

From Bandits to Monte Carlo Tree Search: The optimistic principle applied to Optimization and Planning

Rémi Munos

Sequel project: Sequential Learning
<http://researchers.lille.inria.fr/~munos/>

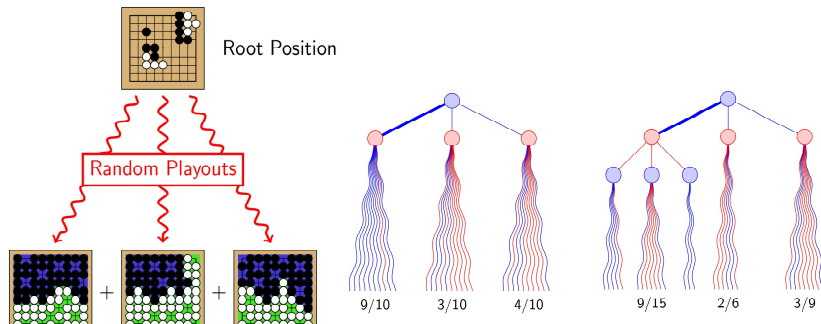
INRIA Lille - Nord Europe / Microsoft-Research NE

AAAI 2013, Bellevue, WA

Jean-Yves Audibert, Sébastien Bubeck, Lucian Buşoniu, Alexandra Carpentier, Pierre-Arnaud Coquelin, Rémi Coulom, Sylvain Gelly, Jean-François Hren, Nathaniel Korda, Odalric-Ambrym Maillard, Gilles Stoltz, Csaba Szepesvári, Olivier Teytaud, Michal Valko, and Yizao Wang.

Initial motivation

Monte-Carlo Tree Search in computer-go



Idea: use bandits at each node of the tree search.

Monte-Carlo Tree Search

- Very efficient in several problems
- Very inefficient in many other problems (even toy problems)
- Not much theoretical guarantee...

We would like

- Understand how the “optimism in the face of uncertainty” principle works in hierarchical problems
- Define classes of problems for which variants of MCTS would be efficient

Outline of this tutorial

- 1 The stochastic multi-armed bandit
 - The K -armed bandit and UCB
 - The many-armed bandit
- 2 Monte Carlo Tree Search
 - Bandits in a hierarchy
 - Computer go and UCT
- 3 Optimistic optimization with known smoothness
 - Deterministic rewards
 - Stochastic rewards (\mathcal{X} -armed bandit)
- 4 Extension to unknown smoothness
- 5 Optimistic planning
 - Deterministic dynamics,
 - Open Loop planning
 - MDPs with a model

100

- applies in a large class of decision making problems in stochastic and deterministic environments
- provides an efficient exploration of the search space by exploring the most promising areas first
- provides a natural transition from global to local search
- Performance depends on the “smoothness” of the function around the maximum w.r.t. some metric,
 - a measure of the quantity of near-optimal solutions,
 - and our knowledge or not of this metric.

The multi-armed bandit problem

Setting:

- Set of K arms (possible actions)
- At each time t , choose an arm $I_t \in \{1, \dots, K\}$ and receive a reward X_t coming from arm I_t .
- **Goal:** find an arm selection policy such as to maximize a function of the rewards.



Exploration-exploitation tradeoff:

- **Exploit:** act optimally according to our current beliefs
- **Explore:** learn more about the environment

Numerous variants

Different settings:

- **Stochastic environments:** the rewards are samples from probability distributions. We compare our strategy to the optimal oracle one.
- **Adversarial environments:** the rewards are chosen by an adversary. We compare our strategy to a class of possible strategies.
- **Action space** can be finite, countably infinite, continuous (function optimization), combinatorial (paths), structured (policies), ...

Different targets:

- maximizing cumulative rewards, returning the best possible solution, estimating the values of all the arms, ...

11

- Clinical trials [Thompson 1933]
- Ads placement on webpages
- Nash equilibria (traffic or communication networks, agent simulation, poker, ...)
- Packet routing, itinerary selection, ...
- Game-playing computers (Go, urban rivals, ...)
- Optimization / planning given a finite numerical budget

Setting:

- Set of K arms, defined by distributions ν_k (with support in $[0, 1]$), whose law is unknown,
- At each time t , choose an arm $I_t \in \{1, \dots, K\}$ and receive reward $X_t \stackrel{i.i.d.}{\sim} \nu_{I_t}$.
- **Goal:** maximize the sum of rewards.

Definitions:

- Let $\mu_k = \mathbb{E}_{X \sim \nu_k}[X]$ be the mean value of arm k ,
- Let $\mu^* = \max_k \mu_k$ the best mean value,

Define the cumulative **regret**:

$$R_n \stackrel{\text{def}}{=} \sum_{t=1}^n (\mu^* - X_t).$$

The cumulative regret

The expected cumulative regret is

$$\begin{aligned}\mathbb{E}R_n &= \mathbb{E}\left[\sum_{t=1}^n \mu^* - X_t\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\sum_{t=1}^n \mu^* - X_t \mid I_t\right]\right] = \mathbb{E}\left[\sum_{t=1}^n \mu^* - \mu_{I_t}\right] \\ &= \mathbb{E}\left[\sum_{k=1}^K (\mu^* - \mu_k) \sum_{t=1}^n \mathbf{1}\{I_t = k\}\right] = \sum_{k=1}^K \Delta_k \mathbb{E}[T_k(n)],\end{aligned}$$

where $\Delta_k \stackrel{\text{def}}{=} \mu^* - \mu_k$, and $T_k(n)$ is the number of times arm k has been pulled up to round n .

Goal: Find an arm selection policy such as to minimize $\mathbb{E}R_n$.

This is an old problem! [Robbins, 1952] Surprisingly, not fully solved yet!

Many proposed strategies:

- **ϵ -greedy exploration:** choose current best action with proba $1 - \epsilon$, or random action with proba ϵ ,
- **Bayesian exploration:** assign prior to the arm distributions and select arm according to the posterior distributions (Gittins index, Thompson sampling, ...)
- **Softmax exploration:** choose arm k with proba $\propto \exp(\beta \hat{\mu}_k)$
- **Follow the perturbed leader:** choose best perturbed arm
- **Optimistic exploration:** select best arm in the most favorable environment compatible with observations

- At time t , for each arm k , use past observations and some probabilistic argument to define high-probability confidence intervals containing the expected reward μ_k
- The most favorable environment for arm k is thus the upper confidence bound (UCB) on μ_k
- Simple implementation: play the arm having the largest UCB!

Intuition of the UCB algorithm

Idea:

- The B-values $B_{k,t}$ are h.p. UCBs on μ_k . Indeed we have:

$$\mathbb{P}\left(|\hat{\mu}_{k,t} - \mu_k| \geq \sqrt{\frac{2 \log(t)}{T_k(t)}}\right) \leq \frac{2}{t^2},$$

using a union bound for all possible values of $T_k(t) \in \{1, \dots, t\}$ together with Chernoff-Hoeffding's inequality:

$$\mathbb{P}\left(\left|\frac{1}{m}\sum_{i=1}^m Y_i - \mu\right| \geq \epsilon\right) \leq 2e^{-2m\epsilon^2}$$

1. *Journal of Management Studies*, 1996, 33, 1, 1-14.

1. *Journal of the American Medical Association*, 1997; 277: 1039-1043.

Could we stay a long time playing a wrong arm?

No, since

- The more we pull an arm k , the smaller the size of the confidence interval and the closer its UCB gets to its mean value μ_k
- But in h.p., it cannot be pulled once its UCB becomes smaller than μ^*

So each sub-optimal arm k can be only pulled a number of times $T_k(n)$ such that the size of its confidence interval $\sqrt{\frac{2 \log n}{T_k(n)}}$ is of order $\Delta_k = \mu^* - \mu_k$.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Let k be a sub-optimal arm, and k^* be an optimal arm. At time t , if arm k is selected, this means that

$$\begin{aligned} B_{k,t} &\geq B_{k^*,t} \\ \hat{\mu}_{k,t} + \sqrt{\frac{2 \log(t)}{T_k(t)}} &\geq \hat{\mu}_{k^*,t} + \sqrt{\frac{2 \log(t)}{T_{k^*}(t)}} \\ \mu_k + 2\sqrt{\frac{2 \log(t)}{T_k(t)}} &\geq \mu^*, \text{ with high probability} \\ T_k(t) &\leq \frac{8 \log(t)}{\Delta_k^2} \end{aligned}$$

Thus, if $T_k(t) > \frac{8 \log(t)}{\Delta_k^2}$, then there is only a small probability that arm k can be selected.

$$\begin{aligned} T_k(n) &\leq u + \sum_{t=u+1}^n \mathbf{1}\{I_t = k; T_k(t) > u\} \\ &\leq u + \sum_{t=u+1}^n \left[\mathbf{1}\{\hat{\mu}_{k,t} - \mu_k \geq \sqrt{\frac{2 \log t}{T_k(t)}}\} + \mathbf{1}\{\hat{\mu}_{k^*,t} - \mu_k \leq -\sqrt{\frac{2 \log t}{T_{k^*}(t)}}\} \right] \end{aligned}$$

Now, taking the expectation of both sides,

$$\begin{aligned}\mathbb{E}[T_k(n)] &\leq u + \sum_{t=u+1}^n 2t^{-2} \\ &\leq \frac{8 \log(n)}{\Delta_k^2} + 1 + \frac{\pi^2}{3}\end{aligned}$$

Better confidence bounds imply smaller regret

- Chernoff-Hoeffding's inequality $1/t$ -confidence bound:

$$\mathbb{E}X \leq \frac{1}{m} \sum_{i=1}^m X_i + \sqrt{\frac{\log t}{2m}}$$

- Bernstein's inequality:

$$\mathbb{E}X \leq \frac{1}{m} \sum_{i=1}^m X_i + \sqrt{\frac{2\sigma^2 \log t}{m}} + \frac{\log t}{3m}$$

- Empirical Bernstein's inequality:

$$\mathbb{E}X \leq \frac{1}{m} \sum_{i=1}^m X_i + \sqrt{\frac{2\hat{\sigma}_t^2 \log(3t)}{m}} + \frac{8 \log(3t)}{3m}$$

[Audibert, Munos, Szepesvári, 2007], [Maurer, Pontil, 2009]

UCB-V

[Audibert, Munos, Szepesvári, 2007]

- UCB using empirical variance estimate:

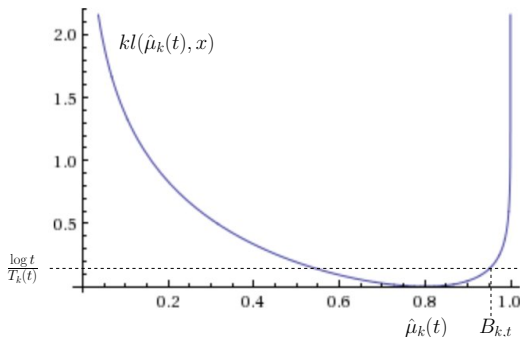
$$B_{k,t} \stackrel{\text{def}}{=} \hat{\mu}_{k,t} + \sqrt{2 \frac{\widehat{\sigma}_{k,t}^2 \log(1.2t)}{T_k(t)}} + \frac{3 \log(1.2t)}{T_k(t)}.$$

Then the expected regret is bounded as:

$$\mathbb{E}R_n \leq 10 \left(\sum_{k: \Delta_k > 0} \frac{\sigma_k^2}{\Delta_k} + 2 \right) \log(n).$$

For Bernoulli distributions, define the kl-UCB

$$B_{k,t} \stackrel{\text{def}}{=} \sup \left\{ x \in [0, 1], \text{kl}(\hat{\mu}_k(t), x) \leq \frac{\log t}{T_k(t)} \right\}$$



(non-asymptotic version of Sanov's theorem)

$$\mathbb{E}R_n = \sum_{k:\Delta_k>0} \frac{\Delta_k}{\text{kl}(\nu_k, \nu^*)} \log n + o(\log n).$$

See also DMED [Honda, Takemura, 2010, 2011] and other related algorithms.

Idea: Use the full empirical distribution to get a refined UCB.

Lower bounds

For single-dimensional distributions [Lai, Robbins, 1985]:

$$\liminf_{n \rightarrow \infty} \frac{\mathbb{E} R_n}{\log n} \geq \sum_{k: \Delta_k > 0} \frac{\Delta_k}{KL(\nu_k, \nu^*)}$$

For larger class of distributions \mathcal{D} [Burnetas, Katehakis, 1996]:

$$\liminf_{n \rightarrow \infty} \frac{\mathbb{E} R_n}{\log n} \geq \sum_{k: \Delta_k > 0} \frac{\Delta_k}{\mathcal{K}_{\text{inf}}(\nu_k, \mu^*)},$$

where

$$\mathcal{K}_{\text{inf}}(\nu, \mu) \stackrel{\text{def}}{=} \inf \left\{ KL(\nu, \nu') : \nu' \in \mathcal{D} \text{ and } \mathbb{E}_{X \sim \nu'}[X] > \mu \right\}.$$

- The smaller the gaps, the harder the problem (for large n)
- For small n , the regret is trivially bounded by $n \max_k \Delta_k$

For small gaps it takes a long time to distinguish which arm is the best.

Proposition 2.

Distribution-independent bounds:

$$\sup_{\text{Distributions}} \mathbb{E} R_n \leq 2\sqrt{2Kn[\log n + 1 + \frac{\pi^2}{3}]}$$

$$\begin{aligned} \mathbb{E}R_n &= \sum_k \Delta_k \mathbb{E}T_k(n) \\ &= \sum_k \Delta_k \sqrt{\mathbb{E}T_k(n)} \sqrt{\mathbb{E}T_k(n)} \\ &\leq \sqrt{\sum_k \Delta_k^2 \mathbb{E}T_k(n)} \sqrt{\sum_k \mathbb{E}T_k(n)} \\ &\leq \sqrt{8Kn \left[\log n + 1 + \frac{\pi^2}{3} \right]} \end{aligned}$$

since $\mathbb{E} T_k(n) \leq 8 \frac{\log n}{\Delta_k^2} + 1 + \frac{\pi^2}{3}$ and $\sum_k T_k(n) = n$.

1. *Journal of Management Studies*, 1997, 34, 1, 1-14.

$$\inf_{\text{Algo}} \sup_{\text{Distributions}} \mathbb{E} R_n = \Omega(\sqrt{Kn})$$

Notice that a refined algorithm (MOSS [Audibert, Bubeck, 2009]) achieves the same order:

$$\sup_{Distributions} \mathbb{E} R_n = O\left(\sqrt{Kn}\right).$$

Unstructured set of actions: The rewards received so far do not tell us anything about the value of unobserved arms.

- among the ones where you have already been
 - because it is good (Exploitation)
 - or not well known (Exploration)
- or choose a new one randomly (Discovery)

Other examples: Marketing (ex: sending catalogues), mining for valuable resources, ...

Many-armed bandits: Assumptions

We make a (probabilistic) assumption about the mean-value of any new arm.

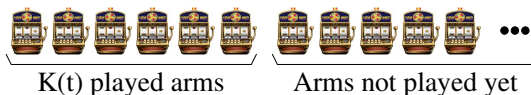
- **Usual assumption:** the distribution of the mean-reward of a new arm is known [Banks, Sundaram, 1992], [Berry, Chen, Zame, Heath, Shepp, 1997].
- **Much weaker assumption:** Assume we know $\beta > 0$ such that

$$\mathbb{P}(\mu(\text{new arm}) > \mu^* - \epsilon) = \Theta(\epsilon^\beta),$$

β characterizes the probability of selecting near-optimal arms

Large $\beta \implies$ small chance of pulling good arm, thus one needs to pull many arms. And vice-versa.

UCB with Arm Increasing Rule [Wang, Audibert, Munos, 2008]



UCB-AIR:

- $K(0) = 0$. At time $t + 1$, pull a new arm if

$$K(t) < \begin{cases} t^{\frac{\beta}{2}} & \text{if } \beta < 1 \text{ and } \mu^* < 1 \\ t^{\frac{\beta}{\beta+1}} & \text{if } \beta \geq 1 \text{ or } \mu^* = 1 \end{cases}$$

- Otherwise, apply UCB-V on the $K(t)$ current arms, i.e., play

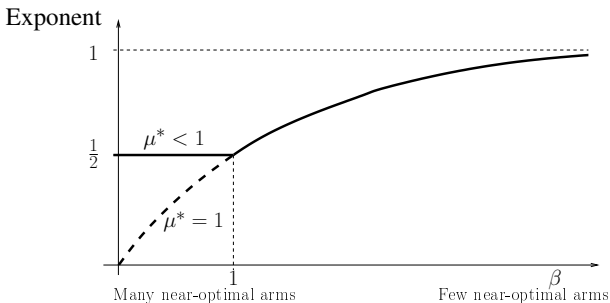
$$\operatorname{argmax}_{1 \leq k \leq K(t)} \underbrace{\hat{\mu}_{k,t}}_{\text{empirical rewards}} + \underbrace{\sqrt{\frac{2\hat{V}_{k,t}\mathcal{E}_t}{T_k(t)} + \frac{3\mathcal{E}_t}{T_k(t)}}_{\text{Confidence interval}},$$

with exploration sequence: $c \log(\log t) \leq \mathcal{E}_t \leq \log t$.

Regret analysis of UCB-AIR

Upper bound on the regret of UCB-AIR:

$$\mathbb{E}R_n = \begin{cases} \tilde{O}(\sqrt{n}) & \text{if } \beta < 1 \text{ and } \mu^* < 1 \\ \tilde{O}(n^{\frac{\beta}{1+\beta}}) & \text{if } \mu^* = 1 \text{ or } \beta \geq 1 \end{cases}$$



Lower bound: $\forall \beta > 0, \mu^* \leq 1$, for any algorithm $\mathbb{E}R_n = \Omega(n^{\frac{\beta}{1+\beta}})$.

- When $\beta > 1$ or $\mu^* = 1$ the upper and lower bounds match (up to logarithmic factor).
- Exploration-Exploitation-Discovery tradeoff:
 - **Exploitation**: Pull a good arm
 - **Exploration**: Pull an uncertain arm
 - **Discovery**: Pull a new arm
- The exploration sequence \mathcal{E}_t can be of order $\log \log t$ (instead of $\log t$): discovery replaces exploration
- **Open question**: similar performance when β is unknown? (i.e. adaptive strategy that estimates β while minimizing regret).

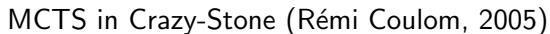
Bandits with a structured set of actions

Optimism in the face of uncertainty extends to:

- **Linear bandits** [Auer, 2002], [Dani, Hayes, Kakade, 2008], [Abbasi-Yadkori, 2009], [Rusmevichientong, Tsitsiklis, 2010], [Filippi, Cappé, Garivier, Szepesvári, 2010]
- **Convex bandits** [Zinkevich, 2003], [Flaxman, Kalai, McMahan, 2005], [Hazan, Agarwal, Kale, 2006], [Bartlett, Hazan, Rakhlin, 2007], [Shalev-Shwartz, 2007], [Abernethy, Bartlett, Rakhlin, Tewari, 2008], [Narayanan, Rakhlin, 2010]
- **Lipschitz bandits** [Agrawal, 1995], [Kleinberg, 2004], [Auer, Ortner, Szepesvári, 2007], [Kleinberg, Slivkins, Upfall, 2008], [Bubeck, Munos, Stoltz, Szepesvári, 2011]
- **Gaussian bandits** [Dorard, Glowacka, Shawe-Taylor, 2009], [Grünewälder, Audibert, Opper, Shawe-Taylor, 2010], [Srinivas, Krause, Kakade, Seeger, 2010]
- **Contextual bandits** [Woodroffe, 1979], [Auer, 2002], [Wang, Kulkarni, Poor, 2005], [Pandey, Agarwal, Chakrabarti, Josifovski, 2007], [Langford, Zhang, 2007], [Hazan, Megiddo, 2007], [Rigollet, Zeevi, 2010], [Chu, Li, Reyzin, Schapire, 2011], [Slivkins, 2011]
- **MDPs** [Burnetas, Katehakis, 1997], [Jaksch, Ortner, Auer, 2010], [Bartlett, Tewari, 2009]
- **Combinatorial bandits** [Cesa-Bianchi, Lugosi, 2009], [Audibert, Bubeck, Lugosi, 2011]

Bandits in a hierarchy

- Bandit = tool to rapidly select the right action, given noisy estimate of their value
- Serve as building block for more complex problems
- Hierarchy of bandits: The reward obtained when pulling an arm is itself the return of another bandit in a hierarchy.
- Illustration: Monte-Carlo Tree Search in computer-go.



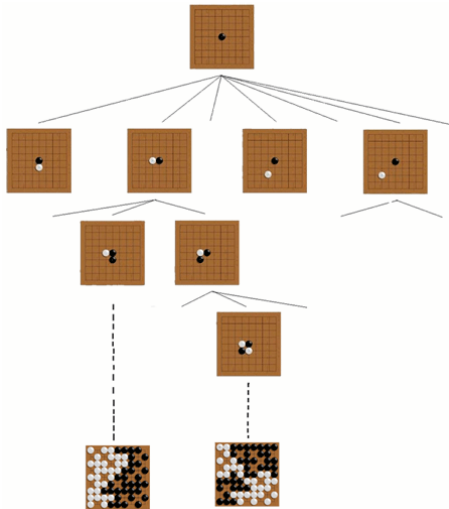
Idea: use bandits at each node of the tree search.

[Gelly, Wang, Munos, Teytaud, 2006] + many others.

Features:

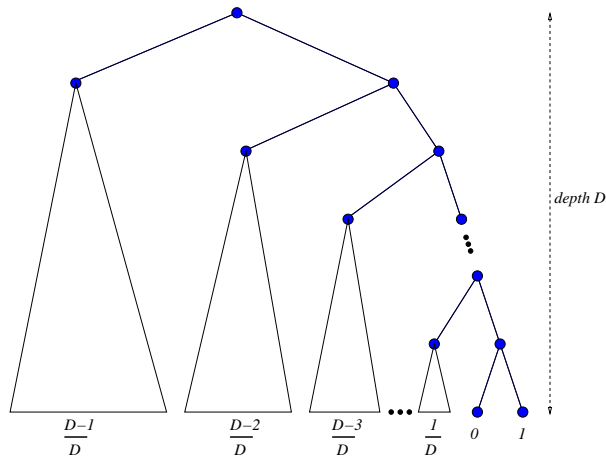
- Explore-Exploit with UCT
- Monte-Carlo evaluation
- Asymmetric tree expansion
- Anytime algo
- Use of features

Among world best programs!



[Kocsis and Szepesvári, 2006]

- In a tree with finite depth, all leaves will be eventually explored an infinite number of times, thus by backward induction, UCT is consistent and the regret is $O(\log n)$.
- However, the constant can be so bad that there is not finite-time guarantee for any reasonable n .



The left branches are explored exponentially more often than the right ones.

Finite-time analysis of UCT

The regret is disastrous: (see [Coquelin and Munos, 2007])

$$\mathbb{E}R_n = \underbrace{\Omega(\exp(\exp(\dots \exp(1)\dots)))}_{D \text{ times}} + O(\log(n)),$$

whereas a uniform exploration of the tree would be “only” exponential in D .

Problem: at each node, the rewards are not i.i.d.

\Rightarrow the B-values are not UCBs.

UCT implicitly makes the assumption that the underlying function is very smooth.

Problems:

- Can we recover the optimistic principle?
- How should we define the smoothness of a function?

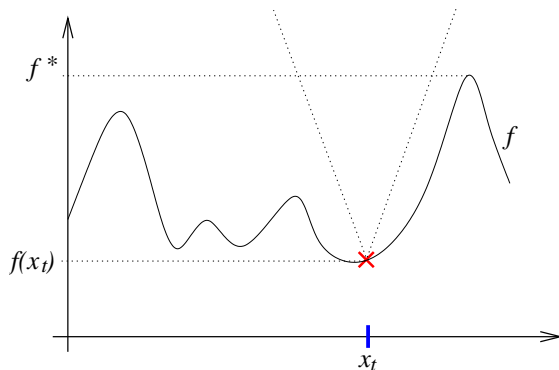
$$|f(x) - f(y)| \leq \ell(x, y).$$

- For each time step $t = 1, 2, \dots, n$ select a state $x_t \in X$
- Observe $f(x_t)$
- Return a state $x(n)$

Performance assessed in terms of the simple regret

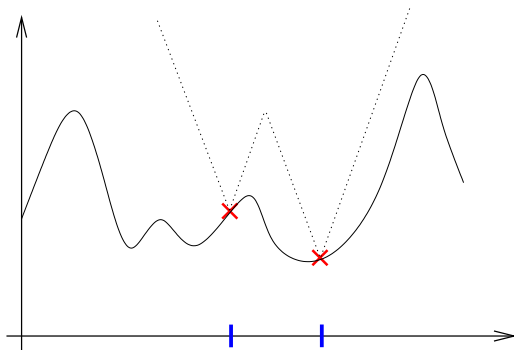
$$r_n = f^* - f(x(n)),$$

where $f^* = \sup_{x \in X} f(x)$.

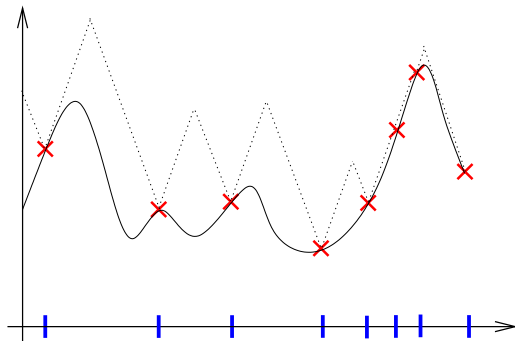


Lipschitz property \rightarrow the evaluation of f at x_t provides a first upper-bound on f .

Example in 1d (continued)



New point \rightarrow refined upper-bound on f .



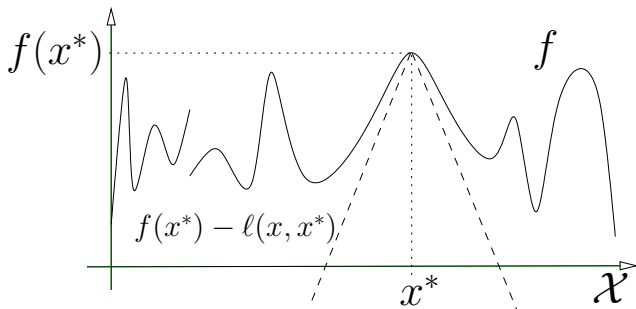
Answer: select the point with highest upper bound!

“Optimism in the face of (partial observation) uncertainty”

- ① \mathcal{X} is equipped with a **semi-metric** ℓ : ℓ is symmetric, and $\ell(x, y) = 0 \Leftrightarrow x = y$.

$$f(x^*) - f(x) \leq \ell(x, x^*).$$

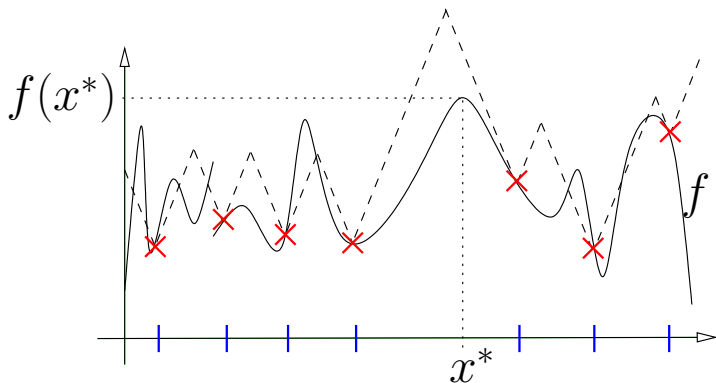
Local smoothness property



For all $x \in \mathcal{X}$,

$$f(x) \geq f(x^*) - \ell(x, x^*).$$

Local smoothness is enough!



Optimistic principle only requires:

- a true bound at the maximum
- the bounds gets refined when adding more points

- For $t = 1$ to n ,
 - Let \mathcal{T}_t be the current partition with cells X_i
 - Define an upper bound for each cell:

where $x_i \in X_i$ and $\text{diam}(X_i) \stackrel{\text{def}}{=} \sup_{x,y \in X_i} \ell(x,y)$

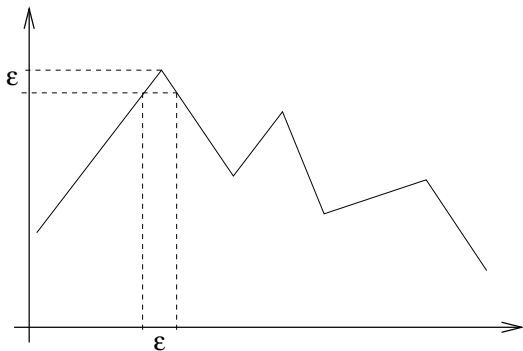
- $$l_t = \underset{j}{\operatorname{argmax}} B_j.$$

- Expand l_t : refine the grid and evaluate f in children cells
- Return $x(n) \stackrel{\text{def}}{=} \operatorname{argmax}_{\{x_t\}_{1 \leq t \leq n}} f(x_t)$

- Thus finite-time performance guarantees can be obtained.

$$X_\varepsilon \stackrel{\text{def}}{=} \{x \in X, f(x) \geq f^* - \varepsilon\}$$

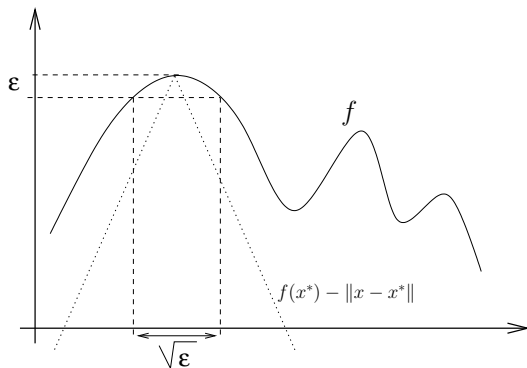
can be covered by $C\varepsilon^{-d}$ ℓ -balls of radius ε .

$$f(x^*) - f(x) = \Theta(\|x^* - x\|).$$


Using $\ell(x, y) = \|x - y\|$, it takes $O(\epsilon^0)$ balls of radius ϵ to cover X_ϵ . Thus $d = 0$.

Assume the function is locally quadratic around its maximum:

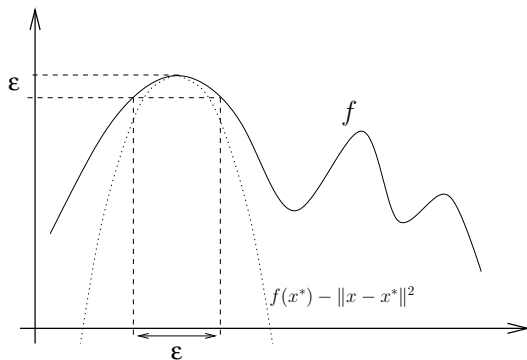
$$f(x^*) - f(x) = \Theta(\|x^* - x\|^2).$$



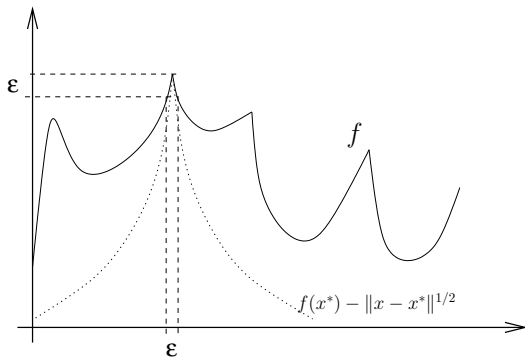
For $\ell(x, y) = \|x - y\|$, it takes $O(\epsilon^{-D/2})$ balls of radius ϵ to cover X_ϵ (of size $O(\epsilon^{D/2})$). Thus $d = D/2$.

Assume the function is locally quadratic around its maximum:

$$f(x^*) - f(x) = \Theta(\|x^* - x\|^2)$$



For $\ell(x, y) = \|x - y\|^2$, it takes $O(\epsilon^0)$ ℓ -balls of radius ϵ to cover X_ϵ . Thus $d = 0$.

$$f(x^*) - f(x) = \Theta(\|x^* - x\|^{1/2})$$


For $\ell(x, y) = \|x - y\|^{1/2}$ we have $d = 0$.

$$f(x^*) - f(x) = \Theta(\|x^* - x\|^\alpha)$$

- If $\alpha = \beta$, then $d = 0$.
- If $\alpha > \beta$, then $d = D(\frac{1}{\beta} - \frac{1}{\alpha}) > 0$.
- If $\alpha < \beta$, then the function is not locally smooth wrt ℓ .

Analysis of DDO (deterministic case)

Assume that the ℓ -diameters of the nodes of depth h decrease exponentially fast with h (i.e., $\text{diam}(h) = c\gamma^h$, for some $c > 0$ and $\gamma < 1$).

This is true for example when $\mathcal{X} = [0, 1]^D$ and $\ell(x, y) = \|x - y\|^\beta$ for some $\beta > 0$.

Theorem 2.

The loss of DOO is

$$r_n = \begin{cases} \left(\frac{C}{1-\gamma^d}\right)^{1/d} n^{-1/d} & \text{for } d > 0, \\ c\gamma^{n/C-1} & \text{for } d = 0. \end{cases}$$

(Remember that $r_n \stackrel{\text{def}}{=} f(x^*) - f(x(n))$).

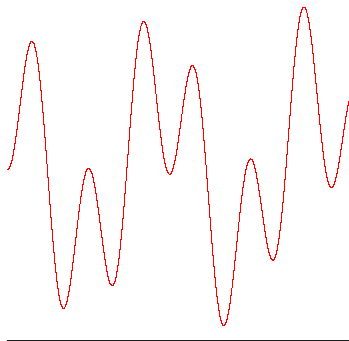
- Only cells X_i of depth h such that $f(x_i) + \text{diam}(h) \geq f(x^*)$ may be expanded by DOO
- From the definition of d , the number of such cells is less than $C \text{diam}(h)^{-d}$
- The number of node expansions $n \leq C \sum_{h=0}^{h_{\max}} \text{diam}(h)^{-d}$
- For $d > 0$, $n = O(\text{diam}(h_{\max})^{-d})$ and the value of the returned point $x(n)$ is at least as good as $f(x_{\max})$ for the deepest expanded node (of depth h_{\max}):

$$f(x(n)) \geq f(x_{\max}) \geq f(x^*) - \text{diam}(h_{\max}) \geq f(x^*) - O(n^{-1/d}).$$

- For $d = 0$, $n = Ch_{\max}$ and $f(x(n)) \geq f(x^*) - O(\gamma^{n/C})$.

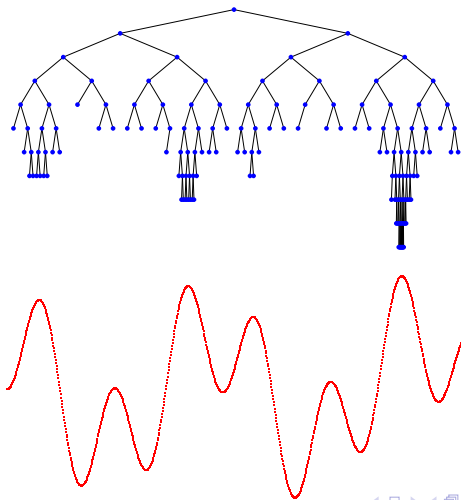
The performance of DOO heavily relies on our knowledge of the true local smoothness.

- $\ell_2(x, y) = 222|x - y|^2$ (i.e., f is locally quadratic).

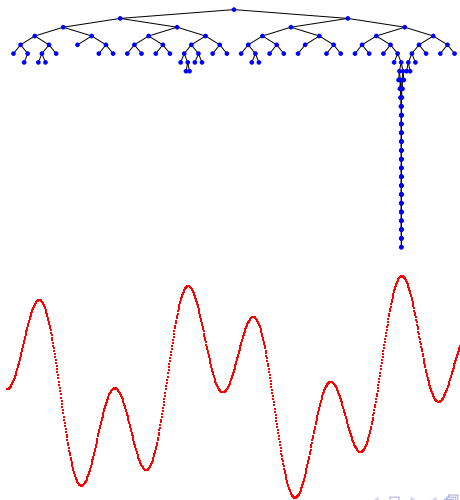


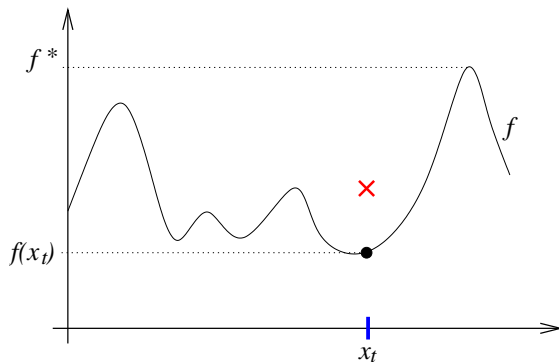
Experiments [2]

Using $\ell_1(x, y) = 14|x - y|$ (i.e., f is globally Lipschitz). $n = 150$.

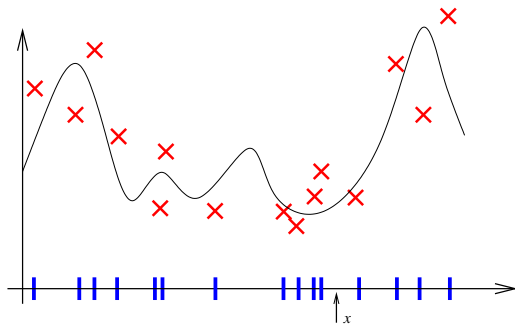


Using $\ell_2(x, y) = 222|x - y|^2$ (i.e., f is locally quadratic). $n = 150$.



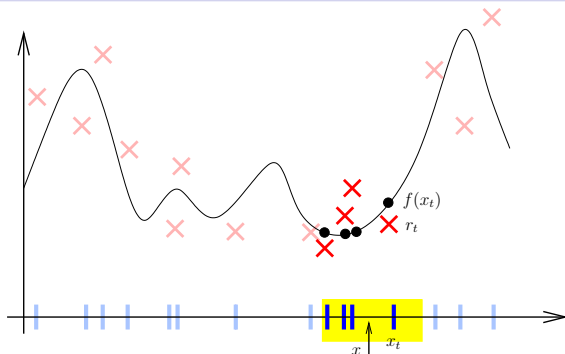
$$r_t = f(x_t) + \epsilon_t, \text{ with } \mathbb{E}[\epsilon_t | x_t] = 0.$$


Where should one sample next?



How to define a high probability upper bound at any state x ?

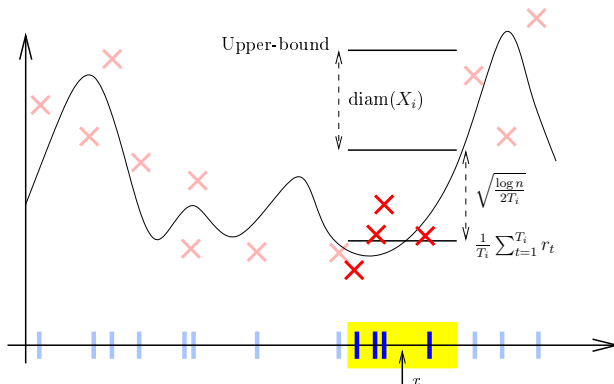
UCB in a given cell



For a fixed domain $X_i \ni x$ containing T_i points $\{x_t\} \in X_i$, we have that $\sum_{t=1}^{T_i} r_t - f(x_t)$ is a Martingale. Thus by Azuma's inequality, with $1/n$ -confidence,

$$\frac{1}{T_i} \sum_{t=1}^{T_i} r_t + \sqrt{\frac{\log n}{2T_i}} \geq \frac{1}{T_i} \sum_{t=1}^{T_i} f(x_t).$$

Upper confidence bound



In any cell X_i define the UCB: $\frac{1}{T_i} \sum_{t=1}^{T_i} r_t + \sqrt{\frac{\log n}{2T_i}} + \text{diam}(X_i)$.

Tradeoff between size of the confidence interval and diameter.

Considering several domains we can derive a tighter UCB

- Consider a series of partitions \mathcal{T}_h of the domain in cells $\{X_{h,i}\}_i$
- Define a UCB for all cells of each partition

$$B_{h,i} = \hat{\mu}_{h,i}(t) + \sqrt{\frac{\log t}{2T_{h,i}(t)}} + \text{diam}(X_{h,i})$$

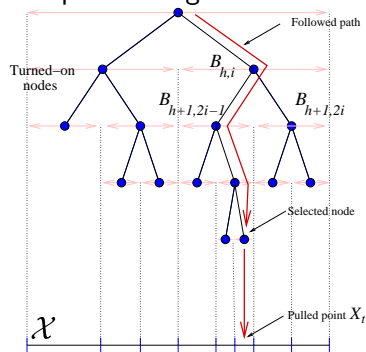
- Define tighter UCB function:

$$B(x) = \min_{X_{h,i} \ni x} B_{h,i},$$

- Select the point with highest UCB:

$$x_{t+1} \in \operatorname{argmax}_x B(x).$$

- Select a leaf J_t of \mathcal{T}_t by following a path from the root that maximizes the B-values,
- Expand J_t :
- Select $x_t \in X_{J_t}$ (arbitrarily)
- Observe reward $r_t = f(x_t) + \epsilon_t$ and update the B-values:



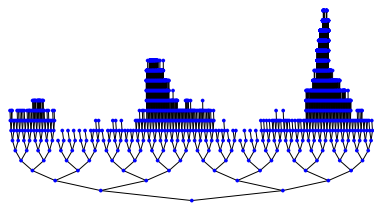
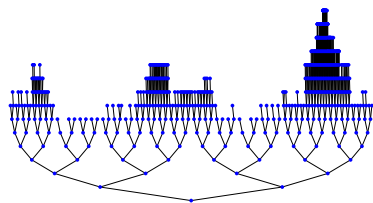
$$B_i(t) \stackrel{\text{def}}{=} \min \left[\hat{\mu}_{i,t} + \sqrt{\frac{2 \log(t)}{T_i(t)}} + \text{diam}(X_i), \max_{j \in \mathcal{C}(i)} B_j(t) \right]$$

- $$J_t \in \operatorname{argmax}_{j \in \mathcal{L}_t} \min_{i \in \mathcal{P}(j)} \left[\hat{\mu}_{i,t} + \sqrt{\frac{2 \log(t)}{T_i(t)}} + \operatorname{diam}(X_i) \right]$$

- For any $X_i \ni x^*$, $B_i(t)$ is a h.p. UCB on $f(x^*)$.
- Thus any suboptimal node X_i such that

$$\sup_{x \in X_i} f(x) + \sqrt{\frac{2 \log(t)}{T_i(t)}} + \text{diam}(X_i) < f(x^*)$$

will not be selected.



Resulting tree at time $n = 1000$ and at $n = 10000$.

$$\forall x, y \in \mathcal{X},$$

$$f(y) - f(x) \leq \max\{f^* - f(y), \ell(x, y)\}$$

Theorem 3.

Assume that the diameters decrease exponentially fast with their depth h . The loss of HOO is

$$r_n = O\left(\left[\frac{n}{\log n}\right]^{-\frac{1}{d+2}}\right)$$

(recall that for deterministic rewards $r_n = O(n^{-1/d})$ for $d > 0$)
(see also the Zooming algorithm [Kleinberg, Slivkins, Upfal, 2008]).

$$f(x^*) - f(x) = \Theta(\|x^* - x\|^\alpha).$$

Choose $\ell(x, y) = \|x - y\|^\beta$.

- If the smoothness of the function is known ($\alpha = \beta$): the loss of HOO is $O(\sqrt{\log n/n})$. The rate is **independent of the dimension**.
- The smoothness is underestimated ($\alpha > \beta$): $d = D(1/\beta - 1/\alpha)$ and the loss is $\tilde{O}(n^{-1/(d+2)})$
- The smoothness is overestimated ($\alpha < \beta$): the weak-Lipschitz assumption is violated, thus there is no finite-time guarantee (e.g., UCT = HOO with $\beta = \infty$)

Assume that the smoothness is unknown

Previous algorithms heavily rely on the knowledge or the local smoothness of the function (i.e. knowledge of the best metric).

Question: When the smoothness is unknown, is it possible to implement the optimistic principle for function optimization?

Some approaches relies on estimating the local or global smoothness of the function [Bubeck, Stoltz, Yu, 2011], [Slivkins, 2011], [Bull, 2013].

DIRECT algorithm

Assume f is Lipschitz but the Lipschitz constant L is unknown.

The DIRECT algorithm [Jones, Perttunen, Stuckman, 1993] expands simultaneously all nodes that may potentially contain the maximum for some value of L .

The sin function and its upper bound for $L = 2$.

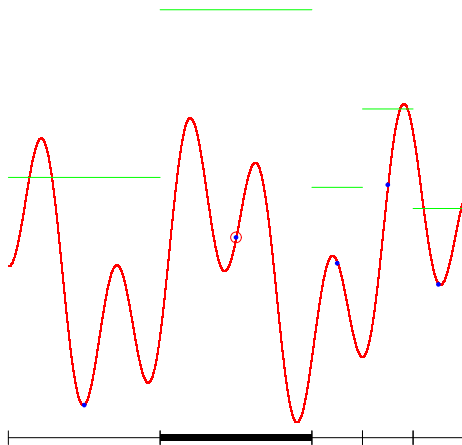
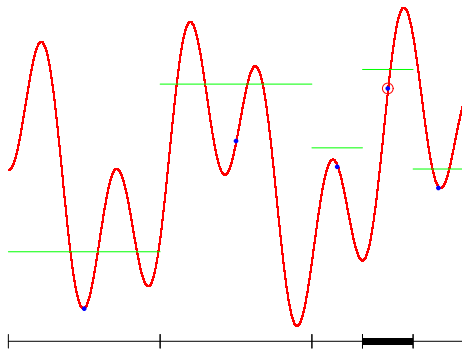
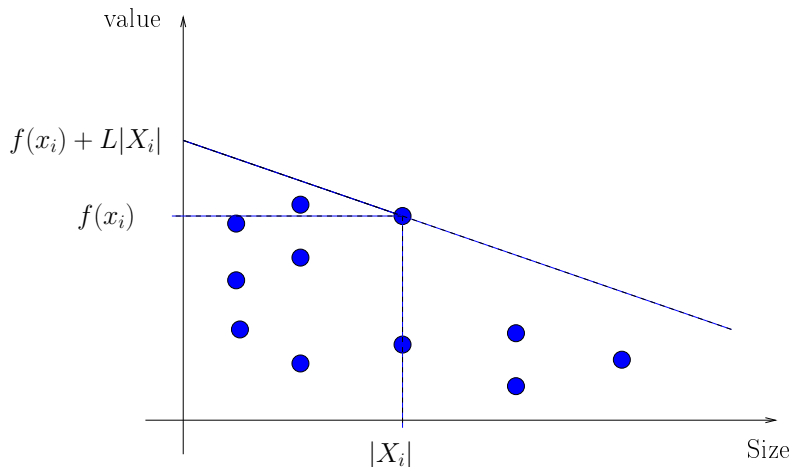


Illustration of DIRECT

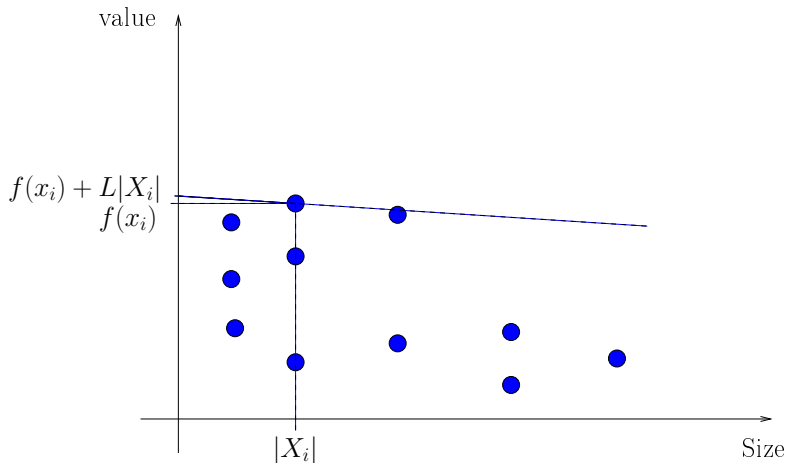
The sin function and its upper bound for $L = 1/2$.



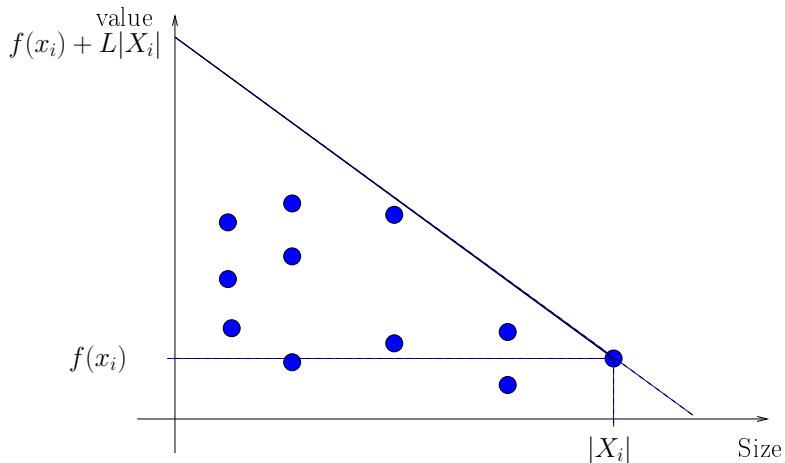
Other representation of DIRECT



Other representation of DIRECT



Other representation of DIRECT



1. *Journal of Management Studies*, 1997, 34, 1, 1-15.

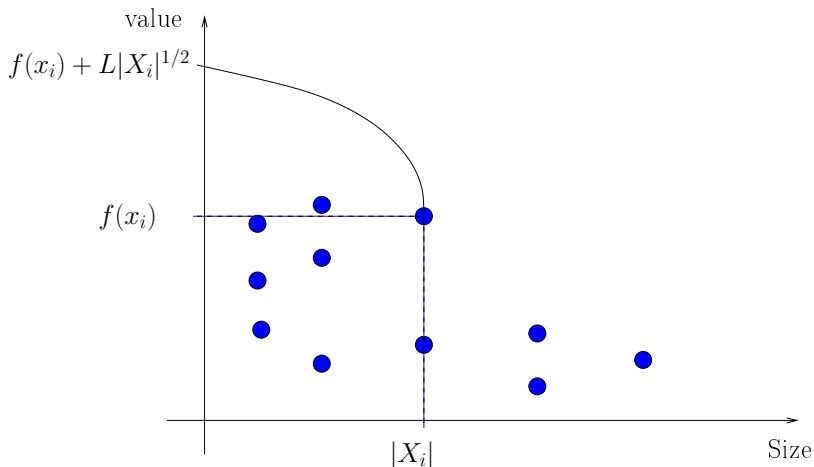


Limitations of DIRECT

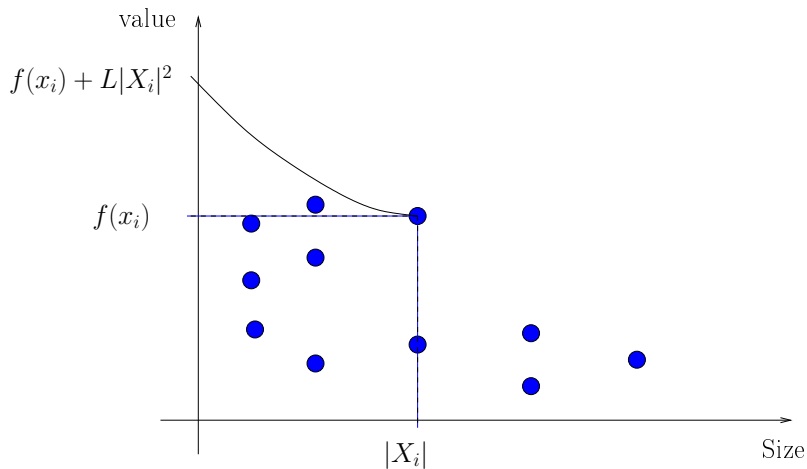
Assuming the function is globally Lipschitz is too restrictive. We would like to handle the general case where:

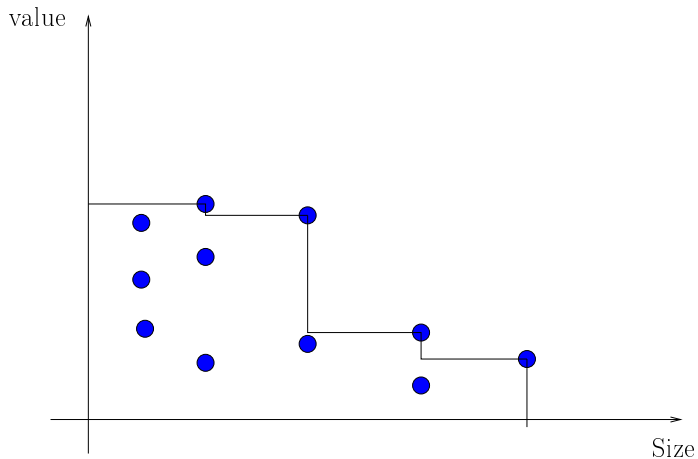
- where the function is only **locally smooth** w.r.t. ℓ
- for **any semi-metric** ℓ

Extension to $\ell(x, y) = L\|x - y\|^{1/2}$



Extension to $\ell(x, y) = L\|x - y\|^2$





- Expand several leaves simultaneously
- SOO expands at most one leaf per depth
- SOO expands a leaf only if its value is larger than the value of all leaves of same or lower depths.
- At round t , SOO does not expand leaves with depth larger than $h_{\max}(t)$

SOO algorithm

Input: the maximum depth function $t \mapsto h_{\max}(t)$

Initialization: $\mathcal{T}_1 = \{(0, 0)\}$ (root node). Set $t = 1$.

while True do

Set $v_{\max} = -\infty$.

for $h = 0$ to $\min(\text{depth}(\mathcal{T}_t), h_{\max}(t))$ **do**

Select the leaf $(h, j) \in \mathcal{L}_t$ of depth h with $\max f(x_{h,j})$ value

if $f(x_{h,i}) > v_{\max}$ **then**

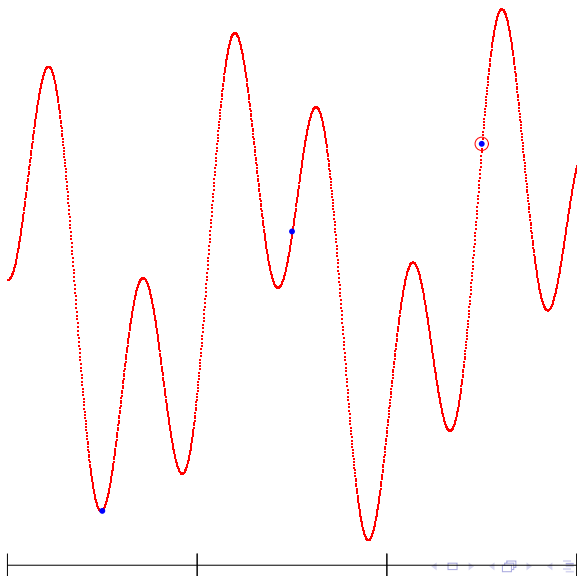
Expand the node (h, i) , Set $v_{\max} = f(x_{h,i})$, Set $t = t + 1$

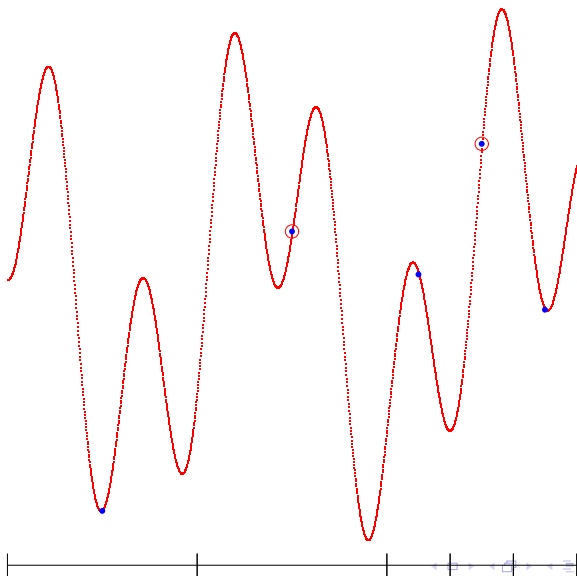
if $t = n$ **then** return $x(n) = \arg \max_{(h,i) \in \mathcal{T}_n} x_{h,i}$

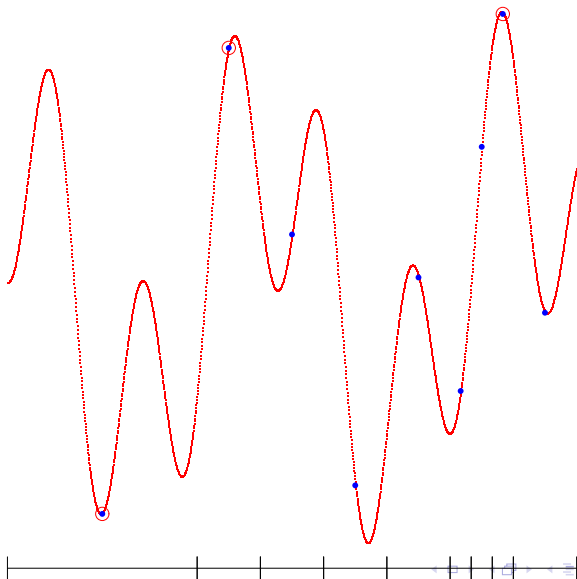
end if

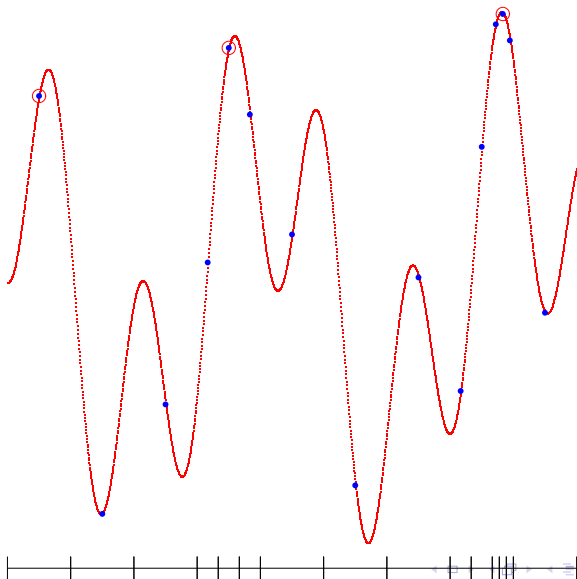
end for

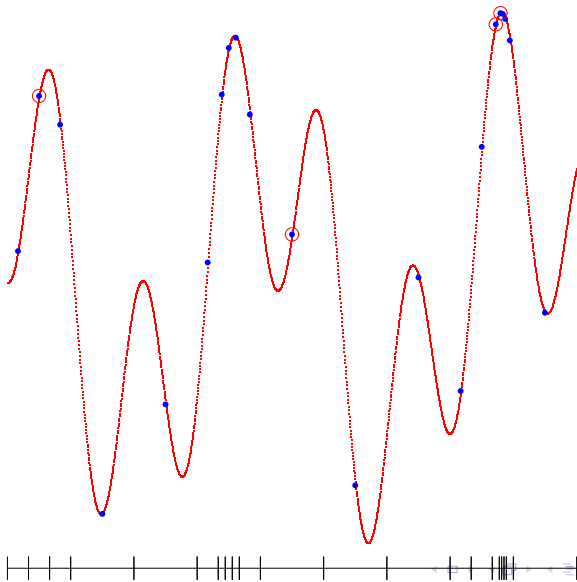
end while.

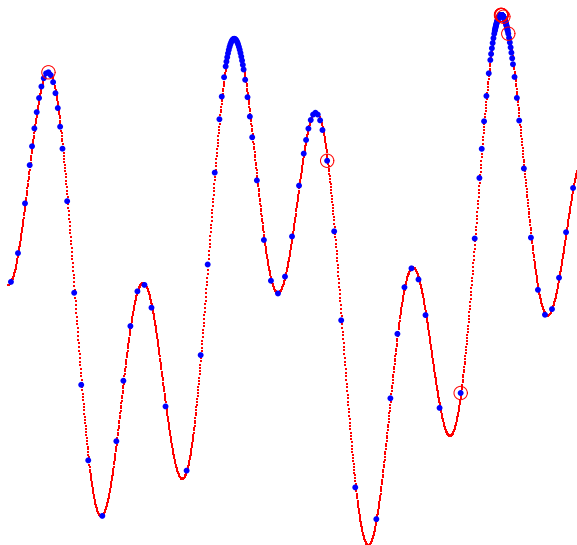


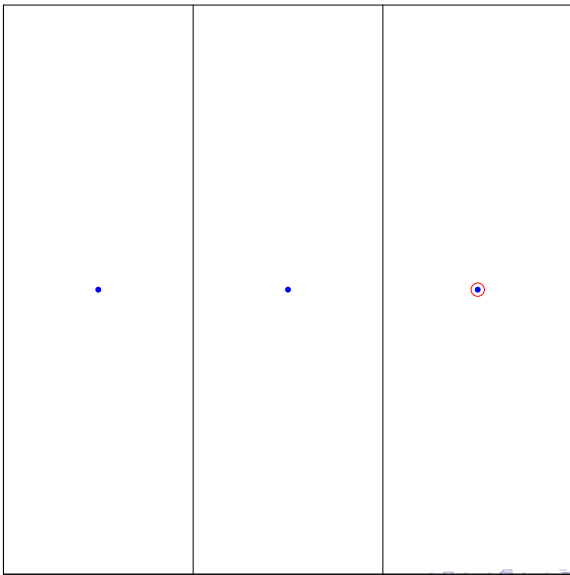






















Performance of S00

Theorem 4.

For any semi-metric ℓ such that

- f is locally smooth w.r.t. ℓ
- The ℓ -diameter of cells of depth h is $c\gamma^h$
- The near-optimality dimension of f w.r.t. ℓ is $d = 0$,

by choosing $h_{\max}(n) = \sqrt{n}$, the expected loss of SOO is

$$r_n \leq c\gamma^{\sqrt{n}/C-1}$$

In the case $d > 0$ a similar statement holds with $\mathbb{E}r_n = \tilde{O}(n^{-1/d})$.

n	SOO	uniform grid	DOO with ℓ_1	DOO with ℓ_2
50	3.56×10^{-4}	1.25×10^{-2}	2.53×10^{-5}	1.20×10^{-2}
100	5.90×10^{-7}	8.31×10^{-3}	2.53×10^{-5}	1.67×10^{-7}
150	1.92×10^{-10}	9.72×10^{-3}	4.93×10^{-6}	4.44×10^{-16}

The case $d = 0$ is non-trivial!

Example:

- f is locally α -smooth around its maximum:

$$f(x^*) - f(x) = \Theta(\|x^* - x\|^\alpha),$$

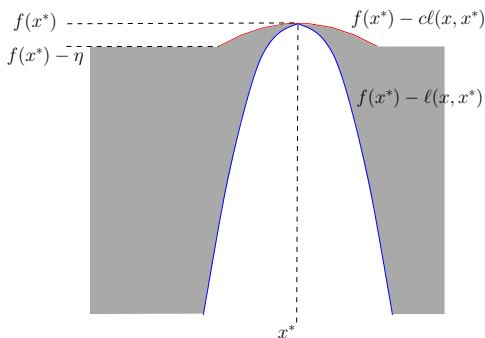
for some $\alpha > 0$.

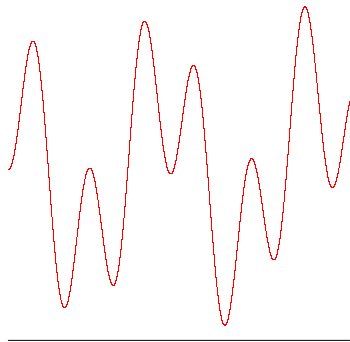
- SOO algorithm does not require the knowledge of ℓ ,
- Using $\ell(x, y) = \|x - y\|^\alpha$ in the analysis, all assumptions are satisfied (with $\gamma = 3^{-\alpha/D}$ and $d = 0$, thus the loss of SOO is $r_n = O(3^{-\sqrt{n}\alpha/(CD)})$ (stretched-exponential loss),
- This is almost as good as DOO optimally fitted!

(Extends to the case $f(x^*) - f(x) \approx \sum_{i=1}^D c_i |x_i^* - x_i|^{\alpha_i}$)

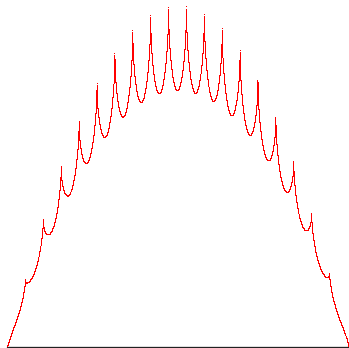
$$\min(\eta, c\ell(x, x^*)) \leq f(x^*) - f(x) \leq \ell(x, x^*), \quad \text{for all } x \in \mathcal{X}.$$

has a near-optimality $d = 0$ (w.r.t. the metric ℓ).

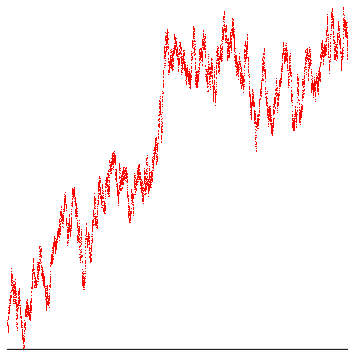




$$\ell(x, y) = c\|x - y\|^2$$

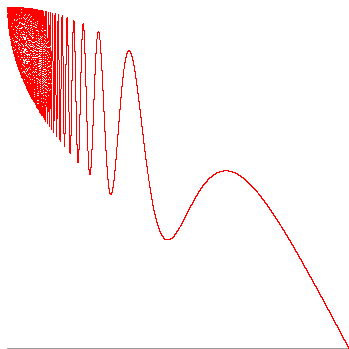


$$\ell(x, y) = c\|x - y\|^{1/2}$$

$d = 0?$ 

$$\ell(x, y) = c\|x - y\|^{1/2}$$

$$f(x) = 1 - \sqrt{x} + (-x^2 + \sqrt{x}) * (\sin(1/x^2) + 1)/2$$



The lower-envelope is of order $1/2$ whereas the upper one is of order 2. We deduce that $d \geq 3/2$.

- For any ℓ , define $I_h = \{ \text{cells } X_i \text{ of depth } h \text{ such that } f(x_i) + \text{diam}(h) \geq f(x^*) \}$
- Once the optimal cell of depth h has been expanded, it takes at most $|I_{h+1}|$ cell expansions of depth $h+1$ before the optimal cell is expanded.
- Thus $n \leq h_{\max}(n) \sum_{h=0}^{\min(h_{\max}(n), h_n^*)} |I_h|$, where h_n^* is the depth of the node containing x^* .
- Assuming $d = 0$, $|I_h| \leq C$, and using $h_{\max}(n) = \sqrt{n}$, we have $\sqrt{n} = C \min(h_{\max}(n), h_n^*) = Ch_n^*$
- Finally the value of the returned point $x(n)$ is at least as good as that of the optimal expanded node i_n^* containing x^* :

$$f(x(n)) \geq f(x_{j_n^*}) \geq f(x^*) - \text{diam}(h_n^*) \geq f(x^*) - c\gamma^{\sqrt{n}/C},$$

where we used that the diameters are $c\gamma^h$.

Comparison SOO versus DIRECT algorithms

- **SOO is much more general than DIRECT:** the function is only **locally smooth** and the space is **semi-metric**.
- **Finite-time analysis of SOO** (whereas only a consistency property $\lim_{n \rightarrow \infty} r_n = 0$ is available for DIRECT in [Finkel and Kelley, 2004])
- **SOO is a rank-based algorithm:** any transformation of the values while preserving their rank will not change anything in the algorithm. Thus extends to the optimization of function givens pair-wise comparisons.
- SOO is easier to implement...

Stochastic S00 (StoS00)

A simple way to extend SOO to the case of stochastic rewards is the following:

- Select a cell i (and sample f at x_i) according to SOO based on the values

$$\hat{\mu}_{i,t} + c\sqrt{\frac{\log n}{T_k(t)}},$$

(where $\hat{\mu}_{i,t}$ is the average rewards received at x_i and $T_i(t)$ is the number of rewards received at state x_i),

- Expand the cell X_i only if $T_i(t) \geq k$, where k is a parameter.

Remark: This really looks like UCT, except that

- several cells are selected at each round,
- a cell is refined only when we received k samples.

Intuition for StoSOO

With high probability, StoSOO acts as a ϵ -perturbed version of SOO where:

- The values $f(x_i)$ are perturbed by ϵ , where $\epsilon = O(\sqrt{\frac{\log n}{k}})$
- There are only $m = n/k$ evaluations to the function.

Thus the loss of StoSOO is

$$\mathbb{E}r_n(\text{StoSOO}) = \mathbb{E}r_m(\text{SOO}) + O(\sqrt{\frac{\log n}{k}})$$

Performance of StoSOO

Theorem 5 (Valko, Carpentier, Munos, 2013).

For any semi-metric ℓ such that

- f is locally smooth w.r.t. ℓ
- The ℓ -diameters of the cells decrease exponentially fast with their depth,
- The near-optimality dimension of f w.r.t. ℓ is $d = 0$,

by choosing $k = \frac{n}{(\log n)^3}$, $h_{\max}(n) = (\log n)^{3/2}$, the expected loss of *StoSOO* is

$$\mathbb{E}r_n = O\left(\frac{(\log n)^2}{\sqrt{n}}\right).$$

Comments about StoS00

In the (rather general) case $d = 0$, StoSOO gives a $\tilde{O}(1/\sqrt{n})$ loss
In comparison to HOO/Zooming:

- HOO optimally fitted gives $O\left(\sqrt{\frac{\log n}{n}}\right)$ loss
- HOO with underestimation of the right smoothness deteriorates
- HOO with overestimation of the right smoothness may not converge.

Thus **StoS00** is almost as good as **H00** optimally fitted.

But there are plenty of other semi-metric spaces:

- Trees, graphs, combinatorial spaces, structured spaces, ...
- Ex: space of policies for MDPs

We only require:

- the search space \mathcal{X} to be equipped with a semi-metric ℓ
- a nested (hierarchical) partitioning of the space, such that the ℓ -diameters of the cells decrease with their depth
- f to satisfy a local smoothness property w.r.t. ℓ
- ℓ may or may not be known.

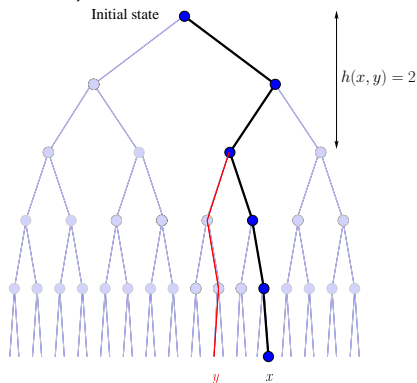
A. J. L. C. J. MDD

1 6 1 1 1 1 1 1

Deterministic transitions and rewards

(infinite time horizon and discounted setting)

- A policy x is an infinite path
- Value $f(x) = \sum_{s \geq 0} \gamma^s r_s(x)$, where the rewards are in $[0, 1]$
- Metric: $\ell(x, y) = \frac{\gamma^{h(x, y)}}{1 - \gamma}$
- Prop: $f(x)$ is Lipschitz w.r.t. ℓ
- \rightarrow Use optimistic search



OPD algorithm

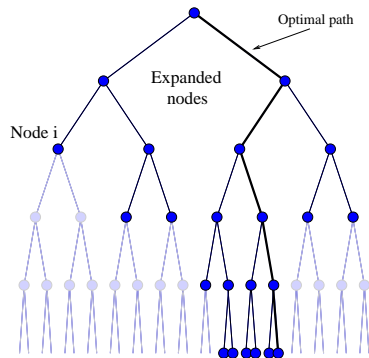
Optimistic Planning in Deterministic systems:

- Define the B-values:

$$B_i \stackrel{\text{def}}{=} \sum_{s=0}^{d(i)} \gamma^s r_s + \frac{\gamma^{d(i)+1}}{1-\gamma}$$

- We have $B_i \geq \max_{x \in i} f(x)$
- For each round $t = 1$ to n ,
expand the node with
highest B-value
- Observe reward, update
B-values

Recommend the first action $a(n)$ of the best policy.



κ -minimax lower bounds

Let \mathcal{M}_κ the set of problems with coefficient κ .

Upper-bound of OPD uniformly over \mathcal{M}_K

$$\sup_{M \in \mathcal{M}_\kappa} r_n(\mathcal{A}_{OPD}, M) = O\left(n^{-\frac{\log 1/\gamma}{\log \kappa}}\right).$$

We can prove that we have a κ -minimax lower-bound:

$$\sup_{\mathcal{A}} \inf_{M \in \mathcal{M}_\kappa} r_n(\mathcal{A}, M) \geq \Omega \left(n^{-\frac{\log 1/\gamma}{\log \kappa}} \right).$$

Sketch of proof:

OPD only expands nodes in I . Reciprocally, I is the set of nodes that need to be expanded in order to find the optimal path.

Using HOO for planning

Apply HOO to the search space \mathcal{X} :

- Assign a B-value to each finite sequence
- At each round t , select a finite sequence x_t maximizing the B-value.
- Observe sample reward $\sum_{s \geq 0} \gamma^s Y_s(x_t)$ of the path x_t and use it to update the B-values.
- The loss is

$$O(n^{-\frac{1}{d+2}}).$$

Problem: HOO does not make full use of the tree structure: It uses the sample reward of the whole path x but not the individual reward samples $Y_s(x)$ collected along the path x .

- At round t , play path x_t (up to depth $h = \frac{1}{2} \frac{\log n}{\log 1/\gamma}$)
- Observe sample rewards $Y_s(x_t)$ of each node along the path x_t
- Compute empirical rewards $\hat{\mu}_t(x_{1:s})$ for each node $x_{1:s}$ of depth $s \leq h$
- Define bound for each path x :

$$B_t(x) = \min_{1 \leq j \leq h} \left[\sum_{s=0}^j \gamma^s \left(\hat{\mu}_t(x_{1:s}) + \sqrt{\frac{2 \log n}{T_{x_{1:s}}(t)}} \right) + \frac{\gamma^{j+1}}{1 - \gamma} \right]$$

- Select path $x_{t+1} = \operatorname{argmax}_x B_t(x)$

This algorithm fully uses the tree structure of the rewards.

Performance of OLOP

Define

- $\beta \geq 0$ such that $\mathbb{P}(\text{Random path is } \epsilon\text{-optimal}) = O(\epsilon^\beta)$.
- or $\kappa \stackrel{\text{def}}{=} K\gamma^\beta \in [1, K]$ the branching factor of the set of near-optimal sequences.
- or the near-optimality dimension, $d = \frac{\log \kappa}{\log 1/\gamma}$.

Theorem 6 (Loss of OLOP).

After n calls to the generative model,

$$\mathbb{E}r_n = f(x^*) - \mathbb{E}f(x(n)) = \begin{cases} \tilde{O}(n^{-1/d}) & \text{if } d \geq 2 \\ \tilde{O}(1/\sqrt{n}) & \text{if } d < 2 \end{cases}$$

Much better than HOO! As good as OPD for $d \geq 2$.

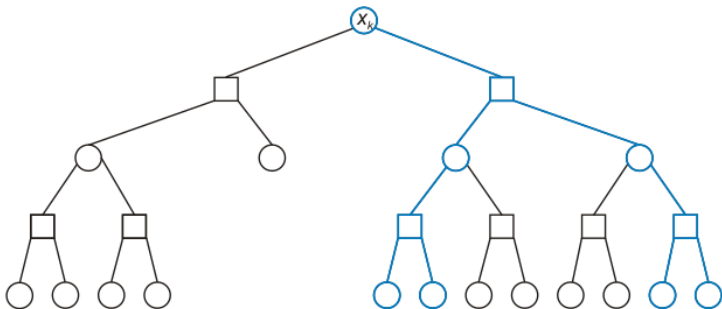


Here a policy is no more a sequence of actions

OP-MDP [Buşoniu and Munos, 2012]:

- The root = current state.
- For $t = 1$ to n :
 - Compute the B-value of all nodes of the current sub-tree
 - Compute the optimistic policy
 - Select a leaf of the optimistic sub-tree and expand it (generates transitions to next states using the model)
- Return first action of the best policy

Illustration of OP-MDP

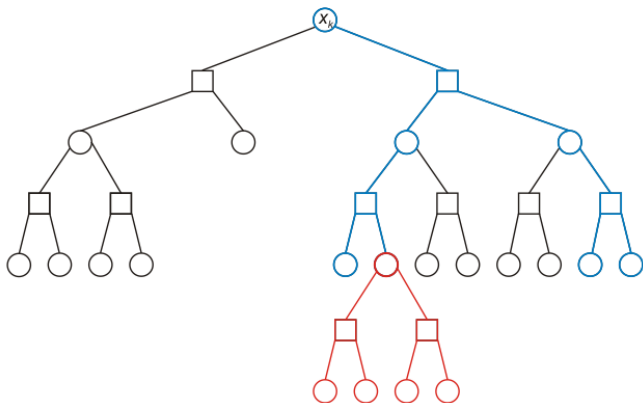


B-values: upper-bounds on the optimal value function $V^*(s)$:

$$B(s) = \frac{1}{1-\gamma} \text{ for leaves}$$

$$B(s) = \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma B(s')]$$

Compute the **optimistic policy** π^+ .



Expand leaf in π^+ with largest contribution: $\arg \max_{s \in \mathcal{L}} P(s)^{\frac{\gamma^{d(s)}}{1-\gamma}}$, where $P(s)$ is the probability to reach s when following π^+ .

- whose “contribution” is at least ϵ
- and that belong to an ϵ -optimal policy

Theorem 7.

The performance of OP-MDP is $r_n = O(n^{-1/d})$.

The performance depends on the quantity of states that contribute significantly to near-optimal policies

Illustration of the performance

Reminder: $r_n = O(n^{-1/d})$.

Uniform rewards and probabilities $d = \frac{\log K + \log N}{\log 1/\gamma}$ (uniform planning)

Structured rewards, uniform probabilities $d = \frac{\log N}{\log 1/\gamma}$ (uniform planning for a single policy)

Uniform rewards, concentrated probabilities $d \rightarrow \frac{\log K}{\log 1/\gamma}$
(planning in deterministic systems)

Structured rewards, concentrated probabilities $d \rightarrow 0$
(exponential rate)

Remarks: d is small when

- Structured rewards
- Heterogeneous transition probabilities

Upper-bound of OP-MDP uniformly over \mathcal{M}_β

$$\sup_{M \in \mathcal{M}_d} r_n(\mathcal{A}_{OP-MDP}, M) \leq O(n^{-1/d}).$$

We conjecture that we have a d -minimax lower-bound:

$$\sup_A \inf_{M \in \mathcal{M}_d} \mathbb{E} r_n(\mathcal{A}, M) \geq \Omega(n^{-1/d}).$$

- Perform optimistic search in policy space.
- In deterministic dynamics, deterministic rewards, can be seen as a direct application of optimistic optimization
- In stochastic rewards, the structure of the reward function can help estimation of paths given samples from other paths
- In MDPs the **near-optimality planning dimension** is a new measure of the quantity of states that need to be expanded (*the set of states that significantly contribute to near-optimal policies*)
- Fast rates when the MDP has structured rewards and heterogeneous transition probabilities.
- Applications to Bayesian Reinforcement learning and planning in POMDPs.

Possible extensions in optimistic planning

- Extends OP-MDP to the case when only a generative model of the dynamics is available
- Extension to a possibly infinite number of next-states
- Apply SOO / StoSOO ideas for planning: Although the value function is Lipschitz w.r.t. the metric $\ell(x, y) = \frac{\gamma^{h(x, y)}}{1 - \gamma}$, it may possess additional smoothness around the maximum with a higher-order semi-metric.

- applies in a large class of decision making problems in stochastic and deterministic environments (in unstructured and structured problems)
- provides an efficient exploration of the search space by exploring the most promising areas first
- provides a natural transition from global to local search
- Performance depends on the “smoothness” of the function around the maximum w.r.t. some metric,
 - a measure of the quantity of near-optimal solutions,
 - and our knowledge or not of this metric.

Thanks !!!

<http://chercheurs.lille.inria.fr/~munos/>