# Computational Learning Theory
## Frequent Substring Mining

Akihiro Yamamoto 山本 章博

http://www.iip.ist.i.kyoto-u.ac.jp/member/akihiro/
akihiro@i.kyoto-u.ac.jp

# Why sequences in ML?(1)

- Sentences are strings(sequences) consisting of characters in an alphabet.



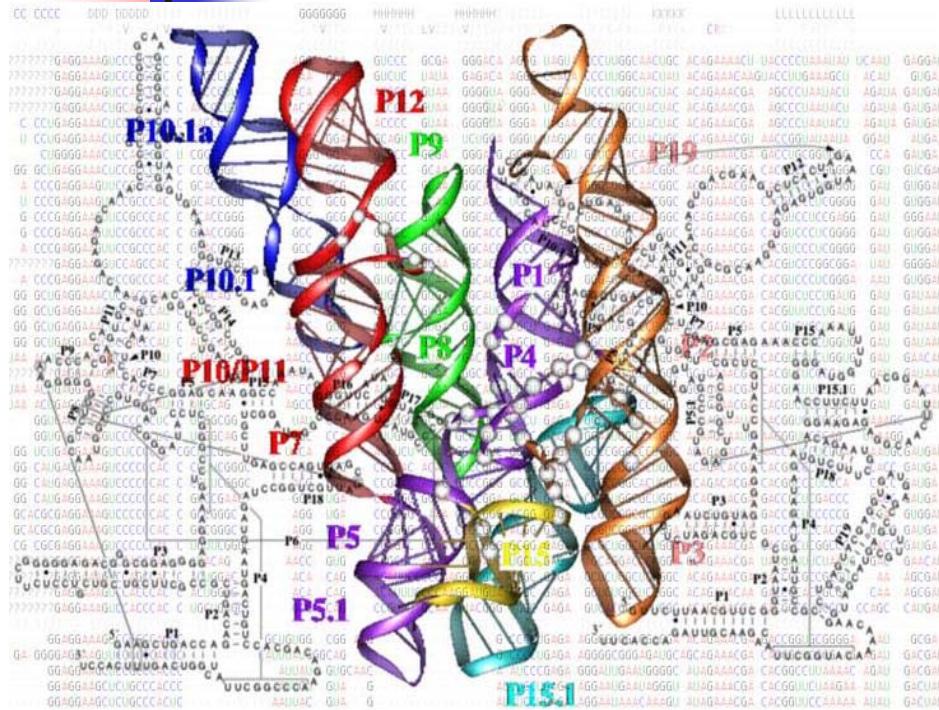ON COMPUTABLE NUMBERS, WITH AN APPLICATION TO THE ENTSCHEIDUNGSPROBLEM

By A. M. TURING.

[Received 28 May, 1936.—Read 12 November, 1936.]

The "computable" numbers may be described briefly as the real numbers whose expressions as a decimal are calculable by finite means. Although the subject of this paper is ostensibly the computable *numbers*, it is almost equally easy to define and investigate computable functions of an integral variable or a real or computable variable, computable predicates, and so forth. The fundamental problems involved are, however, the same in each case, and I have chosen the computable numbers for explicit treatment as involving the least cumbrous technique. I hope shortly to give an account of the relations of the computable numbers, functions, and so forth to one another. This will include a development of the theory of functions of a real variable expressed in terms of com-

[Wikipedia]                    [Davis, M. : *The Unsolvable,* Raven Press]

# Why sequences in ML?(2)



CACAUGUACAAGACUU
5'                          3'

- Many data for academic research is now open. In particular, many string data are provided in the area of bio-informatics.

# FORMALIZATION OF THE PROBLEM

# Formalization of the Task

Inputs : a (long) string $S$

a minimum number of

appearance(threshold) $\sigma \in \mathbf{N}$

Enumerate all patterns $P$ satisfying supp$(P) \geq \sigma$

- Here a pattern $P$ is defined as a substring.
- The support of $P$ means how many times P occurs in $S$.

# Substring and Subsequence

For a string $S = a_0 a_1 \ldots a_n$

- A substring of $S$ is $S[i,\ j] = a_i\, a_{i+1} a_{i+2} \ldots a_j$

  - The terminology "pattern matching" means to find all substrings which is identical to a given sting $T$.

- A subsequence of $S$ is $a_{i_1} a_{i_2} a_{i_3} \ldots a_{i_k}$ $0 \le i_1 < i_2 < \ldots < i_k \le n$

  - just in the terminology "longest common subsequece"

例 For $S = $ abcabcab and $T = $ cccaacaaba,

  cabc is a substring and also a subseqence of $S$, while

  cacb is not a substring of $S$ but a subsequence, and

  a common subsequence of $S$ and $T$.

# Suffix of a String

- A suffix of $S = a_0 a_1 \ldots a_n$ is a subsequence $S[i, n]$ $(i = 1, 2, \ldots, n-1)$ ending with $a_n$.
  - We assume $a_n$ e is a special symbol $

Example $S =$ sakurasaku$

$S[0,10] =$ sakurasaku$      $S[5,10] =$ asaku$

$S[1,10] =$ akurasaku$      $S[6,10] =$ asaku$

$S[2,10] =$ kurasaku$      $S[7,10] =$ aku$

$S[3,10] =$ urasaku$      $S[8,10] =$ ku$

$S[4,10] =$ rasaku$      $S[9,10] =$ u$

# Example (Tsuboi[2003])

$S$ = sakurasaku

minimum support $\sigma = 2$

sakurasaku

sakurasaku          sakurasaku

sakurasaku          sakurasaku          sakurasaku

sakurasaku   sakurasaku   sakurasaku   sakurasaku

sakurasaku

# Tsuboi's Algorithm [Tsuboi 03]

Inputs :

a (long) string $S$ ending with a special symbol $\$$

a minimum number $\sigma \in \mathbf{N}$ of appearance (threshold)

1. Enumerate all suffixes of $S$

2. Make three clusters of the suffixes.

3. For each cluster obtained by Step 2, apply Step 2,

   so we get a tree whose nodes has at most 3 children.

 /* This is hierarchical clustering. */

- Construct the prefix tree (trie) of all suffixes $S[i, n]$ ($i = 1, 2,\ldots, n-1$).

# A Step in Hierarchal Clustering (1)

- Just like the Quick Sort Algorithm, choose a pivot $v$ and make three clusters with the head symbol $S[i]$ of each suffix.
  - If all suffixes have the same head symbol, try to use the second character, third character,…

Example $\quad v = $ k

| $S[i] < v$ | $S[i] = v$ | $S[i] > v$ |
|---|---|---|
| akurasaku$ | kurasaku$ | sakurasaku$ |
| asaku$ | ku$ | urasaku$ |
| aku$ | | rasaku$ |
| | | sa$ |
| | | u$ |

# A Step in Hierarchal Clustering (2)

$v = $ k

akurasaku$
asaku$
aku$

kurasaku$
ku$

sakurasaku$
urasaku$
rasaku$
saku$
u$

$v = $ s

akurasaku$
aku$

asaku$

$v = $ r

kurasaku$
ku$

$v = $ s

rasaku$

sakurasaku$
saku$

urasaku$
u$

# Pruning Clusters (1)

- $\mathrm{supp}(S[i,\ j]) \geq \mathrm{supp}(S[i-k\ ,\ j])\ \ (k > 0)$

- We can delete all clusters which has suffixes less than $\sigma$.

# Pruning Clusters (2)

$v = $ k

| akurasaku$ | kurasaku$ | sakurasaku$ |
| asaku$ | ku$ | urasaku$ |
| aku$ | | rasaku$ |
| | | saku$ |
| | | u$ |

$v = $ s      $v = $ r      $v = $ s

akurasaku$
aku$
asaku$

kurasaku$
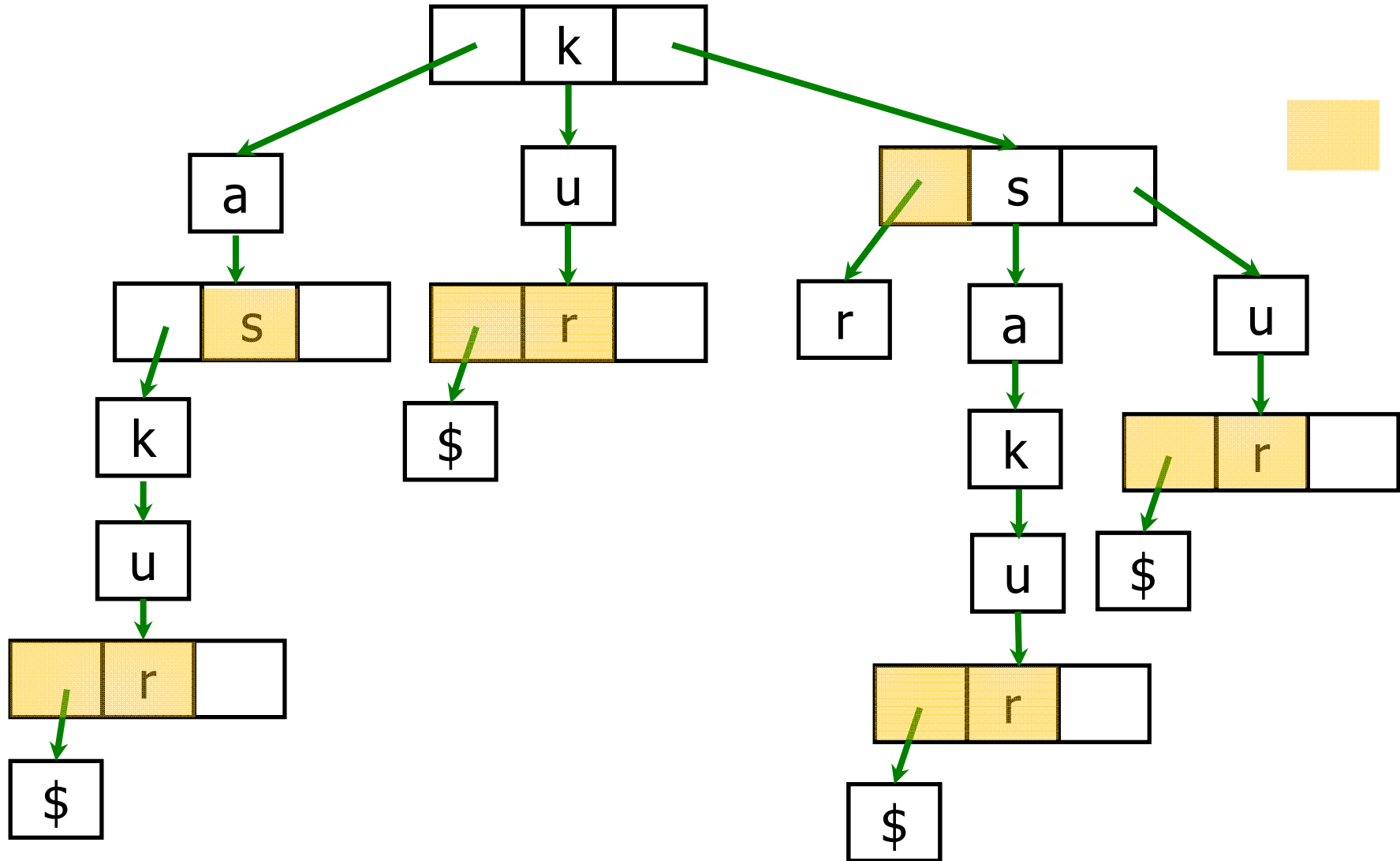ku$

rasaku$
sakurasaku$
saku$
urasaku$
u$

# Implementation with a Tree.

# Implementation with an Array(1)

- Use the suffix array

| s | a | k | u | r | a | s | a | k | u | $ |

| | |
|---|---|
| 7 | aku$ |
| 1 | akurasaku |
| 5 | asaku$ |
| 8 | ku$ |
| 2 | kursasaku$ |
| 4 | rasaku$ |
| 6 | saku$ |
| 0 | sakurasaku$ |
| 9 | urasaku$ |
| 3 | u$ |

# Implementation with an Array(2)

- Divide the suffix array into three recursively.

| s | a | k | u | r | a | s | a | k | u | $ |
|---|---|---|---|---|---|---|---|---|---|---|

| | |
|---|---|
| 7 | aku$ |
| 1 | akurasaku |
| 5 | asaku$ |
| 8 | ku$ |
| 2 | kursasaku$ |
| 4 | rasaku$ |
| 6 | saku$ |
| 0 | sakurasaku$ |
| 9 | urasaku$ |
| 3 | u$ |

# Implementation with an Array(3)

- Divide the suffix array into three recursively.

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| s | a | k | u | r | a | s | a | k | u | $ |

| | |
|---|---|
| 7 | aku$ |
| 1 | akurasaku |
| 5 | asaku$ |
| 8 | ku$ |
| 2 | kursasaku$ |
| 4 | rasaku$ |
| 6 | saku$ |
| 0 | sakurasaku$ |
| 9 | urasaku$ |
| 3 | u$ |