminf_T \frac{\ln(T - T^\alpha)}{\ln(T)} - \frac{\ln\left(\sum_{a' \neq a} \mathbb{E}_{\nu'}[N_T(a')]\right)}{\ln(T)}$$

Using similar steps as for unstructured lower bounds, we get

$\forall a \notin \mathcal{A}^*(\nu), \forall \nu' \in \mathcal{D}$ s.t. $\mathcal{A}^*(\nu') = \{a\}$

$$\liminf_T \frac{\sum_{a' \in \mathcal{A}} \mathbb{E}[N_T(a')] \text{KL}(\nu_{a'}, \nu'_{a'})}{\ln(T)} \geq \liminf_T \frac{\ln(T - T^\alpha)}{\ln(T)} - \overbrace{\frac{\ln\left(\sum_{a' \neq a} \mathbb{E}_{\nu'}[N_T(a')]\right)}{\ln(T)}}^B$$

By uniformly-good assumption, it must be that $B = 0$, hence

$$\liminf_T \sum_{a' \in \mathcal{A}} \frac{\mathbb{E}[N_T(a')]}{\ln(T)} \text{KL}(\nu_{a'}, \nu'_{a'}) = \sum_{a' \in \mathcal{A}} \left(\liminf_T \frac{\mathbb{E}[N_T(a')]}{\ln(T)} \right) \text{KL}(\nu_{a'}, \nu'_{a'}) \geq 1.$$

This holds in particular choosing ν' such that $\forall a' \in \mathcal{A}^*(\nu), \text{KL}(\nu_{a'}, \nu'_{a'}) = 0$. We conclude by remarking that

$$\liminf_{T \rightarrow \infty} \frac{\mathfrak{R}_T}{\ln(T)} = \sum_{a \in \mathcal{A}} \underbrace{\left(\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_T(a)]}{\ln(T)} \right)}_{c_a} (\mu_{\star}(\nu) - \mu_a(\nu)).$$

What is the number of times a sub-optimal arm needs to be pulled?

The fundamental change of measure argument plus a simple reordering gives

$$\mathbb{E}_\nu[N_T(a)] \geq \sup_{\nu' \in \mathcal{D}} \frac{\sup_{\Omega} \text{kl}(\mathbb{P}_{\tilde{\nu}}[\Omega], \mathbb{P}_\nu[\Omega]) - \sum_{a' \in \mathcal{A} \setminus \{a\}} \mathbb{E}_\nu[N_T(a')] \text{KL}(\nu_{a'}, \nu'_{a'})}{\text{KL}(\nu_a, \nu'_a)}.$$

This motivates the following definition:

Definition (Asymptotic price for uniformly-good strategies)

For $\nu \in \mathcal{D}$, $a \notin \mathcal{A}_*(\nu)$, the asymptotic *price* to pay on arm a for *being uniformly-good* on \mathcal{D} is

$$n_T(a, \nu, \mathcal{D}) = \sup_{\nu' \in \mathcal{D}: a \in \mathcal{A}_*(\nu')} \frac{\ln(T) - \sum_{a' \in \mathcal{A} \setminus \{a\}} \mathbb{E}_\nu[N_T(a')] \text{KL}(\nu_{a'}, \nu'_{a'})}{\text{KL}(\nu_a, \nu'_a)}.$$

▷ *No structure* (*most confusing* obtained without changing other arms):

$$\begin{aligned} \mathbb{E}_\nu[N_T(a)] &\geq \sup_{\tilde{\nu} \in \mathcal{D}: \mathcal{A}_*(\tilde{\nu}) = \{a\}} \left\{ \frac{\ln(T)}{\text{KL}(\nu_a, \tilde{\nu}_a)} : \tilde{\nu} = (\nu_1, \dots, \tilde{\nu}_a, \dots, \nu_A) \right\} \\ &= \frac{\ln(T)}{\mathcal{K}_{\mathcal{D}}(\nu_a, \mu^*(\nu))}. \end{aligned}$$

- ▷ **No structure** (*most confusing* obtained without changing other arms):

$$\begin{aligned} \mathbb{E}_\nu[N_T(a)] &\geq \sup_{\tilde{\nu} \in \mathcal{D}: \mathcal{A}_*(\tilde{\nu}) = \{a\}} \left\{ \frac{\ln(T)}{\text{KL}(\nu_a, \tilde{\nu}_a)} : \tilde{\nu} = (\nu_1, \dots, \tilde{\nu}_a, \dots, \nu_A) \right\} \\ &= \frac{\ln(T)}{\mathcal{K}_{\mathcal{D}}(\nu_a, \mu^*(\nu))}. \end{aligned}$$

- ▷ **Structure** (*most confusing* instance requires changing other arms):

$$\mathbb{E}_\nu[N_T(a)] \geq \sup_{\tilde{\nu} \in \mathcal{D}: \mathcal{A}_*(\tilde{\nu}) = \{a\}} \left\{ \frac{\ln(T) - \sum_{a' \in \mathcal{A} \setminus \{a\}} \mathbb{E}_\nu[N_T(a')] \text{KL}(\nu_{a'}, \tilde{\nu}_{a'})}{\text{KL}(\nu_a, \tilde{\nu}_a)} \right\}.$$

How to adapt bandit strategy to handle such structure (ongoing research)?

(*Collections*) $(\mathcal{A}, (\Theta_a)_{a \in \mathcal{A}}, (\mathcal{Y}_a)_{a \in \mathcal{A}}, (\nu_a)_{a \in \mathcal{A}}, (\mu_a)_{a \in \mathcal{A}})$

(*Structure*) $\Theta \subset \prod_{a \in \mathcal{A}} \Theta_a$

(*Parameter*) $\theta \in \Theta$

Finite set \mathcal{A} . For each $a \in \mathcal{A}$:

$$\begin{aligned}
 (\text{Collections}) \quad & (\mathcal{A}, (\Theta_a)_{a \in \mathcal{A}}, (\mathcal{Y}_a)_{a \in \mathcal{A}}, (\nu_a)_{a \in \mathcal{A}}, (\mu_a)_{a \in \mathcal{A}}) \\
 (\text{Structure}) \quad & \Theta \subset \prod_{a \in \mathcal{A}} \Theta_a \\
 (\text{Parameter}) \quad & \theta \in \Theta
 \end{aligned}$$

Finite set \mathcal{A} . For each $a \in \mathcal{A}$:

- ▶ Parameter space Θ_a .

$$\begin{aligned} (\text{Collections}) \quad & (\mathcal{A}, (\Theta_a)_{a \in \mathcal{A}}, (\mathcal{Y}_a)_{a \in \mathcal{A}}, (\nu_a)_{a \in \mathcal{A}}, (\mu_a)_{a \in \mathcal{A}}) \\ (\text{Structure}) \quad & \Theta \subset \prod_{a \in \mathcal{A}} \Theta_a \\ (\text{Parameter}) \quad & \theta \in \Theta \end{aligned}$$

Finite set \mathcal{A} . For each $a \in \mathcal{A}$:

- ▶ Parameter space Θ_a .
- ▶ Observation space \mathcal{Y}_a .

$$\begin{aligned}
 (\text{Collections}) \quad & (\mathcal{A}, (\Theta_a)_{a \in \mathcal{A}}, (\mathcal{Y}_a)_{a \in \mathcal{A}}, (\nu_a)_{a \in \mathcal{A}}, (\mu_a)_{a \in \mathcal{A}}) \\
 (\text{Structure}) \quad & \Theta \subset \prod_{a \in \mathcal{A}} \Theta_a \\
 (\text{Parameter}) \quad & \theta \in \Theta
 \end{aligned}$$

Finite set \mathcal{A} . For each $a \in \mathcal{A}$:

- ▶ Parameter space Θ_a .
- ▶ Observation space \mathcal{Y}_a .
- ▶ Distribution of observations $\nu_a : \Theta_a \rightarrow \mathcal{P}(\mathcal{Y}_a)$

- (*Collections*) $(\mathcal{A}, (\Theta_a)_{a \in \mathcal{A}}, (\mathcal{Y}_a)_{a \in \mathcal{A}}, (\nu_a)_{a \in \mathcal{A}}, (\mu_a)_{a \in \mathcal{A}})$
- (*Structure*) $\Theta \subset \prod_{a \in \mathcal{A}} \Theta_a$
- (*Parameter*) $\theta \in \Theta$

Finite set \mathcal{A} . For each $a \in \mathcal{A}$:

- ▶ Parameter space Θ_a .
- ▶ Observation space \mathcal{Y}_a .
- ▶ Distribution of observations $\nu_a : \Theta_a \rightarrow \mathcal{P}(\mathcal{Y}_a)$
- ▶ Reward: $\mu_a : \Theta \rightarrow \mathbb{R}$ (Θ and not Θ_a !)

- ▶ **Classical Bernoulli MAB:** $\mathcal{A} = \{1, \dots, A\}$, $\Theta_a = [0, 1]$, $\mathcal{Y}_a = \{0, 1\}$, $\nu_a(\theta_a) = \text{Bern}(\theta_a)$, $\Theta = [0, 1]^A$ (unstructured) and $\mu_a(\theta) = \theta_a$.

- ▶ **Classical Bernoulli MAB:** $\mathcal{A} = \{1, \dots, A\}$, $\Theta_a = [0, 1]$, $\mathcal{Y}_a = \{0, 1\}$, $\nu_a(\theta_a) = \text{Bern}(\theta_a)$, $\Theta = [0, 1]^A$ (unstructured) and $\mu_a(\theta) = \theta_a$.
- ▶ **Linear bandits:** $\mathcal{A} \subset \mathbb{R}^d$, $\Theta_a = \{\langle \alpha, a \rangle : \alpha \in \mathbb{R}^d\}$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, 1)$, $\Theta = \{\theta = (\langle \alpha, a \rangle)_{a \in \mathcal{A}}, \alpha \in \mathbb{R}^d\}$, $\mu_a(\theta) = \theta_a$.

- ▶ **Classical Bernoulli MAB:** $\mathcal{A} = \{1, \dots, A\}$, $\Theta_a = [0, 1]$, $\mathcal{Y}_a = \{0, 1\}$, $\nu_a(\theta_a) = \text{Bern}(\theta_a)$, $\Theta = [0, 1]^A$ (unstructured) and $\mu_a(\theta) = \theta_a$.
- ▶ **Linear bandits:** $\mathcal{A} \subset \mathbb{R}^d$, $\Theta_a = \{\langle \alpha, a \rangle : \alpha \in \mathbb{R}^d\}$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, 1)$, $\Theta = \{\theta = (\langle \alpha, a \rangle)_{a \in \mathcal{A}}, \alpha \in \mathbb{R}^d\}$, $\mu_a(\theta) = \theta_a$.
- ▶ **Lipschitz bandits:** $\mathcal{A} \subset \mathcal{X}$, $\Theta_a \subset \mathbb{R}$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, 1)$, $\Theta = \{\theta : \max_{a, a' \in \mathcal{A}} \frac{|\theta_a - \theta_{a'}|}{\ell(a, a')} \leq 1\}$, $\mu_a(\theta) = \theta_a$.

- ▶ **Classical Bernoulli MAB:** $\mathcal{A} = \{1, \dots, A\}$, $\Theta_a = [0, 1]$, $\mathcal{Y}_a = \{0, 1\}$, $\nu_a(\theta_a) = \text{Bern}(\theta_a)$, $\Theta = [0, 1]^A$ (unstructured) and $\mu_a(\theta) = \theta_a$.
- ▶ **Linear bandits:** $\mathcal{A} \subset \mathbb{R}^d$, $\Theta_a = \{\langle \alpha, a \rangle : \alpha \in \mathbb{R}^d\}$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, 1)$, $\Theta = \{\theta = (\langle \alpha, a \rangle)_{a \in \mathcal{A}}, \alpha \in \mathbb{R}^d\}$, $\mu_a(\theta) = \theta_a$.
- ▶ **Lipschitz bandits:** $\mathcal{A} \subset \mathcal{X}$, $\Theta_a \subset \mathbb{R}$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, 1)$, $\Theta = \{\theta : \max_{a, a' \in \mathcal{A}} \frac{|\theta_a - \theta_{a'}|}{\ell(a, a')} \leq 1\}$, $\mu_a(\theta) = \theta_a$.
- ▶ **Combinatorial semi-bandit:** $\mathcal{A} \subset \{0, 1\}^d$, $\Theta_a \subset \mathbb{R}^d$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, I_d)$, $\Theta = \{\theta : \theta_a = (\alpha_1 a_1, \dots, \alpha_d a_d), \alpha \in \mathbb{R}^d\}$, $\mu_a(\theta) = \langle \theta_a, \mathbf{1} \rangle$.

- ▶ **Classical Bernoulli MAB:** $\mathcal{A} = \{1, \dots, A\}$, $\Theta_a = [0, 1]$, $\mathcal{Y}_a = \{0, 1\}$, $\nu_a(\theta_a) = \text{Bern}(\theta_a)$, $\Theta = [0, 1]^A$ (unstructured) and $\mu_a(\theta) = \theta_a$.
- ▶ **Linear bandits:** $\mathcal{A} \subset \mathbb{R}^d$, $\Theta_a = \{\langle \alpha, a \rangle : \alpha \in \mathbb{R}^d\}$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, 1)$, $\Theta = \{\theta = (\langle \alpha, a \rangle)_{a \in \mathcal{A}}, \alpha \in \mathbb{R}^d\}$, $\mu_a(\theta) = \theta_a$.
- ▶ **Lipschitz bandits:** $\mathcal{A} \subset \mathcal{X}$, $\Theta_a \subset \mathbb{R}$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, 1)$, $\Theta = \{\theta : \max_{a, a' \in \mathcal{A}} \frac{|\theta_a - \theta_{a'}|}{\ell(a, a')} \leq 1\}$, $\mu_a(\theta) = \theta_a$.
- ▶ **Combinatorial semi-bandit:** $\mathcal{A} \subset \{0, 1\}^d$, $\Theta_a \subset \mathbb{R}^d$, $\mathcal{Y}_a = \mathbb{R}$, $\nu_a(\theta_a) = \mathcal{N}(\theta_a, I_d)$, $\Theta = \{\theta : \theta_a = (\alpha_1 a_1, \dots, \alpha_d a_d), \alpha \in \mathbb{R}^d\}$, $\mu_a(\theta) = \langle \theta_a, \mathbf{1} \rangle$.
- ▶ **Ranking bandits:** $\mathcal{A} = \{a \in \text{Arr}_N^L\}$, $\Theta_a = [0, 1]^L$, $\mathcal{Y}_a = \{0, 1\}$, $\nu_a(\theta_a) = \text{Fct}\left(\left(\text{Bern}(\theta_{a_\ell})\right)_{\ell \leq L}\right)$, $\Theta = \{\theta : \theta_a = (\alpha_{a_\ell})_{\ell \leq L}, \alpha \in [0, 1]^N\}$, $\mu_a(\theta) = \sum_{\ell=1}^L r(\ell) \theta_{a_\ell} \prod_{i=1}^{\ell-1} (1 - \theta_{a_i})$.

Theorem (Agrawal 1989)

Assume Θ is discrete, $\star(\theta) = \text{Argmax}_{a \in \mathcal{A}} \mu_a(\theta)$ is unique. Then for any uniformly good strategy,

$$\liminf_{T \rightarrow \infty} \frac{R_T(\theta)}{\ln(T)} \geq C(\theta) \quad \text{where}$$

$$C(\theta) = \min \left\{ \frac{\sum_{a \in \mathcal{A} \setminus \star(\theta)} \eta_a (\mu_{\star(\theta)} - \mu_a(\theta))}{\inf_{\lambda \in \Lambda(\theta)} \sum_{a \in \mathcal{A} \setminus \star(\theta)} \eta_a \text{KL}(\nu_a(\theta_a), \nu_a(\lambda_a))} : \eta \in \mathcal{P}(\mathcal{A} \setminus \star(\theta)) \right\}$$

with $\Lambda(\theta) = \{ \lambda \in \Theta : \star(\theta) \neq \star(\lambda), \text{ and } \text{KL}(\nu_a(\theta_a), \nu_a(\lambda_a)) = 0 \text{ for } a = \star(\theta) \}$.

Theorem (Agrawal 1989)

Assume Θ is discrete, $\star(\theta) = \text{Argmax}_{a \in \mathcal{A}} \mu_a(\theta)$ is unique. Then for any uniformly good strategy,

$$\liminf_{T \rightarrow \infty} \frac{R_T(\theta)}{\ln(T)} \geq C(\theta) \quad \text{where}$$

$$C(\theta) = \min \left\{ \frac{\sum_{a \in \mathcal{A} \setminus \star(\theta)} \eta_a (\mu_{\star(\theta)} - \mu_a)}{\inf_{\lambda \in \Lambda(\theta)} \sum_{a \in \mathcal{A} \setminus \star(\theta)} \eta_a \text{KL}(\nu_a(\theta_a), \nu_a(\lambda_a))} : \eta \in \mathcal{P}(\mathcal{A} \setminus \star(\theta)) \right\}$$

with $\Lambda(\theta) = \{\lambda \in \Theta : \star(\theta) \neq \star(\lambda), \text{ and } \text{KL}(\nu_a(\theta_a), \nu_a(\lambda_a)) = 0 \text{ for } a = \star(\theta)\}$.

- Confusing parameters *statistically indistinguishable* from θ when playing only $\star(\theta)$.

Theorem (Graves, Lai 1997)

Assume $\star(\theta) = \operatorname{Argmax}_{a \in \mathcal{A}} \mu_a(\theta)$ is unique. Then for any uniformly good strategy,

$$\liminf_{T \rightarrow \infty} \frac{R_T(\theta)}{\ln(T)} \geq C(\theta) \quad \text{where}$$

$$C(\theta) = \min \left\{ \sum_{a \in \mathcal{A}} n_a (\mu_{\star}(\theta) - \mu_a(\theta)) : \forall a, n_a \geq 0 \right. \\ \left. \text{and } \inf_{\lambda \in \Lambda(\theta)} \sum_{a \in \mathcal{A}} n_a \operatorname{KL}(\nu_a(\theta_a), \nu_a(\lambda_a)) \geq 1 \right\}$$

with $\Lambda(\theta) = \left\{ \lambda \in \Theta : \star(\theta) \neq \star(\lambda), \text{ and } \operatorname{KL}(\nu_a(\theta_a), \nu_a(\lambda_a)) = 0 \text{ for } a = \star(\theta) \right\}$.

Theorem (Graves, Lai 1997)

Assume $\star(\theta) = \text{Argmax}_{a \in \mathcal{A}} \mu_a(\theta)$ is unique. Then for any uniformly good strategy,

$$\liminf_{T \rightarrow \infty} \frac{R_T(\theta)}{\ln(T)} \geq C(\theta) \quad \text{where}$$

$$C(\theta) = \min \left\{ \sum_{a \in \mathcal{A}} n_a (\mu_{\star}(\theta) - \mu_a(\theta)) : \forall a, n_a \geq 0 \right. \\ \left. \text{and } \inf_{\lambda \in \Lambda(\theta)} \sum_{a \in \mathcal{A}} n_a \text{KL}(\nu_a(\theta_a), \nu_a(\lambda_a)) \geq 1 \right\}$$

with $\Lambda(\theta) = \{ \lambda \in \Theta : \star(\theta) \neq \star(\lambda), \text{ and } \text{KL}(\nu_a(\theta_a), \nu_a(\lambda_a)) = 0 \text{ for } a = \star(\theta) \}$.

- Confusing parameters *statistically indistinguishable* from θ when playing only $\star(\theta)$.

STRUCTURES

LINEAR BANDITS

STRUCTURED LOWER BOUNDS

Lower bounds

Lipschitz bandits

Ranking bandits

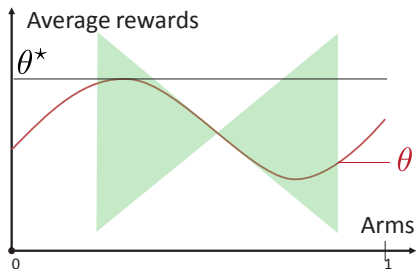
Metric-graph of bandits

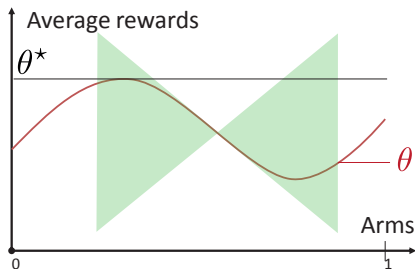
CONCLUSION, PERSPECTIVE

Lipschitz Bandits: Regret Lower Bounds and Optimal Algorithms

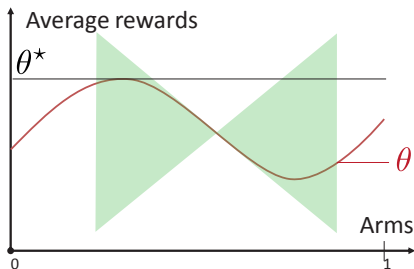
Stefan Magureanu, Richard Combes and Alexandre Proutiere, COLT 2014.

LIPSCHITZ BANDITS - PROBLEM DESCRIPTION

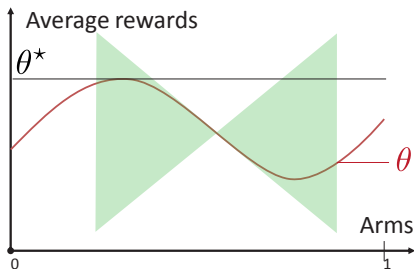




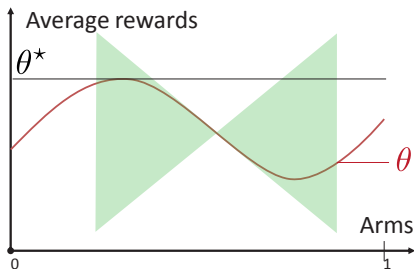
- ▶ The decision maker is given a constant L



- ▶ The decision maker is given a constant L
- ▶ Each $k \in \mathcal{K}$, is assigned a fixed and known coordinate $x_k \in (0, 1)$



- ▶ The decision maker is given a constant L
- ▶ Each $k \in \mathcal{K}$, is assigned a fixed and known coordinate $x_k \in (0, 1)$
- ▶ Then : $\Theta_L = \{\theta \in (0, 1)^K : |\theta_i - \theta_j| \leq L|x_i - x_j|, \forall i, j \leq K\}$



- ▶ The decision maker is given a constant L
- ▶ Each $k \in \mathcal{K}$, is assigned a fixed and known coordinate $x_k \in (0, 1)$
- ▶ Then : $\Theta_L = \{\theta \in (0, 1)^K : |\theta_i - \theta_j| \leq L|x_i - x_j|, \forall i, j \leq K\}$
- ▶ Our goal is to exploit this additional information in order to reduce the achievable regret, relative to that of the classic setting

When $\{x_k : k \in \mathcal{K}\} = (0, 1)$ an efficient algorithm must perform two task:

When $\{x_k : k \in \mathcal{K}\} = (0, 1)$ an efficient algorithm must perform two task:

- ▶ Adaptive discretization (from continuous \mathcal{X} to discrete \mathbb{X})?

When $\{x_k : k \in \mathcal{K}\} = (0, 1)$ an efficient algorithm must perform two task:

- ▶ Adaptive discretization (from continuous \mathcal{X} to discrete \mathbb{X})?
- ▶ Efficient statistical testing:

When $\{x_k : k \in \mathcal{K}\} = (0, 1)$ an efficient algorithm must perform two task:

- ▶ Adaptive discretization (from continuous \mathcal{X} to discrete \mathbb{X})?
- ▶ Efficient statistical testing:
 - ▶ Correctly identify the suboptimal arms by optimally exploiting past observations and structure

When $\{x_k : k \in \mathcal{K}\} = (0, 1)$ an efficient algorithm must perform two task:

- ▶ Adaptive discretization (from continuous \mathcal{X} to discrete \mathbb{X})?
- ▶ Efficient statistical testing:
 - ▶ Correctly identify the suboptimal arms by optimally exploiting past observations and structure
 - ▶ Perform this task optimally: regret lower bounds? algorithms matching this limit? (Magureanu et al., COLT 2014)

LIPSCHITZ BANDITS - REGRET LOWER BOUNDS (PRELIMINARIES)



Let us define the most confusing *bad* parameter λ^k of an arm k :

$$\lambda_j^k = \max(\theta_j, \theta^* - L \times |x_j - x_k|), \forall j \in \mathcal{K}$$

Theorem (Lower bound)

For all $\theta \in \Theta_L$ and uniformly good algorithms π , we have:

$$\liminf_{T \rightarrow \infty} \frac{R^\pi(T)}{\ln(T)} \geq C(\theta)$$

where $C(\theta)$ is the minimal value of the following optimization problem:

$$\begin{aligned} & \min_{c_k > 0; k \in \mathcal{K}^-} \sum_{k \in \mathcal{K}^-} c_k (\theta^* - \theta_k) \\ \text{subject to: } & \sum_{k' \in \mathcal{K}^-} c_{k'} \text{KL}(\theta_{k'}, \lambda_{\theta^*, k'}^k) \geq 1, \quad \forall k \in \mathcal{K}^- \end{aligned}$$

Theorem (Lower bound)

For all $\theta \in \Theta_L$ and uniformly good algorithms π , we have:

$$\liminf_{T \rightarrow \infty} \frac{R^\pi(T)}{\ln(T)} \geq C(\theta)$$

where $C(\theta)$ is the minimal value of the following optimization problem:

$$\begin{aligned} & \min_{c_k > 0; k \in \mathcal{K}^-} \sum_{k \in \mathcal{K}^-} c_k (\theta^* - \theta_k) \\ \text{subject to: } & \sum_{k' \in \mathcal{K}^-} c_{k'} \text{KL}(\theta_{k'}, \lambda_{\theta^*, k'}^k) \geq 1, \quad \forall k \in \mathcal{K}^- \end{aligned}$$

- ▶ Follows result by Graves, Todd L., and Tze Leung Lai. *"Asymptotically efficient adaptive choice of control laws in controlled markov chains."* SIAM journal on control and optimization 35.3 (1997): 715-743

Two algorithms are proposed:

Two algorithms are proposed:

- ▶ OSLB :

Two algorithms are proposed:

- ▶ OSLB :
 - ▶ Asymptotically optimal

Two algorithms are proposed:

- ▶ OSLB :
 - ▶ Asymptotically optimal
 - ▶ Enforces exploration as dictated by the LP in the lower bound

Two algorithms are proposed:

- ▶ OSLB :
 - ▶ Asymptotically optimal
 - ▶ Enforces exploration as dictated by the LP in the lower bound
 - ▶ Computationally complex and performs poorly numerically

Two algorithms are proposed:

- ▶ OSLB :
 - ▶ Asymptotically optimal
 - ▶ Enforces exploration as dictated by the LP in the lower bound
 - ▶ Computationally complex and performs poorly numerically
- ▶ POSLB:

Two algorithms are proposed:

- ▶ OSLB :
 - ▶ Asymptotically optimal
 - ▶ Enforces exploration as dictated by the LP in the lower bound
 - ▶ Computationally complex and performs poorly numerically
- ▶ POSLB:
 - ▶ Asymptotically Pareto-optimal - provably exploits the structure efficiently

Two algorithms are proposed:

- ▶ OSLB :
 - ▶ Asymptotically optimal
 - ▶ Enforces exploration as dictated by the LP in the lower bound
 - ▶ Computationally complex and performs poorly numerically
- ▶ POSLB:
 - ▶ Asymptotically Pareto-optimal - provably exploits the structure efficiently
 - ▶ Computationally light and work well numerically

Two algorithms are proposed:

- ▶ OSLB :
 - ▶ Asymptotically optimal
 - ▶ Enforces exploration as dictated by the LP in the lower bound
 - ▶ Computationally complex and performs poorly numerically
- ▶ POSLB:
 - ▶ Asymptotically Pareto-optimal - provably exploits the structure efficiently
 - ▶ Computationally light and work well numerically
 - ▶ Related to the *UCB* family of algorithms

- ▶ Both algorithms make use of the following index:

$$b_k(n) = \sup \left\{ q \in (\hat{\theta}_k(n), 1) : \sum_{j \in \mathcal{K}} N_j(n) \text{KL}_+(\hat{\theta}_j(n), \lambda_j^{q,k}) \leq f(n) \right\}$$

where $f(n) = \ln(n) + 3K \ln \ln(n)$ and $\text{KL}_+(x, y) = \text{KL}(x, y)$ if $x < y$, and 0 otherwise

ALGORITHM OSLB(ϵ)

- ▶ At each round, OSLB(ε) computes $\hat{c}(n) = c(\hat{\theta}(n))$ - the solution to the LP in the lower bound with θ *replaced by the empirical mean* $\hat{\theta}(n)$

ALGORITHM OSLB(ε)

- ▶ At each round, OSLB(ε) computes $\hat{c}(n) = c(\hat{\theta}(n))$ - the solution to the LP in the lower bound with θ *replaced by the empirical mean* $\hat{\theta}(n)$
- ▶ Let $L(n) = \arg \max_k \hat{\theta}_k(n)$ be the *leader* at round n

ALGORITHM OSLB(ε)

- ▶ At each round, OSLB(ε) computes $\hat{c}(n) = c(\hat{\theta}(n))$ - the solution to the LP in the lower bound with θ *replaced by the empirical mean* $\hat{\theta}(n)$
- ▶ Let $L(n) = \arg \max_k \hat{\theta}_k(n)$ be the *leader* at round n
- ▶ Let $\underline{k}(n) = \arg \min_k N_k(n)$ be the *least played* arm up to time n

- ▶ At each round, OSLB(ε) computes $\hat{c}(n) = c(\hat{\theta}(n))$ - the solution to the LP in the lower bound with θ *replaced by the empirical mean* $\hat{\theta}(n)$
- ▶ Let $L(n) = \arg \max_k \hat{\theta}_k(n)$ be the *leader* at round n
- ▶ Let $\underline{k}(n) = \arg \min_k N_k(n)$ be the *least played* arm up to time n
- ▶ Let $\bar{k}(n) = \arg \min\{N_k(n) : k : \hat{c}_k(n) > N_k(n)/\ln(n)\}$ be the least played arm among the arms played insufficiently many times

Algorithm 1 OSLB(ε)

For all $n \geq 1$, select arm $k(n)$ such that:

If $\hat{\theta}^*(n) \geq \max_{k \neq L(n)} b_k(n)$, then $k(n) = L(n)$;

Else If $N_{\underline{k}(n)}(n) < \frac{\varepsilon}{K} N_{\bar{k}(n)}(n)$, then $k(n) = \underline{k}(n)$; (*Forced Exploration*)

Else $k(n) = \bar{k}(n)$.

Assumption

Assumption

- ▶ The solution of the LP in the lower bound is unique.

Theorem (asymptotic optimality)

For all $\varepsilon > 0$, under the above assumption, the regret achieved under $\pi = \text{OSLB}(\varepsilon)$ satisfies: for all $\theta \in \Theta_L$, for all $\delta > 0$ and $T \geq 1$,

$$R^\pi(T) \leq C^\delta(\theta)(1 + \varepsilon) \ln(T) + C_1 \ln \ln(T) + K^3 \varepsilon^{-1} \delta^{-2} + 3K \delta^{-2}, \quad (3)$$

where $C^\delta(\theta) \rightarrow C(\theta)$, as $\delta \rightarrow 0^+$, and $C_1 > 0$.

- ▶ OSLB(ε) is *computationally expensive* and performs poorly in practice

- ▶ OSLB(ε) is *computationally expensive* and performs poorly in practice
- ▶ Computationally cheaper algorithm: POSLB

- ▶ OSLB(ε) is *computationally expensive* and performs poorly in practice
- ▶ Computationally cheaper algorithm: POSLB
- ▶ POSLB is inspired from the family of *UCB* algorithms

- ▶ OSLB(ε) is *computationally expensive* and performs poorly in practice
- ▶ Computationally cheaper algorithm: POSLB
- ▶ POSLB is inspired from the family of *UCB* algorithms
- ▶ While not optimal it is Pareto optimal :

- ▶ OSLB(ε) is *computationally expensive* and performs poorly in practice
- ▶ Computationally cheaper algorithm: POSLB
- ▶ POSLB is inspired from the family of *UCB* algorithms
- ▶ While not optimal it is Pareto optimal :
 - ▶ Considering $c_k = N_k(T)/\ln(T)$ yields equalities in all constraints in the lower bound LP

Algorithm 2 POSLB

For all $n \geq 1$, select arm $k(n)$ such that:

Algorithm 3 POSLB

For all $n \geq 1$, select arm $k(n)$ such that:

$$q(n) = b_{L(n)}(n);$$

$$k(n) = \arg \max_k f(n) - f_k(n, q(n)) \text{ (ties are broken arbitrarily)}$$

Algorithm 4 POSLB

For all $n \geq 1$, select arm $k(n)$ such that:

$$q(n) = b_{L(n)}(n);$$

$$k(n) = \arg \max_k f(n) - f_k(n, q(n)) \text{ (ties are broken arbitrarily)}$$

$$\text{where } f_k(n, q(n)) = \begin{cases} \sum_{j \in \mathcal{K}} N_j(n) \text{KL}(\hat{\theta}_j(n), \lambda_j^{q(n), k}(n)) & \text{if } k \neq L(n) \\ N_k(n) \text{KL}(\hat{\theta}_k(n), q(n)) & \text{if } k = L(n) \end{cases}.$$

$$\text{and } \lambda_j^{q, k}(n) = \max(q - |k - j|L, \hat{\theta}_j(n)).$$

Theorem (POSLB pulls and pareto optimality)

Under POSLB, for all $\theta \in \Theta_L$, all $T \geq 1$, all $0 < \delta < (\theta^* - \max_{k \neq k^*} \theta_k)/2$, and any suboptimal arm $k \in \mathcal{K}^-$:

$$\mathbb{E}[N_k(T)] \leq \frac{f(T)}{I(\theta_k + \delta, \theta^* - \delta)} + C_1 \ln(\ln(T)) + 2\delta^{-2}.$$

with $C_1 \geq 0$ a constant. Further, under POSLB, for all $\theta \in \Theta_L$ and $k \in \mathcal{K}^-$, we have that:

$$\lim_{T \rightarrow \infty} \frac{\mathbb{E} \left[\sum_{i \in \mathcal{K}^-} N_i(T) \text{KL}_+(\theta_i, \lambda_i^{\theta^*, k}) \right]}{f(T)} = 1.$$

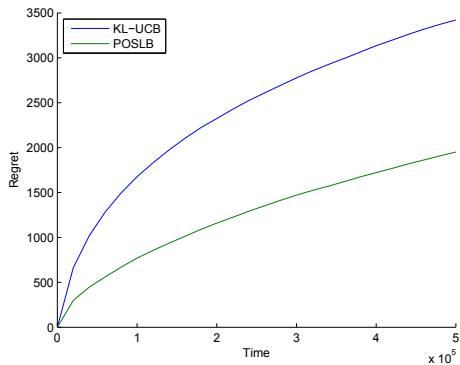
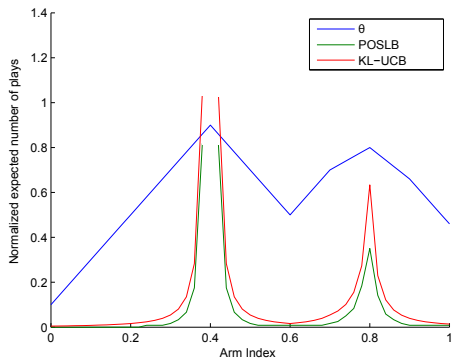


Figure: (Left) The expected rewards and the scaled amount of times suboptimal arms are played under KL-UCB and POSLB as a function of the arm. (Right) Regret under KL-UCB and POSLB as a function of time.

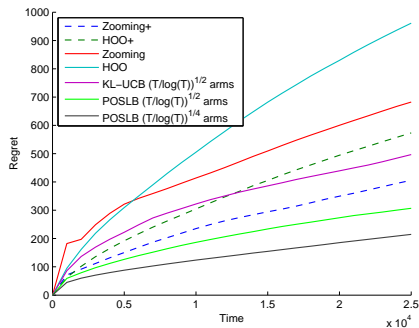
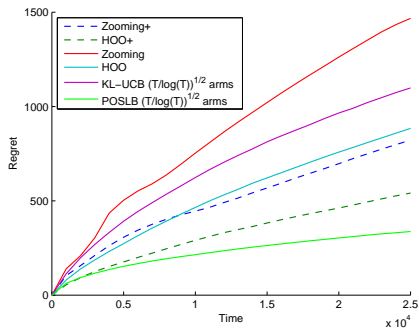


Figure: Expected regret of different algorithms as function of time for a triangular reward function (left) and a quadratic reward function (right).

- ▶ Lower-bound based index that efficiently exploits structure

- ▶ Lower-bound based index that efficiently exploits structure
- ▶ Two algorithms:

- ▶ Lower-bound based index that efficiently exploits structure
- ▶ Two algorithms:
 - ▶ OSLB - asymptotically optimal but complex

- ▶ Lower-bound based index that efficiently exploits structure
- ▶ Two algorithms:
 - ▶ OSLB - asymptotically optimal but complex
 - ▶ POSLB - Pareto-optimal algorithm inspired by the classical UCB

- ▶ Lower-bound based index that efficiently exploits structure
- ▶ Two algorithms:
 - ▶ OSLB - asymptotically optimal but complex
 - ▶ POSLB - Pareto-optimal algorithm inspired by the classical UCB
- ▶ Stepping stone for exploiting structure in generic settings, with more practical applications
- ▶ Tentative generalization to arbitrary structure: OSSB, POSSB (Magureanu 2018, PHD).

STRUCTURES

LINEAR BANDITS

STRUCTURED LOWER BOUNDS

Lower bounds

Lipschitz bandits

Ranking bandits

Metric-graph of bandits

CONCLUSION, PERSPECTIVE

Position in Induced Exploration

Learning to rank: Regret lower bounds and efficient algorithms R Combes, S Magureanu, A Proutiere, C Laroche ACM SIGMETRICS Performance Evaluation Review 43

LEARNING TO RANK : A BANDOT APPROACH

Showing results for still alive.

ARTISTS



Still Alive



Still Alive



Somehow Still Alive



Still Lives



Still Alive
BIGBANG



Still Alive (The Theme from...)



Special Edition 'Still...'

SEE ALL ALBUMS

SEE ALL

SONG

ARTIST

ALBUM



	SONG	ARTIST	ALBUM		
+	Still Alive	BIGBANG	BIGBANG Special Edition Still Alive 1	3:19	
▶	Still Alive	Aperture Science Psychoacoustic Laborat...	Portal 2: Songs to Test By (Collectors Editi...	2:57	
+	Still Alive	Lisa Miskovsky	Mirror's Edge Original Videogame Score	4:34	
+	STILL ALIVE	BIGBANG	Special Edition 'Still Alive'	3:19	
+	Still Alive	The Crash	Melodrama	4:05	
+	Still Alive	Social Distortion	Hard Times And Nursery Rhymes (Deluxe ...	4:06	
+	Still Alive	Nocturnal Rites	Grand Illusion	4:03	
+	Still Alive	Onlap, Charline Max	The Awakening	4:05	
+	Still Alive	Jonathan Coulton	Best. Concert. Ever.	3:05	

Sequential Ranking setup

Sequential Ranking setup

- ▶ N (huge) many given articles

Sequential Ranking setup

- ▶ N (huge) many given articles
- ▶ At each $t = 1, \dots$, a user u_t appears. Choose to display L (ordered) articles.

Sequential Ranking setup

- ▶ N (huge) many given articles
- ▶ At each $t = 1, \dots$, a user u_t appears. Choose to display L (ordered) articles.
- ▶ The user inspects the articles, in order, and clicks on the first *interesting* article then leaves.

Sequential Ranking setup

- ▶ N (huge) many given articles
- ▶ At each $t = 1, \dots$, a user u_t appears. Choose to display L (ordered) articles.
- ▶ The user inspects the articles, in order, and clicks on the first *interesting* article then leaves.
- ▶ The decision maker observes which article was clicked and collects a reward.

- ▶ **Actions:** all combinations of L out of N articles $\mathcal{A} = \{a \in \text{Arr}_N^L\}$

- ▶ **Actions:** all combinations of L out of N articles $\mathcal{A} = \{a \in \text{Arr}_N^L\}$
- ▶ **Feedback X_k for an inspected article k :**
 - 1 if clicked, 0 otherwise; Bernoulli $\mathcal{B}(\theta_k)$

- ▶ **Actions:** all combinations of L out of N articles $\mathcal{A} = \{a \in \text{Arr}_N^L\}$
 - ▶ **Feedback X_k for an inspected article k :**
 - 1 if clicked, 0 otherwise; Bernoulli $\mathcal{B}(\theta_k)$
 - ▶ **Feedback for L displayed articles:**
 - the slot of the clicked article ℓ
 - 0 for each article before ℓ , 1 for the clicked article, *nothing* else
- Click probability on item ℓ in list a : $\theta_{a_\ell} \prod_{i=1}^{\ell} (1 - \theta_{a_i})$.

- ▶ **Actions:** all combinations of L out of N articles $\mathcal{A} = \{a \in \text{Arr}_N^L\}$
- ▶ **Feedback X_k for an inspected article k :**
 - 1 if clicked, 0 otherwise; Bernoulli $\mathcal{B}(\theta_k)$
- ▶ **Feedback for L displayed articles:**
 - the slot of the clicked article ℓ
 - 0 for each article before ℓ , 1 for the clicked article, *nothing* else

Click probability on item ℓ in list a : $\theta_{a_\ell} \prod_{i=1}^{\ell} (1 - \theta_{a_i})$.
- ▶ **Rewards:** $r(\ell)$ - usually decreasing in ℓ .

$$\mu_a(\theta) = \sum_{\ell=1}^L r(\ell) \theta_{a_\ell} \prod_{i=1}^{\ell} (1 - \theta_{a_i}).$$

- ▶ **Actions:** all combinations of L out of N articles $\mathcal{A} = \{a \in \text{Arr}_N^L\}$
- ▶ **Feedback X_k for an inspected article k :**
 - 1 if clicked, 0 otherwise; Bernoulli $\mathcal{B}(\theta_k)$
- ▶ **Feedback for L displayed articles:**
 - the slot of the clicked article ℓ
 - 0 for each article before ℓ , 1 for the clicked article, *nothing* else

Click probability on item ℓ in list a : $\theta_{a_\ell} \prod_{i=1}^{\ell} (1 - \theta_{a_i})$.
- ▶ **Rewards:** $r(\ell)$ - usually decreasing in ℓ .

$$\mu_a(\theta) = \sum_{\ell=1}^L r(\ell) \theta_{a_\ell} \prod_{i=1}^{\ell} (1 - \theta_{a_i}).$$

- ▶ **Goal:** Maximize the cumulative reward over T rounds

$$\mathcal{R}_\theta(T) = T \max_a \mu_a(\theta) - \sum_{t=1}^T \mu_{a_t}(\theta)$$

- ▶ **The set of actions:** Huge $|\mathcal{A}| = N!/(N - L)!$

- ▶ **The set of actions:** Huge $|\mathcal{A}| = N!/(N - L)!$
- ▶ **Feedback for an inspected article:** Random number of observations - depending on the rewards of articles displayed

So?

- ▶ **The set of actions:** Huge $|\mathcal{A}| = N!/(N - L)!$
- ▶ **Feedback for an inspected article:** Random number of observations - depending on the rewards of articles displayed

So?

- ▶ **The set of actions:** We can exploit *structure* to drastically reduce the cost of exploration

- ▶ **The set of actions:** Huge $|\mathcal{A}| = N!/(N - L)!$
- ▶ **Feedback for an inspected article:** Random number of observations - depending on the rewards of articles displayed

So?

- ▶ **The set of actions:** We can exploit *structure* to drastically reduce the cost of exploration
- ▶ **Feedback for an inspected article:** How we explore matters

"Structure":

"Structure":

- ▶ Similarities between users

”Structure”:

- ▶ Similarities between users
- ▶ Similarities between articles

"Structure":

- ▶ Similarities between users
- ▶ Similarities between articles
- ▶ Shape of reward function $r(I)$

Different systems according to the structure that is revealed to the decision maker

Assume $\theta_1 > \theta_2 > \dots > \theta_N$ (item 1 is preferred over 2, etc.)

Let $\Delta_i = r(i) - r(i+1)$, $\Delta_L = r(L)$ and $N_a(t)$ the number of times the set a of articles is displayed until time t

Regret lower bound

If $\Delta_i > \Delta_L > 0$ for all $i < L$, then

$$\liminf_{T \rightarrow \infty} \frac{N_a(T)}{\ln(T)} = \frac{\mathbb{I}\{\exists i : a = \{1, \dots, L-1, i\}\}}{\text{KL}(\mathcal{B}(\theta_i), \mathcal{B}(\theta_L)) \prod_{j < L} (1 - \theta_j)}$$

$$\liminf_{T \rightarrow \infty} \frac{R_{\theta}^{\pi}(T)}{\ln(T)} = r(L) \sum_{i=L+1}^N \frac{\theta_L - \theta_i}{\text{KL}(\mathcal{B}(\theta_i), \mathcal{B}(\theta_L))}$$

\Rightarrow Suggest *exploration* at *last* slot L .

Assume $\theta_1 > \theta_2 > \dots > \theta_N$ (item 1 is preferred over 2, etc.)

Let $\Delta_i = r(i) - r(i+1)$, $\Delta_L = r(L)$ and $N_a(t)$ the number of times the set a of articles is displayed until time t

Regret lower bound

If $r(i) = r(L) > 0$ for all $i < L$:

$$\liminf_{T \rightarrow \infty} \frac{N_a(T)}{\ln(T)} = \frac{\mathbb{I}\{\exists i : u = \{i, 1, \dots, L-1\}\}}{\text{KL}(\mathcal{B}(\theta_i), \mathcal{B}(\theta_L))}$$

$$\liminf_{T \rightarrow \infty} \frac{R_{\theta}^{\pi}(T)}{\ln(T)} = r(L) \prod_{j < L} (1 - \theta_j) \sum_{i=L+1}^N \frac{\theta_L - \theta_i}{\text{KL}(\mathcal{B}(\theta_i), \mathcal{B}(\theta_L))}$$

\Rightarrow Suggest *exploration* at *first* slot 1.

REGRET LOWER BOUNDS - EXPLAINED

Showing results for still alive.

ARTISTS

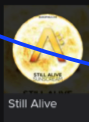
Convex $r(l)$

SEE ALL ALBUMS

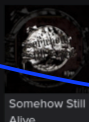
SEE ALL



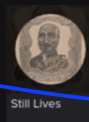
Still Alive



Still Alive



Somehow Still Alive



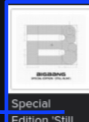
Still Lives



Still Alive
BIGBANG



Still Alive (The Theme from...)



Special Edition 'Still Alive'

SONG

ARTIST

ALBUM

Fixed $r(l)$



	SONG	ARTIST	ALBUM		
+	Still Alive	BIGBANG	BIGBANG Special Edition Still Alive 1	3:19	
▶	Still Alive	Aperture Science Psychoacoustic Laborat...	Portal 2: Songs to Test By (Collectors Editi...	2:57	
+	Still Alive	Lisa Miskovsky	Mirror's Edge Original Videogame Score	4:34	
+	STILL ALIVE	BIGBANG	Special Edition 'Still Alive'	3:19	
+	Still Alive	The Crash	Melodrama	4:05	
+	Still Alive	Social Distortion	Hard Times And Nursery Rhymes (Deluxe ...	4:06	
+	Still Alive	Nocturnal Rites	Grand Illusion	4:03	
+	Still Alive	Onlap, Charline Max	The Awakening	4:05	
+	Still Alive	Jonathan Coulton	Best. Concert. Ever.	3:05	

Theorem (lower bound)

For any uniformly good algorithm π , we have:

$$\liminf_{T \rightarrow \infty} \frac{R^\pi(T)}{\ln(T)} \geq C(\theta),$$

where

$$C(\theta) = \inf_{c_a \geq 0, a \in \mathcal{A}} \sum_{a \in \mathcal{A}} c_a (\mu_{\star}(\theta) - \mu_u(\theta))$$

subject to:

$$\forall i > L, \sum_{a \in \mathcal{A}, i \in a} c_a \text{KL}(\mathcal{B}(\theta_i), \mathcal{B}(\theta_L)) \prod_{s < p_a(i)} (1 - \theta_{a_s}) \geq 1.$$

where $p_a(i) = j$ s.t. $a_j = i$ is the position of i in list a .

Let $j(t) = (j_1(t), \dots, j_N(t))$ be the indices of the items with empirical means sorted in decreasing order and $\mathcal{L}(t) = (j_1(t), \dots, j_L(t))$.

$$\mathcal{E}(t) = \left\{ i \neq \mathcal{L}(t) : \underbrace{\max\{q \in [0, 1] : N_i(t) \text{KL}(\hat{\theta}_i(t), q)\}}_{\text{upper confidence bound}} \leq f(t) \right\} \geq \hat{\theta}_{j_L(t)}(t) \Big\}$$

\Rightarrow *items with high enough upper bound to deserve being explored*

$$U_i^\ell = \{j_1(t), j_2(t), \dots, j_{\ell-1}(t), i, j_\ell(t), \dots, j_{L-1}(t)\}$$

Algorithm 5 Position Induced Exploration(ℓ)

Init: $\mathcal{B}(1) = \emptyset$, $\hat{\theta}_i(1) = 0 = b_i(1) \forall i$, $\mathcal{L}(1) = \{1, \dots, L\}$

For $t \geq 1$:

If $\mathcal{E}(t) = \emptyset$, chooses $a = \mathcal{L}(t)$

Else $\begin{cases} a = \mathcal{L}(t), & w.p. 1/2 \\ a = U_i^\ell(n), i \sim \text{Uniform}(\mathcal{E}(n)) & w.p. 1/2 \end{cases}$

- ▶ Provably asymptotically optimal

- ▶ Provably asymptotically optimal
- ▶ Experiment: compare against

- ▶ Provably asymptotically optimal
- ▶ Experiment: compare against
 - ▶ Slotted-(KL)UCB: top L items in order of their KL-UCB indexes.

- ▶ Provably asymptotically optimal
- ▶ Experiment: compare against
 - ▶ Slotted-(KL)UCB: top L items in order of their KL-UCB indexes.
 - ▶ Ranked Bandit Algorithm: runs L independent instances of KL-UCB on each slot.

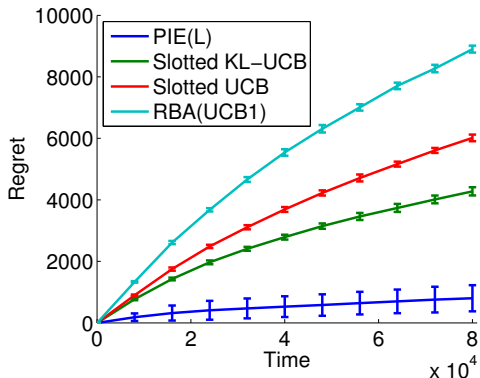
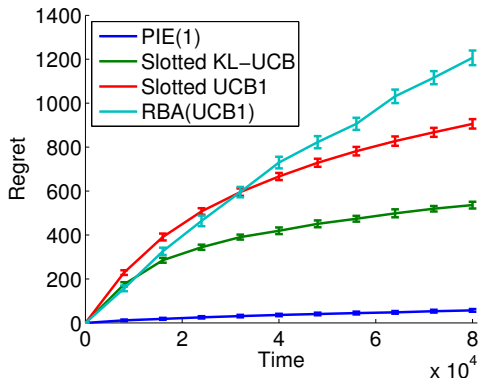
(a) Case 1: $\forall l, r(l) = 2^{1-l}$.(b) Case 2: $\forall l, r(l) = 1$.

Figure: Performance of PIE(1) / PIE(L) and other UCB-based algorithms. A single group of items and users. Error bars represent the standard deviation.

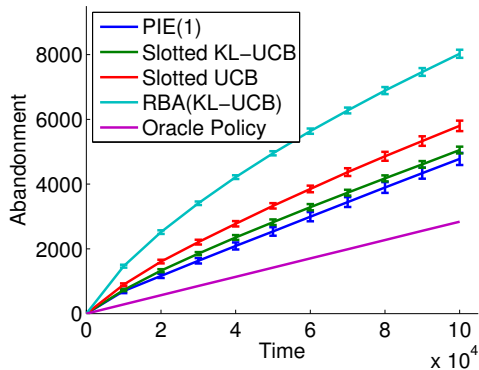


Figure: Performance of PIE(1) on real world data.

- ▶ We consider the Learning to Rank problem as a *Bandit* Optimization problem.

- ▶ We consider the Learning to Rank problem as a *Bandit* Optimization problem.
- ▶ Despite the daunting number of actions, we can *Learn to Rank* with very low cost.

- ▶ We consider the Learning to Rank problem as a *Bandit* Optimization problem.
- ▶ Despite the daunting number of actions, we can *Learn to Rank* with very low cost.
- ▶ Algorithm that optimally exploit structure.

- ▶ We consider the Learning to Rank problem as a *Bandit* Optimization problem.
- ▶ Despite the daunting number of actions, we can *Learn to Rank* with very low cost.
- ▶ Algorithm that optimally exploit structure.
- ▶ plus good empirical performance.

STRUCTURES

LINEAR BANDITS

STRUCTURED LOWER BOUNDS

Lower bounds

Lipschitz bandits

Ranking bandits

Metric-graph of bandits

CONCLUSION, PERSPECTIVE

ACTIVE CONTEXTUAL BANDIT PROBLEM

- ▶ Bandit *configurations*: $\nu = (\nu_{a,b})_{a \in \mathcal{A}, b \in \mathcal{B}}$ with means $(\mu_{a,b})_{a \in \mathcal{A}, b \in \mathcal{B}}$
- ▶ \mathcal{A} : arms, \mathcal{B} : users.
- ▶ *Active contextual* bandit: At time t , learner chooses $b_t \in \mathcal{B}$, then $a_t \in \mathcal{A}$.
- ▶ *Regret*:

$$\mathcal{R}(\nu, T) = \mathbb{E}_\nu \left[\sum_{t=1}^T \max_{a \in \mathcal{A}} \mu_{a,b_t} - X_t \right] = \sum_{a,b \in \mathcal{C}_\nu^-} \Delta_{a,b} \mathbb{E}_\nu [N_{a,b}(T)].$$

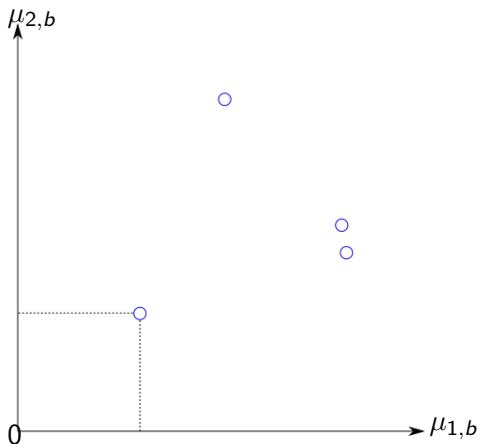
$$\text{where } \mathcal{C}_\nu^- = \left\{ (a, b) \in \mathcal{A} \times \mathcal{B} : \mu_{a,b} < \mu_b^* \right\}.$$

Definition(Uniformly spread strategy)

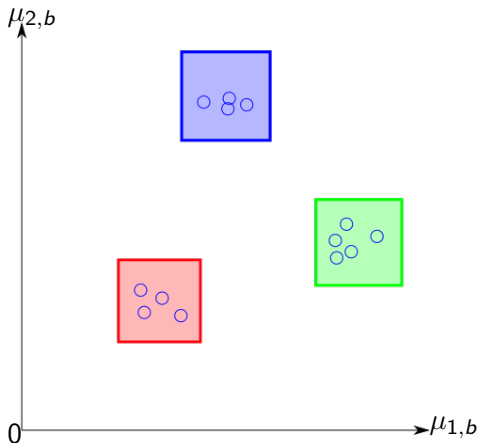
There exists $\gamma_1 > 0$ and a random variable Γ_2 with $\mathbb{E}_\nu[\Gamma_2] < 0$, such that

$$\forall b \in \mathcal{B}, \forall t \in \mathbb{N}, \quad N_b(t) \geq \gamma_1 \cdot t - \Gamma_2.$$

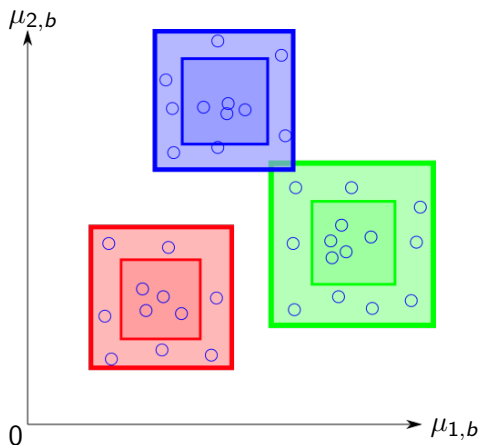
- ▶ Contextual bandits configuration means: $(\mu_{a,b})_{a \in \mathcal{A}, b \in \mathcal{B}}$
- ▶ Set of allowed 2-arm bandits ($\mathcal{A} = \{1, 2\}$):



- ▶ Contextual bandits configuration means: $(\mu_{a,b})_{a \in \mathcal{A}, b \in \mathcal{B}}$
- ▶ Set of allowed 2-arm bandits ($\mathcal{A} = \{1, 2\}$):



- ▶ Contextual bandits configuration means: $(\mu_{a,b})_{a \in \mathcal{A}, b \in \mathcal{B}}$
- ▶ Set of allowed 2-arm bandits $(\mathcal{A} = \{1, 2\})$:



Bandit configurations $(\nu \in \mathcal{P}([0, 1])^{\mathcal{A} \times \mathcal{B}}$ with mean $\mu \in [0, 1]^{\mathcal{A} \times \mathcal{B}}$):

$$\mathcal{D}_\omega = \left\{ \nu : \forall b, b' \in \mathcal{B} \quad \max_{a \in \mathcal{A}} |\mu_{a,b} - \mu_{a,b'}| \leq \omega_{b,b'} \right\},$$

for a known weight matrix $\omega = (\omega_{b,b'})_{b,b' \in \mathcal{B}}$, symmetric, null-diagonal, with positive entries, and satisfying $\omega_{b,b'} \leq \omega_{b,b''} + \omega_{b'',b'}$.

Large values: not structured. Low value: highly structured.

Definition (Consistent strategy)

$$\forall \nu \in \mathcal{D}_\omega, \forall (a, b) \in \mathcal{C}_\nu^-, \forall \alpha \in (0, 1) \quad \lim_{T \rightarrow \infty} \mathbb{E}_\nu \left[\frac{N_{a,b}(T)^\alpha}{N_b(T)} \right] = 0.$$

Proposition (Regret lower bound)

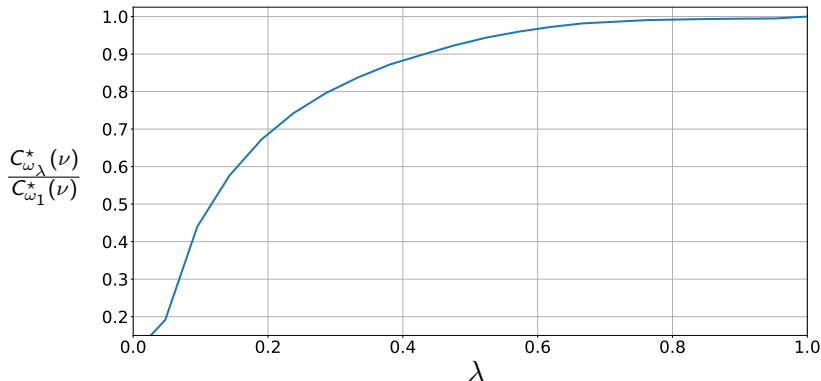
Any uniformly spread and consistent strategy must satisfy

$$\liminf_{T \rightarrow \infty} \frac{\mathcal{R}(\nu, T)}{\ln(T)} \geq C_\omega^*(\nu)$$

where $C_\omega^*(\nu) = \min_{n \in \mathbb{R}_+^{\mathcal{C}^-}} \sum_{a,b \in \mathcal{C}^-} n_{a,b} \Delta_{a,b}$ s.t.

$$\forall (a, b) \in \mathcal{C}^-, \quad \sum_{b' \in \mathcal{B}:(a,b) \in \mathcal{C}^-} \text{kl}^+(\mu_{a,b'} | \mu_b^* - \omega_{b,b'}) n_{a,b'} \geq 1.$$

- ▶ Let ω_λ be a matrix where all the weights are equal to $\lambda \in [0, 1]$ except for the zero diagonal.
- ▶ $\lambda = 1$: *no-structure*, $\lambda = 0$: one unique cluster.
- ▶ We recover that $C_{\omega_1}^*(\nu) = \sum_{a,b \in \mathcal{C}^-} \frac{\Delta_{a,b}}{\mathbb{K}1(\mu_{a,b} | \mu_b^*)}$ (unstructured lower bound)
- ▶ More generally:



- ▶ Explicit lower bound spanning unstructured to highly structured pbs.
- ▶ See (Saber et al., submitted) for an algorithm:
 - ▶ Provably asymptotically optimal.
 - ▶ Computationally cheap
 - ▶ Without explicit forced exploration (still some implicit forcing).

STRUCTURES

LINEAR BANDITS

STRUCTURED LOWER BOUNDS

CONCLUSION, PERSPECTIVE

Confidence bounds in parametric regression: Time and space uniform

$$\forall \delta \in (0, 1), \mathbb{P}\left(\exists t \in \mathbb{N}, x \in \mathcal{X} : |f_{\star}(x) - f_{\theta_t}(x)| \geq \|\varphi(x)\|_{G_{t,\lambda}^{-1}} B_t(\delta)\right) \leq \delta$$

- ▶ Quite tight (Equality everywhere, except Markov inequality and super-martingale).
- ▶ Extends to Kernel regression similarly.
- ▶ Optimal use of it? not quite ("The end of optimism", Lattimore et al.)

Pick your favorite *structured bandit problem*

Study the problem-dependent lower bound

Each arm should be pulled some minimum number of times.

Suggests an algorithm (sometimes optimal) !

- ▶ In Linear bandits:
 - ▶ Features? Representation?
 - ▶ Lower bounds ? Most confusing instances? Optimality?
- ▶ In generic structure:
 - ▶ Generic algorithm (e.g. OSSB)?
 - ▶ Forced exploration?
 - ▶ More informative/Less conservative lower bounds?
 - ▶ Better tracking of information?
- ▶ Beyond structure? No stochastic model?

Habilitation manuscript:

"Mathematics of Statistical Sequential Learning"

<https://hal.archives-ouvertes.fr/tel-02162189>

Open positions:

<http://odalricambrymmaillard.neowordpress.fr/research-projets/open-positions/>

MERCI



Inria Lille - Nord Europe

odalricambrym.maillard@inria.fr

odalricambrymmaillard.wordpress.com