# Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems, Part 2
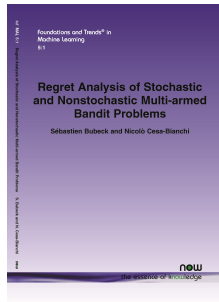
**Sébastien Bubeck**
Theory Group

## The linear bandit problem, Auer [2002]

**Known parameters:** compact action set $\mathcal{A} \subset \mathbb{R}^n$, adversary's action set $\mathcal{L} \subset \mathbb{R}^n$, number of rounds $T$.

# The linear bandit problem, Auer [2002]

**Known parameters:** compact action set $\mathcal{A} \subset \mathbb{R}^n$, adversary's action set $\mathcal{L} \subset \mathbb{R}^n$, number of rounds $T$.

**Protocol:** For each round $t = 1, 2, \ldots, T$, the adversary chooses a loss vector $\ell_t \in \mathcal{L}$ and simultaneously the player chooses $a_t \in \mathcal{A}$ based on past observations and receives a loss/observation $Y_t = \ell_t^\top a_t$.

$$R_T = \mathbb{E} \sum_{t=1}^T \ell_t^\top a_t - \min_{a \in \mathcal{A}} \mathbb{E} \sum_{t=1}^T \ell_t^\top a.$$

# The linear bandit problem, Auer [2002]

**Known parameters:** compact action set $\mathcal{A} \subset \mathbb{R}^n$, adversary's action set $\mathcal{L} \subset \mathbb{R}^n$, number of rounds $T$.

**Protocol:** For each round $t = 1, 2, \ldots, T$, the adversary chooses a loss vector $\ell_t \in \mathcal{L}$ and simultaneously the player chooses $a_t \in \mathcal{A}$ based on past observations and receives a loss/observation $Y_t = \ell_t^\top a_t$.

$$R_T = \mathbb{E} \sum_{t=1}^{T} \ell_t^\top a_t - \min_{a \in \mathcal{A}} \mathbb{E} \sum_{t=1}^{T} \ell_t^\top a.$$

**Other models:** In the i.i.d. model we assume that there is some underlying $\theta \in \mathcal{L}$ such that $\mathbb{E}(Y_t | a_t) = \theta^\top a_t$. In the Bayesian model we assume that we have a prior distribution $\nu$ over the sequence $(\ell_1, \ldots, \ell_T)$ (in this case the expectation in $R_T$ is also over $(\ell_1, \ldots, \ell_T) \sim \nu$). Alternatively we could assume a prior over $\theta$.

# The linear bandit problem, Auer [2002]

**Known parameters:** compact action set $\mathcal{A} \subset \mathbb{R}^n$, adversary's action set $\mathcal{L} \subset \mathbb{R}^n$, number of rounds $T$.

**Protocol:** For each round $t = 1, 2, \ldots, T$, the adversary chooses a loss vector $\ell_t \in \mathcal{L}$ and simultaneously the player chooses $a_t \in \mathcal{A}$ based on past observations and receives a loss/observation $Y_t = \ell_t^\top a_t$.

$$R_T = \mathbb{E} \sum_{t=1}^{T} \ell_t^\top a_t - \min_{a \in \mathcal{A}} \mathbb{E} \sum_{t=1}^{T} \ell_t^\top a.$$

**Other models:** In the i.i.d. model we assume that there is some underlying $\theta \in \mathcal{L}$ such that $\mathbb{E}(Y_t | a_t) = \theta^\top a_t$. In the Bayesian model we assume that we have a prior distribution $\nu$ over the sequence $(\ell_1, \ldots, \ell_T)$ (in this case the expectation in $R_T$ is also over $(\ell_1, \ldots, \ell_T) \sim \nu$). Alternatively we could assume a prior over $\theta$.

**Example:** Part 1 was about $\mathcal{A} = \{e_1, \ldots, e_n\}$ and $\mathcal{L} = [0, 1]^n$.

# The linear bandit problem, Auer [2002]

**Known parameters:** compact action set $\mathcal{A} \subset \mathbb{R}^n$, adversary's action set $\mathcal{L} \subset \mathbb{R}^n$, number of rounds $T$.

**Protocol:** For each round $t = 1, 2, \ldots, T$, the adversary chooses a loss vector $\ell_t \in \mathcal{L}$ and simultaneously the player chooses $a_t \in \mathcal{A}$ based on past observations and receives a loss/observation $Y_t = \ell_t^\top a_t$.
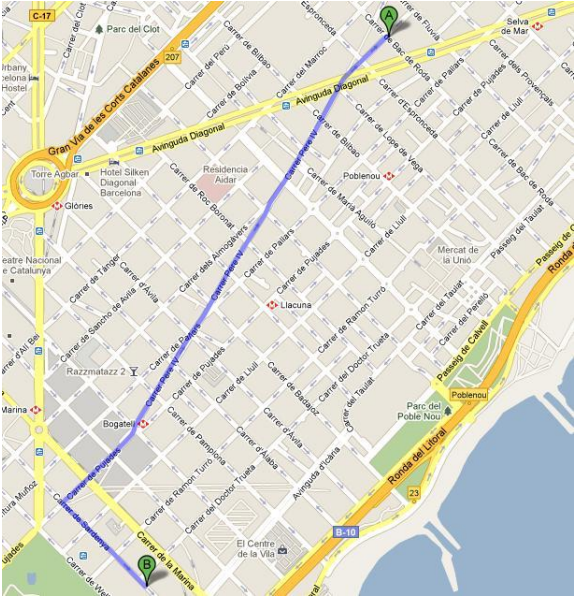
$$R_T = \mathbb{E} \sum_{t=1}^{T} \ell_t^\top a_t - \min_{a \in \mathcal{A}} \mathbb{E} \sum_{t=1}^{T} \ell_t^\top a.$$

**Other models:** In the i.i.d. model we assume that there is some underlying $\theta \in \mathcal{L}$ such that $\mathbb{E}(Y_t | a_t) = \theta^\top a_t$. In the Bayesian model we assume that we have a prior distribution $\nu$ over the sequence $(\ell_1, \ldots, \ell_T)$ (in this case the expectation in $R_T$ is also over $(\ell_1, \ldots, \ell_T) \sim \nu$). Alternatively we could assume a prior over $\theta$.

**Example:** Part 1 was about $\mathcal{A} = \{e_1, \ldots, e_n\}$ and $\mathcal{L} = [0, 1]^n$.
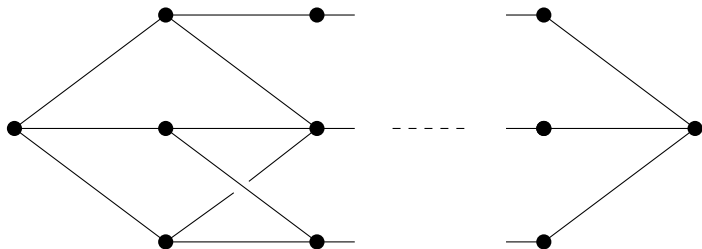
**Assumption:** unless specified otherwise we assume $\mathcal{L} = \mathcal{A}^\circ := \{\ell : \sup_{a \in \mathcal{A}} |\ell^\top a| \leq 1\}$.
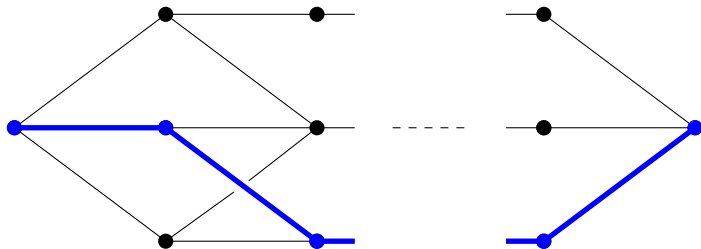
# Example: path planning

# Example: path planning

Adversary



Player

# Example: path planning

Adversary



Player →

# Example: path planning



Adversary ⟶

Player ⟶

# Example: path planning



Adversary $\longrightarrow$

$\ell_1$

$\ell_2$

$\ell_3$

$\ell_4$

$\ell_5$

$\ell_6$

$\ell_7$

$\ell_8$

$\ell_9$

$\ell_{n-2}$

$\ell_{n-1}$

$\ell_n$

Player $\longrightarrow$

# Example: path planning



Adversary →

Player →

loss suffered: $\ell_2 + \ell_7 + \ldots + \ell_n$

# Example: path planning



Adversary →

Feedback: { Full Info: $\ell_1, \ell_2, \ldots, \ell_n$

$\ell_4$

$\ell_1$

$\ell_5$

$\ell_2$

$\ell_6$

$\ell_{n-2}$

$\ell_{n-1}$

$\ell_3$

$\ell_7$

$\ell_n$

$\ell_8$

$\ell_9$

Player →

loss suffered: $\ell_2 + \ell_7 + \ldots + \ell_n$

# Example: path planning



Adversary ⟶ [sign] Feedback: 
{ Full Info: $\ell_1, \ell_2, \ldots, \ell_n$
Semi-Bandit: $\ell_2, \ell_7, \ldots, \ell_n$

$\ell_4$

$\ell_1$

$\ell_5$

$\ell_2$

$\ell_6$

$\ell_{n-2}$

$\ell_3$

$\ell_7$

$\ell_{n-1}$

$\ell_8$

$\ell_9$

$\ell_n$

Player ⟶ [map]

loss suffered: $\ell_2 + \ell_7 + \ldots + \ell_n$

# Example: path planning



Adversary → [road work sign] Feedback:
$\begin{cases} \text{Full Info:} & \ell_1, \ell_2, \ldots, \ell_n \\ \text{Semi-Bandit:} & \ell_2, \ell_7, \ldots, \ell_n \\ \text{Bandit:} & \ell_2 + \ell_7 + \ldots + \ell_n \end{cases}$

Player →

loss suffered: $\ell_2 + \ell_7 + \ldots + \ell_n$

# Thompson Sampling for linear bandit after RVR14

Assume $\mathcal{A} = \{a_1, \ldots, a_{|\mathcal{A}|}\}$. Recall from Part 1 that TS satisfies

$$\sum_i \pi_t(i)(\bar{\ell}_t(i) - \bar{\ell}_t(i, i)) \leq \sqrt{C \sum_{i,j} \pi_t(i)\pi_t(j)(\bar{\ell}_t(i, j) - \bar{\ell}_t(i))^2}$$

$$\Rightarrow R_T \leq \sqrt{C \ T \ \log(|\mathcal{A}|)/2},$$

where $\bar{\ell}_t(i) = \mathbb{E}_t \ell_t(i)$ and $\bar{\ell}_t(i, j) = \mathbb{E}_t(\ell_t(i)|i^* = j)$.

# Thompson Sampling for linear bandit after RVR14

Assume $\mathcal{A} = \{a_1, \ldots, a_{|\mathcal{A}|}\}$. Recall from Part 1 that TS satisfies

$$\sum_i \pi_t(i)(\bar{\ell}_t(i) - \bar{\ell}_t(i,i)) \leq \sqrt{C \sum_{i,j} \pi_t(i)\pi_t(j)(\bar{\ell}_t(i,j) - \bar{\ell}_t(i))^2}$$

$$\Rightarrow R_T \leq \sqrt{C\ T\ \log(|\mathcal{A}|)/2},$$

where $\bar{\ell}_t(i) = \mathbb{E}_t \ell_t(i)$ and $\bar{\ell}_t(i,j) = \mathbb{E}_t(\ell_t(i)|i^* = j)$.

Writing $\bar{\ell}_t(i) = a_i^\top \bar{\ell}_t$, $\bar{\ell}_t(i,j) = a_i^\top \bar{\ell}_t^j$, and
$(M_{i,j}) = \left( \sqrt{\pi_t(i)\pi_t(j)} a_i^\top (\bar{\ell}_t - \bar{\ell}_t^j) \right)$ we want to show that

$$\mathrm{Tr}(M) \leq \sqrt{C} \|M\|_F.$$

## Thompson Sampling for linear bandit after RVR14

Assume $\mathcal{A} = \{a_1, \dots, a_{|\mathcal{A}|}\}$. Recall from Part 1 that TS satisfies

$$\sum_i \pi_t(i)(\bar{\ell}_t(i) - \bar{\ell}_t(i,i)) \leq \sqrt{C \sum_{i,j} \pi_t(i)\pi_t(j)(\bar{\ell}_t(i,j) - \bar{\ell}_t(i))^2}$$

$$\Rightarrow R_T \leq \sqrt{C \ T \ \log(|\mathcal{A}|)/2},$$

where $\bar{\ell}_t(i) = \mathbb{E}_t \ell_t(i)$ and $\bar{\ell}_t(i,j) = \mathbb{E}_t(\ell_t(i)|i^* = j)$.

Writing $\bar{\ell}_t(i) = a_i^\top \bar{\ell}_t$, $\bar{\ell}_t(i,j) = a_i^\top \bar{\ell}_t^j$, and
$(M_{i,j}) = \left( \sqrt{\pi_t(i)\pi_t(j)} a_i^\top (\bar{\ell}_t - \bar{\ell}_t^j) \right)$ we want to show that

$$\mathrm{Tr}(M) \leq \sqrt{C} \|M\|_F.$$

Using the eigenvalue formula for the trace and the Frobenius norm
one can see that $\mathrm{Tr}(M)^2 \leq \mathrm{rank}(M)\|M\|_F^2$.

## Thompson Sampling for linear bandit after RVR14

Assume $\mathcal{A} = \{a_1, \ldots, a_{|\mathcal{A}|}\}$. Recall from Part 1 that TS satisfies

$$\sum_i \pi_t(i)(\bar{\ell}_t(i) - \bar{\ell}_t(i,i)) \leq \sqrt{C \sum_{i,j} \pi_t(i)\pi_t(j)(\bar{\ell}_t(i,j) - \bar{\ell}_t(i))^2}$$

$$\Rightarrow R_T \leq \sqrt{C\ T\ \log(|\mathcal{A}|)/2},$$

where $\bar{\ell}_t(i) = \mathbb{E}_t \ell_t(i)$ and $\bar{\ell}_t(i,j) = \mathbb{E}_t(\ell_t(i)|i^* = j)$.

Writing $\bar{\ell}_t(i) = a_i^\top \bar{\ell}_t$, $\bar{\ell}_t(i,j) = a_i^\top \bar{\ell}_t^j$, and $(M_{i,j}) = \left(\sqrt{\pi_t(i)\pi_t(j)}a_i^\top(\bar{\ell}_t - \bar{\ell}_t^j)\right)$ we want to show that

$$\mathrm{Tr}(M) \leq \sqrt{C}\|M\|_F.$$

Using the eigenvalue formula for the trace and the Frobenius norm one can see that $\mathrm{Tr}(M)^2 \leq \mathrm{rank}(M)\|M\|_F^2$. Moreover the rank of $M$ is at most $n$ since $M = UV^\top$ where $U, V \in \mathbb{R}^{|\mathcal{A}| \times n}$ (the $i^{th}$ row of $U$ is $\sqrt{\pi_t(i)}a_i$ and for $V$ it is $\sqrt{\pi_t(i)}(\bar{\ell}_t - \bar{\ell}_t^i)$).

1. TS satisfies $R_T \leq \sqrt{nT \log(|\mathcal{A}|)}$. To appreciate the improvement recall that without the linear structure one would get a regret of order $\sqrt{|\mathcal{A}| T}$ and that $\mathcal{A}$ can be exponential in the dimension $n$ (think of the path planning example).

# Thompson Sampling for linear bandit after RVR14

1. TS satisfies $R_T \leq \sqrt{nT \log(|\mathcal{A}|)}$. To appreciate the improvement recall that without the linear structure one would get a regret of order $\sqrt{|\mathcal{A}|T}$ and that $\mathcal{A}$ can be exponential in the dimension $n$ (think of the path planning example).

2. Provided that one can efficiently sample from the posterior on $\ell_t$ (or on $\theta$), TS just requires at each step one linear optimization over $\mathcal{A}$.

# Thompson Sampling for linear bandit after RVR14

1. TS satisfies $R_T \leq \sqrt{nT \log(|\mathcal{A}|)}$. To appreciate the improvement recall that without the linear structure one would get a regret of order $\sqrt{|\mathcal{A}| T}$ and that $\mathcal{A}$ can be exponential in the dimension $n$ (think of the path planning example).

2. Provided that one can efficiently sample from the posterior on $\ell_t$ (or on $\theta$), TS just requires at each step one linear optimization over $\mathcal{A}$.

3. TS regret bound is optimal in the following sense. W.l.og. one can assume $|\mathcal{A}| \leq (10T)^n$ and thus TS satisfies $R_T = O(n\sqrt{T \log(T)})$ for any action set. Furthermore one can show that there exists an action set and a prior such that for any strategy one has $R_T = \Omega(n\sqrt{T})$, see Dani, Hayes and Kakade [2008], Rusmevichientong and Tsitsiklis [2010], and Audibert, Bubeck and Lugosi [2011, 2014].

# Adversarial linear bandit after Dani, Hayes, Kakade [2008]

Recall from Part 1 that exponential weights satisfies for any $\widetilde{\ell}_t$ such that $\mathbb{E}\widetilde{\ell}_t(i) = \ell_t(i)$ and $\widetilde{\ell}_t(i) \geq 0$,

$$R_T \leq \frac{\max_i \operatorname{Ent}(\delta_i \| p_1)}{\eta} + \frac{\eta}{2}\mathbb{E}\sum_t \mathbb{E}_{I \sim p_t}\widetilde{\ell}_t(I)^2.$$

# Adversarial linear bandit after Dani, Hayes, Kakade [2008]

Recall from Part 1 that exponential weights satisfies for any $\widetilde{\ell}_t$ such that $\mathbb{E}\widetilde{\ell}_t(i) = \ell_t(i)$ and $\widetilde{\ell}_t(i) \geq 0$,

$$R_T \leq \frac{\max_i \text{Ent}(\delta_i \| p_1)}{\eta} + \frac{\eta}{2} \mathbb{E} \sum_t \mathbb{E}_{I \sim p_t} \widetilde{\ell}_t(I)^2.$$

DHK08 proposed the following (beautiful) unbiased estimator for the linear case:

$$\widetilde{\ell}_t = \Sigma_t^{-1} a_t a_t^\top \ell_t \text{ where } \Sigma_t = \mathbb{E}_{a \sim p_t}(aa^\top).$$

# Adversarial linear bandit after Dani, Hayes, Kakade [2008]

Recall from Part 1 that exponential weights satisfies for any $\widetilde{\ell}_t$ such that $\mathbb{E}\widetilde{\ell}_t(i) = \ell_t(i)$ and $\widetilde{\ell}_t(i) \geq 0$,

$$R_T \leq \frac{\max_i \operatorname{Ent}(\delta_i \| p_1)}{\eta} + \frac{\eta}{2} \mathbb{E} \sum_t \mathbb{E}_{I \sim p_t} \widetilde{\ell}_t(I)^2.$$

DHK08 proposed the following (beautiful) unbiased estimator for the linear case:

$$\widetilde{\ell}_t = \Sigma_t^{-1} a_t a_t^\top \ell_t \text{ where } \Sigma_t = \mathbb{E}_{a \sim p_t}(a a^\top).$$

Again, amazingly, the variance is automatically controlled:

$$\mathbb{E}(\mathbb{E}_{a \sim p_t}(\widetilde{\ell}_t^\top a)^2) = \mathbb{E}\widetilde{\ell}_t^\top \Sigma_t \widetilde{\ell}_t \leq \mathbb{E}a_t^\top \Sigma_t^{-1} a_t = \mathbb{E}\operatorname{Tr}(\Sigma_t^{-1} a_t a_t) = n.$$

# Adversarial linear bandit after Dani, Hayes, Kakade [2008]

Recall from Part 1 that exponential weights satisfies for any $\widetilde{\ell}_t$ such that $\mathbb{E}\widetilde{\ell}_t(i) = \ell_t(i)$ and $\widetilde{\ell}_t(i) \geq 0$,

$$R_T \leq \frac{\max_i \text{Ent}(\delta_i \| p_1)}{\eta} + \frac{\eta}{2} \mathbb{E} \sum_t \mathbb{E}_{I \sim p_t} \widetilde{\ell}_t(I)^2.$$

DHK08 proposed the following (beautiful) unbiased estimator for the linear case:

$$\widetilde{\ell}_t = \Sigma_t^{-1} a_t a_t^\top \ell_t \text{ where } \Sigma_t = \mathbb{E}_{a \sim p_t}(aa^\top).$$

Again, amazingly, the variance is automatically controlled:

$$\mathbb{E}(\mathbb{E}_{a \sim p_t}(\widetilde{\ell}_t^\top a)^2) = \mathbb{E}\widetilde{\ell}_t^\top \Sigma_t \widetilde{\ell}_t \leq \mathbb{E}a_t^\top \Sigma_t^{-1} a_t = \mathbb{E}\text{Tr}(\Sigma_t^{-1} a_t a_t) = n.$$

Up to the issue that $\widetilde{\ell}_t$ can take negative values this suggests the "optimal" $\sqrt{nT \log(|\mathcal{A}|)}$ regret bound.

## Adversarial linear bandit, further development

1. The non-negativity issue of $\widetilde{\ell}_t$ is a manifestation of the need for an added exploration. DHK08 used a suboptimal exploration which led to an additional $\sqrt{n}$ in the regret. This was later improved in Bubeck, Cesa-Bianchi, and Kakade [2012] with an exploration based on the John's ellipsoid (smallest ellipsoid containing $\mathcal{A}$).

## Adversarial linear bandit, further development

1. The non-negativity issue of $\widetilde{\ell}_t$ is a manifestation of the need for an added exploration. DHK08 used a suboptimal exploration which led to an additional $\sqrt{n}$ in the regret. This was later improved in Bubeck, Cesa-Bianchi, and Kakade [2012] with an exploration based on the John's ellipsoid (smallest ellipsoid containing $\mathcal{A}$).

2. Sampling the exp. weights is usually computationally difficult, see Cesa-Bianchi and Lugosi [2009] for some exceptions.

# Adversarial linear bandit, further development

1. The non-negativity issue of $\widetilde{\ell}_t$ is a manifestation of the need for an added exploration. DHK08 used a suboptimal exploration which led to an additional $\sqrt{n}$ in the regret. This was later improved in Bubeck, Cesa-Bianchi, and Kakade [2012] with an exploration based on the John's ellipsoid (smallest ellipsoid containing $\mathcal{A}$).

2. Sampling the exp. weights is usually computationally difficult, see Cesa-Bianchi and Lugosi [2009] for some exceptions.

3. Abernethy, Hazan and Rakhlin [2008] proposed an alternative (beautiful) strategy based on mirror descent. The key idea is to use a $n$-self-concordant barrier for $\mathrm{conv}(\mathcal{A})$ as a mirror map and to sample points uniformly in Dikin ellipses. This method's regret is suboptimal by a factor $\sqrt{n}$ and the computational efficiency depends on the barrier being used.

# Adversarial linear bandit, further development

1. The non-negativity issue of $\widetilde{\ell}_t$ is a manifestation of the need for an added exploration. DHK08 used a suboptimal exploration which led to an additional $\sqrt{n}$ in the regret. This was later improved in Bubeck, Cesa-Bianchi, and Kakade [2012] with an exploration based on the John's ellipsoid (smallest ellipsoid containing $\mathcal{A}$).

2. Sampling the exp. weights is usually computationally difficult, see Cesa-Bianchi and Lugosi [2009] for some exceptions.

3. Abernethy, Hazan and Rakhlin [2008] proposed an alternative (beautiful) strategy based on mirror descent. The key idea is to use a $n$-self-concordant barrier for $\mathrm{conv}(\mathcal{A})$ as a mirror map and to sample points uniformly in Dikin ellipses. This method's regret is suboptimal by a factor $\sqrt{n}$ and the computational efficiency depends on the barrier being used.

4. Bubeck and Eldan [2014]'s entropic barrier allows for a much more information-efficient sampling than AHR08. This gives another strategy with optimal regret which is efficient when $\mathcal{A}$ is convex (and one can do linear optimization on $\mathcal{A}$).

# Adversarial combinatorial bandit after Audibert, Bubeck and Lugosi [2011, 2014]

Combinatorial setting: $\mathcal{A} \subset \{0,1\}^n$, $\max_a \|a\|_1 = m$, $\mathcal{L} = [0,1]^n$.

# Adversarial combinatorial bandit after Audibert, Bubeck and Lugosi [2011, 2014]

Combinatorial setting: $\mathcal{A} \subset \{0,1\}^n$, $\max_a \|a\|_1 = m$, $\mathcal{L} = [0,1]^n$.

1. Full information case goes back to the end of the 90's (Warmuth and co-authors), semi-bandit and bandit were introduced in Audibert, Bubeck and Lugosi [2011] (following several papers that studied specific sets $\mathcal{A}$).

# Adversarial combinatorial bandit after Audibert, Bubeck and Lugosi [2011, 2014]

Combinatorial setting: $\mathcal{A} \subset \{0,1\}^n$, $\max_a \|a\|_1 = m$, $\mathcal{L} = [0,1]^n$.

1. Full information case goes back to the end of the 90's (Warmuth and co-authors), semi-bandit and bandit were introduced in Audibert, Bubeck and Lugosi [2011] (following several papers that studied specific sets $\mathcal{A}$).

2. This is a natural setting to study FPL-type (Follow the Perturbed Leader) strategies, see e.g. Kalai and Vempala [2004] and more recently Devroye, Lugosi and Neu [2013].

# Adversarial combinatorial bandit after Audibert, Bubeck and Lugosi [2011, 2014]

Combinatorial setting: $\mathcal{A} \subset \{0,1\}^n$, $\max_a \|a\|_1 = m$, $\mathcal{L} = [0,1]^n$.

1. Full information case goes back to the end of the 90's (Warmuth and co-authors), semi-bandit and bandit were introduced in Audibert, Bubeck and Lugosi [2011] (following several papers that studied specific sets $\mathcal{A}$).

2. This is a natural setting to study FPL-type (Follow the Perturbed Leader) strategies, see e.g. Kalai and Vempala [2004] and more recently Devroye, Lugosi and Neu [2013].

3. ABL11: Exponential weights is provably suboptimal in this setting! This is in sharp contrast with the case where $\mathcal{L} = \mathcal{A}^\circ$.

# Adversarial combinatorial bandit after Audibert, Bubeck and Lugosi [2011, 2014]

Combinatorial setting: $\mathcal{A} \subset \{0,1\}^n$, $\max_a \|a\|_1 = m$, $\mathcal{L} = [0,1]^n$.

1. Full information case goes back to the end of the 90's (Warmuth and co-authors), semi-bandit and bandit were introduced in Audibert, Bubeck and Lugosi [2011] (following several papers that studied specific sets $\mathcal{A}$).

2. This is a natural setting to study FPL-type (Follow the Perturbed Leader) strategies, see e.g. Kalai and Vempala [2004] and more recently Devroye, Lugosi and Neu [2013].

3. ABL11: Exponential weights is provably suboptimal in this setting! This is in sharp contrast with the case where $\mathcal{L} = \mathcal{A}^\circ$.

4. Optimal regret in the semi-bandit case is $\sqrt{mnT}$ and it can be achieved with mirror descent and the natural unbiased estimator for the semi-bandit situation.

# Adversarial combinatorial bandit after Audibert, Bubeck and Lugosi [2011, 2014]

Combinatorial setting: $\mathcal{A} \subset \{0,1\}^n$, $\max_a \|a\|_1 = m$, $\mathcal{L} = [0,1]^n$.

1. Full information case goes back to the end of the 90's (Warmuth and co-authors), semi-bandit and bandit were introduced in Audibert, Bubeck and Lugosi [2011] (following several papers that studied specific sets $\mathcal{A}$).

2. This is a natural setting to study FPL-type (Follow the Perturbed Leader) strategies, see e.g. Kalai and Vempala [2004] and more recently Devroye, Lugosi and Neu [2013].

3. ABL11: Exponential weights is provably suboptimal in this setting! This is in sharp contrast with the case where $\mathcal{L} = \mathcal{A}^\circ$.

4. Optimal regret in the semi-bandit case is $\sqrt{mnT}$ and it can be achieved with mirror descent and the natural unbiased estimator for the semi-bandit situation.

5. For the bandit case the bound for exponential weights from the previous slides gives $m\sqrt{mnT}$. However the lower bound from ABL14 is $m\sqrt{nT}$, which is conjectured to be tight.

# Preliminaries for the i.i.d. case: a primer on least squares

Assume $Y_t = \theta^\top a_t + \xi_t$ where $(\xi_t)$ is an i.i.d. sequence of centered and sub-Gaussian real-valued random variables. The (regularized) least squares estimator for $\theta$ based on $\mathbb{Y}_t = (Y_1, \ldots, Y_{t-1})^\top$ is, with $\mathbb{A}_t = (a_1 \ldots a_{t-1}) \in \mathbb{R}^{n \times t-1}$ and $\Sigma_t = \lambda I_n + \sum_{s=1}^{t-1} a_s a_s^\top$:

$$\hat{\theta}_t = \Sigma_t^{-1} \mathbb{A}_t \mathbb{Y}_t$$

# Preliminaries for the i.i.d. case: a primer on least squares

Assume $Y_t = \theta^\top a_t + \xi_t$ where $(\xi_t)$ is an i.i.d. sequence of centered and sub-Gaussian real-valued random variables. The (regularized) least squares estimator for $\theta$ based on $\mathbb{Y}_t = (Y_1, \ldots, Y_{t-1})^\top$ is, with $\mathbb{A}_t = (a_1 \ldots a_{t-1}) \in \mathbb{R}^{n \times t-1}$ and $\Sigma_t = \lambda I_n + \sum_{s=1}^{t-1} a_s a_s^\top$:

$$\hat{\theta}_t = \Sigma_t^{-1} \mathbb{A}_t \mathbb{Y}_t$$

Observe that we can also write $\theta = \Sigma_t^{-1}(\mathbb{A}_t(\mathbb{Y}_t + \varepsilon_t) + \lambda\theta)$ where $\varepsilon_t = (\mathbb{E}(Y_1|a_1) - Y_1, \ldots, \mathbb{E}(Y_{t-1}|a_{t-1}) - Y_{t-1})^\top$

# Preliminaries for the i.i.d. case: a primer on least squares

Assume $Y_t = \theta^\top a_t + \xi_t$ where $(\xi_t)$ is an i.i.d. sequence of centered and sub-Gaussian real-valued random variables. The (regularized) least squares estimator for $\theta$ based on $\mathbb{Y}_t = (Y_1, \ldots, Y_{t-1})^\top$ is, with $\mathbb{A}_t = (a_1 \ldots a_{t-1}) \in \mathbb{R}^{n \times t-1}$ and $\Sigma_t = \lambda I_n + \sum_{s=1}^{t-1} a_s a_s^\top$:

$$\hat{\theta}_t = \Sigma_t^{-1} \mathbb{A}_t \mathbb{Y}_t$$

Observe that we can also write $\theta = \Sigma_t^{-1}(\mathbb{A}_t(\mathbb{Y}_t + \varepsilon_t) + \lambda\theta)$ where $\varepsilon_t = (\mathbb{E}(Y_1|a_1) - Y_1, \ldots, \mathbb{E}(Y_{t-1}|a_{t-1}) - Y_{t-1})^\top$ so that

$$\|\theta - \hat{\theta}_t\|_{\Sigma_t} = \|\mathbb{A}_t \varepsilon_t + \lambda\theta\|_{\Sigma_t^{-1}} \leq \|\mathbb{A}_t \varepsilon_t\|_{\Sigma_t^{-1}} + \sqrt{\lambda}\|\theta\|.$$

## Preliminaries for the i.i.d. case: a primer on least squares

Assume $Y_t = \theta^\top a_t + \xi_t$ where $(\xi_t)$ is an i.i.d. sequence of centered and sub-Gaussian real-valued random variables. The (regularized) least squares estimator for $\theta$ based on $\mathbb{Y}_t = (Y_1, \ldots, Y_{t-1})^\top$ is, with $\mathbb{A}_t = (a_1 \ldots a_{t-1}) \in \mathbb{R}^{n \times t-1}$ and $\Sigma_t = \lambda I_n + \sum_{s=1}^{t-1} a_s a_s^\top$:

$$\hat{\theta}_t = \Sigma_t^{-1} \mathbb{A}_t \mathbb{Y}_t$$

Observe that we can also write $\theta = \Sigma_t^{-1}(\mathbb{A}_t(\mathbb{Y}_t + \varepsilon_t) + \lambda\theta)$ where $\varepsilon_t = (\mathbb{E}(Y_1|a_1) - Y_1, \ldots, \mathbb{E}(Y_{t-1}|a_{t-1}) - Y_{t-1})^\top$ so that

$$\|\theta - \hat{\theta}_t\|_{\Sigma_t} = \|\mathbb{A}_t\varepsilon_t + \lambda\theta\|_{\Sigma_t^{-1}} \leq \|\mathbb{A}_t\varepsilon_t\|_{\Sigma_t^{-1}} + \sqrt{\lambda}\|\theta\|.$$

A basic martingale argument (see e.g., Abbasi-Yadkori, Pál and Szepesvári [2011]) shows that w.p. $\geq 1 - \delta$, $\forall t \geq 1$,

$$\|\mathbb{A}_t\varepsilon_t\|_{\Sigma_t^{-1}} \leq \sqrt{\log\det(\Sigma_t) + \log(1/(\delta^2\lambda^n))}.$$

## Preliminaries for the i.i.d. case: a primer on least squares

Assume $Y_t = \theta^\top a_t + \xi_t$ where $(\xi_t)$ is an i.i.d. sequence of centered and sub-Gaussian real-valued random variables. The (regularized) least squares estimator for $\theta$ based on $\mathbb{Y}_t = (Y_1, \ldots, Y_{t-1})^\top$ is, with $\mathbb{A}_t = (a_1 \ldots a_{t-1}) \in \mathbb{R}^{n \times t-1}$ and $\Sigma_t = \lambda \mathrm{I}_n + \sum_{s=1}^{t-1} a_s a_s^\top$:

$$\hat{\theta}_t = \Sigma_t^{-1} \mathbb{A}_t \mathbb{Y}_t$$

Observe that we can also write $\theta = \Sigma_t^{-1}(\mathbb{A}_t(\mathbb{Y}_t + \varepsilon_t) + \lambda\theta)$ where $\varepsilon_t = (\mathbb{E}(Y_1|a_1) - Y_1, \ldots, \mathbb{E}(Y_{t-1}|a_{t-1}) - Y_{t-1})^\top$ so that

$$\|\theta - \hat{\theta}_t\|_{\Sigma_t} = \|\mathbb{A}_t \varepsilon_t + \lambda\theta\|_{\Sigma_t^{-1}} \leq \|\mathbb{A}_t \varepsilon_t\|_{\Sigma_t^{-1}} + \sqrt{\lambda}\|\theta\|.$$

A basic martingale argument (see e.g., Abbasi-Yadkori, Pál and Szepesvári [2011]) shows that w.p. $\geq 1 - \delta$, $\forall t \geq 1$,

$$\|\mathbb{A}_t \varepsilon_t\|_{\Sigma_t^{-1}} \leq \sqrt{\mathrm{logdet}(\Sigma_t) + \log(1/(\delta^2 \lambda^n))}.$$

Note that $\mathrm{logdet}(\Sigma_t) \leq n \log(\mathrm{Tr}(\Sigma_t)/n) \leq n \log(\lambda + t/n)$ (w.l.o.g. we assumed $\|a_t\| \leq 1$).

## i.i.d. linear bandit after DHK08, RT10, AYPS11

Let $\beta = 2\sqrt{n \log(T)}$, and $\mathcal{E}_t = \{\theta' : \|\theta' - \hat{\theta}_t\|_{\Sigma_t} \leq \beta\}$. We showed that w.p. $\geq 1 - 1/T^2$ one has $\theta \in \mathcal{E}_t$ for all $t \in [T]$.

## i.i.d. linear bandit after DHK08, RT10, AYPS11

Let $\beta = 2\sqrt{n\log(T)}$, and $\mathcal{E}_t = \{\theta' : \|\theta' - \hat{\theta}_t\|_{\Sigma_t} \leq \beta\}$. We showed that w.p. $\geq 1 - 1/T^2$ one has $\theta \in \mathcal{E}_t$ for all $t \in [T]$.

The appropriate generalization of UCB is to select:
$(\widetilde{\theta}_t, a_t) = \operatorname{argmin}_{(\theta', a) \in \mathcal{E}_t \times \mathcal{A}} \theta'^\top a$ (this optimization is NP-hard in general, more on that next slide).

## i.i.d. linear bandit after DHK08, RT10, AYPS11

Let $\beta = 2\sqrt{n\log(T)}$, and $\mathcal{E}_t = \{\theta' : \|\theta' - \hat{\theta}_t\|_{\Sigma_t} \leq \beta\}$. We showed that w.p. $\geq 1 - 1/T^2$ one has $\theta \in \mathcal{E}_t$ for all $t \in [T]$.

The appropriate generalization of UCB is to select:
$(\widetilde{\theta}_t, a_t) = \operatorname{argmin}_{(\theta', a) \in \mathcal{E}_t \times \mathcal{A}} \theta'^\top a$ (this optimization is NP-hard in general, more on that next slide). Then one has on the high-probability event:

$$\sum_{t=1}^T \theta^\top (a_t - a^*) \leq \sum_{t=1}^T (\theta - \widetilde{\theta}_t)^\top a_t \leq \beta \sum_{t=1}^T \|a_t\|_{\Sigma_t^{-1}} \leq \beta \sqrt{T \sum_t \|a_t\|_{\Sigma_t^{-1}}^2}.$$

# i.i.d. linear bandit after DHK08, RT10, AYPS11

Let $\beta = 2\sqrt{n \log(T)}$, and $\mathcal{E}_t = \{\theta' : \|\theta' - \hat{\theta}_t\|_{\Sigma_t} \leq \beta\}$. We showed that w.p. $\geq 1 - 1/T^2$ one has $\theta \in \mathcal{E}_t$ for all $t \in [T]$.

The appropriate generalization of UCB is to select:
$(\widetilde{\theta}_t, a_t) = \mathrm{argmin}_{(\theta', a) \in \mathcal{E}_t \times \mathcal{A}} \theta'^{\top} a$ (this optimization is NP-hard in general, more on that next slide). Then one has on the high-probability event:

$$\sum_{t=1}^{T} \theta^{\top}(a_t - a^*) \leq \sum_{t=1}^{T} (\theta - \widetilde{\theta}_t)^{\top} a_t \leq \beta \sum_{t=1}^{T} \|a_t\|_{\Sigma_t^{-1}} \leq \beta \sqrt{T \sum_t \|a_t\|_{\Sigma_t^{-1}}^2}.$$

To control the sum of squares we observe that:

$\det(\Sigma_{t+1}) = \det(\Sigma_t) \det(I_n + \Sigma_t^{-1/2} a_t (\Sigma_t^{-1/2} a_t)^{\top}) = \det(\Sigma_t)(1 + \|a_t\|_{\Sigma_t^{-1}}^2)$

so that (assuming $\lambda \geq 1$)

$\log \det(\Sigma_{T+1}) - \log \det(\Sigma_1) = \sum_t \log(1 + \|a_t\|_{\Sigma_t^{-1}}^2) \geq \frac{1}{2} \sum_t \|a_t\|_{\Sigma_t^{-1}}^2.$

### i.i.d. linear bandit after DHK08, RT10, AYPS11

Let $\beta = 2\sqrt{n \log(T)}$, and $\mathcal{E}_t = \{\theta' : \|\theta' - \hat{\theta}_t\|_{\Sigma_t} \leq \beta\}$. We showed that w.p. $\geq 1 - 1/T^2$ one has $\theta \in \mathcal{E}_t$ for all $t \in [T]$.

The appropriate generalization of UCB is to select:
$(\widetilde{\theta}_t, a_t) = \mathrm{argmin}_{(\theta', a) \in \mathcal{E}_t \times \mathcal{A}} {\theta'}^\top a$ (this optimization is NP-hard in general, more on that next slide). Then one has on the high-probability event:

$$\sum_{t=1}^{T} \theta^\top (a_t - a^*) \leq \sum_{t=1}^{T} (\theta - \widetilde{\theta}_t)^\top a_t \leq \beta \sum_{t=1}^{T} \|a_t\|_{\Sigma_t^{-1}} \leq \beta \sqrt{T \sum_t \|a_t\|_{\Sigma_t^{-1}}^2}.$$

To control the sum of squares we observe that:
$$\det(\Sigma_{t+1}) = \det(\Sigma_t) \det(\mathrm{I}_n + \Sigma_t^{-1/2} a_t (\Sigma_t^{-1/2} a_t)^\top) = \det(\Sigma_t)(1 + \|a_t\|_{\Sigma_t^{-1}}^2)$$

so that (assuming $\lambda \geq 1$)
$$\log \det(\Sigma_{T+1}) - \log \det(\Sigma_1) = \sum_t \log(1 + \|a_t\|_{\Sigma_t^{-1}}^2) \geq \frac{1}{2} \sum_t \|a_t\|_{\Sigma_t^{-1}}^2.$$

Putting things together we see that the regret is $O(n \log(T)\sqrt{T})$.

# What's the point of i.i.d. linear bandit?

So far we did not get any real benefit from the i.i.d. assumption (the regret guarantee we obtained is the same as for the adversarial model). To me the key benefit is in the simplicity of the i.i.d. algorithm which makes it easy to incorporate further assumptions.

# What's the point of i.i.d. linear bandit?

So far we did not get any real benefit from the i.i.d. assumption (the regret guarantee we obtained is the same as for the adversarial model). To me the key benefit is in the simplicity of the i.i.d. algorithm which makes it easy to incorporate further assumptions.

1. Sparsity of $\theta$: instead of regularization with $\ell_2$-norm to define $\hat{\theta}$ one could regularize with $\ell_1$-norm, see e.g., Johnson, Sivakumar and Banerjee [2016].

# What's the point of i.i.d. linear bandit?

So far we did not get any real benefit from the i.i.d. assumption (the regret guarantee we obtained is the same as for the adversarial model). To me the key benefit is in the simplicity of the i.i.d. algorithm which makes it easy to incorporate further assumptions.

1. Sparsity of $\theta$: instead of regularization with $\ell_2$-norm to define $\hat{\theta}$ one could regularize with $\ell_1$-norm, see e.g., Johnson, Sivakumar and Banerjee [2016].

2. Computational constraint: instead of optimizing over $\mathcal{E}_t$ to define $\widetilde{\theta}_t$ one could optimize over a hypercube containing $\mathcal{E}_t$ (this would cost an extra $\sqrt{n}$ in the regret bound).

# What's the point of i.i.d. linear bandit?

So far we did not get any real benefit from the i.i.d. assumption (the regret guarantee we obtained is the same as for the adversarial model). To me the key benefit is in the simplicity of the i.i.d. algorithm which makes it easy to incorporate further assumptions.

1. Sparsity of $\theta$: instead of regularization with $\ell_2$-norm to define $\hat{\theta}$ one could regularize with $\ell_1$-norm, see e.g., Johnson, Sivakumar and Banerjee [2016].

2. Computational constraint: instead of optimizing over $\mathcal{E}_t$ to define $\widetilde{\theta}_t$ one could optimize over a hypercube containing $\mathcal{E}_t$ (this would cost an extra $\sqrt{n}$ in the regret bound).

3. Generalized linear model: $\mathbb{E}(Y_t|a_t) = \sigma(\theta^\top a_t)$ for some known increasing $\sigma : \mathbb{R} \to \mathbb{R}$, see Filippi, Cappe, Garivier and Szepesvari [2011].

# What's the point of i.i.d. linear bandit?

So far we did not get any real benefit from the i.i.d. assumption (the regret guarantee we obtained is the same as for the adversarial model). To me the key benefit is in the simplicity of the i.i.d. algorithm which makes it easy to incorporate further assumptions.
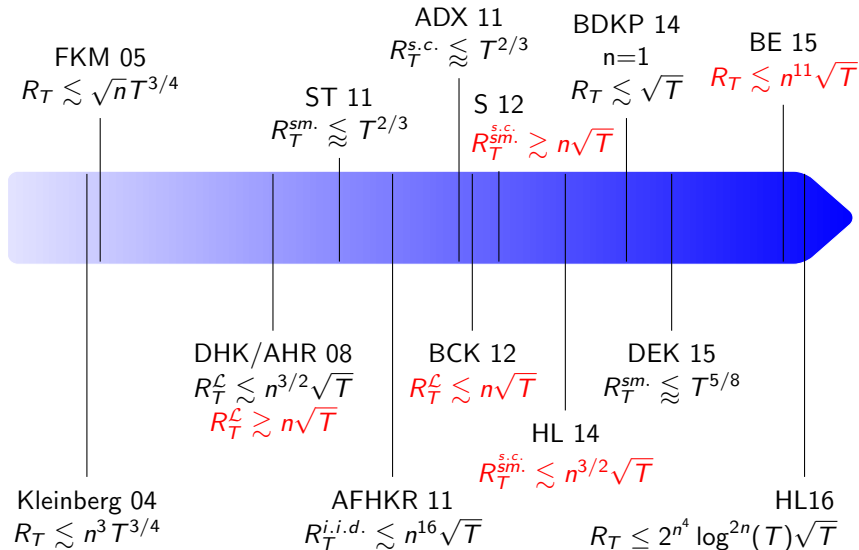
1. Sparsity of $\theta$: instead of regularization with $\ell_2$-norm to define $\hat{\theta}$ one could regularize with $\ell_1$-norm, see e.g., Johnson, Sivakumar and Banerjee [2016].

2. Computational constraint: instead of optimizing over $\mathcal{E}_t$ to define $\widetilde{\theta}_t$ one could optimize over a hypercube containing $\mathcal{E}_t$ (this would cost an extra $\sqrt{n}$ in the regret bound).

3. Generalized linear model: $\mathbb{E}(Y_t|a_t) = \sigma(\theta^\top a_t)$ for some known increasing $\sigma : \mathbb{R} \to \mathbb{R}$, see Filippi, Cappe, Garivier and Szepesvari [2011].

4. $\log(T)$-regime: if $\mathcal{A}$ is finite (note that a polytope is effectively finite for us) one can get $n^2 \log^2(T)/\Delta$ regret:
$$R_T \leq \mathbb{E} \sum_{t=1}^{T} \frac{(\theta^\top(a_t - a^*))^2}{\Delta} \leq \frac{\beta^2}{\Delta} \mathbb{E} \sum_{t=1}^{T} \|a_t\|_{\Sigma_t^{-1}}^2 \lesssim \frac{n^2 \log^2(T)}{\Delta}.$$

# Some non-linear bandit problems

**Lipschitz bandit:** Kleinberg, Slivkins and Upfal [2008, 2016], Bubeck, Munos, Stoltz and Szepesvari [2008, 2011];

**Gaussian process bandit:** Srinivas, Krause, Kakade and Seeger [2010]; and **convex bandit:**

ADX 11
$R_T^{s.c.} \lessapprox T^{2/3}$

BDKP 14
n=1
$R_T \lesssim \sqrt{T}$

BE 15
$R_T \lesssim n^{11}\sqrt{T}$

FKM 05
$R_T \lesssim \sqrt{n} T^{3/4}$

ST 11
$R_T^{sm.} \lessapprox T^{2/3}$

S 12
$R_T^{s.c.} \gtrsim n\sqrt{T}$

DHK/AHR 08
$R_T^{\mathcal{L}} \lesssim n^{3/2}\sqrt{T}$
$R_T^{\mathcal{L}} \gtrsim n\sqrt{T}$

BCK 12
$R_T^{\mathcal{L}} \lesssim n\sqrt{T}$

DEK 15
$R_T^{sm.} \lessapprox T^{5/8}$

HL 14
$R_T^{s.c.} \lesssim n^{3/2}\sqrt{T}$

Kleinberg 04
$R_T \lesssim n^3 T^{3/4}$

AFHKR 11
$R_T^{i.i.d.} \lesssim n^{16}\sqrt{T}$

HL16
$R_T \leq 2^{n^4} \log^{2n}(T)\sqrt{T}$

## Contextual bandit

We now make the game-changing assumption that at the
beginning of each round $t$ a *context* $x_t \in \mathcal{X}$ is revealed to the
player. The ideal notion of regret is now:

$$R_T^{\mathrm{ctx}} = \sum_{t=1}^{T} \ell_t(a_t) - \inf_{\Phi : \mathcal{X} \to \mathcal{A}} \sum_{t=1}^{T} \ell_t(\Phi(x_t)).$$

## Contextual bandit

We now make the game-changing assumption that at the beginning of each round $t$ a *context* $x_t \in \mathcal{X}$ is revealed to the player. The ideal notion of regret is now:

$$R_T^{\mathrm{ctx}} = \sum_{t=1}^{T} \ell_t(a_t) - \inf_{\Phi:\mathcal{X}\to\mathcal{A}} \sum_{t=1}^{T} \ell_t(\Phi(x_t)).$$

## Contextual bandit

We now make the game-changing assumption that at the beginning of each round $t$ a *context* $x_t \in \mathcal{X}$ is revealed to the player. The ideal notion of regret is now:

$$R_T^{\mathrm{ctx}} = \sum_{t=1}^{T} \ell_t(a_t) - \inf_{\Phi:\mathcal{X}\to\mathcal{A}} \sum_{t=1}^{T} \ell_t(\Phi(x_t)).$$

Sometimes it makes sense to restrict the mapping from contexts to actions, so that the infimum is taken over some *policy set* $\Pi \subset \mathcal{A}^{\mathcal{X}}$.

## Contextual bandit

We now make the game-changing assumption that at the beginning of each round $t$ a *context* $x_t \in \mathcal{X}$ is revealed to the player. The ideal notion of regret is now:

$$R_T^{\mathrm{ctx}} = \sum_{t=1}^T \ell_t(a_t) - \inf_{\Phi: \mathcal{X} \to \mathcal{A}} \sum_{t=1}^T \ell_t(\Phi(x_t)).$$

Sometimes it makes sense to restrict the mapping from contexts to actions, so that the infimum is taken over some *policy set* $\Pi \subset \mathcal{A}^{\mathcal{X}}$.

As far as I can tell the contextual bandit problem is an infinite playground and there is no canonical solution (or at least not yet!). Thankfully all we have learned so far can give useful guidance in this challenging problem.

# Linear model after embedding

A natural assumption in several application domains is to suppose linearity in the loss after a correct embedding. Say we know mappings $(\varphi_a)_{a \in \mathcal{A}}$ such that $\mathbb{E}_t(\ell_t(a)) = \varphi_a(x_t)^\top \theta$ for some unknown $\theta \in \mathbb{R}^n$ (or in the adversarial case that $\ell_t(a) = \ell_t^\top \varphi_a(x_t)$).

# Linear model after embedding

A natural assumption in several application domains is to suppose linearity in the loss after a correct embedding. Say we know mappings $(\varphi_a)_{a \in \mathcal{A}}$ such that $\mathbb{E}_t(\ell_t(a)) = \varphi_a(x_t)^\top \theta$ for some unknown $\theta \in \mathbb{R}^n$ (or in the adversarial case that $\ell_t(a) = \ell_t^\top \varphi_a(x_t)$).

This is nothing but a linear bandit problem where the action set is changing over time. All the strategies we described are robust to this modification and thus in this case one can get a regret of $\sqrt{nT \log(|\mathcal{A}|)} \lesssim n\sqrt{T \log(T)}$ (and for the stochastic case one can get efficiently $n^{3/2}\sqrt{T}$).

## Linear model after embedding

A natural assumption in several application domains is to suppose linearity in the loss after a correct embedding. Say we know mappings $(\varphi_a)_{a \in \mathcal{A}}$ such that $\mathbb{E}_t(\ell_t(a)) = \varphi_a(x_t)^\top \theta$ for some unknown $\theta \in \mathbb{R}^n$ (or in the adversarial case that $\ell_t(a) = \ell_t^\top \varphi_a(x_t)$).

This is nothing but a linear bandit problem where the action set is changing over time. All the strategies we described are robust to this modification and thus in this case one can get a regret of $\sqrt{nT \log(|\mathcal{A}|)} \lesssim n\sqrt{T \log(T)}$ (and for the stochastic case one can get efficiently $n^{3/2}\sqrt{T}$).

A much more challenging case is when the correct embedding $\varphi = (\varphi_a)_{a \in \mathcal{A}}$ is only known to belong to some class $\Phi$. Without further assumptions on $\Phi$ we are basically back to the general model. Also note that a natural impulse is to run "bandits on top of bandits", that is first select some $\varphi_t \in \Phi$ and then select $a_t$ based on the assumption that $\varphi_t$ is correct. We won't get into this here, but let us investigate a related idea.

## Exp4, Auer, Cesa-Bianchi, Freund and Schapire [2001]

One can play exponential weights on the set of policies with the following unbiased estimator (obvious notation: $\ell_t(\pi) = \ell_t(\pi(x_t))$, $\pi_t \sim p_t$, and $a_t = \pi_t(x_t)$)

$$\widetilde{\ell}_t(\pi) = \frac{\mathbb{1}\{\pi(x_t) = a_t\}}{\sum_{\pi' : \pi'(x_t) = a_t} p_t(\pi')} \ell_t(a_t).$$

## Exp4, Auer, Cesa-Bianchi, Freund and Schapire [2001]

One can play exponential weights on the set of policies with the following unbiased estimator (obvious notation: $\ell_t(\pi) = \ell_t(\pi(x_t))$, $\pi_t \sim p_t$, and $a_t = \pi_t(x_t)$)

$$\widetilde{\ell}_t(\pi) = \frac{\mathbb{1}\{\pi(x_t) = a_t\}}{\sum_{\pi':\pi'(x_t)=a_t} p_t(\pi')} \ell_t(a_t).$$

Easy exercise: $R_T^{\mathrm{ctx}} \leq \sqrt{2T|\mathcal{A}|\log(|\Pi|)}$ (indeed the relative entropy term is smaller than $\log(|\Pi|)$ while the variance term is exactly $|\mathcal{A}|$).

## Exp4, Auer, Cesa-Bianchi, Freund and Schapire [2001]

One can play exponential weights on the set of policies with the following unbiased estimator (obvious notation: $\ell_t(\pi) = \ell_t(\pi(x_t))$, $\pi_t \sim p_t$, and $a_t = \pi_t(x_t)$)

$$\widetilde{\ell}_t(\pi) = \frac{\mathbb{1}\{\pi(x_t) = a_t\}}{\sum_{\pi' : \pi'(x_t) = a_t} p_t(\pi')} \ell_t(a_t).$$

Easy exercise: $R_T^{\mathrm{ctx}} \leq \sqrt{2T|\mathcal{A}| \log(|\Pi|)}$ (indeed the relative entropy term is smaller than $\log(|\Pi|)$ while the variance term is exactly $|\mathcal{A}|$).

The only issue of this strategy is that the computationally complexity is linear in the policy space, which might be huge. A year and half ago a major paper by Agarwal, Hsu, Kale, Langford, Li and Schapire was posted, with a strategy obtaining the same regret as Exp4 (in the i.i.d. model) but which is also computationally efficient with an oracle for the offline problem (i.e., $\min_{\pi \in \Pi} \sum_{t=1}^{T} \ell_t(\pi(x_t))$). Unfortunately the algorithm is not simple enough yet to be included in these slides.

# The statistician perspective, after Goldenshluger and Zeevi [2009, 2011], Perchet and Rigollet [2011]

Let $\mathcal{X} \subset \mathbb{R}^d$, $\mathcal{A} = [n]$, $(x_t)$ i.i.d. from some $\mu$ absolutely continuous w.r.t. Lebesgue. The reward for playing arm $a$ under context $x$ is drawn from some distribution $\nu_a(x)$ on $[0, 1]$ with mean function $f_a(x)$ which is assumed to be $\beta$-Holder smooth. Let $\Delta(x)$ be the "gap" function.

# The statistician perspective, after Goldenshluger and Zeevi [2009, 2011], Perchet and Rigollet [2011]

Let $\mathcal{X} \subset \mathbb{R}^d$, $\mathcal{A} = [n]$, $(x_t)$ i.i.d. from some $\mu$ absolutely continuous w.r.t. Lebesgue. The reward for playing arm $a$ under context $x$ is drawn from some distribution $\nu_a(x)$ on $[0,1]$ with mean function $f_a(x)$ which is assumed to be $\beta$-Holder smooth. Let $\Delta(x)$ be the "gap" function.

A key parameter is the proportion of contexts with a small gap. The margin assumption is that for some $\alpha > 0$, one has

$$\mu(\{x : \Delta(x) \in (0, \delta)\}) \leq C\delta^\alpha, \forall \delta \in (0, 1].$$

# The statistician perspective, after Goldenshluger and Zeevi [2009, 2011], Perchet and Rigollet [2011]

Let $\mathcal{X} \subset \mathbb{R}^d$, $\mathcal{A} = [n]$, $(x_t)$ i.i.d. from some $\mu$ absolutely continuous w.r.t. Lebesgue. The reward for playing arm $a$ under context $x$ is drawn from some distribution $\nu_a(x)$ on $[0, 1]$ with mean function $f_a(x)$ which is assumed to be $\beta$-Holder smooth. Let $\Delta(x)$ be the "gap" function.

A key parameter is the proportion of contexts with a small gap. The margin assumption is that for some $\alpha > 0$, one has

$$\mu(\{x : \Delta(x) \in (0, \delta)\}) \leq C\delta^\alpha, \forall \delta \in (0, 1].$$

One can achieve a regret of order $T \left( \frac{n \log(n)}{T} \right)^{\frac{\beta(\alpha+1)}{2\beta+d}}$, which is optimal at least in the dependency on $T$. It can be achieved by running Successive Elimination on an adaptively refined partition of the space, see Perchet and Rigollet [2011] for the details.

# The online multi-class classification perspective after Kakade, Shalev-Shwartz, and Tewari [2008]

Here the loss is assumed to be of the following very simple form: $\ell_t(a) = \mathbb{1}\{a \neq a_t^*\}$. In other words using the context $x_t$ one has to predict the best action (which can be interpreted as a *class*) $a_t^* \in [n]$.

# The online multi-class classification perspective after Kakade, Shalev-Shwartz, and Tewari [2008]

Here the loss is assumed to be of the following very simple form: $\ell_t(a) = \mathbb{1}\{a \neq a_t^*\}$. In other words using the context $x_t$ one has to predict the best action (which can be interpreted as a *class*) $a_t^* \in [n]$.

KSST08 introduces the *banditron*, a bandit version of the multi-class perceptron for this problem. While with full information the online multi-class perceptron can be shown to satisfy a "regret" bound on of order $\sqrt{T}$, the banditron attains only a regret of order $T^{2/3}$. See also Chapter 4 in Bubeck and Cesa-Bianchi [2012] for more on this.

# Summary of advanced results

1. The optimal regret for the linear bandit problem is $\widetilde{O}(n\sqrt{T})$. In the Bayesian context Thompson Sampling achieves this bound. In the i.i.d. case one can use an algorithm based on the optimism in face of uncertainty together with concentration properties of the least squares estimator.

2. The i.i.d. algorithm can easily be modified to be computationally efficient, or to deal with sparsity in the unknown vector $\theta$.

3. Extensions/variants: semi-bandit model, non-linear bandit (Lipschitz, Gaussian process, convex).

4. Contextual bandit is still a very active subfield of bandit theory.

5. Many important things were omitted. Example: knapsack bandit, see Badanidiyuru, Kleinberg and Slivkins [2013].

## Some open problems we discussed

1. Prove the lower bound $\mathbb{E}R_T = \Omega(\sqrt{Tn \log(n)})$ for the adversarial $n$-armed bandit with adaptive adversary.

2. Guha and Munagala [2014] conjecture: for product priors, TS is a 2-approximation to the optimal Bayesian strategy for the objective of minimizing the number of pulls on suboptimal arms.

3. Find a "simple" strategy achieving the Bubeck and Slivkins [2012] best of both worlds result.

4. For the combinatorial bandit problem, find a strategy with regret at most $n^{3/2}\sqrt{T}$ (current best is $n^2\sqrt{T}$).

5. Is there a computationally efficient strategy for i.i.d. linear bandit with optimal $n\sqrt{T}$ gap-free regret and with $\log(T)$ gap-based regret?

6. Is there a natural framework to think about "bandits on top of bandits" (while keeping $\sqrt{T}$-regret)?