Bobby Kleinberg Cornell University Joint work with Nicole Immorlica.



24 Oct 2017

Can you construct a dinner schedule that:

- never goes 2 days without macaroni and cheese
- never goes 3 days without pizza
- never goes 5 days without fish?

Answer: Impossible. For $N \ge 60$,

 $\lfloor N/2 \rfloor + \lfloor N/3 \rfloor + \lfloor N/5 \rfloor > N.$

Can you construct a dinner schedule that:

- never goes 2 days without macaroni and cheese
- never goes 4 days without pizza
- never goes 5 days without fish?

Answer: Possible.

Can you construct a dinner schedule that:

- never goes 2 days without macaroni and cheese
- never goes 3 days without pizza
- never goes 100 days without fish?

Answer: Impossible.



Can you construct a dinner schedule that:

- never goes 2 days without macaroni and cheese
- never goes 5 days without pizza
- never goes 100 days without fish
- never goes 7 days without tacos?

Answer: Impossible.

••••

Can you construct a dinner schedule that:

- never goes 2 days without macaroni and cheese
- never goes 5 days without pizza
- never goes 100 days without fish
- never goes 7 days without tacos?

Answer: Impossible.

Can you construct a dinner schedule that:

- never goes 2 days without macaroni and cheese
- never goes 5 days without pizza
- never goes 100 days without fish
- never goes 7 days without tacos?

Answer: Impossible.

••••

The Pinwheel Problem

Given g_1, \ldots, g_n , can \mathbb{Z} be partitioned into S_1, \ldots, S_n such that S_i intersects every interval of length g_i ?

E.g.,
$$(g_1, \ldots, g_5) = (3, 4, 6, 10, 16)$$



The Pinwheel Problem

Given g_1, \ldots, g_n , can \mathbb{Z} be partitioned into S_1, \ldots, S_n such that S_i intersects every interval of length g_i ?

What is the complexity of this decision problem?



The Pinwheel Problem

Given g_1, \ldots, g_n , can \mathbb{Z} be partitioned into S_1, \ldots, S_n such that S_i intersects every interval of length g_i ?

What is the complexity of this decision problem?

It belongs to PSPACE.

No non-trivial lower bounds known.

Later in this talk: PTAS for an optimization version.



The Multi-Armed Bandit Problem



Stochastic Multi-Armed Bandit Problem: A decision-maker ("gambler") chooses one of n actions ("arms") in each time step. Chosen arm yields random payoff from unknown distrib. on [0, 1]. Goal: Maximize expected total payoff.





















Recharging Bandits

- Pulling arm *i* at time *t*, when it was last pulled at time *s*, yields random payoff with expectation $H_i(t s)$.
- H_i is an increasing, concave function; $H_i(t) \leq t$.

Recharging Bandits

- Pulling arm *i* at time *t*, when it was last pulled at time *s*, yields random payoff with expectation $H_i(t s)$.
- H_i is an increasing, concave function; $H_i(t) \leq t$.

Concavity assumption implies free disposal: in step t, pulling i is better than doing nothing because

$$H_i(u-t)+H_i(t-s)\geq H_i(u-s).$$

Recharging Bandits

- Pulling arm *i* at time *t*, when it was last pulled at time *s*, yields random payoff with expectation $H_i(t s)$.
- H_i is an increasing, concave function; $H_i(t) \leq t$.

With known $\{H_i\}$: a special case of deterministic restless bandits. General case is PSPACE-hard [Papadimitriou & Tsitsiklis 1987]. Which reinforcement learning problems have a PTAS?

Recharging Bandits

- Pulling arm *i* at time *t*, when it was last pulled at time *s*, yields random payoff with expectation $H_i(t s)$.
- H_i is an increasing, concave function; $H_i(t) \leq t$.

Plan of attack:

- Analyze optimal play when $\{H_i\}$ are known.
- Use upper confidence bounds + "ironing" to reduce the case when {*H_i*} must be learned to the case when they are known.

Recharging Bandits

- Pulling arm *i* at time *t*, when it was last pulled at time *s*, yields random payoff with expectation $H_i(t s)$.
- H_i is an increasing, concave function; $H_i(t) \leq t$.

Plan of attack:

- Analyze optimal play when $\{H_i\}$ are known.
- Use upper confidence bounds + "ironing" to reduce the case when {*H_i*} must be learned to the case when they are known.

Greedy algorithm: always maximize payoff in current time step.

Greedy/OPT ratio can be arbitrarily close to $1/2\,$

•
$$H_1(t) = 1 - \varepsilon$$
, $H_2(t) = t$.

- Greedy always pulls arm 2.
- "Almost-OPT" pulls arm 1 for $T \gg 1$ time steps, then arm 2.
- Net payoff $(2 \varepsilon)T + 1$ over T + 1 time steps.

Greedy algorithm: always maximize payoff in current time step.

Greedy/OPT is never less than 1/2

- Imagine allowing the algorithm (but not OPT) to pull two arms per time step.
- At each time, supplement the greedy selection with the arm selected by OPT, if they differ.
- This at most doubles the payoff in each time step.
- Net payoff of supplemented schedule \geq OPT. (free disposal property)

For $0 \le x \le 1$, let $R_i(x)$ denote maximum long-run average payoff achievable by playing *i* in at most *x* fraction of time steps.

$$R_i(x) = \sup \left\{ \left. \frac{1}{T} \sum_{j=1}^{\ell} H_i(t_j - t_{j-1}) \right| \begin{array}{l} T < \infty, \ \ell \le x \cdot T, \\ 0 = t_0 < t_1 < \cdots < t_{\ell} \le T \end{array} \right\}.$$

For $0 \le x \le 1$, let $R_i(x)$ denote maximum long-run average payoff achievable by playing *i* in at most *x* fraction of time steps.

$$R_i(x) = \sup \left\{ \left. \frac{1}{T} \sum_{j=1}^{\ell} H_i(t_j - t_{j-1}) \right| \begin{array}{l} T < \infty, \ \ell \le x \cdot T, \\ 0 = t_0 < t_1 < \cdots < t_{\ell} \le T \end{array} \right\}.$$

Fact: R_i is piecewise-linear with breakpoints $R_i(\frac{1}{k}) = \frac{1}{k}H_i(k)$.

Rate of Return Function

For $0 \le x \le 1$, let $R_i(x)$ denote maximum long-run average payoff achievable by playing *i* in at most *x* fraction of time steps.

$$R_i(x) = \sup \left\{ \left. \frac{1}{T} \sum_{j=1}^{\ell} H_i(t_j - t_{j-1}) \right| \left| \begin{array}{c} T < \infty, \ \ell \le x \cdot T, \\ 0 = t_0 < t_1 < \cdots < t_{\ell} \le T \end{array} \right\}.$$

Fact: R_i is piecewise-linear with breakpoints $R_i(\frac{1}{k}) = \frac{1}{k}H_i(k)$.



Rate of Return Function

For $0 \le x \le 1$, let $R_i(x)$ denote maximum long-run average payoff achievable by playing *i* in at most *x* fraction of time steps.

$$\mathcal{R}_i(x) = \sup \left\{ \left. rac{1}{\mathcal{T}} \sum_{j=1}^\ell \mathcal{H}_i(t_j - t_{j-1})
ight| egin{array}{c} \mathcal{T} < \infty, \ \ell \leq x \cdot \mathcal{T}, \ 0 = t_0 < t_1 < \cdots < t_\ell \leq \mathcal{T} \end{array}
ight\}.$$

Fact: R_i is piecewise-linear with breakpoints $R_i(\frac{1}{k}) = \frac{1}{k}H_i(k)$. **Proof sketch:** The optimal sequence $0 = t_0 < \cdots < t_{\ell} \le T$ has

at most two distinct gap sizes, $\lfloor \frac{1}{x} \rfloor$ and $\lceil \frac{1}{x} \rceil$.

Concave Relaxation

The problem

$$\max\left\{\sum_{i=1}^n R_i(x_i) \mid \sum_i x_i \leq 1, \forall i \, x_i \geq 0\right\}$$

specifies an upper bound on the value of the optimal schedule.



Concave Relaxation

The problem

$$\max\left\{\left|\sum_{i=1}^n R_i(x_i)\right| \left|\sum_i x_i \le 1, \forall i \, x_i \ge 0\right\}\right\}$$

specifies an upper bound on the value of the optimal schedule.



Mapping (x_1, \ldots, x_n) to a schedule: pinwheel problem!

Independent Rounding

First idea: every time step, sample arm i with probability x_i .

Then τ_i = delay of arm $i = t_j(i) - t_{j-1}(i)$ is geometrically distributed with expectation $1/x_i$.

Rounding scheme gets $x_i \cdot \mathbb{E}H_i(\tau_i)$ whereas relaxation gets $R_i(x_i) = x_i H_i(1/x_i) = x_i \cdot H_i(\mathbb{E}\tau_i)$.

Fact: if *H* is concave and non-decreasing and *Y* is geometrically distributed then $\mathbb{E}H(Y) \ge (1 - \frac{1}{e}) H(\mathbb{E}Y)$.

Independent Rounding

First idea: every time step, sample arm i with probability x_i .

Then τ_i = delay of arm $i = t_j(i) - t_{j-1}(i)$ is geometrically distributed with expectation $1/x_i$.

Rounding scheme gets $x_i \cdot \mathbb{E}H_i(\tau_i)$ whereas relaxation gets $R_i(x_i) = x_i H_i(1/x_i) = x_i \cdot H_i(\mathbb{E}\tau_i)$.

Fact: if *H* is concave and non-decreasing and *Y* is geometrically distributed then $\mathbb{E}H(Y) \ge (1 - \frac{1}{e}) H(\mathbb{E}Y)$.

To do better, need rounding scheme that reduces variance of τ_i .

Interleaved Arithmetic Progressions

Second idea: round continuous-time schedule to discrete time. In continuous time, pull *i* at $\{\frac{r_i+k}{x_i} \mid k \in \mathbb{N}\}$ where $r_i \sim_{\text{Unif}} [0, 1)$. Map this schedule to discrete time in an order-preserving manner.



Interleaved Arithmetic Progressions

Second idea: round continuous-time schedule to discrete time. In continuous time, pull *i* at $\{\frac{r_i+k}{x_i} \mid k \in \mathbb{N}\}$ where $r_i \sim_{\text{Unif}} [0, 1)$. Map this schedule to discrete time in an order-preserving manner.



Between two pulls of *i*, we pull *j* either $\lfloor x_j/x_i \rfloor$ or $\lceil x_j/x_i \rceil$ times. $\tau_i = 1 + \sum_{i \neq i} Z_i$

 $\{Z_j\}$ are independent, each supported on 2 consecutive integers.

Convex Stochastic Ordering

Definition

If X, Y are random variables, the *convex stochastic ordering* defines $X \leq_{cx} Y$ if and only if $\mathbb{E}\phi(X) \leq \mathbb{E}\phi(Y)$ for every convex function ϕ .

Convex Stochastic Ordering

Definition

If X, Y are random variables, the *convex stochastic ordering* defines $X \leq_{cx} Y$ if and only if $\mathbb{E}\phi(X) \leq \mathbb{E}\phi(Y)$ for every convex function ϕ .

Lemma

If X is a sum of independent Bernoulli random variables and Y is Poisson with $\mathbb{E}Y = \mathbb{E}X$ then $X \leq_{cx} Y$.

Convex Stochastic Ordering

Definition

If X, Y are random variables, the *convex stochastic ordering* defines $X \leq_{cx} Y$ if and only if $\mathbb{E}\phi(X) \leq \mathbb{E}\phi(Y)$ for every convex function ϕ .

Lemma

If X is a sum of independent Bernoulli random variables and Y is Poisson with $\mathbb{E}Y = \mathbb{E}X$ then $X \leq_{cx} Y$.

$$\begin{array}{rcl} \tau_i &=& 1 + \sum_{j \neq i} Z_j &\leq_{\mathrm{cx}} & 1 + \mathsf{Pois}(\frac{1}{x_i} - 1) \\ x_i \cdot \mathbb{E}H_i(\tau_i) &\geq& x_i \cdot \mathbb{E}H_i(1 + \mathsf{Pois}(\frac{1}{x_i} - 1)) \end{array}$$

Approximation Ratio for Interleaved AP Rounding

Fact 1: If *H* is concave and non-decreasing and *Y* is Poisson, then $\mathbb{E}H(1+Y) \ge (1-\frac{1}{2e})H(1+\mathbb{E}Y)$

Fact 2: If *H* is concave and non-decreasing and *Y* is Poisson with $\mathbb{E}Y \ge m$, then

$$\mathbb{E} H(1+Y) \geq \left(1-rac{1}{\sqrt{2\pi m}}
ight) H(1+\mathbb{E} Y)$$

Conclusion: Interleaved AP rounding is

- a $1 \frac{1}{2e} \approx 0.816$ approximation in general
- a 1δ approximation for "small arms" to whom the concave relaxation assigns $x_i < \delta^2$

PTAS for Recharging Bandits

Let $\varepsilon > 0$ be a small constant. Two easy cases . . .

All arms are big. Every arm that gets pulled in the optimal schedule is pulled with frequency ε² or greater.

Then the optimal schedule uses only $1/\varepsilon^2$ arms. Brute-force search takes polynomial time.

All arms are small. If the optimal concave program solution has x_i < ε² for all *i*, then randomly interleaved arithmetic progressions get 1 - ε approximation.

Combine the cases using "partial enumeration". For $p = O_{\varepsilon}(1) \dots$ Outer loop: iterate over *p*-periodic schedules of arms and gaps. Inner loop: fit small arms into gaps using interleaved AP rounding.

- Gaps in the *p*-periodic schedule may not be equally spaced.
 - For each small arm choose just one congruence class (mod *p*) of "eligible gaps."
 - Bin-pack small arms into congruence classes.

- Gaps in the *p*-periodic schedule may not be equally spaced.
 - For each small arm choose just one congruence class (mod *p*) of "eligible gaps."
 - Bin-pack small arms into congruence classes.
 - Works if $x_i < \varepsilon^2/p$ for small arms while $x_i \ge 1/p$ for big arms.

- Gaps in the *p*-periodic schedule may not be equally spaced.
 - For each small arm choose just one congruence class (mod *p*) of "eligible gaps."
 - Bin-pack small arms into congruence classes.
 - Works if $x_i < \varepsilon^2/p$ for small arms while $x_i \ge 1/p$ for big arms.
 - Eliminate intermediate arms by finding k ≤ 1/ε such that arms with x_i ∈ (ε^{4(k+1)}, ε^{4k}] contribute less than ε · OPT.
 - Conclusion: # of big arms $\leq (1/\varepsilon)^{O(1/\varepsilon)}$.

- Gaps in the *p*-periodic schedule may not be equally spaced.
- Why can we assume big arms are scheduled with period $p = O_{\varepsilon}(1)$?
 - We need existence of a *p*-periodic schedule that matches two properties of OPT
 - rate of return from big arms
 - 2 amount of time left over for small arms
 - Existence proof is surprisingly technical; omitted.
 - Conclusion $p = (\#big)/\varepsilon^2$ suffices.

- Gaps in the *p*-periodic schedule may not be equally spaced.
- Why can we assume big arms are scheduled with period $p = O_{\varepsilon}(1)$?
 - We need existence of a *p*-periodic schedule that matches two properties of OPT
 - rate of return from big arms
 - 2 amount of time left over for small arms
 - Existence proof is surprisingly technical; omitted.
 - Conclusion $p = (\#big)/\varepsilon^2$ suffices.
- Grand conclusion: PTAS with running time $n^{(1/\varepsilon)^{(24/\varepsilon)}}$

- Gaps in the *p*-periodic schedule may not be equally spaced.
- Why can we assume big arms are scheduled with period $p = O_{\varepsilon}(1)$?
 - We need existence of a *p*-periodic schedule that matches two properties of OPT
 - rate of return from big arms
 - 2 amount of time left over for small arms
 - Existence proof is surprisingly technical; omitted.
 - Conclusion $p = (\#big)/\varepsilon^2$ suffices.
- Grand conclusion: PTAS with running time $n^{(1/\varepsilon)^{(24/\varepsilon)}}$
- Remark: the 0.816-approximation runs in time $O(n^2 \log n)$.

Now suppose $\{H_i\}$ are not known, must be learned by sampling.

Now suppose $\{H_i\}$ are not known, must be learned by sampling. Idea: divide time into "planning epochs" of length $\phi = O(n/\epsilon)$. In each epoch . . .

- Compute $\overline{H}_i(x)$, an upper confidence bound on $H_i(x)$, $\forall i$.
- **2** Run approx alg. on $\{\overline{H}_i\}$ to schedule arms within epoch.
- **③** Update empirical estimates and confidence radii.

Main challenge: Although H_i is concave, \overline{H}_i may not be.

 R_i

Now suppose $\{H_i\}$ are not known, must be learned by sampling. Idea: divide time into "planning epochs" of length $\phi = O(n/\epsilon)$. In each epoch . . .

- Compute $\overline{H}_i(x)$, an upper confidence bound on $H_i(x)$, $\forall i$.
- **2** Run approx alg. on $\{\overline{H}_i\}$ to schedule arms within epoch.
- **③** Update empirical estimates and confidence radii.

Main challenge: Although H_i is concave, \overline{H}_i may not be.

Solution: Work with \overline{R}_i and "iron" the non-concavity, without disrupting the approximation guarantee.

Now suppose $\{H_i\}$ are not known, must be learned by sampling. Idea: divide time into "planning epochs" of length $\phi = O(n/\epsilon)$. In each epoch . . .

- Compute $\overline{H}_i(x)$, an upper confidence bound on $H_i(x)$, $\forall i$.
- **2** Run approx alg. on $\{\overline{H}_i\}$ to schedule arms within epoch.
- **③** Update empirical estimates and confidence radii.

Approx alg. is almost black box.

Can plug in greedy, interleaved AP rounding, or PTAS.

Approx. factor reduced by $1 - \varepsilon$, plus $O(n \log(n) \sqrt{T \log(nT)})$ regret.

 R_i

Now suppose $\{H_i\}$ are not known, must be learned by sampling. Idea: divide time into "planning epochs" of length $\phi = O(n/\epsilon)$. In each epoch . . .

- Compute $\overline{H}_i(x)$, an upper confidence bound on $H_i(x)$, $\forall i$.
- **2** Run approx alg. on $\{\overline{H}_i\}$ to schedule arms within epoch.
- **③** Update empirical estimates and confidence radii.

Approx alg. is almost black box.

Can plug in greedy, interleaved AP rounding, or PTAS.

Approx. factor reduced by $1 - \varepsilon$, plus $O(n \log(n) \sqrt{T \log(nT)})$ regret.

 R_i

Recharging bandits: A model for learning to schedule recurring tasks (interventions) whose benefit increases with latency.

Approximation algorithms:

- simple greedy $(\frac{1}{2})$;
- rounding concave relaxation using interleaved arithmetic progressions $(1 \frac{1}{2e})$;
- partial enumeration and concave rounding (1ε) .

Nice connections to pinwheel problem in additive combinatorics.

- Pinwheel problem
 - Complexity? (Could be in P. Could be PSPACE-complete.)
 - Is (g_1, \ldots, g_n) always feasible if $\sum_i g_i^{-1} \leq 5/6$?
 - Is $(g_1 + 1, \ldots, g_n + 1)$ always feasible if $\sum_i g_i^{-1} \le 1$?

- Pinwheel problem
 - Complexity? (Could be in P. Could be PSPACE-complete.)
 - 2 Is (g_1, \ldots, g_n) always feasible if $\sum_i g_i^{-1} \leq 5/6$?
 - Is $(g_1 + 1, \dots, g_n + 1)$ always feasible if $\sum_i g_i^{-1} \le 1$? Best result in this direction: increase $g_i + 1$ to $g_i + g_i^{1/2+o(1)}$. [Immorlica-K. 2017]

- Pinwheel problem
 - Complexity? (Could be in P. Could be PSPACE-complete.)
 - **2** Is (g_1, \ldots, g_n) always feasible if $\sum_i g_i^{-1} \le 5/6$?
 - Is (g₁ + 1,..., g_n + 1) always feasible if ∑_i g_i⁻¹ ≤ 1? Best result in this direction: increase g_i + 1 to g_i + g_i^{1/2+o(1)}. [Immorlica-K. 2017]
- Reinforcement learning: What other special cases admit PTAS?

- Applications: extend recharging bandits model to incorporate domain-specific features such as ...
 - (fighting poachers) Strategic arms with endogenous payoffs. [Kempe-Schulman-Tamuz '17]
 - (invasive species removal) Externalities between arms. Movement costs.
 - (education) Payoffs with more complex history-dependency. [Novikoff-Kleinberg-Strogatz '11]