

Analogical proportions in a lattice of sets of alignments built on the common subwords in a finite language

L. Miclet¹, N. Barbot¹, B. Jeudy²

¹IRISA, Lannion & Rennes, France

²Université Jean Monnet, Saint-Étienne, France

ECAI 2012, SAMAI Workshop

- 1 Introduction: Teaching is good for Research
- 2 Subwords, superwords and alignments
- 3 A structured set of sets of alignments for n words
- 4 Conclusion

Roadmap

- 1 Introduction: Teaching is good for Research
- 2 Subwords, superwords and alignments
- 3 A structured set of sets of alignments for n words
- 4 Conclusion

Dynamic Programming and **inspiration**

- Lectures on Dynamic Programming: examples, among which the Longest Common Subword $lcs(u, v)$ problem.
- For the exam, the Shortest Common Superword $SCS(u, v)$.
- **HA-HA !** Notice that $|SCS(u, v)| + |lcs(u, v)| = |u| + |v|$. (N.B.: Not *that* smart: there is a simple demonstration).
- **ANALOGY !** Analog to $LCM(a, b) \times GCD(a, b) = a \times b$.
There is also an analogical proportion between the four words:

$$SCS(u, v) : u :: lcs(u, v) : v$$

Would there be a **lattice of SCS and lcs** generated by a finite language U , with nice analogical proportions ?

This is where **transpiration** begins.

A first glance

$$\begin{array}{rcccccc}
 u = cadb = & c & a & d & \sim & b & \sim \\
 & | & | & | & | & | & | \\
 v = aebf = & \sim & a & \sim & e & b & f \\
 & | & | & | & | & | & | \\
 lcs(u, v) = & \sim & a & \sim & \sim & b & \sim \\
 & | & | & | & | & | & | \\
 SCS(u, v) = & c & a & d & e & b & f \\
 & & & e & d & &
 \end{array}$$

Roadmap

- 1 Introduction: Teaching is good for Research
- 2 **Subwords, superwords and alignments**
- 3 A structured set of sets of alignments for n words
- 4 Conclusion

Basics of stringology

- Σ is an alphabet, *i.e.* a finite set of letters.
- A *word* $u \in \Sigma^*$ is a sequence $u_1 \dots u_n$ of letters in Σ .
- The length of u , denoted $|u|$ is n .
- The empty word, of null length, is ϵ .
- A *language* is a set of words.
- A *subword* of a word u is a word obtained by deleting the letters at some (non necessarily adjacent) positions in u .

non necessarily adjacent	<i>subword</i>	<i>subsequence</i>	<i>sous-mot</i>
adjacent	<i>factor</i>	<i>substring</i>	<i>facteur</i>

The shuffle of words

- We denote $u \bullet v$ or $\bullet(u, v)$ the *shuffle* of two words.
- For example,
 $ab \bullet bc = \{abbc, (abbc), abc b, babc, bac b, bcab\}$
- The set $u \bullet v$ has at most $\binom{|u|+|v|}{|u|}$ elements.
- This operation is associative.

Order relations in Σ^*

- In Σ^* , there is a partial order relation denoted \leq defined by:

$$(u \leq v \Leftrightarrow u \text{ is a subword of } v)$$

- When u is a subword of v , v is called a *superword* of u . For example:

$$abc \leq aab**bc**d$$

$$ab**d** \leq aab**bc**d$$

- Another (total) order relation is that on the lengths of the words.

non necessarily adjacent	<i>superword</i>	<i>supersequence</i>	<i>sur-mot</i>
adjacent	<i>superstring</i>		<i>sur-facteur (?)</i>

SA

For two words

- w is a *common subword* to u and v when $w \leq u$ and $w \leq v$.
- w is a *maximal* common subword to u and v if there not exist any other common subword x to u and v such that $w \leq x$.

Definition

- $\sqcup(u, v)$ is the set of maximal common subwords to u and v
- $\sqcap(u, v)$ is the set of minimal common superwords to u and v

For two words

For example, ab and c are maximal common subwords to $u = cadba$ and $v = fagbhc$, while a is a non maximal common subword.

$u = cadba =$	c	\sim	a	d	\sim	b	a	\sim	\sim
$v = fagbhc =$	\sim	f	a	\sim	g	b	\sim	h	c
$lcs(u, v) =$	\sim	\sim	a	\sim	\sim	b	\sim	\sim	\sim
$SCS(u, v) =$	c	f	a	d	g	b	a	h	c

For two words

For example, ab and c are maximal common subwords to $u = cadba$ and $v = fagbhc$, while a is a non maximal common subword.

$u = cadba =$ $c \sim a \sim d \sim b \sim a \sim \sim$

$v = fagbhc =$ $\sim f \sim a \sim g \sim b \sim h \sim c$

Introduction to locally maximal subwords

a is a subword of ab , but...

$u = cadba =$	c	a	d	b	\sim	a	\sim	\sim	\sim	\sim
$v = fagbhc =$	\sim	\sim	\sim	\sim	f	a	g	b	h	c
$lcs(u, v) =$	\sim	\sim	\sim	\sim	\sim	a	\sim	\sim	\sim	\sim
$SCS(u, v) =$	c	a	d	b	f	a	g	b	h	c

Actually, $SCS(u, v) = (cadb \bullet f) a gbhc$

Call a (second in u , first in v) a *locally maximal subword* of u and v .

Back to SCS and lcs

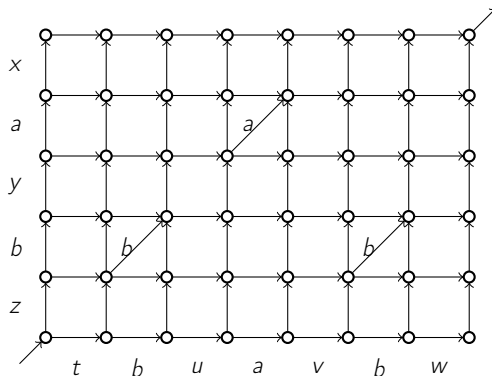
Enumerating $\sqcup(u, v)$ is tricky !

- $u = bac$ and $v = dae$.
- The four words in $(d \bullet b)a(c \bullet e)$ are minimal superwords of length five.
- But u and v have four other minimal superwords of length six, included in $u \bullet v$, namely $ba(c \bullet d)ae$ and $da(e \bullet b)ac$.
- All the other elements of $u \bullet v$ are non minimal.

An example

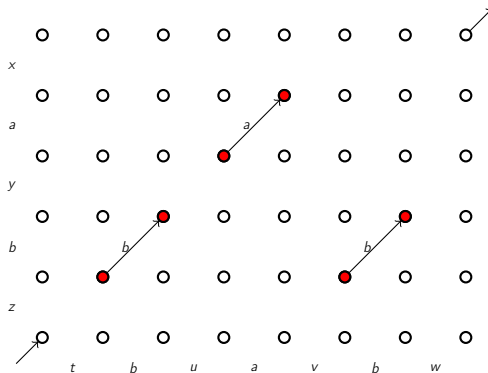
- $\sqcap(cadba, fagbhc) = \{ab, c\}$
- $\sqcup(cadba, fagbhc)$ includes for example
 $cfadgbahc$ and $fagbhcadba$

Graphically speaking, a simpler example



All minimal superwords of $tbuavbw$ and $zbyax$ (and a lot of non minimal) are recognized by this finite automaton. The maximal subword ba is can be read in the diagonals.

Locally maximal subwords



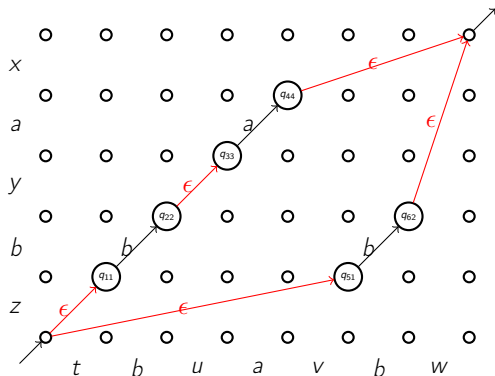
There is a partial order between the diagonal transitions.

The locally maximal subwords are the maximal chains of this order.

Construct an automaton from the Hasse diagram of this order

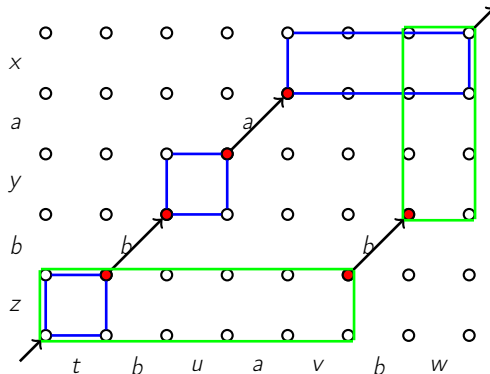
Locally maximal subwords and locally minimal superwords

$\sqcap(u, v)$ and $\sqcup(u, v)$.



The automaton $\mathcal{A}_{\sqcap}(u, v)$ recognizing the locally maximal subwords common to u and v .

Constructing the locally minimal superwords



This automaton recognizes exactly the locally minimal superwords.
A rectangle stands for a shuffle automaton.

Conclusion of the geometric part

- Construct $\sqcap(u, v)$ (maximal subwords)?
Construct the Hasse diagram and eliminate subwords.
 - Construct $\sqcup(u, v)$ (minimal superwords)?
Construct $u \sqcap v$ and test all subwords.
 - Construct $\sqcap(u, v)$ (locally maximal subwords) ?
Construct the Hasse diagram.
 - Construct $\sqcup(u, v)$ (locally minimal superwords) ?
By definition, from the previous.
-
- Generalization to any number of words: quite possible for the locally sub and superwords.
 - We deal in the following only with the locally sub and superwords.

Roadmap

- 1 Introduction: Teaching is good for Research
- 2 Subwords, superwords and alignments
- 3 A structured set of sets of alignments for n words
- 4 Conclusion

Example:

an alignment of three words on a common subword

$$a = \left(\begin{array}{cccccccc} & \boxed{a} & & \boxed{c} & b & d & e & g \\ & \boxed{a} & & \boxed{c} & & & e & h \\ g & \boxed{a} & h & \boxed{c} & & d & & \end{array} \right)$$

$$a = (\boxed{a}\boxed{c}bdeg, \boxed{a}\boxed{c}eh, g\boxed{a}h\boxed{c}d)$$

$$L(a) = (\epsilon \bullet \epsilon \bullet g)a(\epsilon \bullet \epsilon \bullet h)c(bdeg \bullet eh \bullet d)$$

Alignment

Definition

An alignment is a finite set of pairs (w, l) where w is a word and l a set of indices between 1 and $|w|$. The set l defines a subword of w denoted $w[l]$. Moreover, an alignment a must satisfy the following properties for all $(w, l) \in a$ and $(w', l') \in a$:

- ① $w[l] = w'[l']$
- ② $(w = w') \Rightarrow (l = l')$
- ③ $(w \leq w') \Rightarrow (w = w')$

The set of words on which the alignment is defined is called the *support* and is denoted $word(a) = \{w \mid \exists l \text{ with } (w, l) \in a\}$.

Locally maximal alignments and subwords

Definition 1

An alignment $a = \{(w_1, l_1), \dots, (w_n, l_n)\}$ is *locally maximal* if there is no other alignment $b = \{(w_1, l'_1), \dots, (w_n, l'_n)\}$ on the same support such that for all i , $l_i \subset l'_i$.

Definition 2

The set of boxed subwords associated to all locally maximal alignments between a finite set of words $W = \{w_1, \dots, w_n\}$ is called the set of *locally maximal subwords* to W and is denoted $\Box(\{w_1, \dots, w_n\})$.

Example of locally maximal alignment and subword

The set of locally maximal alignments of $W = \{ababc, cabd\}$ is
 $\mathcal{A}(W) =$

$$\left\{ (\boxed{a}\boxed{b}abc, c\boxed{a}\boxed{b}d), (\boxed{a}ba\boxed{b}c, c\boxed{a}\boxed{b}d), \right. \\ \left. (ab\boxed{a}\boxed{b}c, c\boxed{a}\boxed{b}d), (abab\boxed{c}, \boxed{c}abd) \right\}$$

The set of locally maximal subwords is

$$\sqcap(W) = \{ab, c\}$$

Order relation between alignments: an example

$$\left(\begin{array}{c} a \\ d \\ f \end{array} \begin{array}{|c|} \hline b \\ \hline \end{array} \begin{array}{|c|} \hline c \\ \hline \end{array} e \right) \sqsubseteq \left(\begin{array}{c} a \\ a \\ d \\ b \end{array} \begin{array}{|c|} \hline b \\ \hline \end{array} \begin{array}{c} d \\ c \\ a \\ \end{array} c \begin{array}{|c|} \hline b \\ \hline \end{array} e \right)$$

Order relation between alignments

Given two alignments

$a = \{(w_1, l_1), \dots, (w_n, l_n)\}$ and $b = \{(w'_1, l'_1), \dots, (w'_m, l'_m)\}$, we write $a \sqsubseteq b$ if

- For every word w in $word(a)$ there exists a word w' in $word(b)$ such that $u \leq v$.
- and $w_i = w'_j \Rightarrow l'_j \subseteq l_i$.

Homogeneous sets of alignments and order

Homogeneous sets of alignments

A set of alignments is homogeneous when all its elements have the same support.

The family of homogeneous sets of alignments is denoted \mathcal{A}_H .

Order on homogeneous sets of alignments

Let A and B be two homogeneous sets of alignments. We have $A \sqsubseteq B$ if for all $b \in B$, there is $a \in A$ such that $a \sqsubseteq b$.

Property

\sqsubseteq is a partial order on \mathcal{A}_H and the smallest element is $\{\emptyset\}$.

The union operation Υ on \mathcal{A}_H

Let $a \in A_r(\{u_1, \dots, u_n\})$ and $b \in A_s(\{v_1, \dots, v_m\})$, where $a = \{(u_1, l_1), \dots, (u_n, l_n)\}$ and $b = \{(v_1, l'_1), \dots, (v_m, l'_m)\}$.

Firstly, we construct $a + b$, the finite set of alignments $c = \{(w_1, L_1), \dots, (w_p, L_p)\}$ such that

- ① $\{w_1, \dots, w_p\} = \text{word}(a) \cup \text{word}(b)$
- ② for all (i, k) , if $(w_k = u_i)$ then $(L_k \subseteq l_i)$
- ③ for all (j, k) , if $(w_k = v_j)$ then $(L_k \subseteq l'_j)$

Secondly, we denote $a \Upsilon b$ the set of minimal elements of $a + b$ according to \sqsubseteq .

The union operation Υ on \mathcal{A}_H

$$a = \begin{pmatrix} a & \boxed{b} & \boxed{c} & d & & \\ & \boxed{b} & \boxed{c} & & & \\ & & & & a & \end{pmatrix} \quad b = \begin{pmatrix} & \boxed{b} & & \boxed{d} & a & b \\ a & \boxed{b} & c & \boxed{d} & & \\ a & \boxed{b} & c & \boxed{d} & a & \end{pmatrix}$$

$$a + b = \begin{pmatrix} a & \boxed{b} & c & d & & \\ & \boxed{b} & c & & a & \\ & \boxed{b} & & d & a & b \\ a & \boxed{b} & c & d & & \\ a & \boxed{b} & c & d & a & \end{pmatrix} \quad a \Upsilon b = \begin{pmatrix} & \boxed{b} & & d & a & b \\ a & \boxed{b} & c & d & a & \end{pmatrix}$$

Note that $a \sqsubseteq a \Upsilon b$ and $b \sqsubseteq a \Upsilon b$.

The union operation Υ on \mathcal{A}_H

Extension

Let A and B be two homogeneous sets of alignments. We define $A \Upsilon B$ as the set of the minimal elements of $A + B$ according to \sqsubseteq where

$$A + B = \bigcup_{\substack{b \in B \\ a \in A}} (a + b)$$

Property

The operation Υ is internal to \mathcal{A}_H , commutative and idempotent.

An intersection operation

$$a = \{([\boxed{a}]cd, ab[\boxed{a}]c, [\boxed{a}]ba)\}$$

$$b = \{(a[\boxed{c}]d, aba[\boxed{c}], [\boxed{c}]a)\}$$

- The support of $a \wedge b$ is the intersection of the supports of a and b , namely $\{acd, abac\}$.
- The boxed subwords have to be maximal.

$$a \wedge b = \{([\boxed{a}][\boxed{c}]d, [\boxed{a}]ba[\boxed{c}]), ([\boxed{a}][\boxed{c}]d, ab[\boxed{a}][\boxed{c}])\}$$

An intersection operation

Let $a \in A_r(\{u_1, \dots, u_n\})$ and $b \in A_s(\{v_1, \dots, v_m\})$ where $a = \{(u_1, l_1), \dots, (u_n, l_n)\}$ and $b = \{(v_1, l'_1), \dots, (v_m, l'_m)\}$. We construct $a \wedge b$, the finite set of alignments

$c = \{(w_1, L_1), \dots, (w_p, L_p)\}$ such that

- ① $\{w_1, \dots, w_p\} = \text{word}(a) \cap \text{word}(b)$
- ② Either, for all (i, k) such that $w_k = u_i$ we have $l_i \subseteq L_k$, or for all (j, k) such that $w_k = v_j$ we have $l'_j \subseteq L_k$.
- ③ c is a locally maximal alignment.

An intersection operation

Extension

$$A \curlywedge B = \bigcup_{\substack{b \in B \\ a \in A}} (a \curlywedge b)$$

Property

The operation \curlywedge is internal to \mathcal{A}_H , commutative and idempotent.

Finding a lattice structure to \mathcal{A}_H

Definition

We define $\sup(A, B)$ as the minimal set of alignments
 \sqsubseteq
 larger than A and B (if it exists) according to \sqsubseteq .

Similarly, $\inf(A, B)$ is the maximal set of alignments
 \sqsubseteq
 smaller than A and B .

Property

Let A and B be finite homogeneous sets of alignments.

Then $\sup(A, B)$ exists and
 \sqsubseteq

$$\sup(A, B) = A \vee B$$

$$\sqsubseteq$$

Finding a lattice structure to \mathcal{A}_H

Definition

If U is a finite collection of words.

We define the collection of sets of alignments

$$\mathcal{A}(U) = \{A(V) \mid V \subseteq U\}.$$

Property

Let A and B be sets of alignments in $\mathcal{A}(U)$.

Then, in $\mathcal{A}(U)$, $\inf_{\sqsubseteq}(A, B)$ exists and:

$$\inf_{\sqsubseteq}(A, B) = A \sqcap B$$

Finally...

Property

Let $U = \{u_1, u_2, \dots, u_n\}$ be a finite set of words, the operations \wedge and \vee are internal to U .

Final result

Let $U = \{u_1, u_2, \dots, u_n\}$ be a finite set of words, antichain for \leq .

Then $\mathcal{U} = (\mathcal{A}(U), \vee, \wedge)$ is a lattice.

This lattice is said to be built on the finite language U .

Roadmap

- 1 Introduction: Teaching is good for Research
- 2 Subwords, superwords and alignments
- 3 A structured set of sets of alignments for n words
- 4 Conclusion

Conclusion

Summary

- Exploring the structure of subwords and superwords common to a finite language.
- Defining the concept of **locally maximal** common subwords and **minimal** superwords.
- Finding a lattice of alignments on a finite language.

Further work

Use the lattice structure to extend finite languages using

- Analogical proportions
- Generalization