



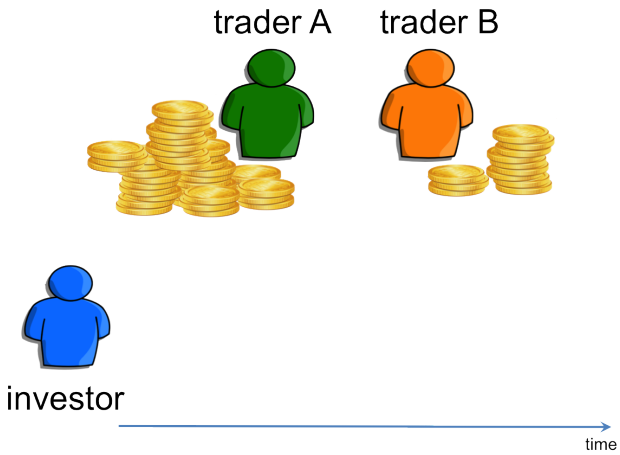
# Optimism and Randomness in Linear Multi-Armed Bandit

Alessandro LAZARIC (*Inria-Lille*)

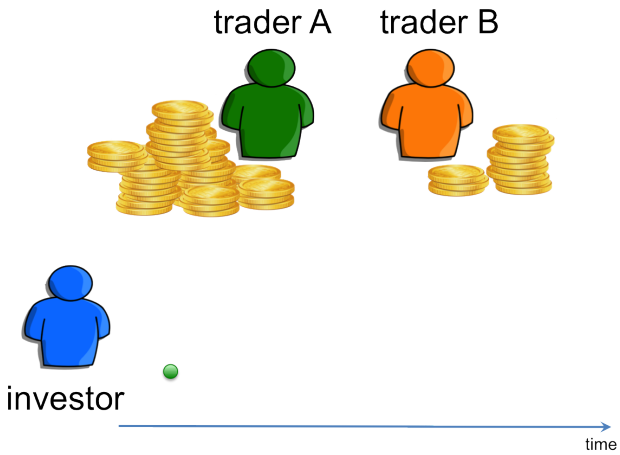
*International Conference on Monte-Carlo Techniques*

*Paris, July 5-8, 2016*

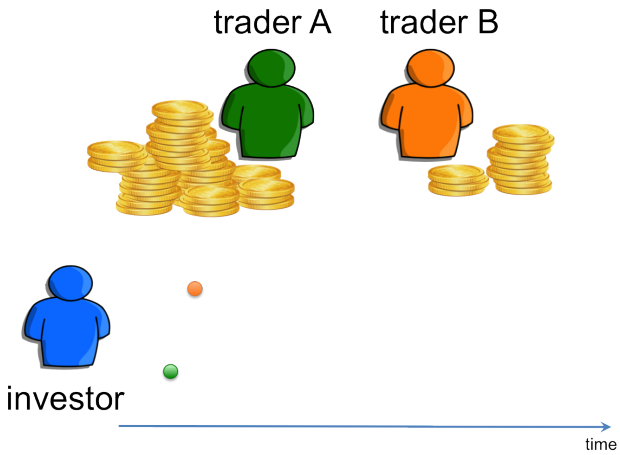
SequeL – Inria Lille



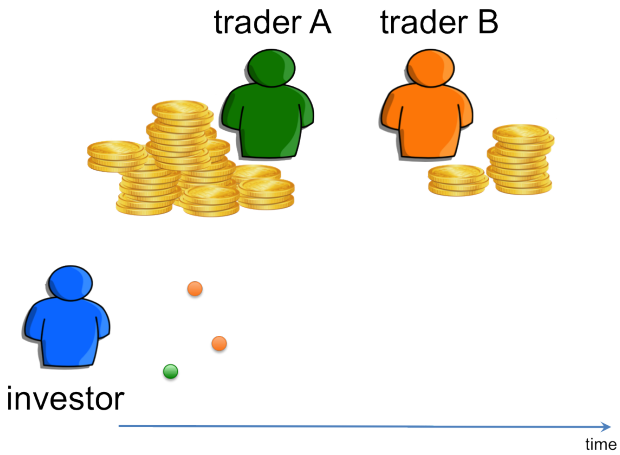
Example from [Lamberton, Pagès, Tarrès, 2004]



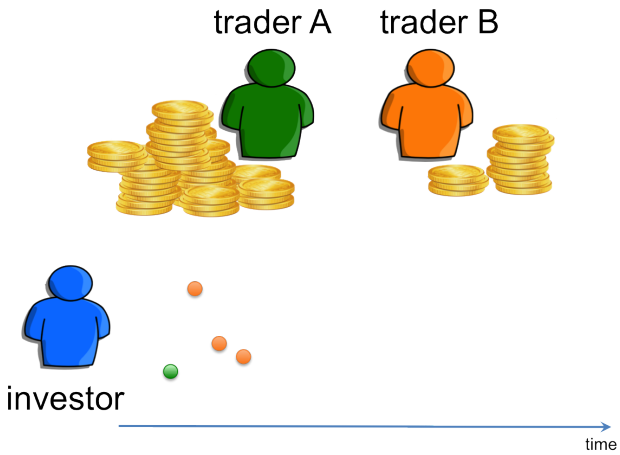
Example from [Lamberton, Pagès, Tarrès, 2004]



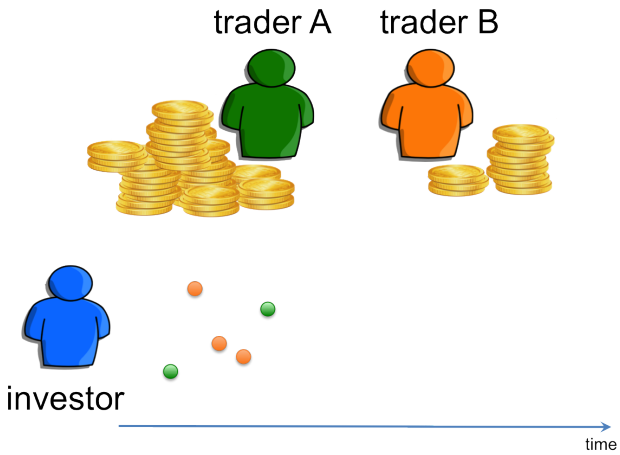
Example from [Lamberton, Pagès, Tarrès, 2004]



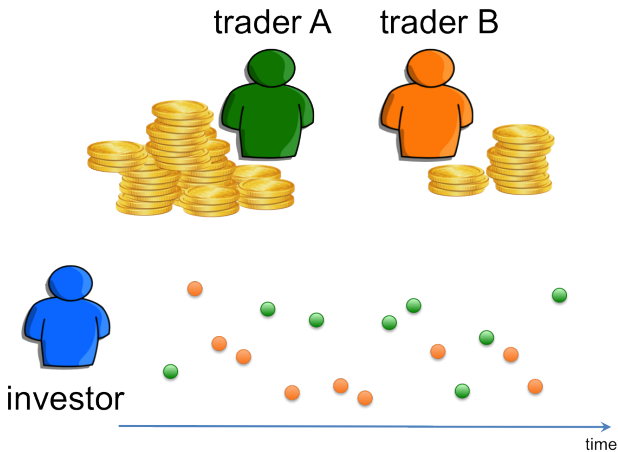
Example from [*Lamberton, Pagès, Tarrès, 2004*]



Example from [*Lamberton, Pagès, Tarrès, 2004*]

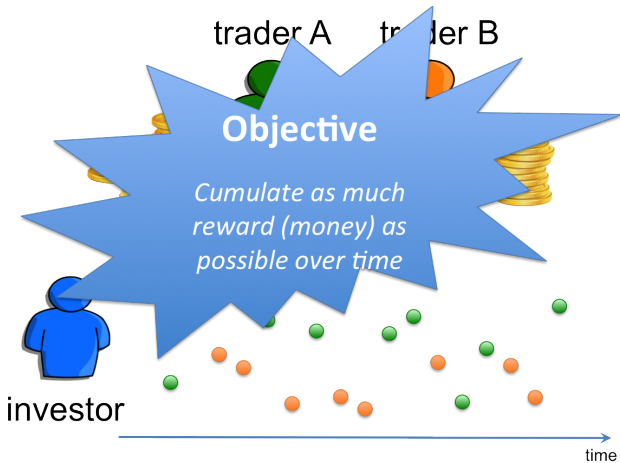


Example from [Lamberton, Pagès, Tarrès, 2004]



Example from [*Lamberton, Pagès, Tarrès, 2004*]





Example from [Lamberton, Pagès, Tarrès, 2004]

# Applications of Multi-armed Bandit

- ▶ Recommendation systems
- ▶ Clinical trials
- ▶ Packet routing, cognitive radios
- ▶ Trading
- ▶ Education
- ▶ ...

# Outline

Linear Bandit Framework

Solving Bandit with Optimism

Solving Bandit with Randomization (and a bit of optimism)

Perspectives

# The Linear Bandit Framework

The setting:

- ▶ Set of arms  $\mathcal{X} \subset \mathbb{R}^d$

# The Linear Bandit Framework

The setting:

- ▶ Set of arms  $\mathcal{X} \subset \mathbb{R}^d$
- ▶ Reward of arm  $x \in \mathcal{X}$

$$r(x) = x^T \theta^* + \xi \quad (\text{standard linear regression model})$$

with  $\theta^* \in \mathbb{R}^d$  unknown and  $\xi$  a zero-mean, sub-Gaussian noise

# The Linear Bandit Framework

The setting:

- ▶ Set of arms  $\mathcal{X} \subset \mathbb{R}^d$

- ▶ Reward of arm  $x \in \mathcal{X}$

$$r(x) = x^T \theta^* + \xi \quad (\text{standard linear regression model})$$

with  $\theta^* \in \mathbb{R}^d$  unknown and  $\xi$  a zero-mean, sub-Gaussian noise

- ▶ Best arm and best value for any parameter  $\theta$

$$x^*(\theta) = \arg \max_{x \in \mathcal{X}} x^T \theta; \quad J(\theta) = \max_{x \in \mathcal{X}} x^T \theta$$

# The Linear Bandit Framework

The setting:

- ▶ Set of arms  $\mathcal{X} \subset \mathbb{R}^d$

- ▶ Reward of arm  $x \in \mathcal{X}$

$$r(x) = x^T \theta^* + \xi \quad (\text{standard linear regression model})$$

with  $\theta^* \in \mathbb{R}^d$  unknown and  $\xi$  a zero-mean, sub-Gaussian noise

- ▶ Best arm and best value for any parameter  $\theta$

$$x^*(\theta) = \arg \max_{x \in \mathcal{X}} x^T \theta; \quad J(\theta) = \max_{x \in \mathcal{X}} x^T \theta$$

- ▶ *Optimal strategy*: select arm  $x^*(\theta^*)$  (*constrained linear optimization*)

# The Linear Bandit Framework

The learning problem:

- ▶ Finite horizon  $T$



# The Linear Bandit Framework

The learning problem:

- ▶ Finite horizon  $T$
- ▶ Select an arm  $x_t$  at each step  $t = 1, \dots, T$

# The Linear Bandit Framework

The learning problem:

- ▶ Finite horizon  $T$
- ▶ Select an arm  $x_t$  at each step  $t = 1, \dots, T$
- ▶ Cumulate as much reward as possible

$$\sum_{t=1}^T x_t^T \theta^* \quad (\text{explore-exploit trade-off})$$

# The Linear Bandit Framework

The learning problem:

- ▶ Finite horizon  $T$
- ▶ Select an arm  $x_t$  at each step  $t = 1, \dots, T$
- ▶ Cumulate as much reward as possible

$$\sum_{t=1}^T x_t^T \theta^* \quad (\text{explore-exploit trade-off})$$

- ▶ *Equivalently*: minimize the regret

$$R(T) = \sum_{t=1}^T (x^*(\theta^*)^T \theta^* - x_t^T \theta^*)$$

# The Linear Bandit Framework

The learning problem:

- ▶ Finite horizon  $T$
- ▶ Select an arm  $x_t$  at each step  $t = 1, \dots, T$
- ▶ Cumulate as much reward as possible

$$\sum_{t=1}^T x_t^T \theta^* \quad (\text{explore-exploit trade-off})$$

- ▶ *Equivalently*: minimize the regret

$$R(T) = \sum_{t=1}^T (x^*(\theta^*)^T \theta^* - x_t^T \theta^*)$$

$\Rightarrow$  a good learning algorithm should have  $o(T)$  regret!

# The Linear Bandit Framework

The core ingredient: *regularized least-squares estimator*

- ▶ Given samples  $\{(x_1, r_1), (x_2, r_2), \dots, (x_{t-1}, r_{t-1})\}$  compute

$$\hat{\theta}_t = \arg \min_{\theta \in \mathbb{R}^d} \sum_{s=1}^{t-1} (r_s - x_s^\top \theta)^2 + \lambda \|\theta\| \quad (\lambda \text{ regularization parameter})$$

# The Linear Bandit Framework

The core ingredient: *regularized least-squares estimator*

- ▶ Given samples  $\{(x_1, r_1), (x_2, r_2), \dots, (x_{t-1}, r_{t-1})\}$  compute

$$\hat{\theta}_t = \arg \min_{\theta \in \mathbb{R}^d} \sum_{s=1}^{t-1} (r_s - x_s^\top \theta)^2 + \lambda \|\theta\| \quad (\lambda \text{ regularization parameter})$$

- ▶ In closed form

$$V_t = \lambda I + \sum_{s=1}^{t-1} x_s x_s^\top \quad (\text{design matrix}) \quad \hat{\theta}_t = V_t^{-1} \sum_{s=1}^{t-1} x_s r_s \quad (\text{RLS estimator})$$

# The Linear Bandit Framework

The core ingredient: *regularized least-squares estimator*

- ▶ Given samples  $\{(x_1, r_1), (x_2, r_2), \dots, (x_{t-1}, r_{t-1})\}$  compute

$$\hat{\theta}_t = \arg \min_{\theta \in \mathbb{R}^d} \sum_{s=1}^{t-1} (r_s - x_s^\top \theta)^2 + \lambda \|\theta\| \quad (\lambda \text{ regularization parameter})$$

- ▶ In closed form

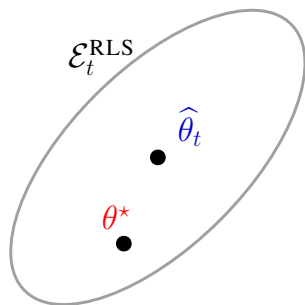
$$V_t = \lambda I + \sum_{s=1}^{t-1} x_s x_s^\top \quad (\text{design matrix}) \quad \hat{\theta}_t = V_t^{-1} \sum_{s=1}^{t-1} x_s r_s \quad (\text{RLS estimator})$$

- ▶ Guarantees (w.h.p.) (*Gauss-Markov confidence interval for martingales*)

- ▶ (estimation)  $\|\hat{\theta}_t - \theta^*\|_{V_t} \leq \sqrt{d \log(t/\delta)}$

- ▶ (prediction)  $|x^\top (\hat{\theta}_t - \theta^*)| \leq \|x\|_{V_t^{-1}} \sqrt{d \log(t/\delta)}$

## Confidence Ellipsoid



$$\mathcal{E}_t^{\text{RLS}} = \{ \theta \in \mathbb{R}^d \mid \| \theta - \hat{\theta}_t \|_{V_t} \leq \sqrt{d \log(1/\delta)} \}$$



# Outline

Linear Bandit Framework

Solving Bandit with Optimism

Solving Bandit with Randomization (and a bit of optimism)

Perspectives

# Optimism in Face of Uncertainty

- ▶ *Exploit*: given past observations, compute  $\hat{\theta}_t$  and confidence ellipsoid  $\mathcal{E}_t^{\text{RLS}}$

# Optimism in Face of Uncertainty

- ▶ *Exploit*: given past observations, compute  $\hat{\theta}_t$  and confidence ellipsoid  $\mathcal{E}_t^{\text{RLS}}$
- ▶ *Explore*: given past observations, any  $\theta \in \mathcal{E}_t^{\text{RLS}}$  could be  $\theta^*$

# Optimism in Face of Uncertainty

- ▶ *Exploit*: given past observations, compute  $\hat{\theta}_t$  and confidence ellipsoid  $\mathcal{E}_t^{\text{RLS}}$
- ▶ *Explore*: given past observations, any  $\theta \in \mathcal{E}_t^{\text{RLS}}$  could be  $\theta^*$
- ▶ *Optimism*: trade-off exploration and exploitation by taking the most “optimistic”  $\theta$  compatible with current estimates

$$\tilde{\theta}_t = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta) = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} \max_{x \in \mathcal{X}} x^\top \theta$$

# Optimism in Face of Uncertainty

- ▶ *Exploit*: given past observations, compute  $\hat{\theta}_t$  and confidence ellipsoid  $\mathcal{E}_t^{\text{RLS}}$
- ▶ *Explore*: given past observations, any  $\theta \in \mathcal{E}_t^{\text{RLS}}$  could be  $\theta^*$
- ▶ *Optimism*: trade-off exploration and exploitation by taking the most “optimistic”  $\theta$  compatible with current estimates

$$\tilde{\theta}_t = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta) = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} \max_{x \in \mathcal{X}} x^T \theta$$

- ▶ Act *as-if*  $\tilde{\theta}_t$  was the true parameter

$$x_t = \arg \max_{x \in \mathcal{X}} x^T \tilde{\theta}_t$$

# Optimism in Face of Uncertainty

- ▶ **Exploit**: given past observations, compute  $\hat{\theta}_t$  and confidence ellipsoid  $\mathcal{E}_t^{\text{RLS}}$
- ▶ **Explore**: given past observations, any  $\theta \in \mathcal{E}_t^{\text{RLS}}$  could be  $\theta^*$
- ▶ **Optimism**: trade-off exploration and exploitation by taking the most “optimistic”  $\theta$  compatible with current estimates

$$\tilde{\theta}_t = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta) = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} \max_{x \in \mathcal{X}} x^\top \theta$$

- ▶ Act *as-if*  $\tilde{\theta}_t$  was the true parameter

$$x_t = \arg \max_{x \in \mathcal{X}} x^\top \tilde{\theta}_t$$

⇒ the resulting algorithm is called LINUCB or OFUL

# How It Works

## High-level intuition

- ▶ The arm choice can be written as (by def. of  $\mathcal{E}_t^{\text{RLS}}$ )

$$x_t = \arg \max_{x \in \mathcal{X}} x^\top \tilde{\theta}_t = \arg \max_{x \in \mathcal{X}} \left( \underbrace{x^\top \hat{\theta}_t}_{\text{exploit}} + \underbrace{\|x\|_{V_t^{-1}} \sqrt{d \log(t/\delta)}}_{\text{explore}} \right)$$

# How It Works

## High-level intuition

- ▶ The arm choice can be written as (by def. of  $\mathcal{E}_t^{\text{RLS}}$ )

$$x_t = \arg \max_{x \in \mathcal{X}} x^T \tilde{\theta}_t = \arg \max_{x \in \mathcal{X}} \left( \underbrace{x^T \hat{\theta}_t}_{\text{exploit}} + \underbrace{\|x\|_{V_t^{-1}} \sqrt{d \log(t/\delta)}}_{\text{explore}} \right)$$

- ▶ *Case 1*:  $x_t = x^*(\theta^*) \Rightarrow$  no regret



# How It Works

## High-level intuition

- ▶ The arm choice can be written as (by def. of  $\mathcal{E}_t^{\text{RLS}}$ )

$$x_t = \arg \max_{x \in \mathcal{X}} x^\top \tilde{\theta}_t = \arg \max_{x \in \mathcal{X}} \left( \underbrace{x^\top \hat{\theta}_t}_{\text{exploit}} + \underbrace{\|x\|_{V_t^{-1}} \sqrt{d \log(t/\delta)}}_{\text{explore}} \right)$$

- ▶ *Case 1*:  $x_t = x^*(\theta^*) \Rightarrow$  no regret
- ▶ *Case 2*:  $x_t \neq x^*(\theta^*) \Rightarrow$  the confidence ellipsoid is *tightened* along the direction whose uncertainty had the largest impact in the decision of  $x_t$

# How It Works

## High-level intuition

- ▶ The arm choice can be written as (by def. of  $\mathcal{E}_t^{\text{RLS}}$ )

$$x_t = \arg \max_{x \in \mathcal{X}} x^\top \tilde{\theta}_t = \arg \max_{x \in \mathcal{X}} \left( \underbrace{x^\top \hat{\theta}_t}_{\text{exploit}} + \underbrace{\|x\|_{V_t^{-1}} \sqrt{d \log(t/\delta)}}_{\text{explore}} \right)$$

- ▶ *Case 1*:  $x_t = x^*(\theta^*) \Rightarrow$  no regret
- ▶ *Case 2*:  $x_t \neq x^*(\theta^*) \Rightarrow$  the confidence ellipsoid is *tightened* along the direction whose uncertainty had the largest impact in the decision of  $x_t$

$\Rightarrow$  either *instantaneous regret is small* or *useful information is obtained* and future regret will be small

# How It Works

Proof sketch

- ▶ Regret decomposition

$$\begin{aligned}
 R(T) &= \sum_{t=1}^T (\mathbf{x}^*(\theta^*)^\top \theta^* - \mathbf{x}_t^\top \theta^*) \\
 &= \sum_{t=1}^T (\mathbf{x}^*(\theta^*)^\top \theta^* - \mathbf{x}_t^\top \tilde{\theta}_t) + \sum_{t=1}^T (\mathbf{x}_t^\top \tilde{\theta}_t - \mathbf{x}_t^\top \theta^*) \\
 &= \underbrace{\sum_{t=1}^T (J(\theta^*) - J(\tilde{\theta}_t))}_{R_1(T)} + \underbrace{\sum_{t=1}^T (\mathbf{x}_t^\top \tilde{\theta}_t - \mathbf{x}_t^\top \theta^*)}_{R_2(T)}
 \end{aligned}$$

# How It Works

Proof sketch

- ▶ Regret decomposition

$$\begin{aligned}
 R(T) &= \sum_{t=1}^T (x^*(\theta^*)^\top \theta^* - x_t^\top \theta^*) \\
 &= \sum_{t=1}^T (x^*(\theta^*)^\top \theta^* - x_t^\top \tilde{\theta}_t) + \sum_{t=1}^T (x_t^\top \tilde{\theta}_t - x_t^\top \theta^*) \\
 &= \underbrace{\sum_{t=1}^T (J(\theta^*) - J(\tilde{\theta}_t))}_{R_1(T)} + \underbrace{\sum_{t=1}^T (x_t^\top \tilde{\theta}_t - x_t^\top \theta^*)}_{R_2(T)}
 \end{aligned}$$

- ▶  $R_1(T) \leq 0$  by construction (recall  $\tilde{\theta}_t = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta)$ )

# How It Works

Proof sketch

- ▶ Regret decomposition

$$\begin{aligned}
 R(T) &= \sum_{t=1}^T (x^*(\theta^*)^\top \theta^* - x_t^\top \theta^*) \\
 &= \sum_{t=1}^T (x^*(\theta^*)^\top \theta^* - x_t^\top \tilde{\theta}_t) + \sum_{t=1}^T (x_t^\top \tilde{\theta}_t - x_t^\top \theta^*) \\
 &= \underbrace{\sum_{t=1}^T (J(\theta^*) - J(\tilde{\theta}_t))}_{R_1(T)} + \underbrace{\sum_{t=1}^T (x_t^\top \tilde{\theta}_t - x_t^\top \theta^*)}_{R_2(T)}
 \end{aligned}$$

- ▶  $R_1(T) \leq 0$  by construction (recall  $\tilde{\theta}_t = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta)$ )
- ▶  $R_2(T)$  is the prediction error on points  $x_t$  used to estimate  $\tilde{\theta}_t$  and  $\tilde{\theta}_t$  is in  $\mathcal{E}_t^{\text{RLS}}$   
 $\Rightarrow$  cumulatively small

# How It Works

## Theorem (Abbasi-Yadkori *et al.*, 2011)

*If OFUL is run over  $T$  steps on arms in  $\mathcal{X} \subset \mathbb{R}^d$ , then it suffers a cumulative regret*

$$R(T) = \tilde{O}(d\sqrt{T})$$

*with high probability.*

# Main Issue

Computing  $\tilde{\theta}_t$  requires solving a doubly-linear optimization problem

$$\tilde{\theta}_t = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta) = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} \max_{x \in \mathcal{X}} x^\top \theta$$

# Main Issue

Computing  $\tilde{\theta}_t$  requires solving a doubly-linear optimization problem

$$\tilde{\theta}_t = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta) = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} \max_{x \in \mathcal{X}} x^T \theta$$

$\Rightarrow$  **computational expensive** for non-trivial arm sets  $\mathcal{X}$



# Outline

Linear Bandit Framework

Solving Bandit with Optimism

Solving Bandit with Randomization (and a bit of optimism)

Perspectives

# Thompson Sampling

A Bayesian algorithm (*dating back to [Thompson, 1933]*)

- ▶ Define a *prior* on parameter  $p(\theta)$  (e.g.,  $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ )

# Thompson Sampling

A Bayesian algorithm (*dating back to [Thompson, 1933]*)

- ▶ Define a *prior* on parameter  $p(\theta)$  (e.g.,  $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ )
- ▶ At each step  $t = 1, \dots, T$ 
  - ▶ Draw  $\tilde{\theta}_t$  from posterior  $p(\theta | x_1, r_1, \dots, x_{t-1}, r_{t-1})$  (e.g.,  $\theta \sim \mathcal{N}(\hat{\theta}_t, V_t^{-1})$ )

# Thompson Sampling

A Bayesian algorithm (*dating back to [Thompson, 1933]*)

- ▶ Define a *prior* on parameter  $p(\theta)$  (e.g.,  $\theta \sim \mathcal{N}(\mathbf{0}, I)$ )
- ▶ At each step  $t = 1, \dots, T$ 
  - ▶ Draw  $\tilde{\theta}_t$  from posterior  $p(\theta | x_1, r_1, \dots, x_{t-1}, r_{t-1})$  (e.g.,  $\theta \sim \mathcal{N}(\hat{\theta}_t, V_t^{-1})$ )
  - ▶ Select arm  $x_t = \arg \max_{x \in \mathcal{X}} x^T \tilde{\theta}_t$

# Thompson Sampling

A Bayesian algorithm (*dating back to [Thompson, 1933]*)

- ▶ Define a *prior* on parameter  $p(\theta)$  (e.g.,  $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ )
- ▶ At each step  $t = 1, \dots, T$ 
  - ▶ Draw  $\tilde{\theta}_t$  from posterior  $p(\theta | x_1, r_1, \dots, x_{t-1}, r_{t-1})$  (e.g.,  $\theta \sim \mathcal{N}(\hat{\theta}_t, V_t^{-1})$ )
  - ▶ Select arm  $x_t = \arg \max_{x \in \mathcal{X}} x^T \tilde{\theta}_t$

$\Rightarrow$  sampling  $\tilde{\theta}_t$  from the posterior implements an *exploration-exploitation* trade-off

# How It Works

► Regret decomposition

$$\begin{aligned}
 R(T) &= \sum_{t=1}^T (x^*(\theta^*)^\top \theta^* - x_t^\top \theta^*) \\
 &= \underbrace{\sum_{t=1}^T (J(\theta^*) - J(\tilde{\theta}_t))}_{R_1(T)} + \underbrace{\sum_{t=1}^T (x_t^\top \tilde{\theta}_t - x_t^\top \theta^*)}_{R_2(T)}
 \end{aligned}$$

# How It Works

- ▶ Regret decomposition

$$\begin{aligned}
 R(T) &= \sum_{t=1}^T (x^*(\theta^*)^\top \theta^* - x_t^\top \theta^*) \\
 &= \underbrace{\sum_{t=1}^T (J(\theta^*) - J(\tilde{\theta}_t))}_{R_1(T)} + \underbrace{\sum_{t=1}^T (x_t^\top \tilde{\theta}_t - x_t^\top \theta^*)}_{R_2(T)}
 \end{aligned}$$

- ▶  $R_2(T)$  is the same as before

# How It Works

- ▶ Regret decomposition

$$\begin{aligned}
 R(T) &= \sum_{t=1}^T (x^*(\theta^*)^\top \theta^* - x_t^\top \theta^*) \\
 &= \underbrace{\sum_{t=1}^T (J(\theta^*) - J(\tilde{\theta}_t))}_{R_1(T)} + \underbrace{\sum_{t=1}^T (x_t^\top \tilde{\theta}_t - x_t^\top \theta^*)}_{R_2(T)}
 \end{aligned}$$

- ▶  $R_2(T)$  is the same as before
  - ▶  $\tilde{\theta}_t^{\text{OFUL}} = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta)$  vs  $\tilde{\theta}_t^{\text{TS}} \sim p(\theta | x_1, r_1, \dots, x_{t-1}, r_{t-1})$
- $\Rightarrow J(\tilde{\theta}_t^{\text{TS}}) = ??$



## How It Works

- ▶ Regret decomposition

$$\begin{aligned}
 R(T) &= \sum_{t=1}^T (x^*(\theta^*)^\top \theta^* - x_t^\top \theta^*) \\
 &= \underbrace{\sum_{t=1}^T (J(\theta^*) - J(\tilde{\theta}_t))}_{R_1(T)} + \underbrace{\sum_{t=1}^T (x_t^\top \tilde{\theta}_t - x_t^\top \theta^*)}_{R_2(T)}
 \end{aligned}$$

- ▶  $R_2(T)$  is the same as before
  - ▶  $\tilde{\theta}_t^{\text{OFUL}} = \arg \max_{\theta \in \mathcal{E}_t^{\text{RLS}}} J(\theta)$  vs  $\tilde{\theta}_t^{\text{TS}} \sim p(\theta | x_1, r_1, \dots, x_{t-1}, r_{t-1})$
- $\Rightarrow J(\tilde{\theta}_t^{\text{TS}}) = ?? \Rightarrow R_1(T) \not\leq 0$

# The Importance of Being Optimistic

Let  $R_t^{\text{TS}} = J(\theta^*) - J(\tilde{\theta}_t)$

- ▶ At step  $\tau$ ,  $\tilde{\theta}_\tau$  is optimistic (i.e.,  $J(\tilde{\theta}_\tau) \geq J(\theta^*)$ ), then  $R_\tau^{\text{TS}} \leq 0$

# The Importance of Being Optimistic

Let  $R_t^{\text{TS}} = J(\theta^*) - J(\tilde{\theta}_t)$

- ▶ At step  $\tau$ ,  $\tilde{\theta}_\tau$  is optimistic (i.e.,  $J(\tilde{\theta}_\tau) \geq J(\theta^*)$ ), then  $R_\tau^{\text{TS}} \leq 0$
- ▶ At any other subsequent (non-optimistic) step  $t$

$$R_t^{\text{TS}} \leq J(\tilde{\theta}_\tau) - J(\hat{\theta}_t) \qquad J(\tilde{\theta}_\tau) \geq J(\theta^*)$$

# The Importance of Being Optimistic

Let  $R_t^{\text{TS}} = J(\theta^*) - J(\tilde{\theta}_t)$

- ▶ At step  $\tau$ ,  $\tilde{\theta}_\tau$  is optimistic (i.e.,  $J(\tilde{\theta}_\tau) \geq J(\theta^*)$ ), then  $R_\tau^{\text{TS}} \leq 0$
- ▶ At any other subsequent (non-optimistic) step  $t$

$$\begin{aligned}
 R_t^{\text{TS}} &\leq J(\tilde{\theta}_\tau) - J(\hat{\theta}_t) && J(\tilde{\theta}_\tau) \geq J(\theta^*) \\
 &\leq \nabla J(\tilde{\theta}_\tau)^\top (\tilde{\theta}_\tau - \hat{\theta}_t) && J(\theta) \text{ is convex}
 \end{aligned}$$

# The Importance of Being Optimistic

Let  $R_t^{\text{TS}} = J(\theta^*) - J(\tilde{\theta}_t)$

- ▶ At step  $\tau$ ,  $\tilde{\theta}_\tau$  is optimistic (i.e.,  $J(\tilde{\theta}_\tau) \geq J(\theta^*)$ ), then  $R_\tau^{\text{TS}} \leq 0$
- ▶ At any other subsequent (non-optimistic) step  $t$

$$\begin{aligned}
 R_t^{\text{TS}} &\leq J(\tilde{\theta}_\tau) - J(\hat{\theta}_t) && J(\tilde{\theta}_\tau) \geq J(\theta^*) \\
 &\leq \nabla J(\tilde{\theta}_\tau)^\top (\tilde{\theta}_\tau - \tilde{\theta}_t) && J(\theta) \text{ is convex} \\
 &\leq x_\tau^\top (\tilde{\theta}_\tau - \tilde{\theta}_t) && \nabla J(\theta) = x^*(\theta) \text{ by def.}
 \end{aligned}$$

# The Importance of Being Optimistic

Let  $R_t^{\text{TS}} = J(\theta^*) - J(\tilde{\theta}_t)$

- ▶ At step  $\tau$ ,  $\tilde{\theta}_\tau$  is optimistic (i.e.,  $J(\tilde{\theta}_\tau) \geq J(\theta^*)$ ), then  $R_\tau^{\text{TS}} \leq 0$
- ▶ At any other subsequent (non-optimistic) step  $t$

$$\begin{aligned}
 R_t^{\text{TS}} &\leq J(\tilde{\theta}_\tau) - J(\hat{\theta}_t) && J(\tilde{\theta}_\tau) \geq J(\theta^*) \\
 &\leq \nabla J(\tilde{\theta}_\tau)^\top (\tilde{\theta}_\tau - \tilde{\theta}_t) && J(\theta) \text{ is convex} \\
 &\leq \mathbf{x}_\tau^\top (\tilde{\theta}_\tau - \tilde{\theta}_t) && \nabla J(\theta) = \mathbf{x}^*(\theta) \text{ by def.} \\
 &\leq \|\mathbf{x}_\tau\|_{\mathbf{V}_\tau^{-1}} \|\tilde{\theta}_\tau - \tilde{\theta}_t\|_{\mathbf{V}_t} && \text{by Cauchy-Schwarz}
 \end{aligned}$$

# The Importance of Being Optimistic

- **Summing up** ( $\nu_k$  time between any two optimistic choices,  $\tau_k$  optimistic times)

$$\sum_{t=1}^T R_t^{\text{TS}} \leq \sqrt{dT} \sum_{k=1}^K \nu_k \|x_{\tau_k}\|_{V_{\tau_k}^{-1}}$$

# The Importance of Being Optimistic

- ▶ **Summing up** ( $\nu_k$  time between any two optimistic choices,  $\tau_k$  optimistic times)

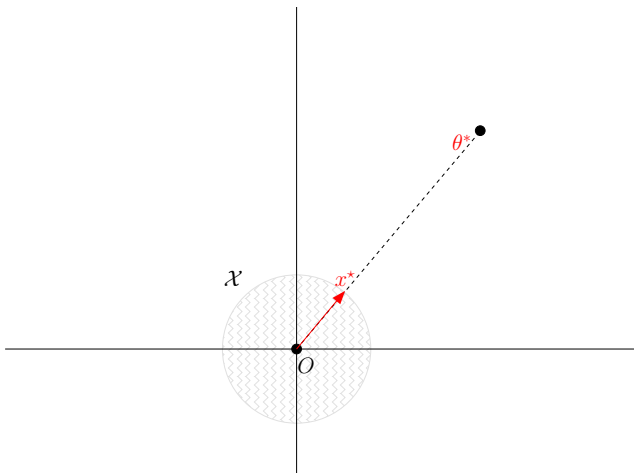
$$\sum_{t=1}^T R_t^{\text{TS}} \leq \sqrt{dT} \sum_{k=1}^K \nu_k \|x_{\tau_k}\|_{V_{\tau_k}^{-1}}$$

- ▶ If  $\tilde{\theta}_t$  is optimistic with probability  $p$ , then  $\mathbb{E}[\nu_k] = 1/p$  and

$$R(T) \leq \tilde{O}(d/p\sqrt{T})$$

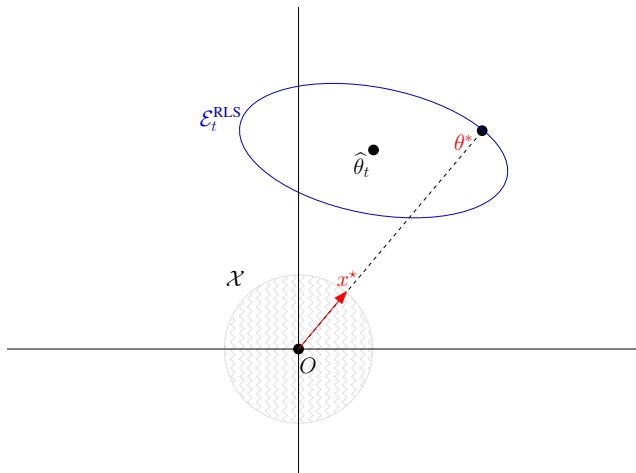


# How It Works



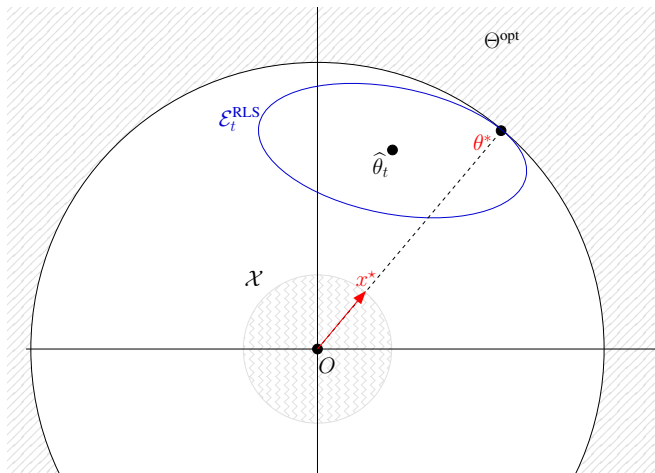
Consider  $\mathcal{X} = \mathcal{B}(0, 1)$ , then  $x^*(\theta) = \theta/\|\theta\|$  and  $J(\theta) = \|\theta\|$

## How It Works



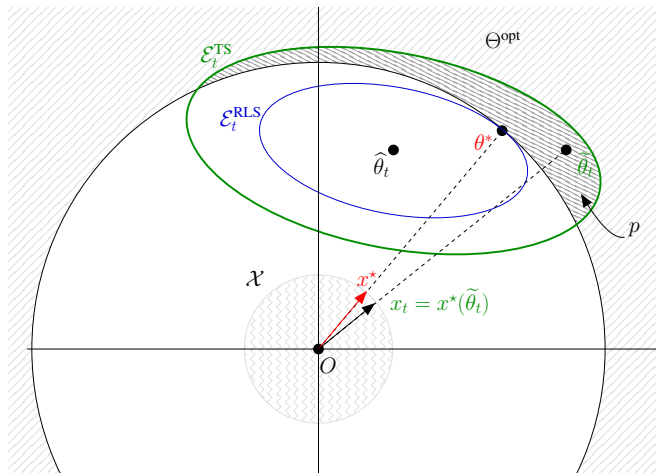
Consider  $\mathcal{X} = \mathcal{B}(0, 1)$ , then  $x^*(\theta) = \theta/\|\theta\|$  and  $J(\theta) = \|\theta\|$

## How It Works



Consider  $\mathcal{X} = \mathcal{B}(0, 1)$ , then  $x^*(\theta) = \theta/\|\theta\|$  and  $J(\theta) = \|\theta\|$

## How It Works



Consider  $\mathcal{X} = \mathcal{B}(0, 1)$ , then  $x^*(\theta) = \theta / \|\theta\|$  and  $J(\theta) = \|\theta\|$

## How It Works

### Theorem (Agrawal & Goyal, 2012; Abeille & L., 2016)

If TS is run over  $T$  steps on arms in  $\mathcal{X} \subset \mathbb{R}^d$ , then it suffers a cumulative regret

$$R(T) = \tilde{O}(d^{3/2}\sqrt{T})$$

with high probability.

## Discussion

- + TS is computationally faster than OFUL
- + TS often performs better than OFUL

## Discussion

- + TS is computationally faster than OFUL
- + TS often performs better than OFUL
- the need for optimism worsens the bound by  $\sqrt{d}$
- the Bayesian design requires choosing appropriate priors

# Outline

Linear Bandit Framework

Solving Bandit with Optimism

Solving Bandit with Randomization (and a bit of optimism)

Perspectives



# Thompson Sampling as a Stochastic Algorithm?

*Probability matching* algorithm

- ▶ Define a *prior* on parameter  $p(\theta)$

# Thompson Sampling as a Stochastic Algorithm?

*Probability matching* algorithm

- ▶ Define a *prior* on parameter  $p(\theta)$
- ▶ At each step  $t = 1, \dots, T$ 
  - ▶ For any arm  $x \in \mathcal{X}$ , compute

$$\pi_t(x) = \mathbb{P}(x = x^* | x_1, r_1, \dots, x_{t-1}, r_{t-1})$$

- ▶ Select arm  $x_t \sim \pi_t$

# Thompson Sampling as a Stochastic Algorithm?

*Probability matching* algorithm

- ▶ Define a *prior* on parameter  $p(\theta)$
- ▶ At each step  $t = 1, \dots, T$ 
  - ▶ For any arm  $x \in \mathcal{X}$ , compute

$$\pi_t(x) = \mathbb{P}(x = x^* | x_1, r_1, \dots, x_{t-1}, r_{t-1})$$

- ▶ Select arm  $x_t \sim \pi_t$

$\Rightarrow$  TS is an *efficient implementation* of probability matching

# Thompson Sampling as a Stochastic Algorithm?

- ▶ In some cases we can explicitly write the Bayesian update as ( $\pi_t$  being a distribution over  $\mathcal{X}$ )

$$\pi_{t+1} = \pi_t + \Delta_t$$

# Thompson Sampling as a Stochastic Algorithm?

- ▶ In some cases we can explicitly write the Bayesian update as ( $\pi_t$  being a distribution over  $\mathcal{X}$ )

$$\pi_{t+1} = \pi_t + \Delta_t$$

- ▶ The *Narendra-Shapiro* algorithm is a *stochastic algorithm* updating a distribution over arms as

$$\pi_{t+1} = \pi_t + \gamma_t \Delta'_t$$

# Thompson Sampling as a Stochastic Algorithm?

- ▶ In some cases we can explicitly write the Bayesian update as ( $\pi_t$  being a distribution over  $\mathcal{X}$ )

$$\pi_{t+1} = \pi_t + \Delta_t$$

- ▶ The *Narendra-Shapiro* algorithm is a *stochastic algorithm* updating a distribution over arms as

$$\pi_{t+1} = \pi_t + \gamma_t \Delta'_t$$

- ▶ The (over-penalized)-*Narendra-Shapiro* is shown [Gadat et al., 2016]
  - ▶ to converge to the stationary distribution of a *piecewise deterministic Markov process*
  - ▶ to suffer from a worst-case regret  $\tilde{O}(\sqrt{T})$  (in the 2-arm independent Bernoulli case)

# Thompson Sampling as a Stochastic Algorithm?

- ▶ In some cases we can explicitly write the Bayesian update as ( $\pi_t$  being a distribution over  $\mathcal{X}$ )

$$\pi_{t+1} = \pi_t + \Delta_t$$

- ▶ The *Narendra-Shapiro* algorithm is a *stochastic algorithm* updating a distribution over arms as

$$\pi_{t+1} = \pi_t + \gamma_t \Delta'_t$$

- ▶ The (over-penalized)-*Narendra-Shapiro* is shown [Gadat et al., 2016]
  - ▶ to converge to the stationary distribution of a *piecewise deterministic Markov process*
  - ▶ to suffer from a worst-case regret  $\tilde{O}(\sqrt{T})$  (in the 2-arm independent Bernoulli case)

⇒ If TS can be seen as a stochastic algorithm, we could have a much better understanding of the dynamics and behavior of bandit algorithms

Thank you!

The Inria logo is displayed in a white rounded square with a teal border. The word "Inria" is written in a red, cursive script font.

*Alessandro Lazaric*

[alessandro.lazaric@inria.fr](mailto:alessandro.lazaric@inria.fr)

[sequel.lille.inria.fr](http://sequel.lille.inria.fr)