

Missing Mass, and Optimal Discovery

based on a joint work with Sébastien Bubeck and Damien Ernst

Aurélien Garivier

Laboratoire de l'Informatique du Parallélisme, le 14 mai 2018

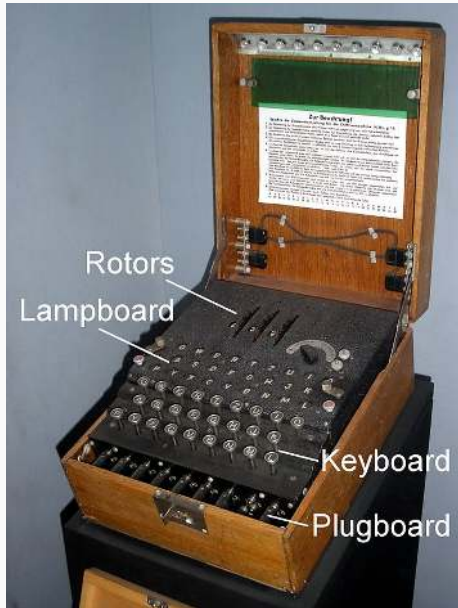
Équipe-projet AOC: Apprentissage, Optimisation, Complexité
Institut de Mathématiques de Toulouse LabeX CIMI
Université Paul Sabatier Toulouse III

Table of contents

1. Estimating the Unseen
2. Discovering dangerous contingencies in electrical systems
3. The Good-UCB algorithm
4. Optimality results

Estimating the Unseen

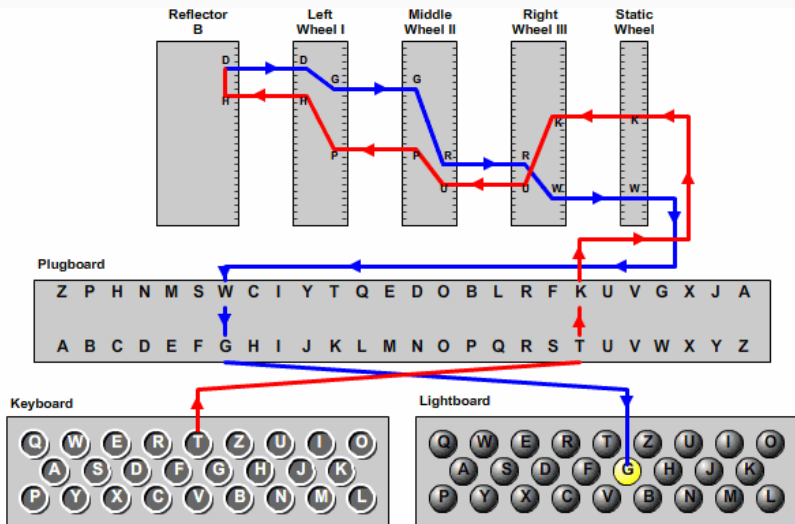
Enigma



- Electro-mechanical rotor cipher machines, 26 characters
- Invented at the end of WW1 by Arthur Scherbius
- Commercial use, then German Army during WW2
- First cracked by Marian Rejewski in the 1930s (Bomb), then improved to $3 \cdot 10^{114}$ configurations
- Read Simon Singh, *The Code Book*



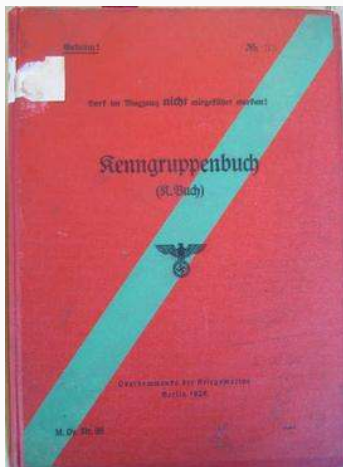
Enigma



© 2006, by Louise Dade

Src: <http://enigma.louisedade.co.uk/>

Battle of the Atlantic



- Massively used by the German Kriegsmarine and Luftwaffe
- **weakness:** 3-letters setting to initiate communication, taken from the *Kenngruppenbuch*
- Government Code and Cypher School: Bletchley Park (on the train line between Cambridge and Oxford)
- Colossus (first programmable computers) in 1943

Estimating probabilities

- Discrete alphabet A .
- Unknown probability P on A
- Sample X_1, \dots, X_n of independent draws of P .
- Goal : use the sample estimate $\hat{P}(a)$ for all $a \in A$.

Natural idea:

$$\hat{P}(a) = \frac{N(a)}{n}, \quad \text{where } N(a) = \#\{i : X_i = a\}$$

Safari preparation

Observe animal sample

1 giraffe, 2 elephants, 3 zebras

Probability estimation?

Empirical frequency

Species	Probability
giraffes	1/6
elephants	2/6
zebras	3/6

Problem?



Learning set:

john read moby dick

mary read a different book

she read a book by cher

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1} w_i)}{\sum_w c(w_{i-1} w)}$$

$$P(s) = \prod_{i=1}^{l+1} p(w_i | w_{i-1})$$

$$\begin{aligned}
 &P(\quad john \quad \quad read \quad \quad a \quad \quad book \quad \quad) \\
 &= P(john | \cdot) \quad P(read | john) \quad P(a | read) \quad P(book | a) \quad P(\cdot | book) \\
 &= \frac{c(\cdot john)}{\sum_w c(\cdot w)} \quad \frac{c(john read)}{\sum_w c(john w)} \quad \frac{c(read a)}{\sum_w c(read w)} \quad \frac{c(a book)}{\sum_w c(a w)} \quad \frac{c(book \cdot)}{\sum_w c(book w)} \\
 &= \frac{1}{3} \quad \frac{1}{1} \quad \frac{2}{3} \quad \frac{1}{2} \quad \frac{1}{2} \\
 &\approx 0.06
 \end{aligned}$$

Learning set:

john read moby dick

mary read a different book

she read a book by cher

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1} w_i)}{\sum_w c(w_{i-1} w)}$$

$$P(s) = \prod_{i=1}^{l+1} p(w_i | w_{i-1})$$

$$\begin{aligned} & P(\text{ cher read a book }) \\ = & P(\text{ cher} | \cdot) P(\text{ read} | \text{ john}) P(\text{ a} | \text{ read}) P(\text{ book} | \text{ a}) P(\cdot | \text{ book}) \\ = & \frac{c(\cdot \text{ cher})}{\sum_w c(\cdot w)} \frac{c(\text{ cher read})}{\sum_w c(\text{ cher } w)} \frac{c(\text{ read a})}{\sum_w c(\text{ read } w)} \frac{c(\text{ a book})}{\sum_w c(\text{ a } w)} \frac{c(\text{ book } \cdot)}{\sum_w c(\text{ book } w)} \\ = & \frac{0}{3} \quad \frac{0}{1} \quad \frac{2}{3} \quad \frac{1}{2} \quad \frac{1}{2} \\ = & \mathbf{0} \end{aligned}$$

⇒ useless, the unseen **must** be treated correctly.

Bayesian Approach: Laplace Estimator

Pierre-Simon de Laplace (1749-1827), Thomas Bayes (1702-1761)

Will the sun rise tomorrow?

$$\hat{P}(a) = \frac{N(a) + 1}{n + |A|}$$

- good for small alphabets and many samples
- very bad when lots of items seen once (ex: DNA sequences)
- $|A|$ can be very large (or even infinite), but P concentrated on few items

⇒ not a satisfying solution to the problem

Alan Turing



1912-1954
student of Godfrey Harold Hardy
in Cambridge
PhD from Princeton with Alonzo
Church

Irving John Good

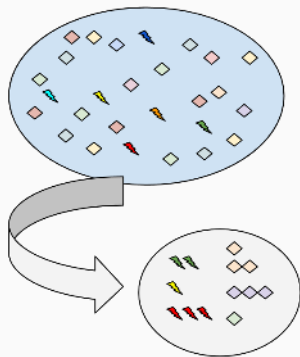


1916-2009
Graduated in Cambridge
Academic career in Bayesian statistics
in Manchester and then in the
University of Virginia (USA)

Missing mass estimation

X_1, \dots, X_n independent draws of $P \in \mathfrak{M}_1(A)$.

$$O_n(x) = \sum_{m=1}^n \mathbb{1}\{X_m = x\}$$



How to 'estimate' the **total mass of the *unseen*** items

$$R_n = \sum_{x \in A} P(x) \mathbb{1}\{O_n(x) = 0\} ?$$

The Good-Turing Estimator

See [I.J. Good, 1953], credits idea to A. Turing

Idea: in order to estimate the mass of the unseen

$$R_n = \sum_{x \in A} P(x) \mathbb{1}\{O_n(x) = 0\},$$

use the number of **hapaxes** = items seen only once (linguistic)

$$\hat{R}_n = \frac{U_n}{n}, \quad \text{where } U_n = \sum_{x \in A} \mathbb{1}\{O_n(x) = 1\}$$

Lemma [Good '53]: For every distribution P ,

$$0 \leq \mathbb{E}[\hat{R}_n] - \mathbb{E}[R_n] \leq \frac{1}{n}$$

Completely non-parametric: no assumption on P

Bias of the Good-Turing Estimator

$$\begin{aligned}\mathbb{E}[\hat{R}_n] - \mathbb{E}[R_n] &= \frac{1}{n} \sum_{x \in A} \mathbb{P}(O_n(x) = 1) - \sum_{x \in A} P(x) \mathbb{P}(O_n(x) = 0) \\ &= \frac{1}{n} \sum_{x \in A} n P(x) (1 - P(x))^{n-1} - \sum_{x \in A} P(x) (1 - P(x))^n \\ &= \sum_{x \in A} P(x) (1 - P(x))^{n-1} (1 - (1 - P(x))) \\ &= \frac{1}{n} \sum_{x \in A} P(x) \times n P(x) (1 - P(x))^{n-1} \\ &= \frac{1}{n} \sum_{x \in A} P(x) \mathbb{P}(O_n(x) = 1) \\ &= \frac{1}{n} \mathbb{E} \left[\sum_{x \in A} P(x) \mathbb{1}\{O_n(x) = 1\} \right] \in \left[0, \frac{1}{n} \right]\end{aligned}$$

Jackknife interpretation

If we had additional samples, we would estimate R_n by the proportion of unseen elements in X_{n+1}, X_{n+2}, \dots

We have no additional samples, **but** we keep every observation as a "test", pretending that the samples was made of everything else:

$$\begin{aligned}\hat{R}_n &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \notin \{x_j : j \neq i\}\} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{O_n(x_i) = 1\} \\ &= \frac{1}{n} \sum_{x \in A} \mathbb{1}\{O_n(x) = 1\}\end{aligned}$$

Remark: jackknife is a **resampling method**, related to **bootstrap** and **crossvalidation** (of great use in Machine Learning).

Deviation Bounds

Proposition: With probability at least $1 - \delta$ for every P ,

$$\hat{R}_n - \frac{1}{n} - (1 + \sqrt{2})\sqrt{\frac{\log(4/\delta)}{n}} \leq R_n \leq \hat{R}_n + (1 + \sqrt{2})\sqrt{\frac{\log(4/\delta)}{n}}$$

See [McAllester and Schapire '00, McAllester and Ortiz '03]:

- deviations of \hat{R}_n : **McDiarmid's inequality**
- deviations of R_n : **negative association**

Other tool: Poissonization [see Optimal Probability Estimation with Applications to Prediction and Classification, by Acharya, Jafarpour, Orlitsky Suresh, Colt 2013]

Application to Classification: minimax optimality

[Optimal Probability Estimation with Applications to Prediction and Classification, by Acharya, Jafarpour, Orlicy Suresh, Colt 2013]

- P_1, P_2 probability distributions on A
- Given: two samples (X_1^1, \dots, X_n^1) of P_1 and (X_1^2, \dots, X_n^2) of P_2
- Goal: if $I = 1, 2$ with probability $1/2$ and if $X \sim P_I$, build a classifier $\phi_n : A \rightarrow \{1, 2\}$ so that $P(\phi_n(X) = I)$ is as large as possible
- Maximal risk :

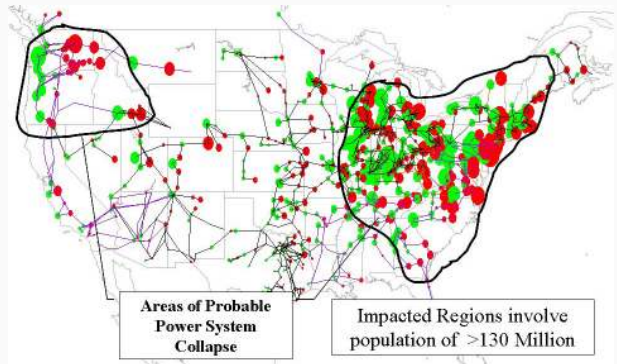
$$\bar{R}_n(\phi) = \max_{P_1, P_2} \mathbb{P}(\phi(X) \neq I)$$

- **Prop:** if $\phi_n^{\text{ML}}(x) = \arg \max_i \#\{j : X_j^i = x\}$ then there exists $c > 0$ such that for all $n \geq 1$, $\bar{R}_n(\phi_n^{\text{ML}}) \geq \min_{\phi} R_n(\phi) + c$.
- **Theorem:** there exists a Good-Turing based classifier ϕ_n^{GT} such that for all $n \geq 1$, $\bar{R}_n(\phi_n^{\text{GT}}) \leq \min_{\phi} R_n(\phi) + O(n^{-1/5})$.

Discovering dangerous contingencies in electrical systems

The problem

Power system security assessment

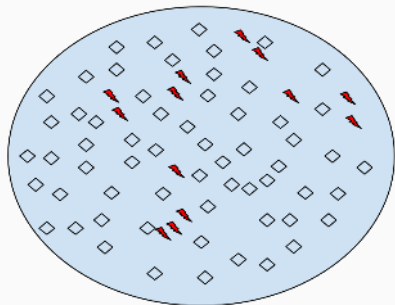


By Mark MacAlester, Federal Emergency Management Agency [Public domain], via Wikimedia Commons

Damien Ernst (Electrical Engineering, Liège): How to **identify quickly contingencies/scenarios** that could lead to unacceptable operating conditions (dangerous contingencies) if no preventive actions were taken?

The model

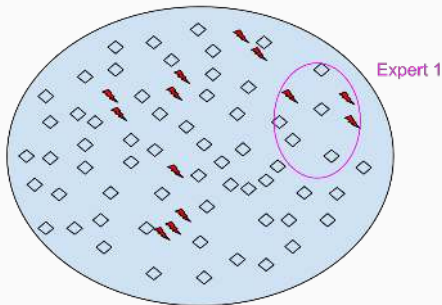
- Subset $A \subset \mathcal{X}$ of important items
- $|\mathcal{X}| \gg 1$, $|A| \ll |\mathcal{X}|$
- Access to \mathcal{X} only by probabilistic experts $(P_i)_{1 \leq i \leq K}$: sequential independent draws



Goal: discover rapidly the elements of A

The model

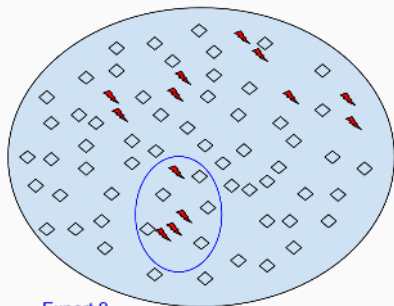
- Subset $A \subset \mathcal{X}$ of important items
- $|\mathcal{X}| \gg 1$, $|A| \ll |\mathcal{X}|$
- Access to \mathcal{X} only by probabilistic experts $(P_i)_{1 \leq i \leq K}$: sequential independent draws



Goal: discover rapidly the elements of A

The model

- Subset $A \subset \mathcal{X}$ of important items
- $|\mathcal{X}| \gg 1$, $|A| \ll |\mathcal{X}|$
- Access to \mathcal{X} only by probabilistic experts $(P_i)_{1 \leq i \leq K}$: sequential independent draws

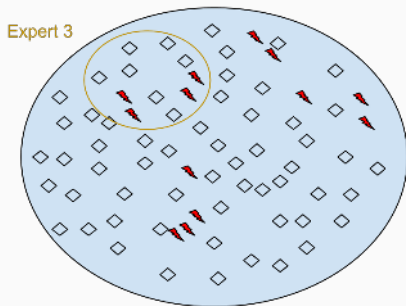


Expert 2

Goal: discover rapidly the elements of A

The model

- Subset $A \subset \mathcal{X}$ of important items
- $|\mathcal{X}| \gg 1$, $|A| \ll |\mathcal{X}|$
- Access to \mathcal{X} only by probabilistic experts $(P_i)_{1 \leq i \leq K}$: sequential independent draws



Goal: discover rapidly the elements of A

At each time step $t = 1, 2, \dots$:

- pick an index $I_t = \pi_t(I_1, Y_1, \dots, I_{s-1}, Y_{s-1}) \in \{1, \dots, K\}$ according to past observations
- observe $Y_t = X_{I_t, n_{I_t, t}} \sim P_{I_t}$, where

$$n_{i,t} = \sum_{s \leq t} \mathbb{1}\{I_s = i\}$$

Goal: design the strategy $\pi = (\pi_t)_t$ so as to **maximize the number of important items found** after t requests

$$F^\pi(t) = \left| A \cap \{Y_1, \dots, Y_t\} \right|$$

Assumption: non-intersecting supports

$$A \cap \text{supp}(P_i) \cap \text{supp}(P_j) = \emptyset \text{ for } i \neq j$$

Is it a Bandit Problem ?

It looks like a bandit problem...

- sequential choices among K options
- want to maximize cumulative rewards
- exploration vs exploitation dilemma

... but it is **not a bandit problem** !

- rewards are not i.i.d.
- **destructive rewards**: no interest to observe twice the same important item
- all strategies eventually equivalent

The oracle strategy

Proposition: Under the non-intersecting support hypothesis, the greedy oracle strategy

$$I_t^* \in \arg \max_{1 \leq i \leq K} P_i(A \setminus \{Y_1, \dots, Y_t\})$$

is optimal: for every possible strategy π , $\mathbb{E}[F^\pi(t)] \leq \mathbb{E}[F^*(t)]$.

Remark: the proposition is false if the supports may intersect

\implies estimate the “*missing mass of important items*”!

The Good-UCB algorithm

Optimal Discovery with Probabilistic Expert Advice: Finite Time Analysis and Macroscopic Optimality

Sébastien Bubeck

Department of Operations Research and Financial Engineering
MIT Operations Research Center
77 Massachusetts Ave
Cambridge, MA 02139, USA

sbubeck@mit.edu

Thaman Ernst

Department of Electrical Engineering and Computer Science
University of Texas at Austin, 78712
3101 Speedway
Austin, TX 78712, USA

ernst@utexas.edu

Aurélien Garivier

UMR 5175 - Laboratoire de Mathématiques
Université Paul Sabatier
118 route de Narbonne
F-31062 Toulouse Cedex 9, France

garivier@math.ups-tlse.fr

Editor: Shiroko Kobayashi

Abstract

We consider an optimal problem where a learner chooses the best of several experts at a point in space and the true state (optimal discovery) will probably be a random vector. We address this in a general framework of independent, parametric and/or the Good-Turner energy mass estimation. We prove novel lower bounds on the performance of the algorithm under such conditions of its probabilistic aspects. These lower bounds are tight, as we also prove a matching upper bound. Finally, we compare the algorithm here with a naive energy and with another existing results as pseudo-regret and regret in the literature. These theoretical findings

Keywords: optimal discovery, multi-armed bandit, stochastic geometry, Good-Turner estimator, UCB



The Good-UCB algorithm

Estimator of the missing important mass for expert i :

$$\hat{R}_{i, n_{i, t-1}} = \frac{1}{n_{i, t-1}} \sum_{x \in A} \mathbb{1} \left\{ \sum_{s=1}^{n_{i, t-1}} \mathbb{1} \{X_{i, s} = x\} = 1 \right.$$
$$\left. \text{and } \sum_{j=1}^K \sum_{s=1}^{n_{j, t-1}} \mathbb{1} \{X_{j, s} = x\} = 1 \right\}$$

Good-UCB algorithm:

- 1: For $1 \leq t \leq K$ choose $I_t = t$.
- 2: **for** $t \geq K + 1$ **do**
- 3: Choose $I_t = \arg \max_{1 \leq i \leq K} \left\{ \hat{R}_{i, n_{i, t-1}} + C \sqrt{\frac{\log(4t)}{n_{i, t-1}}} \right\}$
- 4: Observe Y_t distributed as P_{I_t}
- 5: Update the missing mass estimates accordingly
- 6: **end for**

Optimality results

Theorem: For any $t \geq 1$, under the non-intersecting support assumption, Good-UCB (with constant $C = (1 + \sqrt{2})\sqrt{3}$) satisfies

$$\mathbb{E} [F^*(t) - F^{UCB}(t)] \leq 17\sqrt{Kt \log(t)} + 20\sqrt{Kt} + K + K \log(t/K)$$

Remark: Usual result for bandit problem, but not-so-simple analysis

Sketch of proof

1. On a set $\tilde{\Omega}$ of probability at least $1 - \sqrt{\frac{K}{t}}$, the “confidence intervals” hold true simultaneously all $u \geq \sqrt{Kt}$
2. Let $\bar{l}_u = \arg \max_{1 \leq i \leq K} R_{i, n_{i, u-1}}$. On $\tilde{\Omega}$,

$$R_{l_u, n_{l_u, u-1}} \geq R_{\bar{l}_u, n_{\bar{l}_u, u-1}} - \frac{1}{n_{l_u, u-1}} - 2(1 + \sqrt{2}) \sqrt{\frac{3 \log(4u)}{n_{l_u, u-1}}}$$

3. But one shows that $\mathbb{E} F^*(t) \leq \sum_{u=1}^t \mathbb{E} R_{\bar{l}_u, n_{\bar{l}_u, u-1}}^\pi$

4. Thus

$$\begin{aligned} & \mathbb{E} [F^*(t) - F^{UCB}(t)] \\ & \leq \sqrt{Kt} + \mathbb{E} \left[\sum_{u=1}^t \frac{1}{n_{l_u, u-1}} + 2(1 + \sqrt{2}) \sqrt{\frac{3 \log(4t)}{n_{l_u, u-1}}} \right] \\ & \leq \sqrt{Kt} + K + K \log(t/K) + 4(1 + \sqrt{2}) \sqrt{3Kt \log(4t)} \end{aligned}$$

Experiment: restoring property

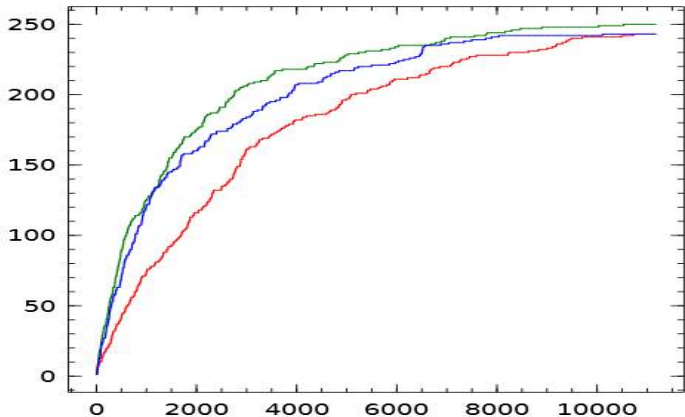


Figure 1: green: oracle, blue: Good-UCB, red: uniform sampling

Another analysis of Good-UCB

For $\lambda \in (0, 1)$, $T(\lambda)$ = time at which missing mass of important items is smaller than λ on all experts:

$$T(\lambda) = \inf \left\{ t : \forall i \in \{1, \dots, K\}, P_i(A \setminus \{Y_1, \dots, Y_t\}) \leq \lambda \right\}$$

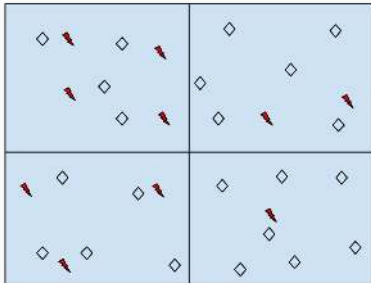
Theorem: Let $c > 0$ and $S \geq 1$. Under the non-intersecting support assumption, for Good-UCB with $C = (1 + \sqrt{2})\sqrt{c+2}$, with probability at least $1 - \frac{K}{cS^c}$, for any $\lambda \in (0, 1)$,

$$T_{UCB}(\lambda) \leq T^* + KS \log(8T^* + 16KS \log(KS)),$$

$$\text{where } T^* = T^* \left(\lambda - \frac{3}{S} - 2(1 + \sqrt{2})\sqrt{\frac{c+2}{S}} \right)$$

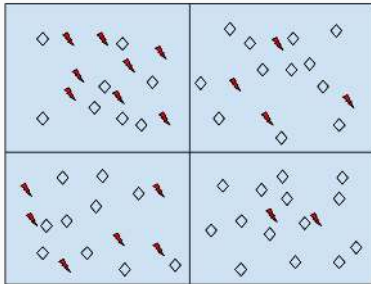
The macroscopic limit

- Restricted framework: $P_i = \mathcal{U}\{1, \dots, N\}$
- $N \rightarrow \infty$
- $|A \cap \text{supp}(P_i)|/N \rightarrow q_i \in (0, 1)$, $q = \sum_i q_i$



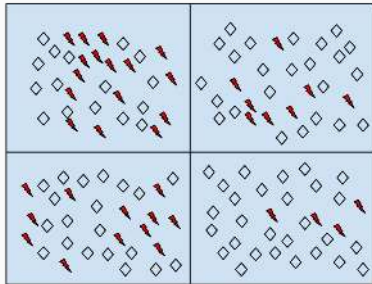
The macroscopic limit

- Restricted framework: $P_i = \mathcal{U}\{1, \dots, N\}$
- $N \rightarrow \infty$
- $|A \cap \text{supp}(P_i)|/N \rightarrow q_i \in (0, 1)$, $q = \sum_i q_i$



The macroscopic limit

- Restricted framework: $P_i = \mathcal{U}\{1, \dots, N\}$
- $N \rightarrow \infty$
- $|A \cap \text{supp}(P_i)|/N \rightarrow q_i \in (0, 1)$, $q = \sum_i q_i$



The limiting discovery process of the Oracle strategy is *deterministic*

Proposition: For every $\lambda \in (0, q_1)$, for every sequence $(\lambda^N)_N$ converging to λ as N goes to infinity, almost surely

$$\lim_{N \rightarrow \infty} \frac{T_*^N(\lambda^N)}{N} = \sum_i \left(\log \frac{q_i}{\lambda} \right)_+$$

Oracle vs. uniform sampling

Oracle: The proportion of important items not found after Nt draws tends to

$$q - F^*(t) = I(t) \underline{q}_{I(t)} \exp(-t/I(t)) \leq K \underline{q}_K \exp(-t/K)$$

with $\underline{q}_K = \left(\prod_{i=1}^K q_i\right)^{1/K}$ the geometric mean of the $(q_i)_i$.

Uniform: The proportion of important items not found after Nt draws tends to $K \bar{q}_K \exp(-t/K)$

⇒ Asymptotic ratio of efficiency

$$\rho(q) = \frac{\bar{q}_K}{\underline{q}_K} = \frac{\frac{1}{K} \sum_{i=1}^K q_i}{\left(\prod_{i=1}^K q_i\right)^{1/K}} \geq 1$$

larger if the $(q_i)_i$ are unbalanced

Theorem: Take $C = (1 + \sqrt{2})\sqrt{c + 2}$ with $c > 3/2$ in the Good-UCB algorithm.

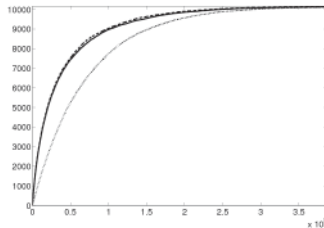
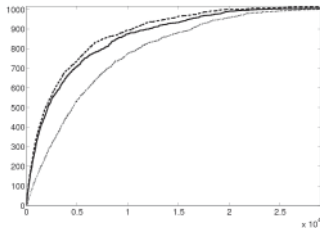
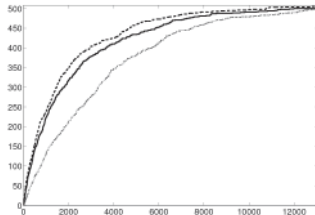
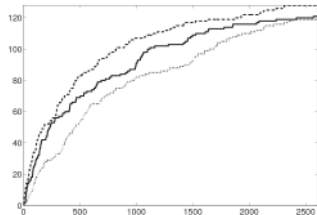
- For every sequence $(\lambda^N)_N$ converging to λ as N goes to infinity, almost surely

$$\limsup_{N \rightarrow +\infty} \frac{T_{UCB}^N(\lambda^N)}{N} \leq \sum_i \left(\log \frac{q_i}{\lambda} \right)_+$$

- The proportion of items found after Nt steps $F^{GUCB}(Nt)$ converges *uniformly* to $F^*(Nt)$ as N goes to infinity

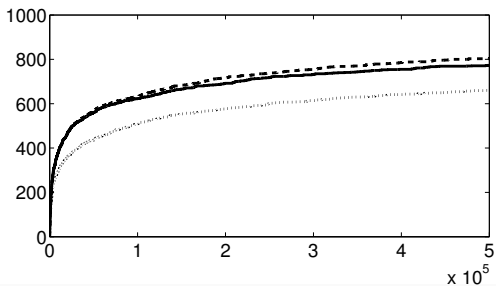
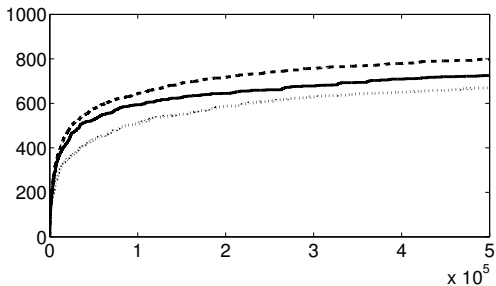
Experiment

Number of items found by Good-UCB (solid), the OCL (dashed), and uniform sampling (dotted) as a function of time for sizes $N = 128$, $N = 500$, $N = 1000$ and $N = 10000$ in a 7-experts setting.



And when the assumptions are not satisfied?

Number of primes found by **Good-UCB** (solid), the **oracle** (dashed) and **uniform** sampling (dotted) using geometric experts with means 100, 300, 500, 700, 900, for $C = 0.1$ (top) and $C = 0.02$ (bottom).



Conclusion and perspectives

- We propose an algorithm for the optimal discovery with probabilistic expert advice
- We give a standard regret analysis under the only assumption that the supports of the experts are non-overlapping
- We propose a different optimality result, which permits a macroscopic analysis in the uniform case
- Another interesting limit to consider is when the number of important items to find is fixed, but the total number of items tends to infinity (Poisson regime)
- Then, the behavior of the algorithm is not very good: too large confidence bonus because no tight deviations bounds for the Good-Turing estimator when the proportion of important items tends to 0. Improvement by better deviation bounds?

Thank you for your attention!