## Ewens-like distributions and Analysis of Algorithms

Nicolas Auger, Mathilde Bouvel, Cyril Nicaud, Carine Pivoteau

July 6, 2016

## Notion of presortedness

- In practice, data are often presorted.
  - No reasons to be uniformly distributed.
  - Few alterations in databases.
- First intuition in [Knuth73] and formalized in [Mannila86].





- In practice :
  - Used in standard libraries
- 🍨 python"
- Java's developers benchmarks, using spies
- TimSort

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2 If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \le |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2) If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \le |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2 If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \le |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2 If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \leq |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2 If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \leq |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2) If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \le |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2 If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \leq |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2 If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \leq |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

Let  $X = (x_1, ..., x_n)$  and  $Y = (y_1, ..., y_\ell)$  two sequences of elements from a set E;  $m : E^+ \to \mathbb{Z}^+$  is a **measure of presortedness** iff

• 
$$m(X) = 0$$
 if X is sorted.

2 If 
$$n = \ell$$
 and  $x_i < x_j \iff y_i < y_j$ , then  $m(X) = m(Y)$ .

**3** If Y is a subsequence of X, then  $m(Y) \le m(X)$ .

• If 
$$X < Y$$
, then  $m(XY) \le m(X) + m(Y)$ .

So For any element a,  $m(aX) \leq |X| + m(X)$ .

- number of Runs -1, Runs(41536827) = 4
- number of Inversions, Inv(41536827) = 9

# Adaptiveness of sorting algorithms

#### Theorem

Let X be a sequence s.t. m(X) = k. Any algorithm uses at least C(n, k) comparisons to sort X, with  $C(n, k) \in \Theta(n + \log(\|below_m(n, k)\|)$  and  $below_m(n, k) = \{\sigma \in \mathfrak{S}_n : m(\sigma) \le k\}.$ 

### Definition

A sorting algorithm is **m-optimal** if it reaches this bound.

- Natural Merge Sort
  [Knuth73]
- $\mathcal{O}(n \log r)$ , where *r* is the number of runs
- Runs-optimal



## Records as a measure of presortedness

Let  $X = (x_1, ..., x_n)$  be a sequence;  $x_i$  is a **record**(left-to-right maximum) iff  $x_j < x_i$  whenever j < i.

Lemma

For any sequence X of size n,  $m_{rec}(X) = n - record(X)$  is a measure of presortedness.

Example : For X = 32418567, record(X) = 3 and  $m_{rec}(X) = 5$ .



Complexity  $\mathcal{O}(n + k \log k)$ 

 $\|below_{m_{rec}}(n,k)\| \ge k!$ 

Under the uniform distribution, for most measures m:

- $\|below_m(n, \mathbb{E}[m])\| = \Theta(n!).$
- $\mathcal{O}(n \log n)$  in average.

### Questions

- How to define a probabilistic framework well-suited for presortedness measures ?
- Analysis of algorithms ?

Any permutation can be seen as a composition of cycles. Example : 145263 is composed of 3 cycles : (1), (563) and (42).

We denote  $cycle(\sigma)$  the number of cycles of  $\sigma$ .

### Definition (Ewens distribution)

[Ewens72]

To any σ ∈ 𝔅<sub>n</sub>, we associate a weight w(σ) = θ<sup>cycle(σ)</sup>, where θ is an arbitrary positive real number.

• Total weight : 
$$\sum_{\sigma\in\mathfrak{S}_n}\mathsf{w}(\sigma)= heta^{(n)}$$
 .

• 
$$\mathbb{P}(\sigma) = \frac{\theta^{\operatorname{cycle}(\sigma)}}{\theta^{(n)}}.$$

Notation :  $\theta^{(n)} = \theta(\theta + 1) \dots (\theta + n - 1)$ 

## Generalizing the distribution

## Definition (Ewens-like distribution)

- Let  $\chi$  be any statistic on  $\sigma \in \mathfrak{S}_n$ .
- To any  $\sigma \in \mathfrak{S}_n$ , we associate a weight  $w(\sigma) = \theta^{\chi(\sigma)}$ .
- Let  $W_n = \sum_{\sigma \in \mathfrak{S}_n} w(\sigma)$  and  $\mathbb{P}(\sigma) = \frac{w(\sigma)}{W_n}$ .

## Generalizing the distribution

## Definition (Ewens-like distribution)

- Let  $\chi$  be any statistic on  $\sigma \in \mathfrak{S}_n$ .
- To any  $\sigma \in \mathfrak{S}_n$ , we associate a weight  $w(\sigma) = \theta^{\chi(\sigma)}$ .

• Let 
$$W_n = \sum_{\sigma \in \mathfrak{S}_n} w(\sigma)$$
 and  $\mathbb{P}(\sigma) = rac{w(\sigma)}{W_n}$  .

#### Analytic combinatorics

Let 
$$F(z, u) = \sum f_{n,k} z^n u^k$$
, where  $f_{n,k} = \|\{\sigma \in \mathfrak{S}_n : \chi(\sigma) = k\}\|$ .

$$W_n = n![z^n]F(z,\theta)$$
 and  $\mathbb{E}_n[\chi] = \frac{\theta[z^n] \left.\frac{\mathrm{d}F(z,u)}{\mathrm{d}u}\right|_{u=\theta}}{[z^n]F(z,\theta)}$ 

But can be difficult when  $\theta$  depends on n.

#### Recall

For any sequence X of size n,  $m_{rec}(X) = n - record(X)$  is a measure of presortedness.

### Definition (Ewens-like distribution for records)

In the following, we focus on this distribution.

# Some probabilities

### Results



# Some probabilities

## Results

$$\begin{array}{c|c} \mathbb{P}_n(\operatorname{Record at position } i) & \frac{\theta}{\theta+i-1} \\ \hline \mathbb{P}_n(\sigma(i-1) > \sigma(i)) & \frac{(i-1)(2\theta+i-2)}{2(\theta+i-1)(\theta+i-2)} \\ \hline \mathbb{P}_n(\sigma(1) = k) & \frac{(n-1)!\theta^{(n-k)}\theta}{(n-k)!\theta^{(n)}} \\ \hline \mathbb{P}_n(inv_j(\sigma) = k) & \begin{cases} \frac{\theta}{\theta+j-1} & \text{if } k = 0 \\ \frac{1}{\theta+j-1} & \text{otherwise} \end{cases} \end{array}$$

$$\mathbb{P}_n(\text{Record at position } i) = \frac{\theta^{(i-1)}\theta}{\theta^{(i)}} = \frac{\theta}{\theta + i - 1}$$



## Asymptotic equivalents

## Results

	$\theta = 1$	fixed $\theta > 0$	$\theta := n^{\epsilon},$	$\theta := \lambda n$ ,	$\theta := n^{\delta}$
	(uniform)		$0<\epsilon<1$	$\lambda > 0$	$\delta > 1$
$\mathbb{E}_n$ [record]	log n	$\theta \cdot \log n$	$(1-\epsilon)\cdot n^\epsilon\log n$	$\lambda \log(1+1/\lambda) \cdot n$	n
$\mathbb{E}_n[desc]$	n/2	n/2	n/2	$n/2(\lambda+1)$	$n^{2-\delta}/2$
$\mathbb{E}_n[\sigma(1)]$	n/2	n/( heta+1)	$n^{1-\epsilon}$	$(\lambda+1)/\lambda$	1
$\mathbb{E}_n[inv]$	<i>n</i> <sup>2</sup> /4	<i>n</i> <sup>2</sup> /4	n <sup>2</sup> /4	$n^2/4 \cdot f(\lambda)$	$n^{3-\delta}/6$

With  $f(\lambda) = 1 - 2\lambda + 2\lambda^2 \log (1 + 1/\lambda)$ .

$$\mathbb{P}_n(\text{Record at position } i) = rac{ heta^{(i-1)} heta}{ heta^{(i)}} = rac{ heta}{ heta+i-1}$$













134567829

123456789





- Adapts to the number of *inversions*.
- Sorts a sequence X in  $\Theta(Inv(X))$  comparisons.

Recall								
	$\theta = 1$	fixed $\theta > 0$	$\theta := n^{\epsilon}$ ,	$\theta := \lambda n$ ,	$\theta := n^{\delta}$			
	(uniform)		$0 < \epsilon < 1$	$\lambda > 0$	$\delta > 1$			
$\mathbb{E}_n[inv]$	<i>n</i> <sup>2</sup> /4	<i>n</i> <sup>2</sup> /4	<i>n</i> <sup>2</sup> /4	$n^2/4 \cdot f(\lambda)$	$n^{3-\delta}/6$			
With $f(\lambda) = 1 - 2\lambda + 2\lambda^2 \log (1 + 1/\lambda)$ .								

Unless  $\theta \gg n$ , InsertSort remains in  $\Theta(n^2)$  on average.

## Introduction to min/max search

NAIVEMINMAX(T, n)

return min, max

2n comparisons

3/2-MINMAX(T, n)

 $\begin{array}{c} \min, \max \leftarrow T[n], T[n] \\ \text{for } i \leftarrow 2 \text{ to } n \text{ by } 2 \text{ do} \\ & \text{if } T[i-1] < T[i] \text{ do} \\ & \begin{tabular}{ll} & pMin, pMax \leftarrow T[i-1], T[i] \\ & \text{else} \\ & \begin{tabular}{ll} & pMin, pMax \leftarrow T[i], T[i-1] \\ & \text{if } pMin < \min \text{ do } \min \leftarrow pMin \\ & \text{if } pMax > \max \text{ do } \max \leftarrow pMax \end{array}$ 

3n/2 comparisons

NAIVEMINMAX is faster than 3/2-MINMAX, when the data are uniformly distributed in [0, 1]. [Auger,Nicaud,Pivoteau,STACS 2016]

## Introduction to modern architecture



Prediction rule : for each if, same as last time

# Analysis under the uniform distribution

NAIVEMINMAX(T, n)

return min, max

3/2-MINMAX(T, n)

 $\ensuremath{\textit{misprediction}}\xspace = \ensuremath{\textit{alternation}}\xspace$  of record and non-record

### Results

- NAIVEMINMAX generates  $\Theta(\log(n))$  mispredictions.
- 3/2-MINMAX generates  $\Theta(n)$  mispredictions.

What happen if the number of records increases significantly ?

## Average analysis of the number of mispredictions

When  $\theta = \lambda n$  for some real  $\lambda$  and for our prediction rule, we have :

- $\mu$  number of mispredictions of NAIVEMINMAX.
- $\nu$  number of misprediction of 3/2-MINMAX.



## Discussion

## Questions

What's next ?

- Ewens-like distribution for other meaningful statistics that take part in (sorting) algorithms.
- For example, the runs for the analysis of TimSort.
- Explain the asymptotic shape of the diagrams below.(done)

