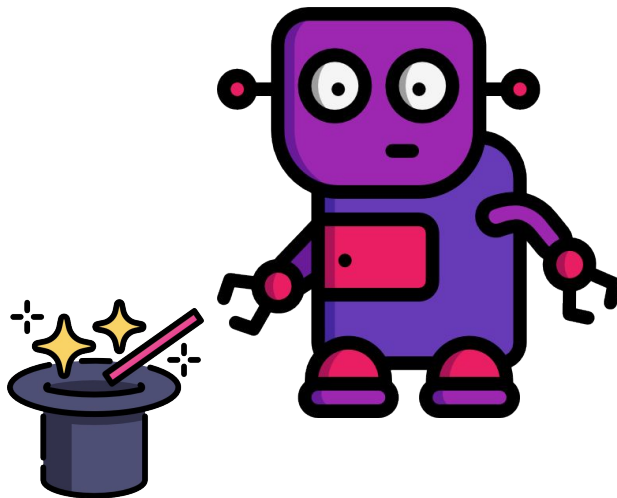




Programa de Pós Graduação em Ciências da Computação
Análise de Algoritmos e Estruturas de Dados - Profa. Dra. Lilian Berton

Análise: Learned Sort

Willian Dihanster Gomes de Oliveira



12 de Abril de 2021
São José dos Campos - SP

Introdução

- Uso de **algoritmos inteligentes** cresceu muito nos últimos anos.
- **Resultados interessantes** para diversas aplicações.
- Podemos **aplicar em ordenação**?

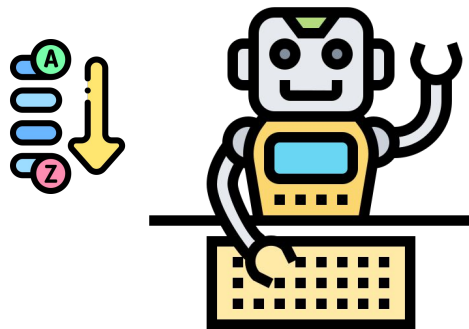


Figura 1: Ilustração de um robô fazendo uma tarefa de ordenação.

Ordenação + Algoritmos Inteligentes?

- Modelo **sabe** onde **colocar** o elemento atual.

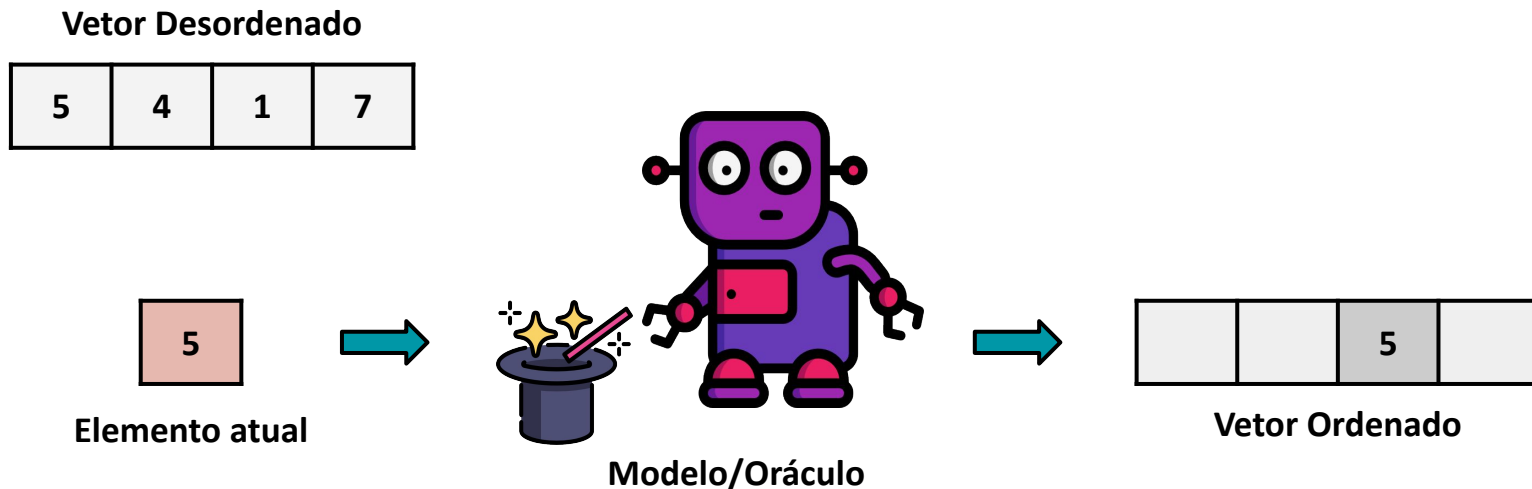


Figura 2: Esquema prático do uso de um algoritmo inteligente em ordenação.

Learned Sort v1.0 (KRISTO, Ani et al., 2020) [1]

- Supõe **padrão na distribuição** dos valores de entrada.
- Modelo **CDF (Cumulative Density Function)**.
 - Treina apenas com amostra para gerar aproximação;
 - Inferência para cada elemento $P(X \leq x)$.

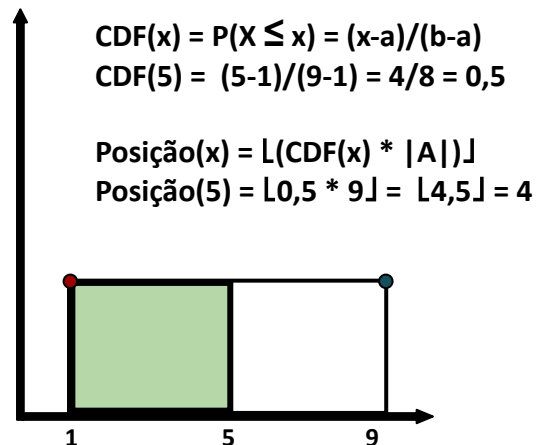
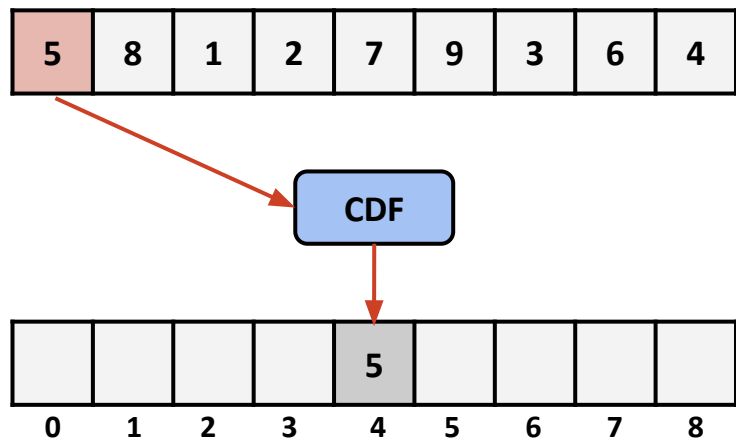


Figura 3: Esquema prático de aplicação do modelo CDF e cálculo de uma CDF.

Learned Sort v1.0 - Colisão

- **Linear:** coloca na **próxima posição livre**;
- **Encadeamento:** cria um **encadeamento na posição** que gerou colisão;
- **Spill Bucket:** aloca os valores em um **bucket especial** separado.

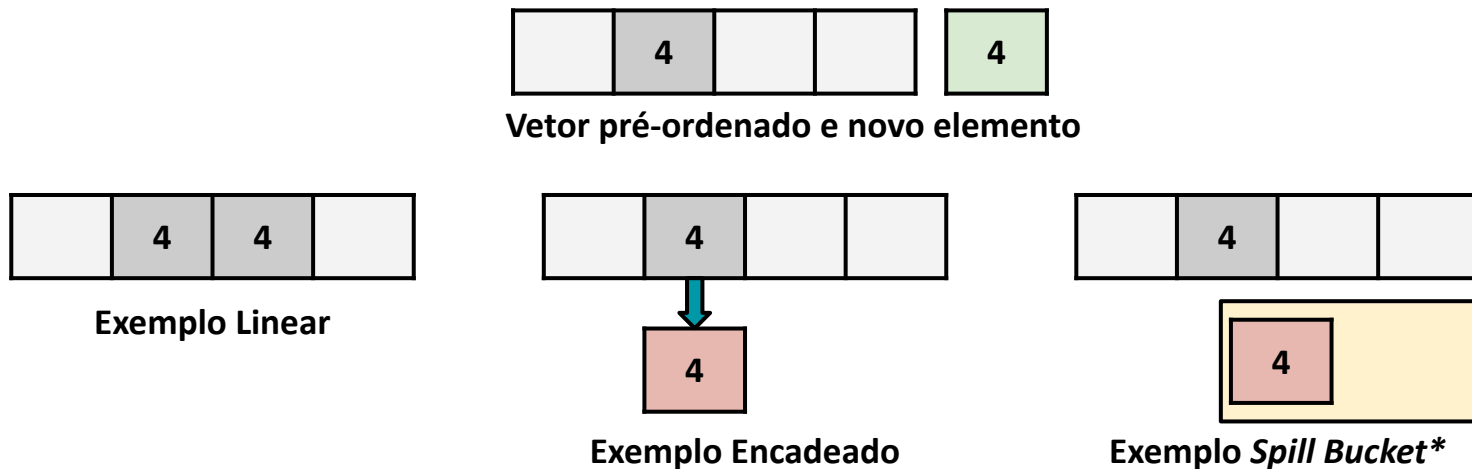


Figura 4: Esquema de tratamentos de colisões.

Learned Sort v1.0 - Imprecisão do Modelo

- Modelo pode **não acertar tudo**.
- Passada final com **Insertion Sort**.
- É esperado **pouca confusão**.



Figura 5: Esquema da passada final do algoritmo.

Learned Sort v1.0 - Algoritmo e Complexidade

Algorithm 1 A first Learned Sort

Input A - the array to be sorted

Input F_A - the CDF model for the distribution of A

Input o - the over-allocation rate. Default=1

Output A' - the sorted version of array A

```
1: procedure LEARNED-SORT( $A, F_A, o$ )
2:    $N \leftarrow A.length$ 
3:    $A' \leftarrow$  empty array of size  $(N \cdot o)$ 
4:   for  $x$  in  $A$  do
5:      $pos \leftarrow \lfloor F_A(x) \cdot N \cdot o \rfloor$ 
6:     if  $EMPTY(A'[pos])$  then  $A'[pos] \leftarrow x$ 
7:     else COLLISION-HANDLER( $x$ )
8:   if  $o > 1$  then COMPACT( $A'$ )
9:   if NON-MONOTONIC then INSERTION-SORT( $A'$ )
10:  return  $A'$ 
```

$\left. \begin{array}{l} O(n) \\ O(1) \\ O(1) \\ O(m) \end{array} \right\}$

$O(n)$

$O(n)$ ou $O(n^2)$

Depende da qualidade do modelo!

Sem colisões e ordenação
quase perfeita $\approx O(n)$

$O(n \cdot m + n + n) = O(n \cdot m)$
ou
 $O(n \cdot m + n + n^2) = O(n^2)$

Figura 6: Pseudocódigo do 1º algoritmo [1].

Learned Sort v2.0 - Cache-Otimizado

- Evoluções:
 - Prevê **bucket** em vez de posição exata;
 - Ordenação baseada em **Counting** e novo modelo CDF.

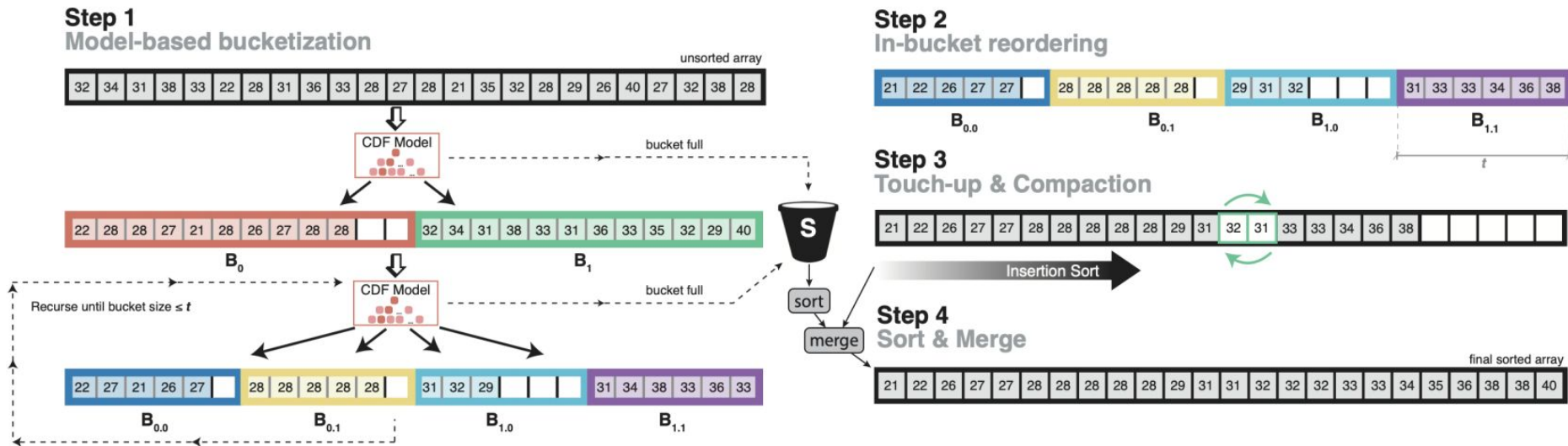


Figura 7: Esquema prático da versão 2.0 do algoritmo [1].

Learned Sort v2.0 - Novo Modelo

- **Modelo Hierárquico**
 - “Sabe” qual modelo pode lidar melhor com a chave.

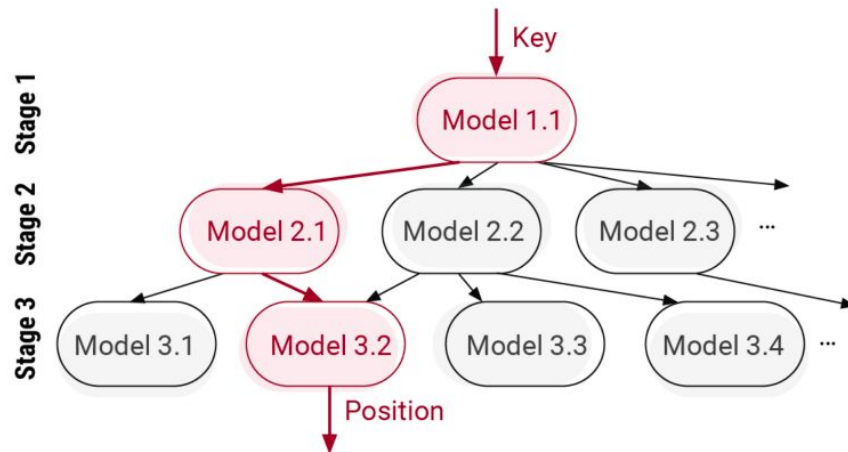


Figura 8: Esquema do 2º estimador [1].

Learned Sort v2.0 - Algoritmo e Complexidade

Algorithm 2 The Learned Sort algorithm

Input A - the array to be sorted
Input F_A - the CDF model for the distribution of A
Input f - fan-out of the algorithm
Input t - threshold for bucket size
Output A' - the sorted version of array A

```

1: procedure LEARNED-SORT( $A, F_A, f, t$ )
2:    $N \leftarrow |A|$                                  $\triangleright$  Size of the input array
3:    $n \leftarrow f$                                  $\triangleright n$  represents the number of buckets
4:    $b \leftarrow \lfloor N/f \rfloor$                          $\triangleright b$  represents the bucket capacity
5:    $B \leftarrow [] \times N$                            $\triangleright$  Empty array of size  $N$ 
6:    $I \leftarrow [0] \times n$                          $\triangleright$  Records bucket sizes
7:    $S \leftarrow []$                                  $\triangleright$  Spill bucket
8:    $\text{read\_arr} \leftarrow$  pointer to  $A$ 
9:    $\text{write\_arr} \leftarrow$  pointer to  $B$ 

10:  // Stage 1: Model-based bucketization
11:  while  $b \geq t$  do                                 $\triangleright$  Until bucket capacity reaches the threshold  $t$ 
12:     $I \leftarrow [0] \times n$                              $\triangleright$  Reset array  $I$ 
13:    for  $x \in \text{read\_arr}$  do
14:       $\text{pos} \leftarrow \lfloor \text{INFER}(F_A, x) \cdot n \rfloor$ 
15:      if  $I[\text{pos}] \geq b$  then                             $\triangleright$  Bucket is full
16:         $S.\text{append}(x)$                                  $\triangleright$  Add to spill bucket
17:      else                                             $\triangleright$  Write into the predicted bucket
18:         $\text{write\_arr}[\text{pos} \cdot b + I[\text{pos}]] \leftarrow x$ 
19:         $\text{INCREMENT } I[\text{pos}]$ 
20:       $b \leftarrow \lfloor b/f \rfloor$                              $\triangleright$  Update bucket capacity
21:       $n \leftarrow \lfloor N/b \rfloor$                              $\triangleright$  Update the number of buckets
22:       $\text{PtrSWP}(\text{read\_arr}, \text{write\_arr})$                  $\triangleright$  Pointer swap to reuse memory

```

$O(I \cdot n) =$
 $O(n)$

```

23:  // Stage 2: In-bucket reordering
24:   $\text{offset} \leftarrow 0$ 
25:  for  $i \leftarrow 0$  up to  $n$  do                                 $\triangleright$  Process each bucket
26:     $K \leftarrow [0] \times b$                                  $\triangleright$  Array of counts

27:    for  $j \leftarrow 0$  up to  $I[i]$  do  $\triangleright$  Record the counts of the predicted positions
28:       $\text{pos} \leftarrow \lfloor \text{INFER}(F_A, \text{read\_arr}[\text{offset} + j]) \cdot N \rfloor$ 
29:       $\text{INCREMENT } K[\text{pos} - \text{offset}]$ 

30:    for  $j \leftarrow 1$  up to  $|K|$  do                                 $\triangleright$  Calculate the running total
31:       $K[j] \leftarrow K[j] + K[j - 1]$ 

32:     $T \leftarrow []$                                              $\triangleright$  Temporary auxiliary memory
33:    for  $j \leftarrow 0$  up to  $I[i]$  do                                 $\triangleright$  Order keys w.r.t. the cumulative counts
34:       $\text{pos} \leftarrow \lfloor \text{INFER}(F_A, \text{read\_arr}[\text{offset} + j]) \cdot N \rfloor$ 
35:       $T[j] \leftarrow \text{read\_arr}[\text{offset} + K[\text{pos} - \text{offset}]]$ 
36:       $\text{DECREMENT } K[\text{pos} - \text{offset}]$ 
37:    Copy  $T$  back to  $\text{read\_arr}[\text{offset}]$ 
38:     $\text{offset} \leftarrow \text{offset} + b$ 

39:  // Stage 3: Touch-up
40:   $\text{INSERTION-SORT-AND-COMPACT}(\text{read\_arr})$ 

41:  // Stage 4: Sort & Merge
42:   $\text{SORT}(S)$ 
43:   $A' \leftarrow \text{MERGE}(\text{read\_arr}, S)$ 

44:  return  $A'$ 

```

$O(n \cdot t)$

$O(n \cdot t)$ ou $O(n^2)$

$O(\text{slogs}) + O(n)$

$O(n)$ ou $O(n^2)$

Figura 9: Pseudocódigo do 2º algoritmo [1].

Resultados

- Learned Sort foi melhor que Radix Sort somente em ordem crescente.
- Mas melhor que `std::sort` em todos os casos.

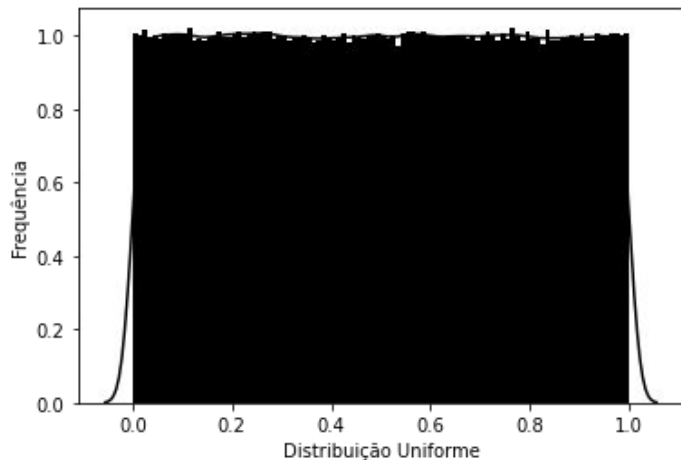


Figura 10: Distribuição Uniforme utilizada 1M com $a=0$ e $b=1$.

Algoritmo	Aleatório	Crescente	Decrescente
Learned Sort	212096	11285	230521
Radix Sort	38690	40929	38392
<code>std::sort</code>	377897	339338	342604

Tabela 1: Resultados obtidos (em μs).

Resultados

- Learned Sort foi melhor que Radix Sort somente em ordem crescente.
- Mas melhor que `std::sort` em todos os casos.

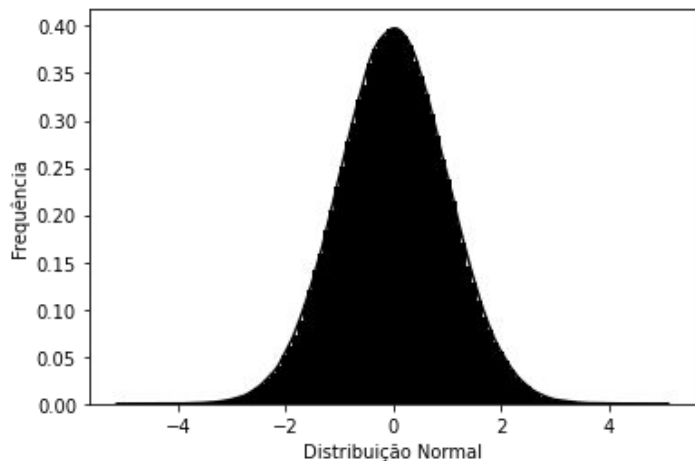


Figura 10: Distribuição Normal utilizada
1M com $\mu=0$ e $\sigma=1$.

Algoritmo	Aleatório	Crescente	Decrescente
Learned Sort	218696	11710	219092
Radix Sort	44780	41175	39319
<code>std::sort</code>	347259	339811	342700

Tabela 2: Resultados obtidos (em μs).

Conclusões

- Mais **experimentos**.
 - Não foi obtido os resultados esperados. Parâmetros e ambiente computacional?
 - Autores conseguiram resultados melhores nos seguintes cenários:
 - Número de elementos $\lll K$ e quando o número de elementos é tão grande que não cabe na memória cache L3.
 - No algoritmo 1 mesmo com **modelo prevendo 100% não bateu RadixSort***, pois RadixSort faz **acesso sequencial** e Learned Sort faz **acesso aleatório**.
- Pode ser uma **boa opção** quando:
 - A **distribuição** dos dados é **conhecida**;
 - **Não** há muitas **duplicatas**.

Referências

- KRISTO, Ani et al. The case for a learned sorting algorithm. In: **Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data**. 2020. p. 1001-1016.

Código disponível em:

<https://github.com/learnedsystems/LearnedSort>

Obrigado!