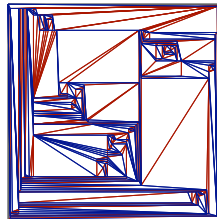
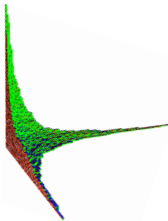
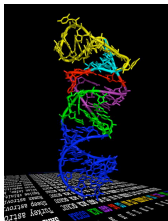


Effective Boltzmann sampling

Carine Pivoteau

based on work by P. Duchon, P. Flajolet, E. Fusy,
G. Louchard, C. Pivoteau and G. Schaeffer



Plan of the talk

- 1 Boltzmann model
- 2 Algorithms
- 3 Technical matters

Context

Random sampling under Boltzmann model

- approximate size sampling,
 - size distribution spread over the whole combinatorial class, but uniform for a sub-class of objects of the same size,
 - control parameter,
 - automatized sampling: the sampler is compiled from specification automatically,
 - very large objects can be sampled.
1. *Boltzmann samplers for the random generation of combinatorial structures.* P. Duchon, P. Flajolet, G. Louchard, G. Schaeffer. Combinatorics, Probability and Computing, 13(4-5):577-625, 2004. Special issue on Analysis of Algorithms.
 2. *Boltzmann sampling of unlabelled structures.* Ph. Flajolet, E. Fusy, C. Pivoteau. Proceedings of ANALCO07, january 2007.

Boltzmann model

Framework: constructible classes

| | specification | ordinary g.f. (unlabelled) | exponential g.f. (labelled) |
|----------------------|--|--|--|
| ε / atom | 1 / \mathcal{Z} | 1 / x | 1 / x |
| Union | $\mathcal{C} = \mathcal{A} \cup \mathcal{B}$ | $C(x) = A(x) + B(x)$ | $\hat{C}(x) = A(x) + B(x)$ |
| Product | $\mathcal{C} = \mathcal{A} \times \mathcal{B}$ | $C(x) = A(x) \times B(x)$ | $\hat{C}(x) = A(x) \times B(x)$ |
| Sequence | $\mathcal{C} = \text{SEQ}(\mathcal{A})$ | $C(x) = \frac{1}{1-A(x)}$ | $\hat{C}(x) = \frac{1}{1-A(x)}$ |
| Set | $\mathcal{C} = \text{SET}(\mathcal{A})$ | $\exp\left(\sum_{k=1}^{\infty} \frac{(-1)^{k-1}}{k} A(z^k)\right)$ | $\hat{C}(x) = \exp(A(x))$ |
| Multiset | $\mathcal{C} = \text{MSET}(\mathcal{A})$ | $\exp\left(\sum_{k=1}^{\infty} \frac{1}{k} A(z^k)\right)$ | — |
| Cycle | $\mathcal{C} = \text{CYC}(\mathcal{A})$ | $\sum_{k=1}^{\infty} \frac{1}{k} \log \frac{1}{1-A(z^k)}$ | $\hat{C}(x) = \log \frac{1}{1-A(z^k)}$ |

Model definition

Definition

In the **unlabelled** case, Boltzmann model assigns to any object $c \in \mathcal{C}$ the following probability:

$$\mathbb{P}_x(c) = \frac{x^{|c|}}{C(x)}$$

In the **labelled** case, this probability becomes:

$$\mathbb{P}_x(c) = \frac{1}{\hat{C}(x)} \frac{x^{|c|}}{|c|!}$$

A Boltzmann sampler $\Gamma C(x)$ for the class \mathcal{C} is a process that produces objects from \mathcal{C} according to this model.

→ 2 object of the same size will be drawn with the same probability.

Size distribution

The probability of drawing an object of size N is then:

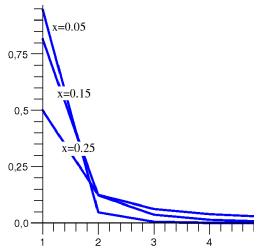
$$\mathbb{P}_x(N = n) = \sum_{|c|=n} \mathbb{P}_x(c) = \frac{C_n x^n}{C(x)} \quad \text{or} \quad \frac{C_n x^n}{n! \hat{C}(x)}$$

Then, the expected size of an object drawn by a generator with parameter x is:

$$\mathbb{E}_x(N) = x \frac{C'(x)}{C(x)} \quad \text{or} \quad x \frac{\hat{C}'(x)}{\hat{C}(x)}$$

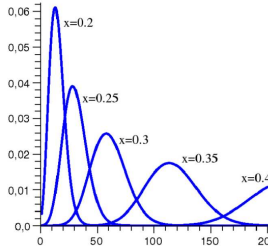
Size distribution – 2

General plane trees



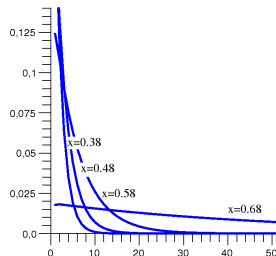
Peaked

Set partitions



Bumpy

Surjections



Flat

Algorithms

1 Boltzmann model

2 Algorithms

- Basic unlabelled constructions
- Labelled classes
- Back to unlabelled

3 Technical matters

Unions, products

Disjoint unions

Boltzmann sampler ΓC for $\mathcal{C} = \mathcal{A} \cup \mathcal{B}$:

With probability $\frac{A(x)}{C(x)}$ do $\Gamma A(x)$ else do $\Gamma B(x)$ \rightarrow Bernoulli.

$$\text{if } \gamma \in \mathcal{A} \text{ then } \mathbb{P}_x(\gamma) = \frac{x^{|\gamma|}}{A(x)} \cdot \frac{A(x)}{C(x)} = \frac{x^{|\gamma|}}{C(x)}.$$

Products

Boltzmann sampler ΓC for $\mathcal{C} = \mathcal{A} \times \mathcal{B}$:

Generate a pair $\langle \Gamma A(x), \Gamma B(x) \rangle \rightarrow$ independent calls.

$$\text{if } \gamma = \alpha \times \beta \in \mathcal{C} \text{ then } \mathbb{P}_x(\gamma) = \frac{x^{|\alpha|}}{A(x)} \cdot \frac{x^{|\beta|}}{B(x)} = \frac{x^{|\gamma|}}{C(x)}.$$

Remark: $A(x)$, $B(x)$ and $C(x)$ are given by an *oracle*.

Sequences

Sequences

Boltzmann sampler ΓC for $\mathcal{C} = \text{SEQ}(\mathcal{A})$:

Generate k according to a **geometric law** of parameter $A(x)$

Generate a k -tuple $\langle \Gamma A(x), \dots, \Gamma A(x) \rangle \rightarrow$ **independent** calls.

Proof.

Recursive equation: $\mathcal{C} = \mathbf{1} + \mathcal{AC}$ with $+$, \times constructions.

With proba. $A(x)$, call $\Gamma A(x)$ and call recursively $\Gamma C(x)$; else stop.

Number of successful trials of $\text{Bern}(A(x))$ is $\text{Geom}(A(x))$.

Remark: $A(x)$ is given by an **oracle**.

Examples of specifications with $\{\cup, \times, \text{Seq}\}$

Regular specifications (non recursive).

- integer **compositions**,
- **polyominoes** that have rational g.f.: column-convex,



- **regular languages**,
- binary **words** without long runs:

$$\mathcal{L} = \text{SEQ}_m(b) \times \text{SEQ}(a \times \text{SEQ}_m(a) \times b \times \text{SEQ}_m(b)) \times \text{SEQ}_m(a)$$

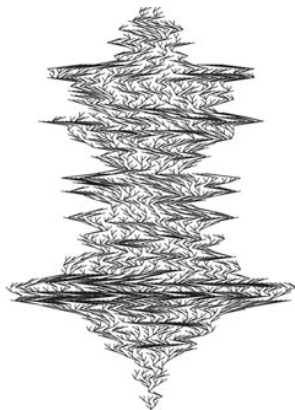
where $\text{SEQ}_m(a) = a + aa + \dots + a^m$.

- ...

Specifications with $\{\cup, \times, \text{Seq}\}$ and recursion

Context-free specifications.

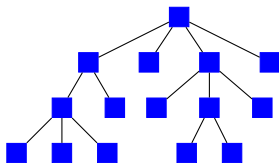
- binary, k -ary, 2–3–4 trees, ...,
- triangulations,
- general planar rooted trees,
- noncrossing graphs,
- ...



General planar rooted trees

$$\mathcal{T} = \mathcal{Z} \times \text{SEQ}(\mathcal{T})$$

$$T(z) = \frac{z}{1 - T(z)} = \frac{1 - \sqrt{1 - 4z}}{2}$$



Algorithm: $\Gamma B(x)$

$k \leftarrow \text{Geom}(T(x));$

$children \leftarrow \text{the } k\text{-tuple } \langle \Gamma T(x), \dots, \Gamma T(x) \rangle;$

Return $\blacksquare \times children$;

1 Boltzmann model

2 Algorithms

- Basic unlabelled constructions
- Labelled classes
- Back to unlabelled

3 Technical matters

Labelled classes

Same algorithms, with exponential g.f.

| construction | sampler |
|---|--|
| $\mathcal{C} = \emptyset \text{ or } \mathcal{Z}$ | $\Gamma C(x) := \varepsilon \text{ or atom}$ |
| $\mathcal{C} = \mathcal{A} + \mathcal{B}$ | $\Gamma C(x) := \text{Bern } \frac{\hat{A}(x)}{\hat{C}(x)} \longrightarrow \Gamma A(x) \mid \Gamma B(x)$ |
| $\mathcal{C} = \mathcal{A} \times \mathcal{B}$ | $\Gamma C(x) := \langle \Gamma A(x) ; \Gamma B(x) \rangle$ |
| $\mathcal{C} = \text{SEQ}(\mathcal{A})$ | $\Gamma C(x) := \text{Geom } \hat{A}(x) \Longrightarrow \Gamma A(x)$ |

Put the labels at the end !

Labelled sets and cycles

Sets

Boltzmann sampler ΓC for $\mathcal{C} = \text{SET}(\mathcal{A})$:

Generate k according to a **Poisson law** of parameter $A(x)$

Generate a k -tuple $\langle \Gamma A(x), \dots, \Gamma A(x) \rangle$

Poisson law: $\mathbb{P}(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}$

Cycles

Boltzmann sampler ΓC for $\mathcal{C} = \text{CYC}(\mathcal{A})$:

Generate k according to a **logarithmic law** of parameter $A(x)$

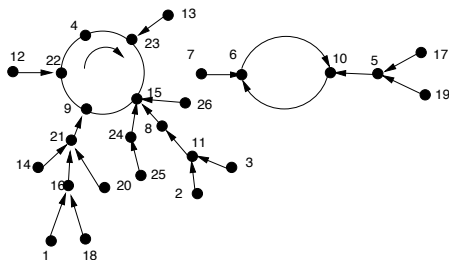
Generate a k -tuple $\langle \Gamma A(x), \dots, \Gamma A(x) \rangle$

Logarithmic law: $\mathbb{P}(X = k) = \frac{1}{\log(1 - \lambda)^{-1}} \frac{\lambda^k}{k}$

Remark: the laws are given by simple sequential algorithms

Examples of possible labelled classes

- permutations, derangements, involutions,
- surjections,
- set partitions,
- necklaces,
- labelled (planar) trees,
- functional graphs,
- ...



Theorem (Labelled free Boltzmann samplers)

For a labelled class \mathcal{C} specified (poss. recursively) using the following constructions:

$$\varepsilon, \mathcal{Z}, +, \times, \text{SEQ}, \text{SET}, \text{CYC}$$

*The **free Boltzmann sampler** $\Gamma_{\mathcal{C}}(x)$ operates in **linear time** in the size of the object produced.*

- oracle complexity is not involved,
- size is not controlled (yet).

1 Boltzmann model

2 Algorithms

- Basic unlabelled constructions
- Labelled classes
- Back to unlabelled

3 Technical matters

To begin: MSet_2

$\text{MSet}_2(\mathcal{A}) \cong$ unordered set of **two** objects of \mathcal{A}

$$\begin{aligned} \mathcal{C} &= \text{MSet}_2(\mathcal{A}) \\ C(z) &= \frac{1}{2}A^2(z) + \frac{1}{2}A(z^2) \quad \left[2\text{MSet}_2(\mathcal{A}) = \mathcal{A}^2 + \Delta\mathcal{A}^2 \right] \end{aligned}$$

Algorithm: $\Gamma C(x)$

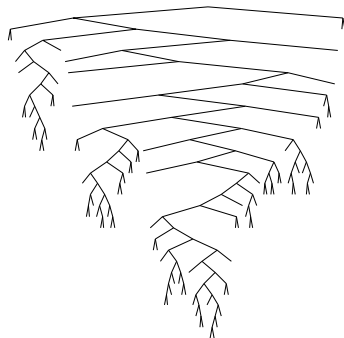
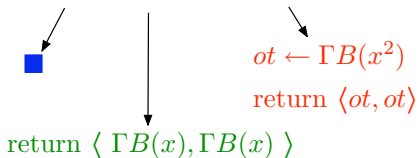
```
if  $\text{Bern}\left(\frac{1}{2} \frac{A^2(x)}{C(x)}\right) = 1$  then
  Return  $\langle \Gamma A(x), \Gamma A(x) \rangle$ 
else
   $a \leftarrow \Gamma A(x^2)$ ;
  Return  $\langle a, a \rangle$ ;
end if
```

A simple example

Unlabelled binary trees
(Otter trees)

$$\mathcal{B} = \mathcal{Z} + \text{MSET}_2(\mathcal{B})$$

$$B(z) = z + \frac{1}{2}B(z)^2 + \frac{1}{2}B(z^2)$$



MSet: the general case

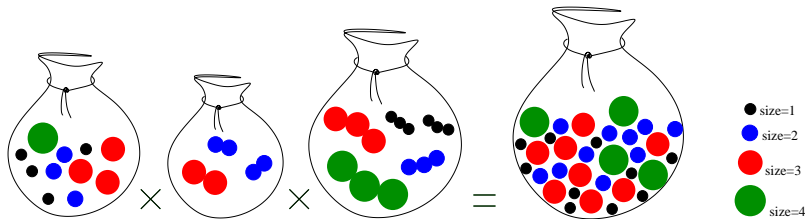
$$\mathcal{M} = \text{MSET}(\mathcal{A}) \cong \prod_{\gamma \in \mathcal{A}} \text{SEQ}(\gamma) \Rightarrow M(z) = \prod_{\gamma \in \mathcal{A}} (1 - z^{|\gamma|})^{-1}$$

$$M(z) = \exp \left(\sum_{k=1}^{\infty} \frac{1}{k} A(z^k) \right) = \prod_{k=1}^{\infty} \exp \left(\frac{1}{k} A(z^k) \right)$$

MSet: the general case

$$\mathcal{M} = \text{MSET}(\mathcal{A}) \cong \prod_{\gamma \in \mathcal{A}} \text{SEQ}(\gamma) \Rightarrow M(z) = \prod_{\gamma \in \mathcal{A}} (1 - z^{|\gamma|})^{-1}$$

$$M(z) = \exp \left(\sum_{k=1}^{\infty} \frac{1}{k} A(z^k) \right) = \prod_{k=1}^{\infty} \exp \left(\frac{1}{k} A(z^k) \right)$$



MSet

Algorithm $\Gamma MSet[\mathcal{A}](x)$

- Draw k , the **max. index** of a subset, depending on x ;
- For each index i of a subset until $k - 1$
 - Draw the **number p of elements to sample**, according to a Poisson law of parameter $\frac{1}{i}A(x^i)$.
 - Call $\Gamma A(x^i)$ p times, and each time, add **i copies** of the result to the multiset.
- for index k , draw the number p of elements to generate, according to a **non zero** Poisson law.

Proof: begin with the product sampler and use $\text{Geom}(\lambda) \equiv \sum_k \text{Pois}(\frac{\lambda^k}{k})$.

How to compute the max. index for $\mathcal{C} = \text{MSet}(\mathcal{A})$

Draw U , uniformly at random in $[0, 1]$ and find k such that:

$$\frac{\exp\left(\sum_{i=1}^{k-1} \frac{1}{i} A(x^i)\right)}{C(x)} < U \leq \frac{\exp\left(\sum_{i=1}^k \frac{1}{i} A(x^i)\right)}{C(x)}$$

with i.e.:

$$\sum_{i=k+1}^{\infty} \frac{1}{i} A(x^i) < \log \frac{1}{U} \leq \sum_{i=k}^{\infty} \frac{1}{i} A(x^i)$$

Algorithm

$U \leftarrow \text{rand}(0, 1); V \leftarrow \log \frac{1}{U};$

$p \leftarrow \log C(x); k \leftarrow 0;$

while $V < p$ **do**

$k \leftarrow k + 1; p \leftarrow p - \frac{A(x^k)}{k};$

Return $k;$

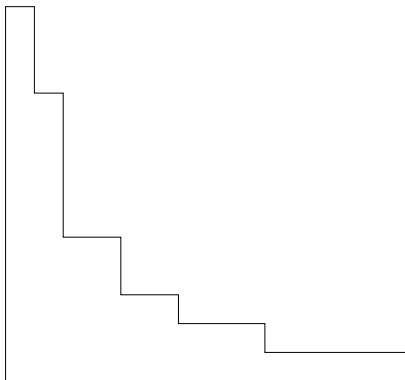
Sampling of a partition $\text{MSet}(\mathcal{Z} \times \text{Seq}(\mathcal{Z}))$

$k = 3;$

$p_1 = 5;$

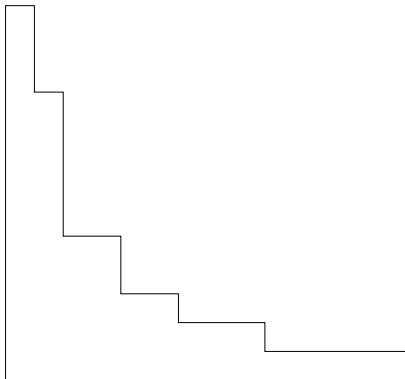
$p_2 = 3;$

$p_3 = 1;$



Sampling of a partition $\text{MSet}(\mathcal{Z} \times \text{Seq}(\mathcal{Z}))$

$k = 3;$
 $p_1 = 5;$
 $p_2 = 3;$
 $p_3 = 1;$



Sampling of a partition $\text{MSet}(\mathcal{Z} \times \text{Seq}(\mathcal{Z}))$

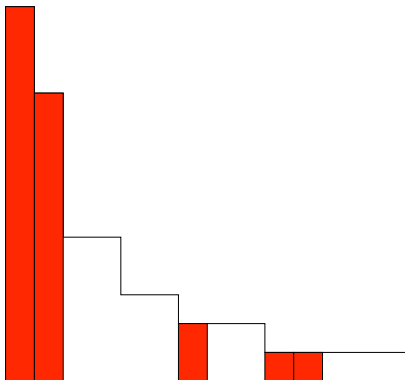
$k = 3;$

$p_1 = 5;$

$p_2 = 3;$

$p_3 = 1;$

$\rightarrow 13, 10, 2, 1, 1$



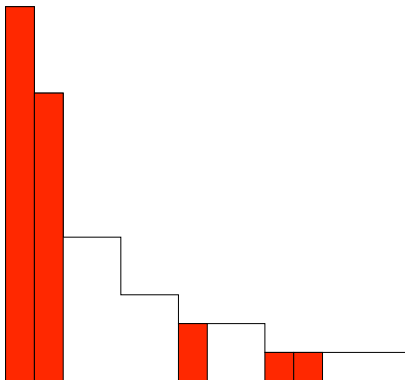
Sampling of a partition $\text{MSet}(\mathcal{Z} \times \text{Seq}(\mathcal{Z}))$

$k = 3;$

$p_1 = 5;$

$p_2 = 3;$

$p_3 = 1;$



Sampling of a partition $\text{MSet}(\mathcal{Z} \times \text{Seq}(\mathcal{Z}))$

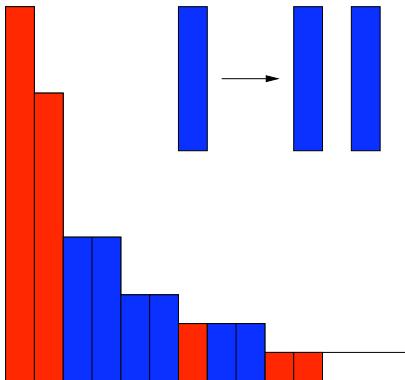
$k = 3;$

$p_1 = 5;$

$p_2 = 3;$

$p_3 = 1;$

$\rightarrow 5, 3, 2$



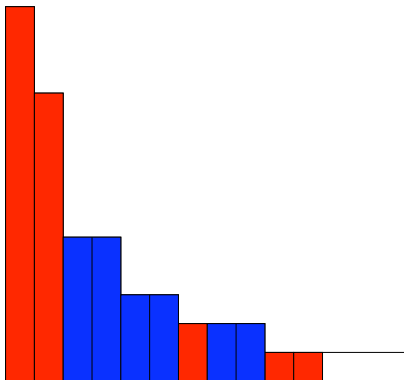
Sampling of a partition $\text{MSet}(\mathcal{Z} \times \text{Seq}(\mathcal{Z}))$

$k = 3$;

$p_1 = 5$;

$p_2 = 3$;

$p_3 = 1$;



Sampling of a partition $\text{MSet}(\mathcal{Z} \times \text{Seq}(\mathcal{Z}))$

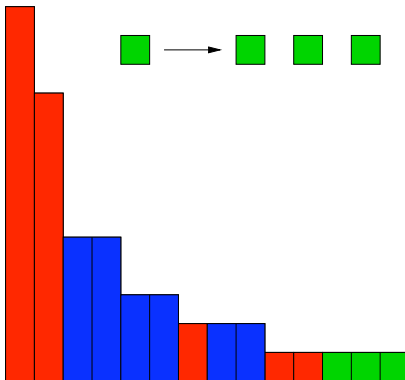
$k = 3;$

$p_1 = 5;$

$p_2 = 3;$

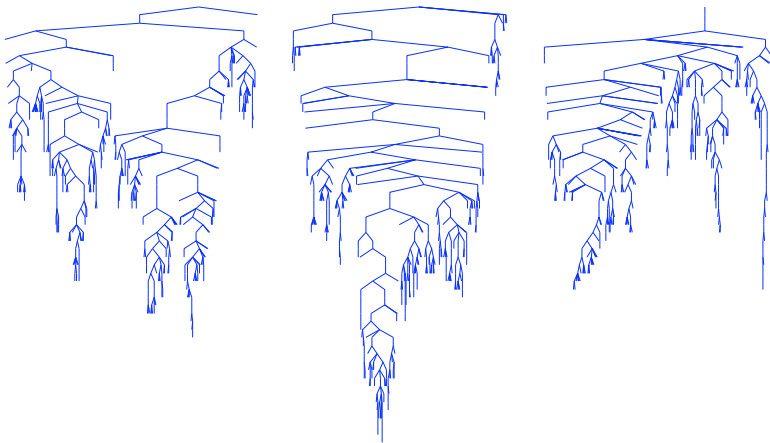
$p_3 = 1;$

$\rightarrow 1$



Cayley trees

$$\mathcal{T} = \mathcal{Z} \times \text{MSET}(\mathcal{T})$$



From MSet to PSet

Principle : Use the following non ambiguous decomposition:

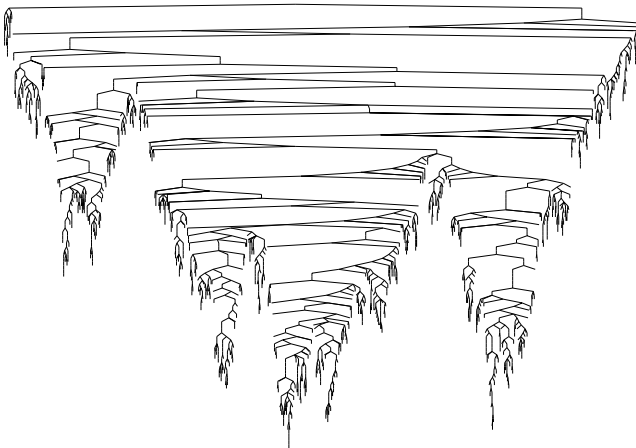
$$\begin{aligned} \text{MSET}(\mathcal{A}) &= \text{PSET}(\mathcal{A}) \times \text{MSET}(\mathcal{A}^{(2)}) \\ \prod_{\gamma} \frac{1}{1-z^{|\gamma|}} &= \prod_{\gamma} \frac{1+z^{|\gamma|}}{1-z^{2|\gamma|}} = \prod_{\gamma} (1+z^{|\gamma|}) \prod_{\gamma} \frac{1}{1-z^{2|\gamma|}} \end{aligned}$$

The algo. $\Gamma \text{PSet}[\mathcal{A}](x)$ to sample a powerset of objects of \mathcal{A} is:

- Sample a multiset with $\Gamma \text{MSet}[\mathcal{A}](x)$,
- Extract the **corresponding powerset** :
 - by removing objects with even multiplicity,
 - and keeping **only one occurrence** of objects with **odd multiplicity**.

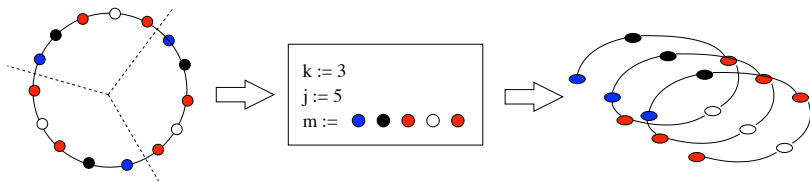
Trees without twins

$$\mathcal{T} = \mathcal{Z} \times \text{PSet}(\mathcal{T})$$



Cycles

$$\mathcal{C} = \text{Cyc}(\mathcal{A}) \Rightarrow C(z) = \sum_{k \geq 1} \frac{\varphi(k)}{k} \log \frac{1}{1 - A(z^k)}$$



Algorithm

$\Gamma\text{Cyc}[\mathcal{A}](x)$

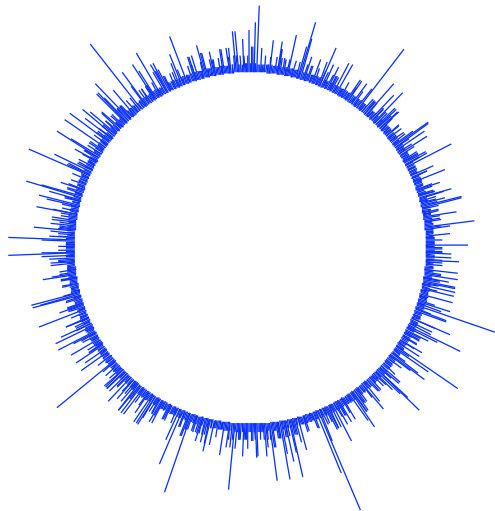
- Draw the replication order k of the cycle.
- Draw the length j of the pattern according to a logarithmic law of parameter $A(x^k)$.
- Draw the pattern m , calling $\Gamma A(x^k)$ j times.
- Return a cycle composed of k copies of m .

Proof: find $\mathbb{P}_x(\gamma)$ for any $\gamma \in \text{Cyc}(\mathcal{A})$, using $\mathbb{P}(k, j) = \frac{\varphi(k)}{kj} \frac{A(x^k)^j}{C(x)}$ and $\sum_{k|m} \varphi(k) = m$.

Cyclic compositions

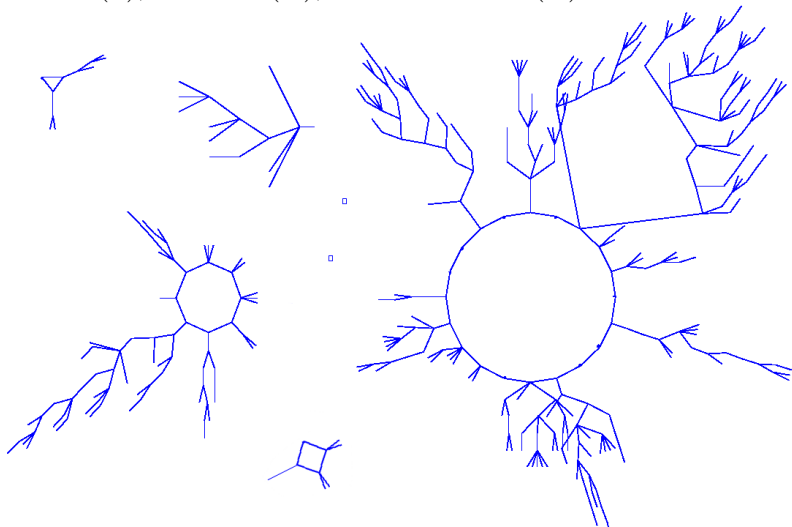
$$\mathcal{C} = \text{CYC}(\mathcal{Z} \times \text{SEQ}(\mathcal{Z}))$$

$$C(z) = \sum_{k=1}^{\infty} \frac{\varphi(k)}{k} \log \frac{1}{1 - \frac{z^k}{1-z^k}}$$



Functional graphs

$$\mathcal{G} = \text{SET}(\mathcal{C}), \mathcal{C} = \text{CYC}(\mathcal{T}), \mathcal{T} = \mathcal{Z} \times \text{MSET}(\mathcal{T})$$



MSet_k(\mathcal{A}) et Cyc_k(\mathcal{A})

Samplers are deduced from the generating functions:

- MSET_k(\mathcal{A})

$$M_k(z) = [u^k] \exp \left(uA(z) + \frac{u^2}{2} A(z^2) + \frac{u^3}{3} A(z^3) + \dots \right)$$

Example: $M_3(z) = \frac{1}{6}A(z)^3 + \frac{1}{2}A(z)A(z^2) + \frac{1}{3}A(z^3)$

- CYC_k(\mathcal{A})

$$C_k(z) = [u^k] \sum_{l=1}^{\infty} \frac{\varphi(l)}{l} \log \frac{1}{1 - u^l A(z^l)}$$

Example: $C_4(z) = \frac{1}{4}A(z)^4 + \frac{1}{4}A(z^2)^2 + \frac{1}{2}A(z^4)$

Other constraints

- $\text{MSET}_{<k}$, $\text{PSET}_{<k}$, $\text{CYC}_{<k}$

Union over $i < k$ of sets or cycles with i elements.

Example: $\text{CYC}_{<k}(\mathcal{A}) = \bigcup_{i=1}^{k-1} \text{CYC}_i(\mathcal{A})$

- PSET_k , $\text{MSET}_{>k}$, $\text{PSET}_{>k}$, $\text{CYC}_{>k}$

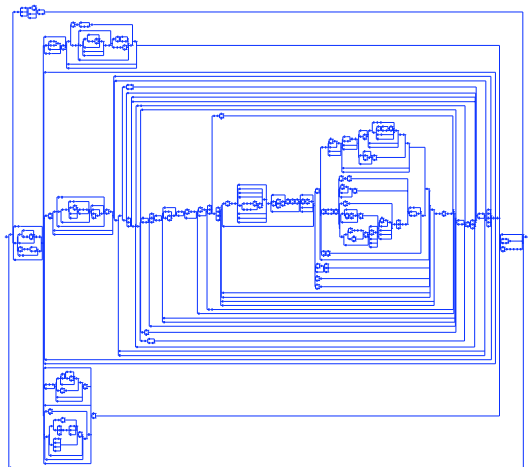
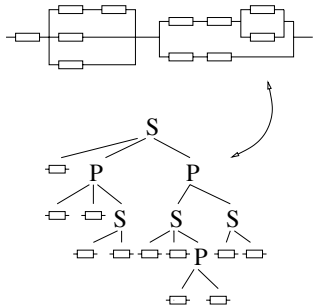
Rejection.

Series-parallel circuits

$$\mathcal{C} = \mathcal{P} + \mathcal{S} + \mathcal{Z}$$

$$\mathcal{S} = \text{SEQ}_{\geq 2}(\mathcal{P} + \mathcal{Z})$$

$$\mathcal{P} = \text{MSET}_{\geq 2}(\mathcal{S} + \mathcal{Z})$$



Theorem (Unlabelled free Boltzmann samplers)

For an unlabelled class \mathcal{C} specified (poss. recursively) using the following constructions:

$$\varepsilon, \mathcal{Z}, +, \times, \text{SEQ}, \text{SEQ}_k, \text{MSET}, \text{MSET}_k, \text{CYC}, \text{CYC}_k$$

The free Boltzmann sampler $\Gamma_{\mathcal{C}}(x)$ operates in linear time in the size of the object produced.

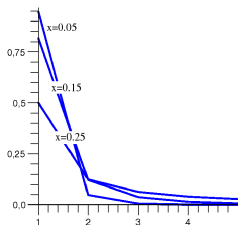
SET: not so bad!

if $\rho < 1$ then the *overhead* (total size of the discarded elements) is bounded by a constant.

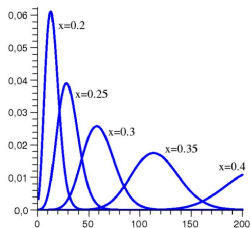
Technical matters

Size control – parameter tuning

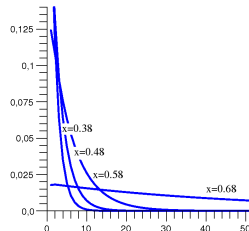
- Free samplers: produce objects with randomly varying sizes!
- Tuned samplers: choose x so that expected size is n .
- Size distribution of free sampler determines complexity.



Peaked



Bumpy



Flat

Size control – 2

Cost of tuned Boltzmann samplers for a labelled or unlabelled class \mathcal{C} specified (poss. recursively) using the following constructions:

ε , \mathcal{Z} , $+$, \times , SEQ, SEQ_k, MSET, MSET_k, CYC, CYC_k

Theorem (Bumpy type)

Approximate-size complexity = $\mathcal{O}(n)$. *Exact size* = $o(n^2)$.

Theorem (Flat type)

Approximate-size complexity = $\mathcal{O}(n)$. *Exact-size* = $o(n^2)$.

Theorem (Peaked type)

Approximate-size complexity = $\mathcal{O}(n^2)$ \rightarrow *pointing*.

Pointing

If \mathcal{A} is a class, then $\mathcal{C} = \mathcal{A}^\bullet$ is the set of objects with one atom pointed, and

$$C_n = nA_n, \quad C(z) = z \frac{d}{dz} A(z).$$

Uniformity at given size is preserved, but size profile is altered.

Transforms peaked (inefficient) to flat (efficient).

For example, **binary trees** \mathcal{B} :

$$\mathcal{B} = \mathcal{Z} + \mathcal{B} \times \mathcal{B} \quad \Longrightarrow \quad \begin{cases} \mathcal{B} &= \mathcal{Z} + \mathcal{B} \times \mathcal{B} \\ \mathcal{B}^\bullet &= \mathcal{Z}^\bullet + \mathcal{B}^\bullet \times \mathcal{B} + \mathcal{B} \times \mathcal{B}^\bullet. \end{cases}$$

It works for all simple families of trees.

Singular samplers

- critical sequences



exact-size: $\mathcal{O}(n)$.

- regular languages

if the automaton recognizing \mathcal{L} is strongly connected:

exact-size: $\mathcal{O}(n)$.

- trees

with early interrupted execution:

approx-size: $\mathcal{O}(n)$, exact-size: $\mathcal{O}(n^2)$.

Concluding remarks

Already done

- BaNi06 *Accessible and deterministic automata: enumeration and Boltzmann samplers*, by F. Bassino C. Nicaud. In *Fourth Colloquium on Mathematics and Computer Science*.
- BoFuPi06 *Random sampling of plane partitions*, by O. Bodini, E. Fusy, and C. Pivoteau. In *GASCOM-2006*.
- DaSo07 *Degree distribution of random Apollonian network structures and Boltzmann sampling*, by A. Darrasse and M. Soria. In *International Conference on Analysis of Algorithms*, 2007, *DIMACS*.
- Fusy05 *Quadratic exact-size and linear approximate-size random sampling of planar graphs*, by E. Fusy. In *International Conference on Analysis of Algorithms*, 2005, *DMTCS Conference Volume AD (2005)*, pp. 125-138.
- PaWe07 *Properties of Random Graphs via Boltzmann Samplers*, by K. Panagiotou and A. Weiß. In *International Conference on Analysis of Algorithms*, 2007, *DIMACS*.
- Ponty06 *Modélisation de séquences génomiques structurées, génération aléatoire et application*, by Yann Ponty, PhD Thesis, Université Paris-Sud, 2006.

...

Coming soon...?

- other constructions: box operator, shuffle, ...
- multivariate Boltzmann samplers,
- automatic oracle, singularities,
- discrete samplers,
- specialized samplers,
- new applications,
- ...