

Bandits on Graphs

Exploiting smoothness and side observations

Michal Valko (Sequel INRIA)

joint work with

Shipra Agrawal (MSR India)

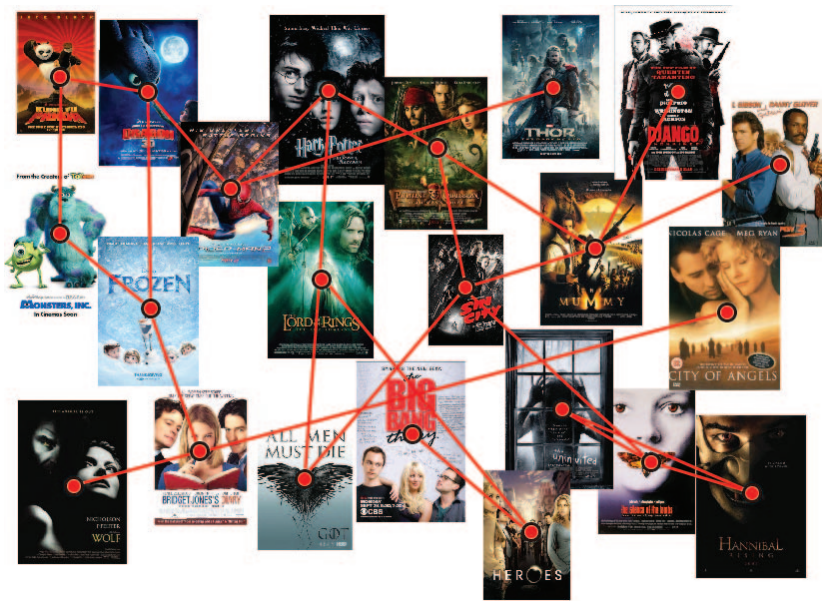
Tomáš Kocák (Sequel INRIA)

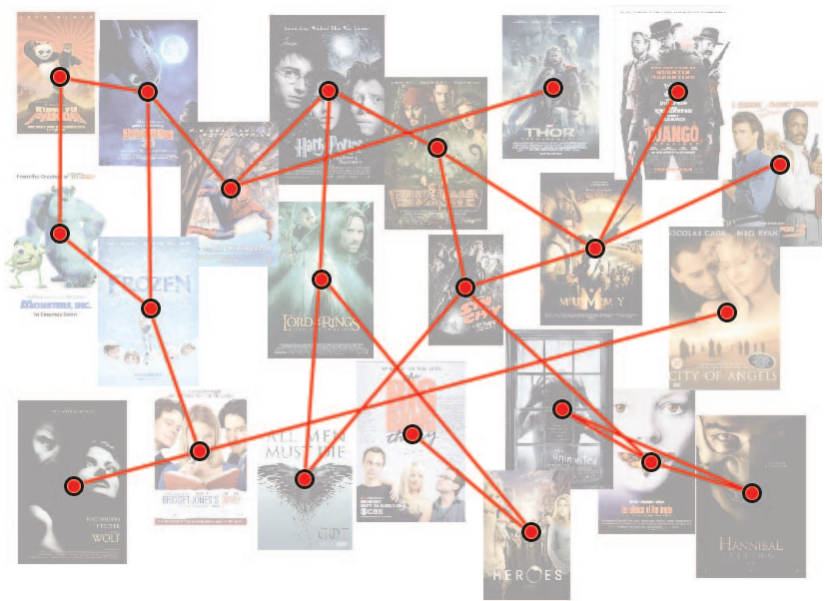
Branislav Kveton (Technicolor → Adobe)

Rémi Munos (Sequel INRIA/Google Deepmind)

Gergely Neu (Sequel INRIA)







Overview

- ▶ Sequential decision making in structured settings

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a graph

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a graph
 - ▶ the graph encodes some **structural property** of the setting

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ **goal: maximize the sum of the outcomes**

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ **application: recommender systems**

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ application: recommender systems
- ▶ Exploiting **smoothness**

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ application: recommender systems
- ▶ Exploiting **smoothness**
 - ▶ **fixed graph**

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ application: recommender systems
- ▶ Exploiting **smoothness**
 - ▶ fixed graph
 - ▶ **iid outcomes**

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ application: recommender systems
- ▶ Exploiting **smoothness**
 - ▶ fixed graph
 - ▶ iid outcomes
 - ▶ **neighboring nodes have similar outcomes**

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ application: recommender systems
- ▶ Exploiting **smoothness**
 - ▶ fixed graph
 - ▶ iid outcomes
 - ▶ neighboring nodes have similar outcomes
- ▶ Exploiting **side observations**

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ application: recommender systems
- ▶ Exploiting **smoothness**
 - ▶ fixed graph
 - ▶ iid outcomes
 - ▶ neighboring nodes have similar outcomes
- ▶ Exploiting **side observations**
 - ▶ **changing graph**

Overview

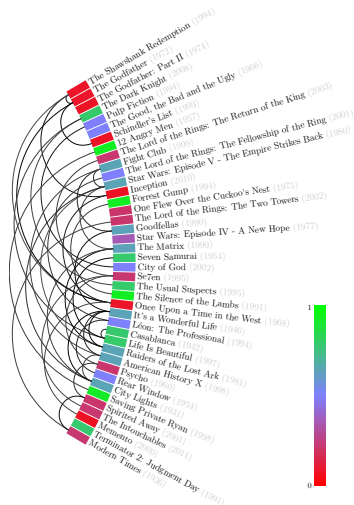
- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ application: recommender systems
- ▶ Exploiting **smoothness**
 - ▶ fixed graph
 - ▶ iid outcomes
 - ▶ neighboring nodes have similar outcomes
- ▶ Exploiting **side observations**
 - ▶ changing graph
 - ▶ **non-stochastic outcomes**

Overview

- ▶ Sequential decision making in structured settings
 - ▶ we are asked to pick a node (or a few nodes) in a **graph**
 - ▶ the **graph** encodes some **structural property** of the setting
 - ▶ goal: maximize the sum of the outcomes
 - ▶ application: recommender systems
- ▶ Exploiting **smoothness**
 - ▶ fixed graph
 - ▶ iid outcomes
 - ▶ neighboring nodes have similar outcomes
- ▶ Exploiting **side observations**
 - ▶ changing graph
 - ▶ non-stochastic outcomes
 - ▶ **side observations**

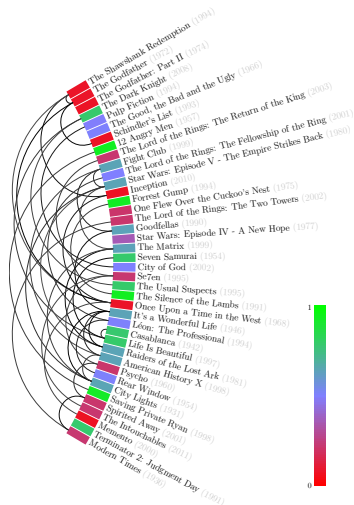
Movie recommendation: (in each time step)

- ▶ Recommend movies to a **single user**.



Movie recommendation: (in each time step)

- ▶ Recommend movies to a **single user**.
- ▶ Good prediction after a few steps ($T \ll N$).

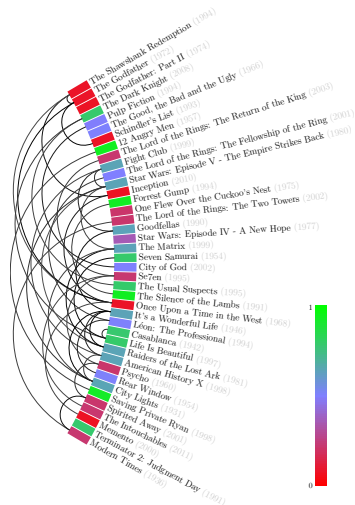


Movie recommendation: (in each time step)

- ▶ Recommend movies to a **single user**.
- ▶ Good prediction after a few steps ($T \ll N$).

Goal:

- ▶ Maximize overall reward (sum of ratings).



Movie recommendation: (in each time step)

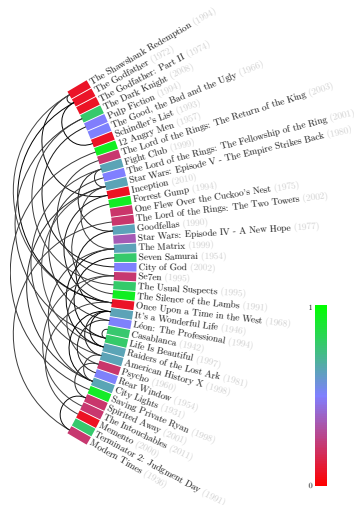
- ▶ Recommend movies to a **single user**.
- ▶ Good prediction after a few steps ($T \ll N$).

Goal:

- ▶ Maximize overall reward (sum of ratings).

Assumptions:

- ▶ Unknown reward function $f : V(G) \rightarrow \mathbb{R}$.
- ▶ Function f is **smooth** on a graph.
- ▶ Neighboring movies \Rightarrow similar preferences.
- ▶ Similar preferences \nRightarrow neighboring movies.



Smooth graph function

- ▶ Graph G with vertex set $V(G) = \{1, \dots, N\}$ and edge set $E(G)$.

Smooth graph function

- ▶ Graph G with vertex set $V(G) = \{1, \dots, N\}$ and edge set $E(G)$.
- ▶ f_1, \dots, f_N : Values of the function on the vertices of the graph.

Smooth graph function

- ▶ Graph G with vertex set $V(G) = \{1, \dots, N\}$ and edge set $E(G)$.
- ▶ f_1, \dots, f_N : Values of the function on the vertices of the graph.
- ▶ $w_{i,j}$: Weight of the edge connecting nodes i and j .

Smooth graph function

- ▶ Graph G with vertex set $V(G) = \{1, \dots, N\}$ and edge set $E(G)$.
- ▶ f_1, \dots, f_N : Values of the function on the vertices of the graph.
- ▶ $w_{i,j}$: Weight of the edge connecting nodes i and j .

Smooth graph function

- ▶ Graph G with vertex set $V(G) = \{1, \dots, N\}$ and edge set $E(G)$.
- ▶ f_1, \dots, f_N : Values of the function on the vertices of the graph.
- ▶ $w_{i,j}$: Weight of the edge connecting nodes i and j .

Smoothness of the function:

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

Smaller value of $S_G(f)$, smoother the function f is.

Smooth graph function

- ▶ Graph G with vertex set $V(G) = \{1, \dots, N\}$ and edge set $E(G)$.
- ▶ f_1, \dots, f_N : Values of the function on the vertices of the graph.
- ▶ $w_{i,j}$: Weight of the edge connecting nodes i and j .

Smoothness of the function:

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

Smaller value of $S_G(f)$, smoother the function f is.

Examples:

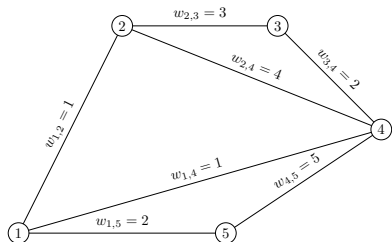
- ▶ **Complete graph:** Only constant function has smoothness 0.
- ▶ **Edgeless graph:** Every function has smoothness 0.
- ▶ **Constant function:** Smoothness 0 for every graph.

Graph Laplacian

- ▶ \mathcal{W} : $N \times N$ matrix of the edge weights $w_{i,j}$.
- ▶ \mathcal{D} : Diagonal matrix with the entries $d_i = \sum_j w_{i,j}$.
- ▶ $\mathcal{L} = \mathcal{D} - \mathcal{W}$: Graph Laplacian.
 - ▶ Positive semidefinite matrix.
 - ▶ Diagonally dominant matrix.

Example:

$$\mathcal{L} = \begin{pmatrix} 4 & -1 & 0 & -1 & -2 \\ -1 & 8 & -3 & -4 & 0 \\ 0 & -3 & 5 & -2 & 0 \\ -1 & -4 & -2 & 12 & -5 \\ -2 & 0 & 0 & -5 & 7 \end{pmatrix}$$



Smoothness of the function and Laplacian

- ▶ $\mathbf{f} = (f_1, \dots, f_N)^T$: Vector of function values.

Smoothness of the function and Laplacian

- ▶ $\mathbf{f} = (f_1, \dots, f_N)^\top$: Vector of function values.
- ▶ Let $\mathcal{L} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ be the eigendecomposition of the Laplacian.

Smoothness of the function and Laplacian

- ▶ $\mathbf{f} = (f_1, \dots, f_N)^\top$: Vector of function values.
- ▶ Let $\mathcal{L} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ be the eigendecomposition of the Laplacian.
 - ▶ Diagonal matrix $\mathbf{\Lambda}$ whose diagonal entries are eigenvalues of \mathcal{L} .

Smoothness of the function and Laplacian

- ▶ $\mathbf{f} = (f_1, \dots, f_N)^\top$: Vector of function values.
- ▶ Let $\mathcal{L} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ be the eigendecomposition of the Laplacian.
 - ▶ Diagonal matrix $\mathbf{\Lambda}$ whose diagonal entries are eigenvalues of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} are eigenvectors of \mathcal{L} .

Smoothness of the function and Laplacian

- ▶ $\mathbf{f} = (f_1, \dots, f_N)^T$: Vector of function values.
- ▶ Let $\mathcal{L} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ be the eigendecomposition of the Laplacian.
 - ▶ Diagonal matrix $\mathbf{\Lambda}$ whose diagonal entries are eigenvalues of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} are eigenvectors of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} form a basis.

Smoothness of the function and Laplacian

- ▶ $\mathbf{f} = (f_1, \dots, f_N)^\top$: Vector of function values.
- ▶ Let $\mathcal{L} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ be the eigendecomposition of the Laplacian.
 - ▶ Diagonal matrix $\mathbf{\Lambda}$ whose diagonal entries are eigenvalues of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} are eigenvectors of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} form a basis.
- ▶ $\boldsymbol{\alpha}^*$: Unique vector such that $\mathbf{Q}\boldsymbol{\alpha}^* = \mathbf{f}$ Note: $\mathbf{Q}^\top \mathbf{f} = \boldsymbol{\alpha}^*$

Smoothness of the function and Laplacian

- ▶ $\mathbf{f} = (f_1, \dots, f_N)^\top$: Vector of function values.
- ▶ Let $\mathcal{L} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ be the eigendecomposition of the Laplacian.
 - ▶ Diagonal matrix $\mathbf{\Lambda}$ whose diagonal entries are eigenvalues of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} are eigenvectors of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} form a basis.
- ▶ α^* : Unique vector such that $\mathbf{Q}\alpha^* = \mathbf{f}$ Note: $\mathbf{Q}^\top \mathbf{f} = \alpha^*$

Smoothness of the function and Laplacian

- ▶ $\mathbf{f} = (f_1, \dots, f_N)^T$: Vector of function values.
- ▶ Let $\mathcal{L} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ be the eigendecomposition of the Laplacian.
 - ▶ Diagonal matrix $\mathbf{\Lambda}$ whose diagonal entries are eigenvalues of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} are eigenvectors of \mathcal{L} .
 - ▶ Columns of \mathbf{Q} form a basis.
- ▶ α^* : Unique vector such that $\mathbf{Q}\alpha^* = \mathbf{f}$ Note: $\mathbf{Q}^T\mathbf{f} = \alpha^*$

$$S_G(f) = \mathbf{f}^T \mathcal{L} \mathbf{f} = \mathbf{f}^T \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \mathbf{f} = \alpha^{*\top} \mathbf{\Lambda} \alpha^* = \|\alpha^*\|_{\mathbf{\Lambda}}^2 = \sum_{i=1}^N \lambda_i (\alpha_i^*)^2$$

Smoothness and regularization: Small value of

(a) $S_G(f)$ (b) $\mathbf{\Lambda}$ norm of α^* (c) α_i^* for large λ_i

Setting

Learning setting for a bandit algorithm π

- ▶ In each time t step choose a node $\pi(t)$.

Setting

Learning setting for a bandit algorithm π

- ▶ In each time t step choose a node $\pi(t)$.
- ▶ the $\pi(t)$ -th row $\mathbf{x}_{\pi(t)}$ of the matrix \mathbf{Q} corresponds to the arm $\pi(t)$.

Setting

Learning setting for a bandit algorithm π

- ▶ In each time t step choose a node $\pi(t)$.
- ▶ the $\pi(t)$ -th row $\mathbf{x}_{\pi(t)}$ of the matrix \mathbf{Q} corresponds to the arm $\pi(t)$.
- ▶ Obtain noisy reward $r_t = \mathbf{x}_{\pi(t)}^\top \boldsymbol{\alpha}^* + \varepsilon_t$. Note: $\mathbf{x}_{\pi(t)}^\top \boldsymbol{\alpha}^* = f_{\pi(t)}$
 - ▶ ε_t is R -sub-Gaussian noise. $\forall \xi \in \mathbb{R}, \mathbb{E}[e^{\xi \varepsilon_t}] \leq \exp(\xi^2 R^2 / 2)$
- ▶ Minimize cumulative regret

$$R_T = T \max_a (\mathbf{x}_a^\top \boldsymbol{\alpha}^*) - \sum_{t=1}^T \mathbf{x}_{\pi(t)}^\top \boldsymbol{\alpha}^*.$$

- ▶ Can't we just use *linear bandits*?

Solutions

- ▶ **Linear bandit algorithms**

(Existing solutions)

- ▶ **LinUCB**

(Li et al., 2010)

- ▶ Regret bound $\approx D\sqrt{T \ln T}$

- ▶ **LinearTS**

(Agrawal and Goyal, 2013)

- ▶ Regret bound $\approx D\sqrt{T \ln N}$

Note: D is ambient dimension, in our case N , length of x_i .

Number of actions, e.g., all possible movies → **HUGE!**

Solutions

▶ Linear bandit algorithms

(Existing solutions)

▶ LinUCB

(Li et al., 2010)

▶ Regret bound $\approx D\sqrt{T \ln T}$

▶ LinearTS

(Agrawal and Goyal, 2013)

▶ Regret bound $\approx D\sqrt{T \ln N}$

Note: D is ambient dimension, in our case N , length of x_i .

Number of actions, e.g., all possible movies → **HUGE!**

▶ Spectral bandit algorithms

(Our solutions)

▶ SpectralUCB

(Valko et al., ICML 2014)

▶ Regret bound $\approx d\sqrt{T \ln T}$

▶ SpectralTS

(Kocák et al., AAAI 2014)

▶ Regret bound $\approx d\sqrt{T \ln N}$

Note: d is **effective dimension**, usually much smaller than D .

Solutions

▶ Linear bandit algorithms

(Existing solutions)

▶ LinUCB

(Li et al., 2010)

▶ Regret bound $\approx D\sqrt{T \ln T}$

▶ LinearTS

(Agrawal and Goyal, 2013)

▶ Regret bound $\approx D\sqrt{T \ln N}$

Note: D is ambient dimension, in our case N , length of x_i .

Number of actions, e.g., all possible movies → **HUGE!**

▶ Spectral bandit algorithms

(Our solutions)

▶ SpectralUCB

(Valko et al., ICML 2014)

▶ Regret bound $\approx d\sqrt{T \ln T}$

▶ Operations per step: $D^2 N$

▶ SpectralTS

(Kocák et al., AAAI 2014)

▶ Regret bound $\approx d\sqrt{T \ln N}$

▶ Operations per step: $D^2 + DN$

Note: d is **effective dimension**, usually much smaller than D .

Effective dimension

- ▶ **Effective dimension:** Largest d such that

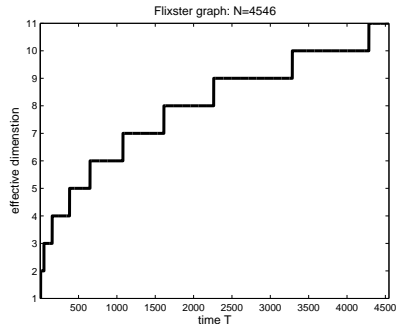
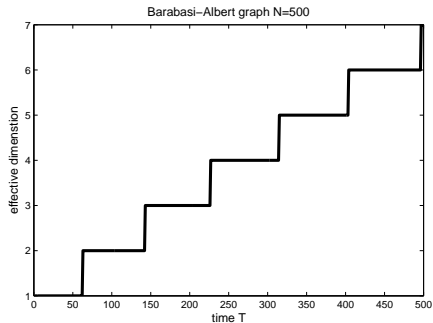
$$(d - 1)\lambda_d \leq \frac{T}{\log(1 + T/\lambda)}.$$

- ▶ Function of time horizon and graph properties
- ▶ λ_i : i -th smallest eigenvalue of $\mathbf{\Lambda}$.
- ▶ λ : Regularization parameter of the algorithm.

Properties:

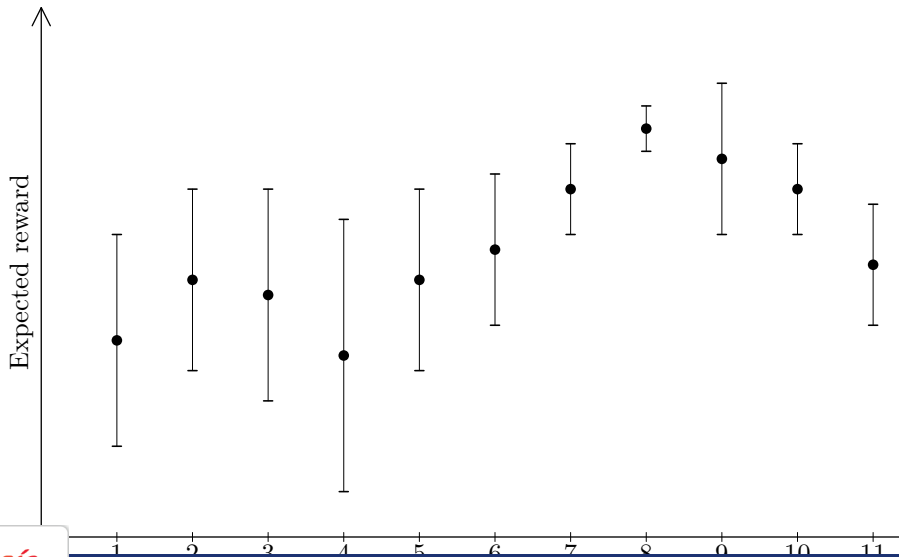
- ▶ d is small when the coefficients λ_i grow rapidly above time.
- ▶ d is related to the number of “non-negligible” dimensions.
- ▶ Usually d is much smaller than D in real world graphs.
- ▶ Can be computed beforehand.

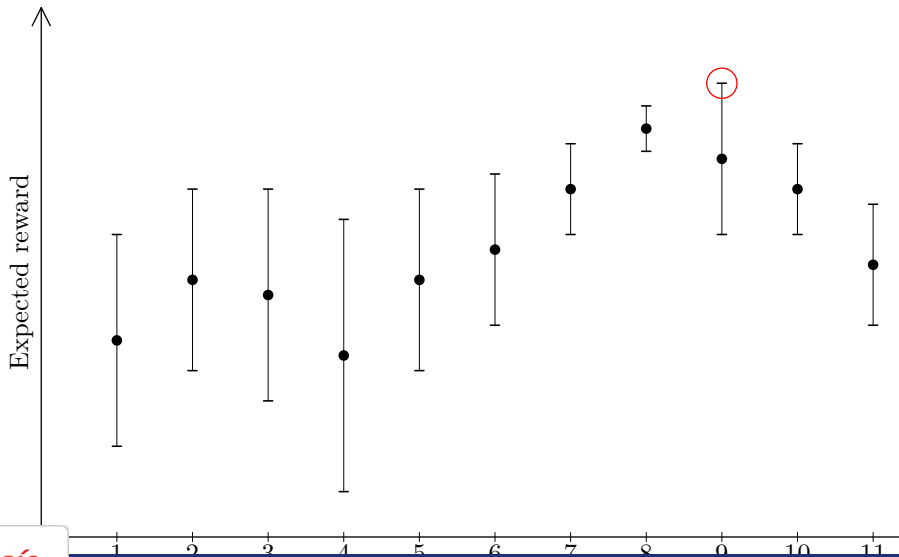
Effective dimension vs. Ambient dimension

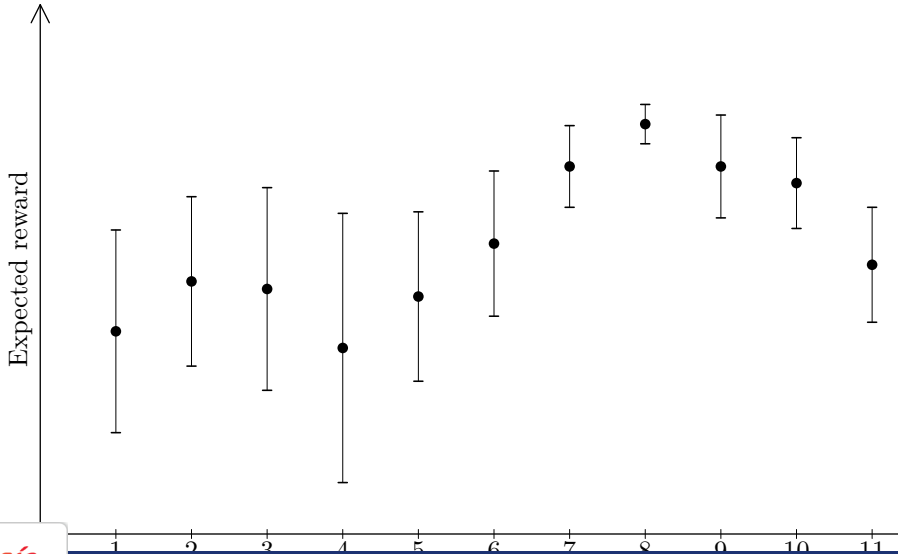


$$d \ll D$$

Note: In our setting $T < N = D$.

UCB style algorithms: **Estimate**

UCB style algorithms: **Sample**

UCB style algorithms: **Estimate ...**

SpectralUCB

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Run:**
- 4: $\Lambda \leftarrow \Lambda_{\mathcal{L}} + \lambda \mathbf{I}$
- 5: $d \leftarrow \max\{d : (d-1)\lambda_d \leq T/\ln(1+T/\lambda)\}$
- 6: **for** $t = 1$ **to** T **do**
- 7: Update the basis coefficients $\hat{\alpha}$:
- 8: $\mathbf{X}_t \leftarrow [\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(t-1)}]^\top$
- 9: $\mathbf{r} \leftarrow [r_1, \dots, r_{t-1}]^\top$
- 10: $\mathbf{V}_t \leftarrow \mathbf{X}_t \mathbf{X}_t^\top + \Lambda$
- 11: $\hat{\alpha}_t \leftarrow \mathbf{V}_t^{-1} \mathbf{X}_t^\top \mathbf{r}$
- 12: $c_t \leftarrow 2R\sqrt{d \ln(1+t/\lambda) + 2 \ln(1/\delta)} + C$
- 13: $\pi(t) \leftarrow \arg \max_a (\mathbf{x}_a^\top \hat{\alpha}_t + c_t \|\mathbf{x}_a\|_{\mathbf{V}_t^{-1}})$
- 14: Observe the reward r_t
- 15: **end for**

SpectralUCB

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C, \mathcal{L}$
- 3: **Run:**
- 4: $\Lambda \leftarrow \Lambda_{\mathcal{L}} + \lambda \mathbf{I}$
- 5: $d \leftarrow \max\{d : (d-1)\lambda_d \leq T/\ln(1+T/\lambda)\}$
- 6: **for** $t = 1$ **to** T **do**
- 7: Update the basis coefficients $\hat{\alpha}$:
- 8: $\mathbf{X}_t \leftarrow [\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(t-1)}]^\top$
- 9: $\mathbf{r} \leftarrow [r_1, \dots, r_{t-1}]^\top$
- 10: $\mathbf{V}_t \leftarrow \mathbf{X}_t \mathbf{X}_t^\top + \Lambda$
- 11: $\hat{\alpha}_t \leftarrow \mathbf{V}_t^{-1} \mathbf{X}_t^\top \mathbf{r}$
- 12: $c_t \leftarrow 2R\sqrt{d \ln(1+t/\lambda) + 2 \ln(1/\delta)} + C$
- 13: $\pi(t) \leftarrow \arg \max_a (\mathbf{x}_a^\top \hat{\alpha}_t + c_t \|\mathbf{x}_a\|_{\mathbf{V}_t^{-1}})$
- 14: Observe the reward r_t
- 15: **end for**

SpectralUCB

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C, \mathcal{L}$
- 3: **Run:**
- 4: $\Lambda \leftarrow \Lambda_{\mathcal{L}} + \lambda \mathbf{I}$
- 5: $d \leftarrow \max\{d : (d-1)\lambda_d \leq T/\ln(1+T/\lambda)\}$
- 6: **for** $t = 1$ **to** T **do**
- 7: Update the basis coefficients $\hat{\alpha}$:
- 8: $\mathbf{X}_t \leftarrow [\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(t-1)}]^\top$
- 9: $\mathbf{r} \leftarrow [r_1, \dots, r_{t-1}]^\top$
- 10: $\mathbf{V}_t \leftarrow \mathbf{X}_t \mathbf{X}_t^\top + \Lambda$
- 11: $\hat{\alpha}_t \leftarrow \mathbf{V}_t^{-1} \mathbf{X}_t^\top \mathbf{r}$
- 12: $c_t \leftarrow 2R\sqrt{d \ln(1+t/\lambda) + 2 \ln(1/\delta)} + C$
- 13: $\pi(t) \leftarrow \arg \max_a \left(\mathbf{x}_a^\top \hat{\alpha}_t + c_t \|\mathbf{x}_a\|_{\mathbf{V}_t^{-1}} \right)$
- 14: Observe the reward r_t
- 15: **end for**

SpectralUCB

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Run:**
- 4: $\Lambda \leftarrow \Lambda_{\mathcal{L}} + \lambda \mathbf{I}$
- 5: $d \leftarrow \max\{d : (d-1)\lambda_d \leq T/\ln(1+T/\lambda)\}$
- 6: **for** $t = 1$ **to** T **do**
- 7: Update the basis coefficients $\hat{\alpha}$:
- 8: $\mathbf{X}_t \leftarrow [\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(t-1)}]^\top$
- 9: $\mathbf{r} \leftarrow [r_1, \dots, r_{t-1}]^\top$
- 10: $\mathbf{V}_t \leftarrow \mathbf{X}_t \mathbf{X}_t^\top + \Lambda$
- 11: $\hat{\alpha}_t \leftarrow \mathbf{V}_t^{-1} \mathbf{X}_t^\top \mathbf{r}$
- 12: $c_t \leftarrow 2R\sqrt{d \ln(1+t/\lambda) + 2 \ln(1/\delta)} + C$
- 13: $\pi(t) \leftarrow \arg \max_a \left(\mathbf{x}_a^\top \hat{\alpha} + c_t \|\mathbf{x}_a\|_{\mathbf{V}_t^{-1}} \right)$
- 14: Observe the reward r_t
- 15: **end for**

SpectralUCB

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Run:**
- 4: $\Lambda \leftarrow \Lambda_{\mathcal{L}} + \lambda \mathbf{I}$
- 5: $d \leftarrow \max\{d : (d-1)\lambda_d \leq T/\ln(1+T/\lambda)\}$
- 6: **for** $t = 1$ **to** T **do**
- 7: Update the basis coefficients $\hat{\alpha}$:
- 8: $\mathbf{X}_t \leftarrow [\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(t-1)}]^\top$
- 9: $\mathbf{r} \leftarrow [r_1, \dots, r_{t-1}]^\top$
- 10: $\mathbf{V}_t \leftarrow \mathbf{X}_t \mathbf{X}_t^\top + \Lambda$
- 11: $\hat{\alpha}_t \leftarrow \mathbf{V}_t^{-1} \mathbf{X}_t^\top \mathbf{r}$
- 12: $c_t \leftarrow 2R\sqrt{d \ln(1+t/\lambda) + 2 \ln(1/\delta)} + C$
- 13: $\pi(t) \leftarrow \arg \max_a \left(\mathbf{x}_a^\top \hat{\alpha}_t + c_t \|\mathbf{x}_a\|_{\mathbf{V}_t^{-1}} \right)$
- 14: Observe the reward r_t
- 15: **end for**

SpectralUCB

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Run:**
- 4: $\Lambda \leftarrow \Lambda_{\mathcal{L}} + \lambda \mathbf{I}$
- 5: $d \leftarrow \max\{d : (d-1)\lambda_d \leq T/\ln(1+T/\lambda)\}$
- 6: **for** $t = 1$ **to** T **do**
- 7: Update the basis coefficients $\hat{\alpha}$:
- 8: $\mathbf{X}_t \leftarrow [\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(t-1)}]^\top$
- 9: $\mathbf{r} \leftarrow [r_1, \dots, r_{t-1}]^\top$
- 10: $\mathbf{V}_t \leftarrow \mathbf{X}_t \mathbf{X}_t^\top + \Lambda$
- 11: $\hat{\alpha}_t \leftarrow \mathbf{V}_t^{-1} \mathbf{X}_t^\top \mathbf{r}$
- 12: $c_t \leftarrow 2R\sqrt{d \ln(1+t/\lambda) + 2 \ln(1/\delta)} + C$
- 13: $\pi(t) \leftarrow \arg \max_a \left(\mathbf{x}_a^\top \hat{\alpha}_t + c_t \|\mathbf{x}_a\|_{\mathbf{V}_t^{-1}} \right)$
- 14: Observe the reward r_t
- 15: **end for**

SpectralUCB regret bound

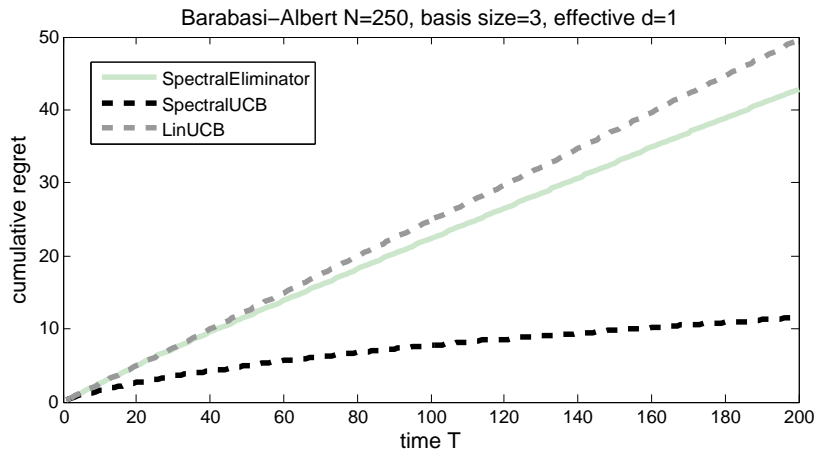
- ▶ d : Effective dimension.
- ▶ λ : Minimal eigenvalue of $\mathbf{\Lambda} = \mathbf{\Lambda}_{\mathcal{L}} + \lambda \mathbf{I}$.
- ▶ C : Smoothness upper bound, $\|\boldsymbol{\alpha}^*\|_{\mathbf{\Lambda}} \leq C$.
- ▶ $\mathbf{x}_i^T \boldsymbol{\alpha}^* \in [-1, 1]$ for all i .

The **cumulative regret** R_T of **SpectralUCB** is with probability $1 - \delta$ bounded as

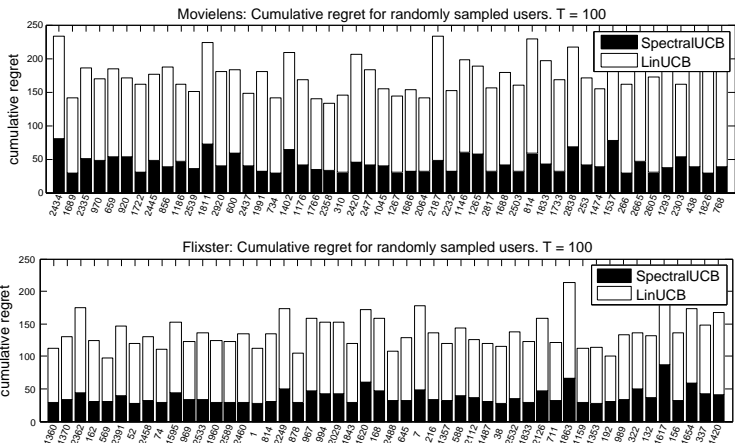
$$R_T \leq \left(8R \sqrt{d \ln \frac{\lambda + T}{\lambda} + 2 \ln \frac{1}{\delta} + 4C + 4} \right) \sqrt{dT \ln \frac{\lambda + T}{\lambda}}.$$

$$R_T \approx d \sqrt{T \ln T}$$

Synthetic experiment

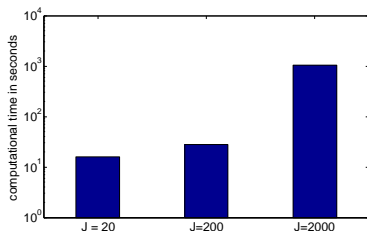
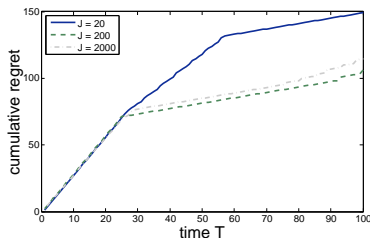


Real world experiment



Improving the running time: reduced eigenbasis

- ▶ **Reduced basis:** We only need first few eigenvectors.
- ▶ **Getting J eigenvectors:** $\mathcal{O}(Jm \log m)$ time for m edges
- ▶ Computationally less expensive, comparable performance.



How to make it faster?

- ▶ UCB-style algorithms need to (re)-compute UCBs every t

How to make it faster?

- ▶ UCB-style algorithms need to (re)-compute UCBs every t
- ▶ Can be a problem for large set of arms $\rightarrow D^2N \rightarrow N^3$

How to make it faster?

- ▶ UCB-style algorithms need to (re)-compute UCBs every t
- ▶ Can be a problem for large set of arms $\rightarrow D^2N \rightarrow N^3$
- ▶ Optimistic (UCB) approach vs. Thompson Sampling

How to make it faster?

- ▶ UCB-style algorithms need to (re)-compute UCBs every t
- ▶ Can be a problem for large set of arms $\rightarrow D^2N \rightarrow N^3$
- ▶ Optimistic (UCB) approach vs. Thompson Sampling
 - ▶ Play the arm maximizing probability of being the best

How to make it faster?

- ▶ UCB-style algorithms need to (re)-compute UCBs every t
- ▶ Can be a problem for large set of arms $\rightarrow D^2 N \rightarrow N^3$
- ▶ Optimistic (UCB) approach vs. Thompson Sampling
 - ▶ Play the arm maximizing probability of being the best
 - ▶ Sample $\tilde{\mu}$ from the distribution $\mathcal{N}(\hat{\mu}, v^2 \mathbf{B}^{-1})$

How to make it faster?

- ▶ UCB-style algorithms need to (re)-compute UCBs every t
- ▶ Can be a problem for large set of arms $\rightarrow D^2 N \rightarrow N^3$
- ▶ Optimistic (UCB) approach vs. Thompson Sampling
 - ▶ Play the arm maximizing probability of being the best
 - ▶ Sample $\tilde{\boldsymbol{\mu}}$ from the distribution $\mathcal{N}(\hat{\boldsymbol{\mu}}, v^2 \mathbf{B}^{-1})$
 - ▶ Play arm which maximizes $\mathbf{b}^\top \tilde{\boldsymbol{\mu}}$ and observe reward

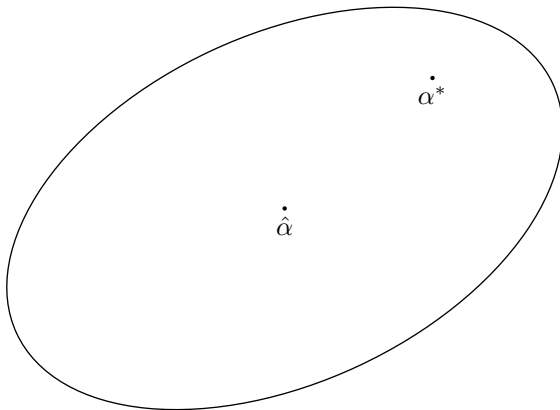
How to make it faster?

- ▶ UCB-style algorithms need to (re)-compute UCBs every t
- ▶ Can be a problem for large set of arms $\rightarrow D^2 N \rightarrow N^3$
- ▶ Optimistic (UCB) approach vs. Thompson Sampling
 - ▶ Play the arm maximizing probability of being the best
 - ▶ Sample $\tilde{\boldsymbol{\mu}}$ from the distribution $\mathcal{N}(\hat{\boldsymbol{\mu}}, v^2 \mathbf{B}^{-1})$
 - ▶ Play arm which maximizes $\mathbf{b}^\top \tilde{\boldsymbol{\mu}}$ and observe reward
 - ▶ Compute posterior distribution according to reward received

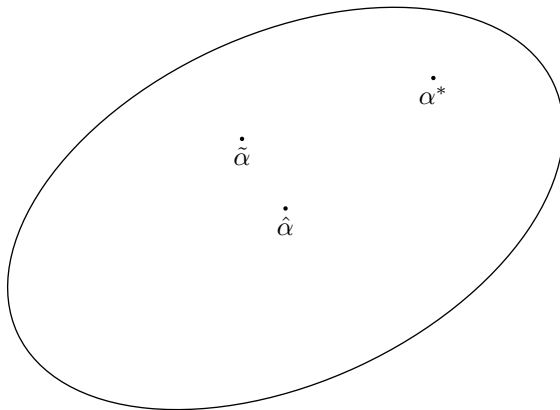
How to make it faster?

- ▶ UCB-style algorithms need to (re)-compute UCBs every t
- ▶ Can be a problem for large set of arms $\rightarrow D^2 N \rightarrow N^3$
- ▶ Optimistic (UCB) approach vs. Thompson Sampling
 - ▶ Play the arm maximizing probability of being the best
 - ▶ Sample $\tilde{\boldsymbol{\mu}}$ from the distribution $\mathcal{N}(\hat{\boldsymbol{\mu}}, v^2 \mathbf{B}^{-1})$
 - ▶ Play arm which maximizes $\mathbf{b}^\top \tilde{\boldsymbol{\mu}}$ and observe reward
 - ▶ Compute posterior distribution according to reward received
- ▶ Only requires $D^2 + DN \rightarrow N^2$ per step update

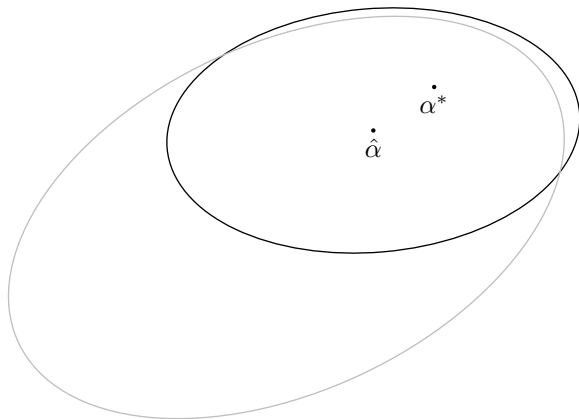
Thomson Sampling: Estimate



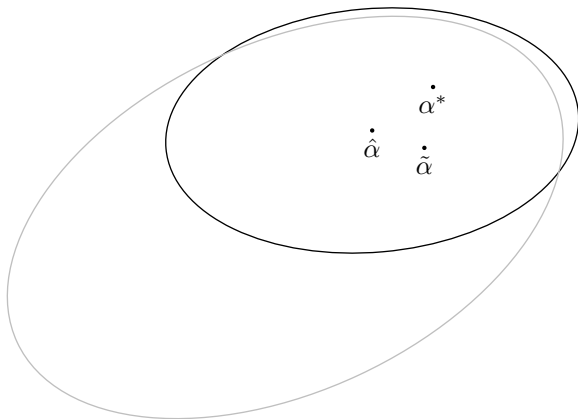
Thomson Sampling: **Sample**

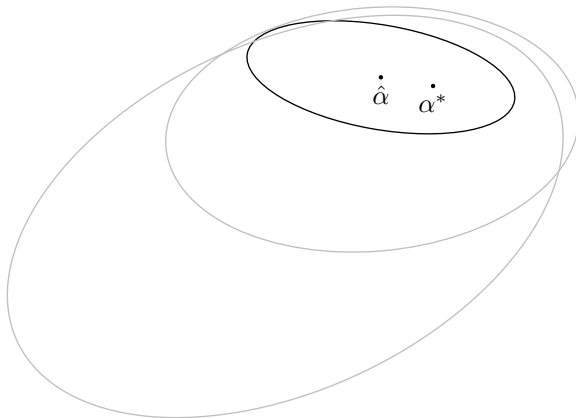


Thomson Sampling: Estimate



Thomson Sampling: **Sample**



Thomson Sampling: **Estimate ...**

SpectralTS algorithm

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Initialization:**
- 4: $v = R\sqrt{6d \log((\lambda + T)/\delta\lambda)} + C$
- 5: $\hat{\alpha} = 0_N$
- 6: $\mathbf{f} = 0_N$
- 7: $\mathbf{V} = \Lambda_{\mathcal{L}} + \lambda \mathbf{I}_N$
- 8: **Run:**
- 9: **for** $t = 1$ **to** T **do**
- 10: Sample $\tilde{\alpha} \sim \mathcal{N}(\hat{\alpha}, v^2 \mathbf{V}^{-1})$
- 11: $\pi(t) \leftarrow \arg \max_a \mathbf{x}_a^T \tilde{\alpha}$
- 12: Observe a noisy reward $r(t) = \mathbf{x}_{\pi(t)}^T \boldsymbol{\alpha}^* + \varepsilon_t$
- 13: $\mathbf{f} \leftarrow \mathbf{f} + \mathbf{x}_{\pi(t)} r(t)$
- 14: Update $\mathbf{V} \leftarrow \mathbf{V} + \mathbf{x}_{\pi(t)} \mathbf{x}_{\pi(t)}^T$
- 15: Update $\hat{\alpha} \leftarrow \mathbf{V}^{-1} \mathbf{f}$
- 16: **end for**

SpectralTS algorithm

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Initialization:**
- 4: $v = R\sqrt{6d \log((\lambda + T)/\delta\lambda)} + C$
- 5: $\hat{\alpha} = 0_N$
- 6: $\mathbf{f} = 0_N$
- 7: $\mathbf{V} = \Lambda_{\mathcal{L}} + \lambda \mathbf{I}_N$
- 8: **Run:**
- 9: **for** $t = 1$ **to** T **do**
- 10: Sample $\tilde{\alpha} \sim \mathcal{N}(\hat{\alpha}, v^2 \mathbf{V}^{-1})$
- 11: $\pi(t) \leftarrow \arg \max_a \mathbf{x}_a^T \tilde{\alpha}$
- 12: Observe a noisy reward $r(t) = \mathbf{x}_{\pi(t)}^T \boldsymbol{\alpha}^* + \varepsilon_t$
- 13: $\mathbf{f} \leftarrow \mathbf{f} + \mathbf{x}_{\pi(t)} r(t)$
- 14: Update $\mathbf{V} \leftarrow \mathbf{V} + \mathbf{x}_{\pi(t)} \mathbf{x}_{\pi(t)}^T$
- 15: Update $\hat{\alpha} \leftarrow \mathbf{V}^{-1} \mathbf{f}$
- 16: **end for**

SpectralTS algorithm

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Initialization:**
- 4: $v = R\sqrt{6d \log((\lambda + T)/\delta\lambda)} + C$
- 5: $\hat{\alpha} = 0_N$
- 6: $\mathbf{f} = 0_N$
- 7: $\mathbf{V} = \Lambda_{\mathcal{L}} + \lambda \mathbf{I}_N$
- 8: **Run:**
- 9: **for** $t = 1$ **to** T **do**
- 10: **Sample** $\tilde{\alpha} \sim \mathcal{N}(\hat{\alpha}, v^2 \mathbf{V}^{-1})$
- 11: $\pi(t) \leftarrow \arg \max_a \mathbf{x}_a^T \tilde{\alpha}$
- 12: Observe a noisy reward $r(t) = \mathbf{x}_{\pi(t)}^T \boldsymbol{\alpha}^* + \varepsilon_t$
- 13: $\mathbf{f} \leftarrow \mathbf{f} + \mathbf{x}_{\pi(t)} r(t)$
- 14: Update $\mathbf{V} \leftarrow \mathbf{V} + \mathbf{x}_{\pi(t)} \mathbf{x}_{\pi(t)}^T$
- 15: Update $\hat{\alpha} \leftarrow \mathbf{V}^{-1} \mathbf{f}$
- 16: **end for**

SpectralTS algorithm

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Initialization:**
- 4: $v = R\sqrt{6d \log((\lambda + T)/\delta\lambda)} + C$
- 5: $\hat{\alpha} = 0_N$
- 6: $\mathbf{f} = 0_N$
- 7: $\mathbf{V} = \Lambda_{\mathcal{L}} + \lambda \mathbf{I}_N$
- 8: **Run:**
- 9: **for** $t = 1$ **to** T **do**
- 10: Sample $\tilde{\alpha} \sim \mathcal{N}(\hat{\alpha}, v^2 \mathbf{V}^{-1})$
- 11: $\pi(t) \leftarrow \arg \max_a \mathbf{x}_a^T \tilde{\alpha}$
- 12: Observe a noisy reward $r(t) = \mathbf{x}_{\pi(t)}^T \alpha^* + \varepsilon_t$
- 13: $\mathbf{f} \leftarrow \mathbf{f} + \mathbf{x}_{\pi(t)} r(t)$
- 14: Update $\mathbf{V} \leftarrow \mathbf{V} + \mathbf{x}_{\pi(t)} \mathbf{x}_{\pi(t)}^T$
- 15: Update $\hat{\alpha} \leftarrow \mathbf{V}^{-1} \mathbf{f}$
- 16: **end for**

SpectralTS algorithm

- 1: **Input:**
- 2: $N, T, \{\Lambda_{\mathcal{L}}, \mathbf{Q}\}, \lambda, \delta, R, C$
- 3: **Initialization:**
- 4: $v = R\sqrt{6d \log((\lambda + T)/\delta\lambda)} + C$
- 5: $\hat{\alpha} = 0_N$
- 6: $\mathbf{f} = 0_N$
- 7: $\mathbf{V} = \Lambda_{\mathcal{L}} + \lambda \mathbf{I}_N$
- 8: **Run:**
- 9: **for** $t = 1$ **to** T **do**
- 10: Sample $\tilde{\alpha} \sim \mathcal{N}(\hat{\alpha}, v^2 \mathbf{V}^{-1})$
- 11: $\pi(t) \leftarrow \arg \max_a \mathbf{x}_a^T \tilde{\alpha}$
- 12: Observe a noisy reward $r(t) = \mathbf{x}_{\pi(t)}^T \boldsymbol{\alpha}^* + \varepsilon_t$
- 13: $\mathbf{f} \leftarrow \mathbf{f} + \mathbf{x}_{\pi(t)} r(t)$
- 14: Update $\mathbf{V} \leftarrow \mathbf{V} + \mathbf{x}_{\pi(t)} \mathbf{x}_{\pi(t)}^T$
- 15: Update $\hat{\alpha} \leftarrow \mathbf{V}^{-1} \mathbf{f}$
- 16: **end for**

SpectralTS regret bound

- ▶ d : Effective dimension.
- ▶ λ : Minimal eigenvalue of $\mathbf{\Lambda} = \mathbf{\Lambda}_{\mathcal{L}} + \lambda \mathbf{I}$.
- ▶ C : Smoothness upper bound, $\|\boldsymbol{\alpha}^*\|_{\mathbf{\Lambda}} \leq C$.
- ▶ $\mathbf{x}_i^T \boldsymbol{\alpha}^* \in [-1, 1]$ for all i .

The **cumulative regret** R_T of **SpectralTS** is with probability $1 - \delta$ bounded as

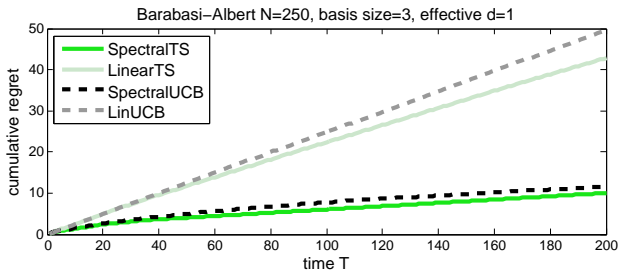
$$R_T \leq \frac{11g}{p} \sqrt{\frac{4 + 4\lambda}{\lambda} d T \log \frac{\lambda + T}{\lambda}} + \frac{1}{T} + \frac{g}{p} \left(\frac{11}{\sqrt{\lambda}} + 2 \right) \sqrt{2T \log \frac{2}{\delta}},$$

where $p = 1/(4e\sqrt{\pi})$ and

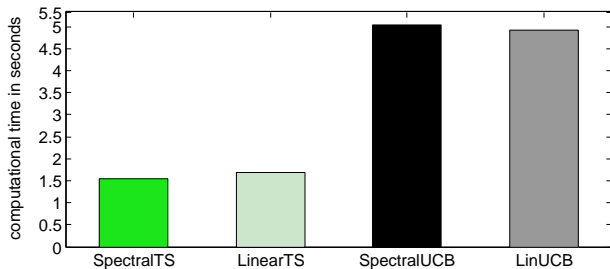
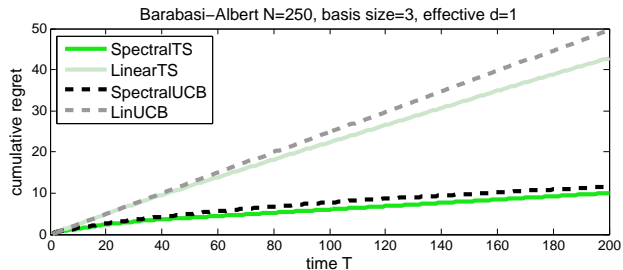
$$g = \sqrt{4 \log TN} \left(R \sqrt{6d \log \left(\frac{\lambda + T}{\delta \lambda} \right)} + C \right) + R \sqrt{2d \log \left(\frac{(\lambda + T) T^2}{\delta \lambda} \right)} + C.$$

$$R_T \approx d \sqrt{T \log N}$$

Synthetic experiment

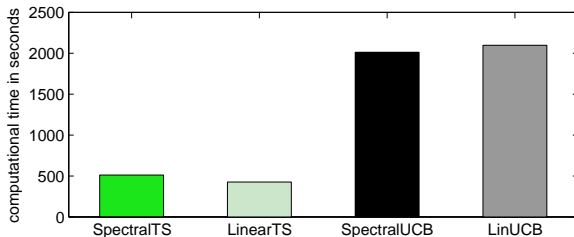
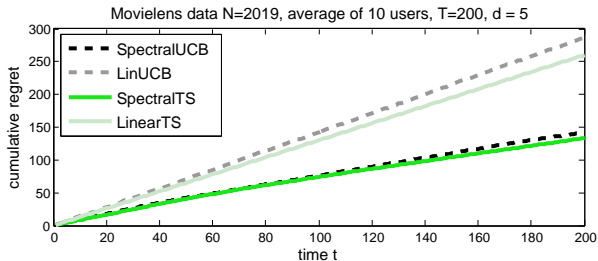


Synthetic experiment



Real world experiment

MovieLens dataset of 6k users who rated one million movies.



Spectral Bandits Summary

- ▶ New spectral bandit setting (**for smooth graph functions**).

Spectral Bandits Summary

- ▶ New spectral bandit setting (**for smooth graph functions**).
- ▶ **SpectralUCB**
 - ▶ Regret bound $\approx d\sqrt{T \ln T}$

Spectral Bandits Summary

- ▶ New spectral bandit setting (**for smooth graph functions**).
- ▶ **SpectralUCB**
 - ▶ Regret bound $\approx d\sqrt{T \ln T}$
- ▶ **SpectralTS**
 - ▶ Regret bound $\approx d\sqrt{T \ln N}$
 - ▶ Computationally more efficient.

Spectral Bandits Summary

- ▶ New spectral bandit setting (**for smooth graph functions**).
- ▶ **SpectralUCB**
 - ▶ Regret bound $\approx d\sqrt{T \ln T}$
- ▶ **SpectralTS**
 - ▶ Regret bound $\approx d\sqrt{T \ln N}$
 - ▶ Computationally more efficient.
- ▶ **SpectralEliminator**
 - ▶ Regret bound $\approx \sqrt{dT \ln T}$
 - ▶ Side result: **LinearEliminator** with $\mathcal{O}(\sqrt{DT \ln T})$ regret for (contextual) linear bandits.

Spectral Bandits Summary

- ▶ New spectral bandit setting (**for smooth graph functions**).
- ▶ **SpectralUCB**
 - ▶ Regret bound $\approx d\sqrt{T \ln T}$
- ▶ **SpectralTS**
 - ▶ Regret bound $\approx d\sqrt{T \ln N}$
 - ▶ Computationally more efficient.
- ▶ **SpectralEliminator**
 - ▶ Regret bound $\approx \sqrt{dT \ln T}$
 - ▶ Side result: **LinearEliminator** with $\mathcal{O}(\sqrt{DT \ln T})$ regret for (contextual) linear bandits.
- ▶ Bounds scale with **effective dimension** $d \ll D$.

Spectral Bandits Summary

- ▶ New spectral bandit setting (**for smooth graph functions**).
- ▶ **SpectralUCB**
 - ▶ Regret bound $\approx d\sqrt{T \ln T}$
- ▶ **SpectralTS**
 - ▶ Regret bound $\approx d\sqrt{T \ln N}$
 - ▶ Computationally more efficient.
- ▶ **SpectralEliminator**
 - ▶ Regret bound $\approx \sqrt{dT \ln T}$
 - ▶ Side result: **LinearEliminator** with $\mathcal{O}(\sqrt{DT \ln T})$ regret for (contextual) linear bandits.
- ▶ Bounds scale with **effective dimension** $d \ll D$.

Exploiting side observations

Example 1: undirected observations



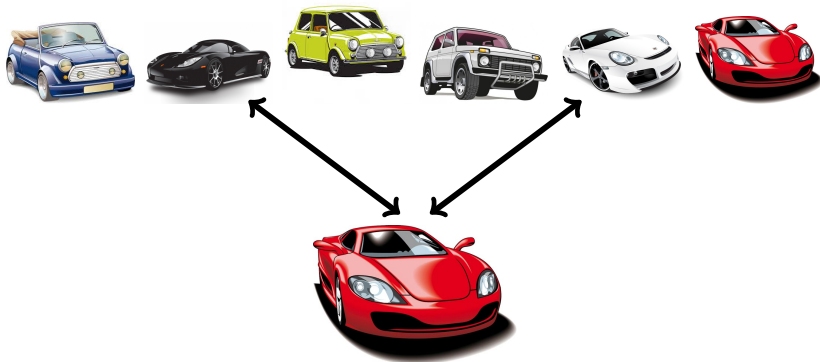
Exploiting side observations

Example 1: undirected observations

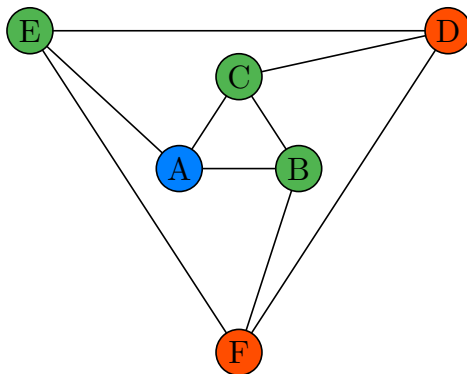


Exploiting side observations

Example 1: undirected observations



Example 1: Graph Representation



Example 2: Directed observation



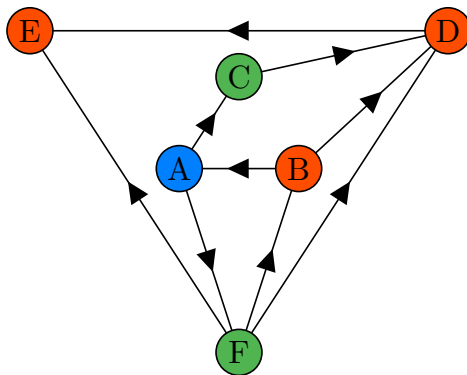
Example 2: Directed observation



Example 2: Directed observation



Example 2



Learning setting

In each time step $t = 1, \dots, T$

- ▶ **Environment (adversary):**
 - ▶ Privately assigns losses to actions
 - ▶ Generates an observation graph

Learning setting

In each time step $t = 1, \dots, T$

- ▶ **Environment (adversary):**
 - ▶ Privately assigns losses to actions
 - ▶ Generates an observation graph
 - ▶ Undirected / Directed

Learning setting

In each time step $t = 1, \dots, T$

- ▶ **Environment (adversary):**
 - ▶ Privately assigns losses to actions
 - ▶ Generates an observation graph
 - ▶ Undirected / Directed
 - ▶ Disclosed / Not disclosed

Learning setting

In each time step $t = 1, \dots, T$

▶ **Environment (adversary):**

- ▶ Privately assigns losses to actions
- ▶ Generates an observation graph
 - ▶ Undirected / Directed
 - ▶ Disclosed / Not disclosed

▶ **Learner:**

- ▶ Plays action $I_t \in [N]$
- ▶ Obtain loss ℓ_{t,I_t} of action played
- ▶ Observe losses of neighbors of I_t

Learning setting

In each time step $t = 1, \dots, T$

▶ **Environment (adversary):**

- ▶ Privately assigns losses to actions
- ▶ Generates an observation graph
 - ▶ Undirected / Directed
 - ▶ Disclosed / Not disclosed

▶ **Learner:**

- ▶ Plays action $I_t \in [N]$
- ▶ Obtain loss ℓ_{t,I_t} of action played
- ▶ Observe losses of neighbors of I_t
 - ▶ **Graph: disclosed**

Learning setting

In each time step $t = 1, \dots, T$

▶ **Environment (adversary):**

- ▶ Privately assigns losses to actions
- ▶ Generates an observation graph
 - ▶ Undirected / Directed
 - ▶ Disclosed / Not disclosed

▶ **Learner:**

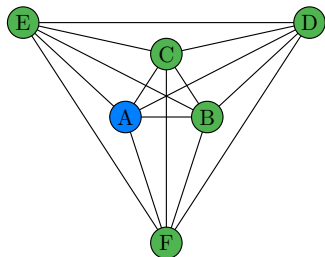
- ▶ Plays action $I_t \in [N]$
- ▶ Obtain loss ℓ_{t,I_t} of action played
- ▶ Observe losses of neighbors of I_t
 - ▶ **Graph: disclosed**

▶ **Performance measure:** Total expected regret

$$R_T = \max_{i \in [N]} \mathbb{E} \left[\sum_{t=1}^T (\ell_{t,I_t} - \ell_{t,i}) \right]$$

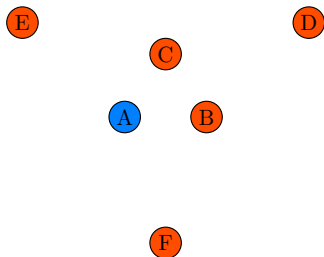
Full Information setting

- ▶ Pick an action (e.g. action A)
- ▶ Observe losses of all actions
- ▶ $R_T = \tilde{O}(\sqrt{T})$



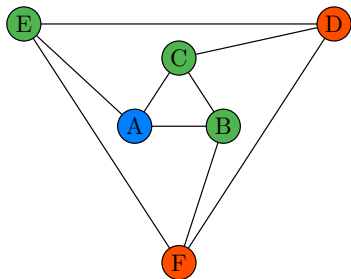
Bandit setting

- ▶ Pick an action (e.g. action A)
- ▶ Observe loss of a chosen action
- ▶ $R_T = \tilde{O}(\sqrt{NT})$



Side observation (Undirected case)

- ▶ Pick an action (e.g. action A)
- ▶ Observe losses of neighbors

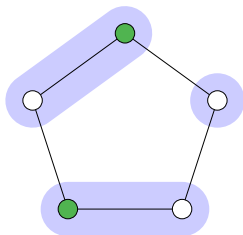
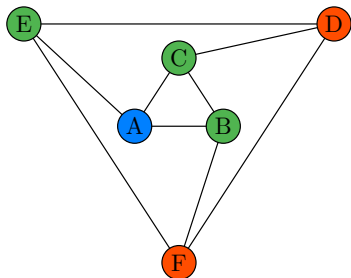


Side observation (Undirected case)

- ▶ Pick an action (e.g. action A)
- ▶ Observe losses of neighbors

Mannor and Shamir (ELP algorithm)

- ▶ Need to know graph
- ▶ Clique decomposition (c cliques)
- ▶ $R_T = \tilde{O}(\sqrt{cT})$



Side observation (Undirected case)

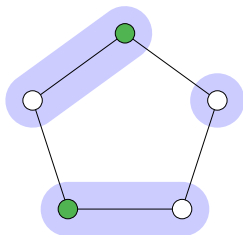
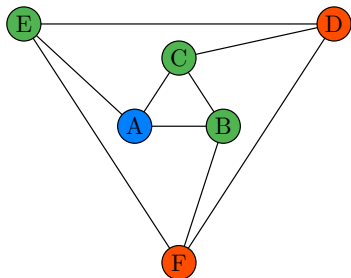
- ▶ Pick an action (e.g. action A)
- ▶ Observe losses of neighbors

Mannor and Shamir (ELP algorithm)

- ▶ Need to know graph
- ▶ Clique decomposition (c cliques)
- ▶ $R_T = \tilde{O}(\sqrt{cT})$

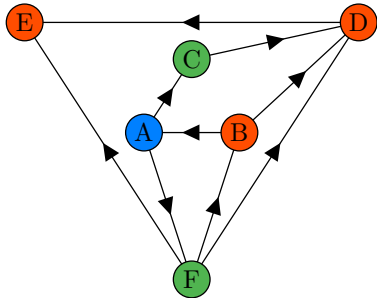
Alon, Cesa-Bianchi, Gentile, Mansour

- ▶ No need to know graph
- ▶ Independence set of α actions
- ▶ $R_T = \tilde{O}(\sqrt{\alpha T})$



Side observation (Directed case)

- ▶ Pick an action (e.g. action A)
- ▶ Observe losses of neighbors

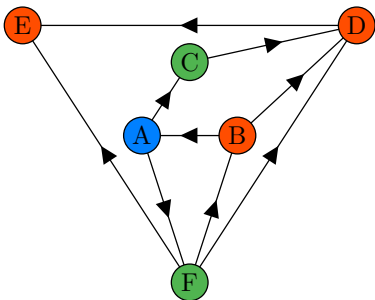


Side observation (Directed case)

- ▶ Pick an action (e.g. action A)
- ▶ Observe losses of neighbors

Alon, Cesa-Bianchi, Gentile, Mansour

- ▶ **Exp3-DOM**
- ▶ Need to know graph
- ▶ Need to find dominating set
- ▶ $R_T = \tilde{O}(\sqrt{\alpha T})$



Side observation (Directed case)

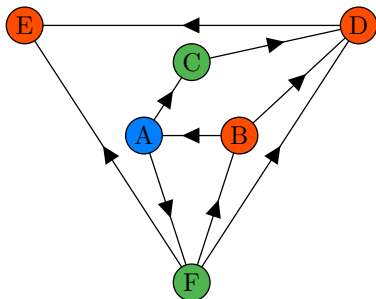
- ▶ Pick an action (e.g. action A)
- ▶ Observe losses of neighbors

Alon, Cesa-Bianchi, Gentile, Mansour

- ▶ **Exp3-DOM**
- ▶ Need to know graph
- ▶ Need to find dominating set
- ▶ $R_T = \tilde{O}(\sqrt{\alpha T})$

Our solution: Exp3-IX

- ▶ No need to know graph
- ▶ $R_T = \tilde{O}(\sqrt{\alpha T})$



Exp3 algorithms in general

- ▶ **Compute weights** using loss estimates $\hat{\ell}_{t,i}$.

$$w_{t,i} = \exp \left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_{s,i} \right)$$

- ▶ **Play action** I_t such that

$$\mathbb{P}(I_t = i) = p_{t,i} = \frac{w_{t,i}}{W_t} = \frac{w_{t,i}}{\sum_{j=1}^N w_{t,j}}$$

- ▶ **Update loss estimates** (using observability graph)

Exp3 algorithms in general

- ▶ **Compute weights** using loss estimates $\hat{\ell}_{t,i}$.

$$w_{t,i} = \exp \left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_{s,i} \right)$$

- ▶ **Play action** I_t such that

$$\mathbb{P}(I_t = i) = p_{t,i} = \frac{w_{t,i}}{W_t} = \frac{w_{t,i}}{\sum_{j=1}^N w_{t,j}}$$

- ▶ **Update loss estimates** (using observability graph)

How the algorithms approach to bias variance tradeoff?

Bias variance tradeoff approaches

- ▶ Approach of previous algorithms – **Mixing**
 - ▶ Bias sampling distribution \mathbf{p}_t over actions
 - ▶ $\mathbf{p}'_t = (1 - \gamma)\mathbf{p}_t + \gamma\mathbf{s}_t$ – mixed distribution
 - ▶ \mathbf{s}_t – probability distribution which supports exploration
 - ▶ Loss estimates $\hat{\ell}_{t,i}$ are unbiased

- ▶ Approach of our algorithm – **Implicit eXploration (IX)**
 - ▶ Bias loss estimates $\hat{\ell}_{t,i}$
 - ▶ Biased loss estimates \implies biased weights
 - ▶ Biased weights \implies biased probability distribution
 - ▶ No need for mixing

Mannor and Shamir - ELP algorithm

- ▶ $\mathbb{E}[\hat{\ell}_{t,i}] = \ell_{t,i}$ – unbiased loss estimates
- ▶ $p'_{t,i} = (1 - \gamma)p_{t,i} + \gamma s_{t,i}$ – bias by mixing
- ▶ $\mathbf{s}_t = \{s_{t,1}, \dots, s_{t,N}\}$ – probability distribution over the action set

$$\mathbf{s}_t = \arg \max_{\mathbf{s}_t} \left[\min_{j \in [N]} \left(s_{t,j} + \sum_{k \in N_{t,j}} s_{t,k} \right) \right] = \arg \max_{\mathbf{s}_t} \left[\min_{j \in [N]} q_{t,j} \right]$$

- ▶ $q_{t,j}$ – probability that loss of j is observed according to \mathbf{s}_t

Mannor and Shamir - ELP algorithm

- ▶ $\mathbb{E}[\hat{\ell}_{t,i}] = \ell_{t,i}$ – unbiased loss estimates
- ▶ $p'_{t,i} = (1 - \gamma)p_{t,i} + \gamma s_{t,i}$ – bias by mixing
- ▶ $\mathbf{s}_t = \{s_{t,1}, \dots, s_{t,N}\}$ – probability distribution over the action set

$$\mathbf{s}_t = \arg \max_{\mathbf{s}_t} \left[\min_{j \in [N]} \left(s_{t,j} + \sum_{k \in N_{t,j}} s_{t,k} \right) \right] = \arg \max_{\mathbf{s}_t} \left[\min_{j \in [N]} q_{t,j} \right]$$

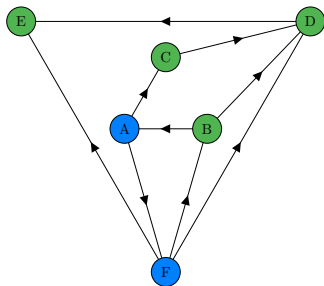
- ▶ $q_{t,j}$ – probability that loss of j is observed according to \mathbf{s}_t
- ▶ **Computation of \mathbf{s}_t**
 - ▶ Graph needs to be disclosed
 - ▶ Solving simple linear program
- ▶ Needs to know graph before playing an action
- ▶ Graphs can be only undirected

Alon, Cesa-Bianchi, Gentile, Mansour - Exp3-DOM

- ▶ $\mathbb{E}[\hat{\ell}_{t,i}] = \ell_{t,i}$ – unbiased loss estimates
- ▶ $p'_{t,i} = (1 - \gamma)p_{t,i} + \gamma s_{t,i}$ – bias by mixing
- ▶ $\mathbf{s}_t = \{s_{t,1}, \dots, s_{t,N}\}$ – probability distribution over the action set

$$s_{t,i} = \begin{cases} \frac{1}{r} & \text{if } i \in R; |R| = r \\ 0 & \text{otherwise.} \end{cases}$$

- ▶ R – dominating set of r elements
- ▶ \mathbf{s}_t – uniform distribution over R
- ▶ Needs to know graph beforehand
- ▶ Graphs can be directed



Previous algorithms - loss estimates

$$\hat{\ell}_{t,i} = \begin{cases} \ell_{t,i}/o_{t,i} & \text{if } \ell_{t,i} \text{ is observed} \\ 0 & \text{otherwise.} \end{cases}$$

$$\mathbb{E}[\hat{\ell}_{t,i}] = \frac{\ell_{t,i}}{o_{t,i}} o_{t,i} + 0(1 - o_{t,i}) = \ell_{t,i}$$

Previous algorithms - loss estimates

$$\hat{\ell}_{t,i} = \begin{cases} \ell_{t,i}/o_{t,i} & \text{if } \ell_{t,i} \text{ is observed} \\ 0 & \text{otherwise.} \end{cases}$$

$$\mathbb{E}[\hat{\ell}_{t,i}] = \frac{\ell_{t,i}}{o_{t,i}} o_{t,i} + 0(1 - o_{t,i}) = \ell_{t,i}$$

Exp3-IX - loss estimates

$$\hat{\ell}_{t,i} = \begin{cases} \ell_{t,i}/(o_{t,i} + \gamma) & \text{if } \ell_{t,i} \text{ is observed} \\ 0 & \text{otherwise.} \end{cases}$$

$$\mathbb{E}[\hat{\ell}_{t,i}] = \frac{\ell_{t,i}}{o_{t,i} + \gamma} o_{t,i} + 0(1 - o_{t,i}) = \ell_{t,i} - \ell_{t,i} \frac{\gamma}{o_{t,i} + \gamma} \leq \ell_{t,i}$$

► No mixing!

Analysis of Exp3 algorithms in general

- Evolution of W_{t+1}/W_t

$$\frac{1}{\eta} \log \frac{W_{t+1}}{W_t} = \frac{1}{\eta} \log \left(1 - \eta \sum_{i=1}^N p_{t,i} \hat{\ell}_{t,i} + \frac{\eta^2}{2} \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right),$$

$$\sum_{i=1}^N p_{t,i} \hat{\ell}_{t,i} \leq \left[\frac{\log W_t}{\eta} - \frac{\log W_{t+1}}{\eta} \right] + \frac{\eta}{2} \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2$$

- Taking expectation and summing over time

$$\mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} \hat{\ell}_{t,i} \right] - \mathbb{E} \left[\sum_{t=1}^T \hat{\ell}_{t,k} \right] \leq \mathbb{E} \left[\frac{\log N}{\eta} \right] + \mathbb{E} \left[\frac{\eta}{2} \sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right]$$

Regret bound of Exp3-IX

$$\underbrace{\mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} \hat{\ell}_{t,i} \right]}_A - \underbrace{\mathbb{E} \left[\sum_{t=1}^T \hat{\ell}_{t,k} \right]}_B \leq \mathbb{E} \left[\frac{\log N}{\eta} \right] + \underbrace{\mathbb{E} \left[\frac{\eta}{2} \sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right]}_C$$

Lower bound of A (using definition of loss estimates)

$$\mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} \hat{\ell}_{t,i} \right] \geq \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} \ell_{t,i} \right] - \mathbb{E} \left[\gamma \sum_{t=1}^T \sum_{i=1}^N \frac{p_{t,i}}{\sigma_{t,i} + \gamma} \right]$$

Lower bound of B (optimistic loss estimates: $\mathbb{E}[\hat{\ell}] < \mathbb{E}[\ell]$)

$$-\mathbb{E} \left[\sum_{t=1}^T \hat{\ell}_{t,k} \right] \geq -\mathbb{E} \left[\sum_{t=1}^T \ell_{t,k} \right]$$

Upper bound of C (using definition of loss estimates)

$$\mathbb{E} \left[\frac{\eta}{2} \sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right] \leq \mathbb{E} \left[\frac{\eta}{2} \sum_{t=1}^T \sum_{i=1}^N \frac{p_{t,i}}{\sigma_{t,i} + \gamma} \right]$$

Regret bound of Exp3-IX

$$R_T \leq \frac{\log N}{\eta} + \left(\frac{\eta}{2} + \gamma\right) \sum_{t=1}^T \mathbb{E} \left[\sum_{i=1}^N \frac{p_{t,i}}{o_{t,i} + \gamma} \right]$$

$$R_T \approx \mathcal{O} \left(\sqrt{\log N \sum_{t=1}^T \mathbb{E} \left[\sum_{i=1}^N \frac{p_{t,i}}{o_{t,i} + \gamma} \right]} \right)$$

Regret bound of Exp3-IX

$$R_T \leq \frac{\log N}{\eta} + \left(\frac{\eta}{2} + \gamma\right) \sum_{t=1}^T \mathbb{E} \left[\sum_{i=1}^N \frac{p_{t,i}}{o_{t,i} + \gamma} \right]$$

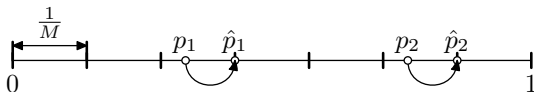
$$R_T \approx \mathcal{O} \left(\sqrt{\log N \sum_{t=1}^T \mathbb{E} \left[\sum_{i=1}^N \frac{p_{t,i}}{o_{t,i} + \gamma} \right]} \right)$$

Graph lemma

- ▶ Graph G with $V(G) = \{1, \dots, N\}$
- ▶ d_i^- – in-degree of vertex i
- ▶ α – independence set of G
- ▶ Turán's Theorem + induction

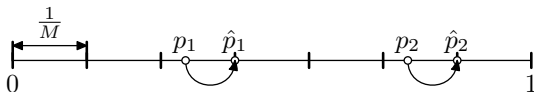
$$\sum_{i=1}^N \frac{1}{1 + d_i^-} \leq 2\alpha \log \left(1 + \frac{N}{\alpha} \right)$$

Discretization



$$\sum_{i=1}^N \frac{p_{t,i}}{o_{t,i} + \gamma} = \sum_{i=1}^N \frac{p_{t,i}}{p_{t,i} + \sum_{j \in N_i^-} p_{t,j} + \gamma} \leq \sum_{i=1}^N \frac{\hat{p}_{t,i}}{\hat{p}_{t,i} + \sum_{j \in N_i^-} \hat{p}_{t,j}} + 2$$

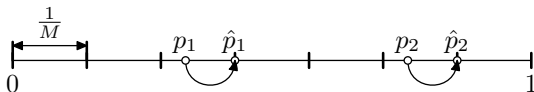
Discretization



$$\sum_{i=1}^N \frac{p_{t,i}}{o_{t,i} + \gamma} = \sum_{i=1}^N \frac{p_{t,i}}{p_{t,i} + \sum_{j \in N_i^-} p_{t,j} + \gamma} \leq \sum_{i=1}^N \frac{\hat{p}_{t,i}}{\hat{p}_{t,i} + \sum_{j \in N_i^-} \hat{p}_{t,j}} + 2$$

Note: we set $M = \lceil N^2/\gamma \rceil$

Discretization



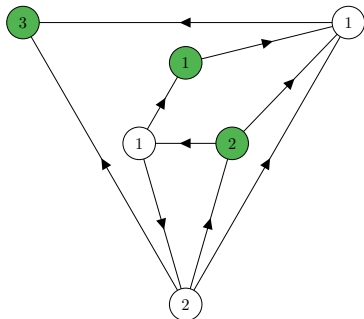
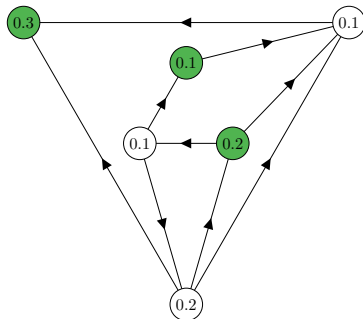
$$\sum_{i=1}^N \frac{p_{t,i}}{o_{t,i} + \gamma} = \sum_{i=1}^N \frac{p_{t,i}}{p_{t,i} + \sum_{j \in N_i^-} p_{t,j} + \gamma} \leq \sum_{i=1}^N \frac{\hat{p}_{t,i}}{\hat{p}_{t,i} + \sum_{j \in N_i^-} \hat{p}_{t,j}} + 2$$

Note: we set $M = \lceil N^2/\gamma \rceil$

$$\sum_{i=1}^N \frac{\hat{p}_{t,i}}{\hat{p}_{t,i} + \sum_{j \in N_i^-} \hat{p}_{t,j}}$$

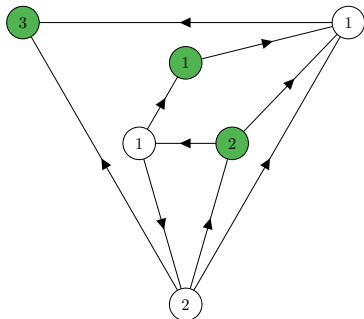
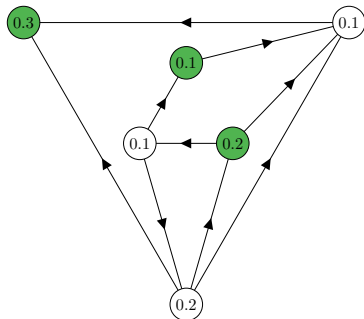
$$\sum_{i=1}^N \frac{\hat{p}_{t,i}}{\hat{p}_{t,i} + \sum_{j \in N_i^-} \hat{p}_{t,j}}$$

Example: let $M = 10$



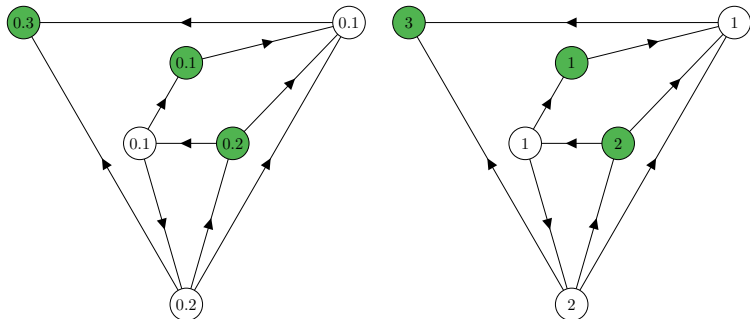
$$\sum_{i=1}^N \frac{M \hat{p}_{t,i}}{M \hat{p}_{t,i} + \sum_{j \in N_i^-} M \hat{p}_{t,j}}$$

Example: let $M = 10$



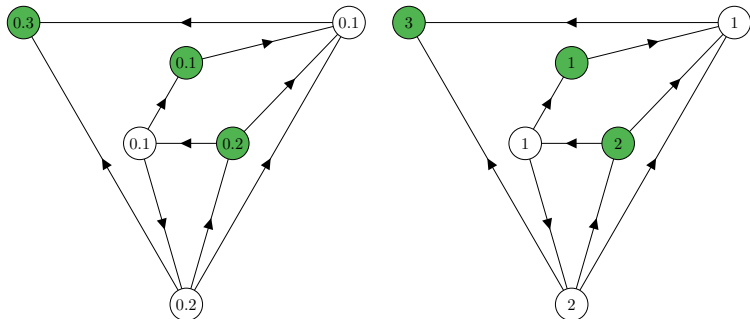
$$\sum_{i=1}^N \frac{M \hat{p}_{t,i}}{M \hat{p}_{t,i} + \sum_{j \in N_i^-} M \hat{p}_{t,j}} = \sum_{i=1}^N \sum_{k \in C_i} \frac{1}{1 + d_k^-}$$

Example: let $M = 10$



$$\sum_{i=1}^N \frac{M \hat{p}_{t,i}}{M \hat{p}_{t,i} + \sum_{j \in N_i^-} M \hat{p}_{t,j}} = \sum_{i=1}^N \sum_{k \in C_i} \frac{1}{1 + d_k^-} \leq 2\alpha \log \left(1 + \frac{M + N}{\alpha} \right)$$

Example: let $M = 10$



Exp3-IX regret bound

$$R_T \leq \frac{\log N}{\eta} + \left(\frac{\eta}{2} + \gamma\right) \sum_{t=1}^T \mathbb{E} \left[2\alpha_t \log \left(1 + \frac{\lceil N^2/\gamma \rceil + N}{\alpha_t} \right) + 2 \right]$$

$$R_T = \tilde{O} \left(\sqrt{\bar{\alpha} T \log(N)} \right)$$

Exp3-IX regret bound

$$R_T \leq \frac{\log N}{\eta} + \left(\frac{\eta}{2} + \gamma\right) \sum_{t=1}^T \mathbb{E} \left[2\alpha_t \log \left(1 + \frac{\lceil N^2/\gamma \rceil + N}{\alpha_t} \right) + 2 \right]$$

$$R_T = \tilde{O} \left(\sqrt{\bar{\alpha} T \log(N)} \right)$$

Next step

Exp3-IX regret bound

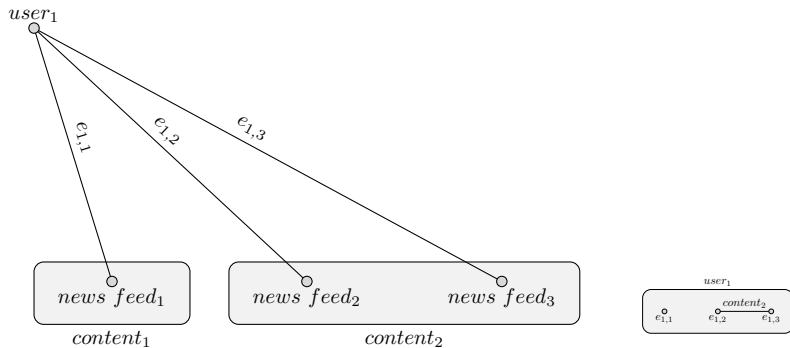
$$R_T \leq \frac{\log N}{\eta} + \left(\frac{\eta}{2} + \gamma\right) \sum_{t=1}^T \mathbb{E} \left[2\alpha_t \log \left(1 + \frac{\lceil N^2/\gamma \rceil + N}{\alpha_t} \right) + 2 \right]$$

$$R_T = \tilde{O} \left(\sqrt{\bar{\alpha} T \log(N)} \right)$$

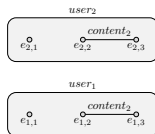
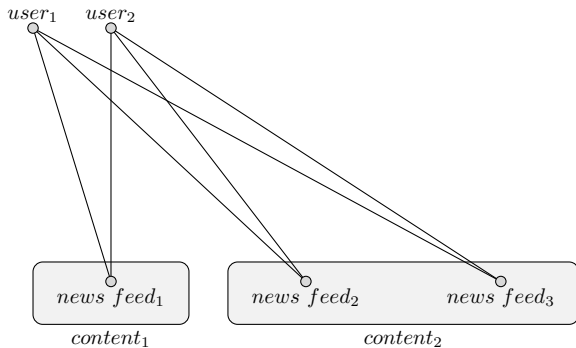
Next step

Generalization of the setting to **combinatorial actions**

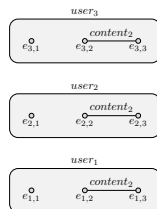
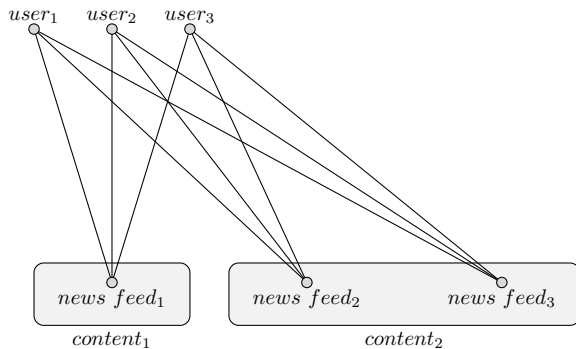
Example



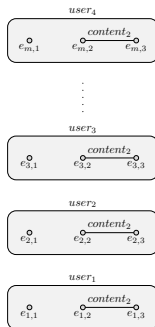
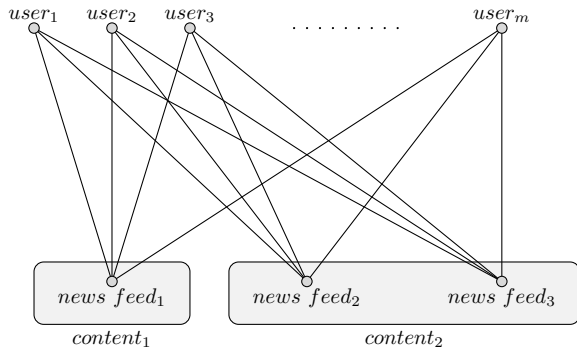
Example



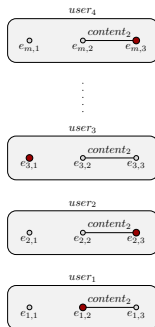
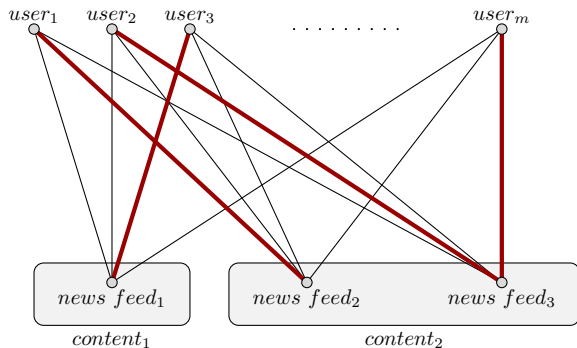
Example



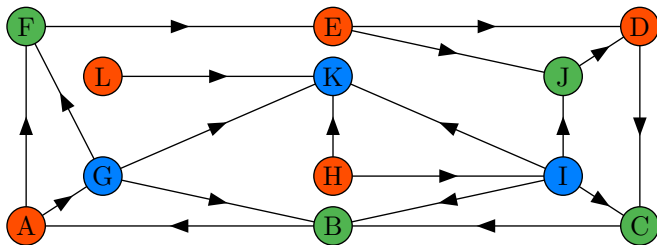
Example



Example



- ▶ Play m out of N nodes (combinatorial structure)
- ▶ Obtain losses of all played nodes
- ▶ Observe losses of all neighbors of played nodes



- ▶ Play action $\mathbf{V}_t \in S \subset \{0, 1\}^N$, $\|\mathbf{v}\|_1 \leq m$ for all $\mathbf{v} \in S$
- ▶ Obtain losses $\mathbf{V}_t^\top \ell_t$
- ▶ Observe additional losses according to the graph

FPL-IX algorithm

- ▶ Draw perturbation $Z_{t,i} \sim \text{Exp}(1)$ for all $i \in [M]$
- ▶ Play “the best” action \mathbf{V}_t according to total loss estimate $\hat{\mathbf{L}}_{t-1}$ and perturbation \mathbf{Z}_t

$$\mathbf{V}_t = \arg \min_{\mathbf{v} \in \mathcal{S}} \mathbf{v}^\top \left(\eta_t \hat{\mathbf{L}}_{t-1} - \mathbf{Z}_t \right)$$

- ▶ Compute loss estimates

$$\hat{\ell}_{t,i} = \ell_{t,i} K_{t,i} \mathbb{1}\{\ell_{t,i} \text{ is observed}\}$$

- ▶ $K_{t,i}$: geometric random variable with

$$\mathbb{E}[K_{t,i}] = \frac{1}{o_{t,i} + (1 - o_{t,i})\gamma}$$

FPL-IX - regret bound

$$R_T = \tilde{O} \left(m^{3/2} \sqrt{\sum_{t=1}^T \alpha_t} \right) = \tilde{O} \left(m^{3/2} \sqrt{\bar{\alpha} T} \right)$$

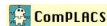
Side Observation Summary

- ▶ **Implicit eXploration** idea
- ▶ **New algorithm for simple actions - Exp3-IX**
 - ▶ Using implicit exploration idea
 - ▶ Same regret bound as previous algorithm
 - ▶ No need to know graph before an action is played
 - ▶ Computationally efficient
- ▶ **New combinatorial setting with side observations**
- ▶ **Algorithm for combinatorial setting - FPL-IX**
- ▶ **Future directions**
 - ▶ No need to know graph after an action is played
 - ▶ Stochastic side observations
 - ▶ Random graph models
 - ▶ Exploiting the communities



Microsoft
Research

technicolor



Michal Valko

michal.valko@inria.fr

sequel.lille.inria.fr

Sylvester's determinant theorem:

$$|\mathbf{A} + \mathbf{x}\mathbf{x}^\top| = |\mathbf{A}||\mathbf{I} + \mathbf{A}^{-1}\mathbf{x}\mathbf{x}^\top| = |\mathbf{A}|(1 + \mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x})$$

Goal:

- ▶ Upperbound determinant $|\mathbf{A} + \mathbf{x}\mathbf{x}^\top|$ for $\|\mathbf{x}\|_2 \leq 1$
- ▶ Upperbound $\mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x}$

Sylvester's determinant theorem:

$$|\mathbf{A} + \mathbf{x}\mathbf{x}^\top| = |\mathbf{A}||\mathbf{I} + \mathbf{A}^{-1}\mathbf{x}\mathbf{x}^\top| = |\mathbf{A}|(1 + \mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x})$$

Goal:

- ▶ Upperbound determinant $|\mathbf{A} + \mathbf{x}\mathbf{x}^\top|$ for $\|\mathbf{x}\|_2 \leq 1$
- ▶ Upperbound $\mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x}$

$$\mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x} = \mathbf{x}^\top\mathbf{Q}\mathbf{\Lambda}^{-1}\mathbf{Q}^\top\mathbf{x} = \mathbf{y}^\top\mathbf{\Lambda}^{-1}\mathbf{y} = \sum_{i=1}^N \lambda_i y_i^2$$

Sylvester's determinant theorem:

$$|\mathbf{A} + \mathbf{x}\mathbf{x}^\top| = |\mathbf{A}||\mathbf{I} + \mathbf{A}^{-1}\mathbf{x}\mathbf{x}^\top| = |\mathbf{A}|(1 + \mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x})$$

Goal:

- ▶ Upperbound determinant $|\mathbf{A} + \mathbf{x}\mathbf{x}^\top|$ for $\|\mathbf{x}\|_2 \leq 1$
- ▶ Upperbound $\mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x}$

$$\mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x} = \mathbf{x}^\top\mathbf{Q}\mathbf{\Lambda}^{-1}\mathbf{Q}^\top\mathbf{x} = \mathbf{y}^\top\mathbf{\Lambda}^{-1}\mathbf{y} = \sum_{i=1}^N \lambda_i y_i^2$$

- ▶ $\|\mathbf{y}\|_2 \leq 1$.
- ▶ \mathbf{y} is a canonical vector.
- ▶ $\mathbf{x} = \mathbf{Q}\mathbf{y}$ is an eigenvector of \mathbf{A} .

Corollary:

Determinant $|\mathbf{V}_T|$ of $\mathbf{V}_T = \mathbf{\Lambda} + \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top$ is maximized when all \mathbf{x}_t are aligned with axes.

$$\begin{aligned}
 |\mathbf{V}_T| &\leq \max_{\sum t_i = T} \prod (\lambda_i + t_i) \\
 \ln \frac{|\mathbf{V}_T|}{|\mathbf{\Lambda}|} &\leq \max_{\sum t_i = T} \sum \ln \left(1 + \frac{t_i}{\lambda_i} \right) \\
 \ln \frac{|\mathbf{V}_T|}{|\mathbf{\Lambda}|} &\leq \sum_{i=1}^d \ln \left(1 + \frac{T}{\lambda} \right) + \sum_{i=d+1}^N \ln \left(1 + \frac{t_i}{\lambda_{d+1}} \right) \\
 &\leq d \ln \left(1 + \frac{T}{\lambda} \right) + \frac{T}{\lambda_{d+1}} \\
 &\leq 2d \ln \left(1 + \frac{T}{\lambda} \right)
 \end{aligned}$$

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

Proof:

$$\begin{aligned} \mathbf{f}^\top \mathcal{L} \mathbf{f} &= \mathbf{f}^\top \mathcal{D} \mathbf{f} - \mathbf{f}^\top \mathcal{W} \mathbf{f} = \sum_{i=1}^N d_i f_i^2 - \sum_{i,j \leq N} w_{i,j} f_i f_j \\ &= \frac{1}{2} \left(\sum_{i=1}^N d_i f_i^2 - 2 \sum_{i,j \leq N} w_{i,j} f_i f_j + \sum_{j=1}^N d_j f_j^2 \right) = \frac{1}{2} \sum_{i,j \leq N} w_{i,j} (f_i - f_j)^2 \end{aligned}$$

SpectralUCB analysis sketch

- ▶ Derivation of the confidence ellipsoid for $\hat{\alpha}$ with probability $1 - \delta$.
 - ▶ Using analysis of OFUL (Abbasi-Yadkori et al., 2011)

$$|\mathbf{x}^\top(\hat{\alpha} - \alpha^*)| \leq \|\mathbf{x}\|_{\mathbf{V}_t^{-1}} \left(R \sqrt{2 \ln \left(\frac{|\mathbf{V}_t|^{1/2}}{\delta |\mathbf{\Lambda}|^{1/2}} \right)} + C \right)$$

- ▶ Regret in one time step: $r_t = \mathbf{x}_*^\top \alpha^* - \mathbf{x}_{\pi(t)}^\top \alpha^* \leq 2c_t \|\mathbf{x}_{\pi(t)}\|_{\mathbf{V}_t^{-1}}$
- ▶ Cumulative regret:

$$R_T = \sum_{t=1}^T r_t \leq \sqrt{T \sum_{t=1}^T r_t^2} \leq 2(c_T + 1) \sqrt{2T \ln \frac{|\mathbf{V}_T|}{|\mathbf{\Lambda}|}}$$

SpectralUCB analysis sketch

- ▶ Derivation of the confidence ellipsoid for $\hat{\alpha}$ with probability $1 - \delta$.
 - ▶ Using analysis of OFUL (Abbasi-Yadkori et al., 2011)

$$|\mathbf{x}^\top(\hat{\alpha} - \alpha^*)| \leq \|\mathbf{x}\|_{\mathbf{V}_t^{-1}} \left(R \sqrt{2 \ln \left(\frac{|\mathbf{V}_t|^{1/2}}{\delta |\mathbf{\Lambda}|^{1/2}} \right)} + C \right)$$

- ▶ Regret in one time step: $r_t = \mathbf{x}_*^\top \alpha^* - \mathbf{x}_{\pi(t)}^\top \alpha^* \leq 2c_t \|\mathbf{x}_{\pi(t)}\|_{\mathbf{V}_t^{-1}}$
- ▶ Cumulative regret:

$$R_T = \sum_{t=1}^T r_t \leq \sqrt{T \sum_{t=1}^T r_t^2} \leq 2(c_T + 1) \sqrt{2T \ln \frac{|\mathbf{V}_T|}{|\mathbf{\Lambda}|}}$$

SpectralUCB analysis sketch

- ▶ Derivation of the confidence ellipsoid for $\hat{\alpha}$ with probability $1 - \delta$.
 - ▶ Using analysis of OFUL (Abbasi-Yadkori et al., 2011)

$$|\mathbf{x}^\top(\hat{\alpha} - \alpha^*)| \leq \|\mathbf{x}\|_{\mathbf{V}_t^{-1}} \left(R \sqrt{2 \ln \left(\frac{|\mathbf{V}_t|^{1/2}}{\delta |\mathbf{\Lambda}|^{1/2}} \right)} + C \right)$$

- ▶ Regret in one time step: $r_t = \mathbf{x}_*^\top \alpha^* - \mathbf{x}_{\pi(t)}^\top \alpha^* \leq 2c_t \|\mathbf{x}_{\pi(t)}\|_{\mathbf{V}_t^{-1}}$
- ▶ Cumulative regret:

$$R_T = \sum_{t=1}^T r_t \leq \sqrt{T \sum_{t=1}^T r_t^2} \leq 2(c_T + 1) \sqrt{2T \ln \frac{|\mathbf{V}_T|}{|\mathbf{\Lambda}|}}$$

- ▶ Upperbound for $\ln(|\mathbf{V}_t|/|\mathbf{\Lambda}|)$

$$\ln \frac{|\mathbf{V}_t|}{|\mathbf{\Lambda}|} \leq \ln \frac{|\mathbf{V}_T|}{|\mathbf{\Lambda}|} \leq 2d \ln \left(\frac{\lambda + T}{\lambda} \right)$$

SpectralTS analysis sketch

Divide arms into two groups

- ▶ $\Delta_i = \mathbf{b}_*^T \boldsymbol{\mu} - \mathbf{b}_i^T \boldsymbol{\mu} \leq g \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$ arm i is **unsaturated**
- ▶ $\Delta_i = \mathbf{b}_*^T \boldsymbol{\mu} - \mathbf{b}_i^T \boldsymbol{\mu} > g \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$ arm i is **saturated**

SpectralTS analysis sketch

Divide arms into two groups

- ▶ $\Delta_i = \mathbf{b}_*^T \boldsymbol{\mu} - \mathbf{b}_i^T \boldsymbol{\mu} \leq g \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$ arm i is **unsaturated**
- ▶ $\Delta_i = \mathbf{b}_*^T \boldsymbol{\mu} - \mathbf{b}_i^T \boldsymbol{\mu} > g \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$ arm i is **saturated**

Saturated arm

- ▶ Small standard deviation \rightarrow accurate regret estimate.
- ▶ **High regret** on playing the arm \rightarrow **Low probability** of picking

SpectralTS analysis sketch

Divide arms into two groups

- ▶ $\Delta_i = \mathbf{b}_*^T \boldsymbol{\mu} - \mathbf{b}_i^T \boldsymbol{\mu} \leq g \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$ arm i is **unsaturated**
- ▶ $\Delta_i = \mathbf{b}_*^T \boldsymbol{\mu} - \mathbf{b}_i^T \boldsymbol{\mu} > g \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$ arm i is **saturated**

Saturated arm

- ▶ Small standard deviation \rightarrow accurate regret estimate.
- ▶ **High regret** on playing the arm \rightarrow **Low probability** of picking

Unsaturated arm

- ▶ **Low regret** bounded by a factor of standard deviation
- ▶ **High probability** of picking

SpectralTS analysis sketch

- ▶ Confidence ellipsoid for estimate $\hat{\boldsymbol{\mu}}$ of $\boldsymbol{\mu}$ (with probability $1 - \delta/T^2$)
 - ▶ Using analysis of OFUL algorithm (Abbasi-Yadkori et al., 2011)

$$|\mathbf{b}_i^\top \hat{\boldsymbol{\mu}} - \mathbf{b}_i^\top \boldsymbol{\mu}| \leq \left(R \sqrt{2 \log \left(\frac{|\mathbf{B}_T|^{1/2} T^2}{|\boldsymbol{\Lambda}|^{1/2} \delta} \right)} + C \right) \|\mathbf{b}_i\|_{\mathbf{B}_T^{-1}}$$

SpectralTS analysis sketch

- ▶ Confidence ellipsoid for estimate $\hat{\boldsymbol{\mu}}$ of $\boldsymbol{\mu}$ (with probability $1 - \delta/T^2$)
 - ▶ Using analysis of OFUL algorithm (Abbasi-Yadkori et al., 2011)

$$|\mathbf{b}_i^\top \hat{\boldsymbol{\mu}} - \mathbf{b}_i^\top \boldsymbol{\mu}| \leq \left(R \sqrt{2 \log \left(\frac{|\mathbf{B}_T|^{1/2} T^2}{|\boldsymbol{\Lambda}|^{1/2} \delta} \right)} + C \right) \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$$

- ▶ Our key result coming from spectral properties of \mathbf{B}_t .

$$\log \frac{|\mathbf{B}_t|}{|\boldsymbol{\Lambda}|} \leq 2d \log \left(1 + \frac{T}{\lambda} \right)$$

SpectralTS analysis sketch

- ▶ Confidence ellipsoid for estimate $\hat{\boldsymbol{\mu}}$ of $\boldsymbol{\mu}$ (with probability $1 - \delta/T^2$)
 - ▶ Using analysis of OFUL algorithm (Abbasi-Yadkori et al., 2011)

$$|\mathbf{b}_i^\top \hat{\boldsymbol{\mu}} - \mathbf{b}_i^\top \boldsymbol{\mu}| \leq \left(R \sqrt{2d \log \left(\frac{(\lambda + T)T^2}{\delta\lambda} \right)} + C \right) \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}} = \ell \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$$

- ▶ Our key result coming from spectral properties of \mathbf{B}_t .

$$\log \frac{|\mathbf{B}_t|}{|\boldsymbol{\Lambda}|} \leq 2d \log \left(1 + \frac{T}{\lambda} \right)$$

SpectralTS analysis sketch

- ▶ Confidence ellipsoid for estimate $\hat{\boldsymbol{\mu}}$ of $\boldsymbol{\mu}$ (with probability $1 - \delta/T^2$)
 - ▶ Using analysis of OFUL algorithm (Abbasi-Yadkori et al., 2011)

$$|\mathbf{b}_i^\top \hat{\boldsymbol{\mu}} - \mathbf{b}_i^\top \boldsymbol{\mu}| \leq \left(R \sqrt{2d \log \left(\frac{(\lambda + T)T^2}{\delta\lambda} \right)} + C \right) \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}} = \ell \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}$$

- ▶ Our key result coming from spectral properties of \mathbf{B}_t .

$$\log \frac{|\mathbf{B}_t|}{|\boldsymbol{\Lambda}|} \leq 2d \log \left(1 + \frac{T}{\lambda} \right)$$

- ▶ Concentration of sample $\tilde{\boldsymbol{\mu}}$ around mean $\hat{\boldsymbol{\mu}}$ (with probability $1 - 1/T^2$)
 - ▶ Using concentration inequality for Gaussian random variable.

$$|\mathbf{b}_i^\top \tilde{\boldsymbol{\mu}} - \mathbf{b}_i^\top \hat{\boldsymbol{\mu}}| \leq \left(R \sqrt{6d \log \left(\frac{\lambda + T}{\delta\lambda} \right)} + C \right) \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}} \sqrt{4 \log(TN)} = v \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}} \sqrt{4 \log(TN)}$$

SpectralTS analysis sketch

Define $\text{regret}'(t) = \text{regret}(t) \cdot \mathbb{1}\{|\mathbf{b}_i^\top \hat{\boldsymbol{\mu}}(t) - \mathbf{b}_i^\top \boldsymbol{\mu}| \leq \ell \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}\}$

$$\text{regret}'(t) \leq \frac{11g}{\rho} \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} + \frac{1}{T^2}$$

SpectralTS analysis sketch

Define $\text{regret}'(t) = \text{regret}(t) \cdot \mathbb{1}\{\|\mathbf{b}_i^\top \hat{\boldsymbol{\mu}}(t) - \mathbf{b}_i^\top \boldsymbol{\mu}\| \leq \ell \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}\}$

$$\text{regret}'(t) \leq \frac{11g}{\rho} \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} + \frac{1}{T^2}$$

Super-martingale (i.e. $\mathbb{E}[Y_t - Y_{t-1} | \mathcal{F}_{t-1}] \leq 0$)

$$X_t = \text{regret}'(t) - \frac{11g}{\rho} \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} - \frac{1}{T^2}$$

$$Y_t = \sum_{w=1}^t X_w.$$

$(Y_t; t = 0, \dots, T)$ is a **super-martingale** process w.r.t. history \mathcal{F}_t .

SpectralTS analysis sketch

Define $\text{regret}'(t) = \text{regret}(t) \cdot \mathbb{1}\{\|\mathbf{b}_i^\top \hat{\boldsymbol{\mu}}(t) - \mathbf{b}_i^\top \boldsymbol{\mu}\| \leq \ell \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}\}$

$$\text{regret}'(t) \leq \frac{11g}{\rho} \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} + \frac{1}{T^2}$$

Super-martingale (i.e. $\mathbb{E}[Y_t - Y_{t-1} | \mathcal{F}_{t-1}] \leq 0$)

$$X_t = \text{regret}'(t) - \frac{11g}{\rho} \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} - \frac{1}{T^2}$$

$$Y_t = \sum_{w=1}^t X_w.$$

$(Y_t; t = 0, \dots, T)$ is a **super-martingale** process w.r.t. history \mathcal{F}_t .

Azuma-Hoeffding inequality for super-martingale, w. p. $1 - \delta/2$:

$$\sum_{t=1}^T \text{regret}'(t) \leq \frac{11g}{\rho} \sum_{t=1}^T \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} + \frac{1}{T} + \frac{g}{\rho} \left(\frac{11}{\sqrt{\lambda}} + 2 \right) \sqrt{2T \ln \frac{2}{\delta}}$$

SpectralTS analysis sketch

Define $\text{regret}'(t) = \text{regret}(t) \cdot \mathbb{1}\{\|\mathbf{b}_i^\top \hat{\boldsymbol{\mu}}(t) - \mathbf{b}_i^\top \boldsymbol{\mu}\| \leq \ell \|\mathbf{b}_i\|_{\mathbf{B}_t^{-1}}\}$

$$\text{regret}'(t) \leq \frac{11g}{\rho} \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} + \frac{1}{T^2}$$

Super-martingale (i.e. $\mathbb{E}[Y_t - Y_{t-1} | \mathcal{F}_{t-1}] \leq 0$)

$$X_t = \text{regret}'(t) - \frac{11g}{\rho} \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} - \frac{1}{T^2}$$

$$Y_t = \sum_{w=1}^t X_w.$$

$(Y_t; t = 0, \dots, T)$ is a **super-martingale** process w.r.t. history \mathcal{F}_t .

Azuma-Hoeffding inequality for super-martingale, w. p. $1 - \delta/2$:

$$\sum_{t=1}^T \text{regret}'(t) \leq \frac{11g}{\rho} \sum_{t=1}^T \|\mathbf{b}_{a(t)}\|_{\mathbf{B}_t^{-1}} + \frac{1}{T} + \frac{g}{\rho} \left(\frac{11}{\sqrt{\lambda}} + 2 \right) \sqrt{2T \ln \frac{2}{\delta}}$$

Backup: SpectralEliminator pseudocode

Input:

N : the number of nodes, T : the number of pulls

$\{\Lambda_{\mathcal{L}}, \mathbf{Q}\}$ spectral basis of \mathcal{L}

λ : regularization parameter

$\beta, \{t_j\}_j^J$ parameters of the elimination and phases

$A_1 \leftarrow \{\mathbf{x}_1, \dots, \mathbf{x}_K\}$.

for $j = 1$ **to** J **do**

$\mathbf{V}_{t_j} \leftarrow \gamma \Lambda_{\mathcal{L}} + \lambda \mathbf{I}$

for $t = t_j$ **to** $\min(t_{j+1} - 1, T)$ **do**

Play $\mathbf{x}_t \in A_j$ with the largest width to observe r_t :

$\mathbf{x}_t \leftarrow \arg \max_{\mathbf{x} \in A_j} \|\mathbf{x}\|_{\mathbf{V}_t^{-1}}$

$\mathbf{V}_{t+1} \leftarrow \mathbf{V}_t + \mathbf{x}_t \mathbf{x}_t^\top$

end for

Eliminate the arms that are not promising:

$\hat{\alpha}_t \leftarrow \mathbf{V}_t^{-1} [\mathbf{x}_{t_j}, \dots, \mathbf{x}_t] [r_{t_j}, \dots, r_t]^\top$

$A_{j+1} \leftarrow \left\{ \mathbf{x} \in A_j, \langle \hat{\alpha}_t, \mathbf{x} \rangle + \|\mathbf{x}\|_{\mathbf{V}_t^{-1}} \beta \geq \max_{\mathbf{x} \in A_j} \left[\langle \hat{\alpha}_t, \mathbf{x} \rangle - \|\mathbf{x}\|_{\mathbf{V}_t^{-1}} \beta \right] \right\}$

end for

Backup: SpectralEliminator analysis

SpectralEliminator

- ▶ Divide time into sets ($t_1 = 1 \leq t_2 \leq \dots$) to introduce independence for Azuma-Hoeffding inequality and observe

$$R_T \leq \sum_{j=0}^J (t_{j+1} - t_j) [\langle \mathbf{x}^* - \mathbf{x}_t, \hat{\alpha}_j \rangle + (\|\mathbf{x}^*\|_{\mathbf{V}_j^{-1}} + \|\mathbf{x}_t\|_{\mathbf{V}_j^{-1}})\beta]$$
- ▶ Bound $\langle \mathbf{x}^* - \mathbf{x}_t, \hat{\alpha}_j \rangle$ for each phase
- ▶ No bad arms: $\langle \mathbf{x}^* - \mathbf{x}_t, \hat{\alpha}_j \rangle \leq (\|\mathbf{x}^*\|_{\mathbf{V}_j^{-1}} + \|\mathbf{x}_t\|_{\mathbf{V}_j^{-1}})\beta$
- ▶ By algorithm: $\|\mathbf{x}\|_{\mathbf{V}_j^{-1}}^2 \leq \frac{1}{t_j - t_{j-1}} \sum_{s=t_{j-1}+1}^{t_j} \|\mathbf{x}_s\|_{\mathbf{V}_{s-1}^{-1}}^2$
- ▶ $\sum_{s=t_{j-1}+1}^{t_j} \min\left(1, \|\mathbf{x}_s\|_{\mathbf{V}_{s-1}^{-1}}^2\right) \leq \log \frac{|\mathbf{V}_j|}{|\Lambda|}$