

Periodicity of Degenerate Strings

Pengfei Wang¹

Joint with Estéban Gabory², Eric Rivals¹, Michelle Sweering² and
Hilde Verbeek²

¹LIRMM, Univ. Montpellier, CNRS, Montpellier, France.

²CWI, Amsterdam, The Netherlands.

Prague Stringology Conference
Aug.2023



Multiple Sequence Alignment

Multiple Sequence Alignment (MSA)

```
CA - - AGCGCTAA - - - TT
C - - - AGCCGAAGT - - AT
CA - CAAGTCAAG - - - - T
```

Local Gapless Alignment (LGA)

```
CAAGCGCTAATT
CAGCCGAAGTAT
CACAAAGTCAAGT
```

Degenerate string from LGA

Local Gapless Alignment

C	A	A	G	C	G	C	T	A	A	T	T
C	A	G	C	C	G	A	A	G	T	A	T
C	A	C	A	A	G	T	C	A	A	G	T

Degenerate String

$$\{C\} \cdot \{A\} \cdot \begin{Bmatrix} A \\ G \\ C \end{Bmatrix} \cdot \begin{Bmatrix} G \\ C \\ A \end{Bmatrix} \cdot \begin{Bmatrix} C \\ C \\ A \end{Bmatrix} \cdot \{G\} \cdot \begin{Bmatrix} C \\ A \\ T \end{Bmatrix} \cdot \begin{Bmatrix} T \\ A \\ C \end{Bmatrix} \cdot \begin{Bmatrix} A \\ G \\ A \end{Bmatrix} \cdot \begin{Bmatrix} A \\ T \\ A \end{Bmatrix} \cdot \begin{Bmatrix} T \\ A \\ G \end{Bmatrix} \cdot \{T\}$$

Period of classical string

Example: String $u = aabbaa$ has periods 0, 4 and 5

	0	1	2	3	4	5	6	7	8	9	10
u	a	a	b	b	a	a	-	-	-	-	-
u	a	a	b	b	a	a	-	-	-	-	-
	-	a	a	b	b	a	a	-	-	-	-
	-	-	a	a	b	b	a	a	-	-	-
	-	-	-	a	a	b	b	a	a	-	-
	-	-	-	-	a	a	b	b	a	a	-
	-	-	-	-	-	a	a	b	b	a	a

- Period set = $\{0, 4, 5\}$ (and corresponding autocorrelation 100011)

Known results about periods of classical strings

- Characterization of period sets [JCTA, Guibas and Odlyzko,1981].
 - (Auto)correlations
 - Lower bound on the number of period sets of a given length
 - Populations (i.e. how many strings share a given period set)

- The combinatorics of periods [JCTA, Rivals and Rahmann,2003].

- The convergence of the number of period sets for strings of given length [ICALP , Rivals, Sweering and Wang, 2022].

How to define periodicity of DS?

0	1	2	3	4	5	6	7	8
$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-	-	-	-
$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-	-	-	-
-	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-	-	-
-	-	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-	-
-	-	-	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-
-	-	-	-	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$

Table of Contents

- 1 Degenerate strings, three types of periodicity
- 2 Characterization of periodicities
- 3 Counting the number of period sets

Table of Contents

1 Degenerate strings, three types of periodicity

2 Characterization of periodicities

3 Counting the number of period sets

Degenerate string

Σ : a finite alphabet of size σ and n an integer.

Definition (Degenerate string)

$\hat{w} = \hat{w}[0 \dots n-1] \in (\mathcal{P}(\Sigma) \setminus \emptyset)^n$ is a string of length n over $\mathcal{P}(\Sigma) \setminus \emptyset$.

- $\hat{w}[i]$ is called an undetermined symbol
- We say two degenerate strings \hat{x} and \hat{y} of length n match, if for all $i \in \{0, \dots, n-1\}$ the intersection $\hat{x}[i] \cap \hat{y}[i]$ is non-empty.
- A hollow string \hat{w} is a degenerate string such that $\hat{w}[i] = \emptyset$ for at least one $i \in \{0, \dots, n-1\}$

Definition (Language)

The language of a degenerate string \hat{w} of length n is set of all classical strings that match with it:

$$\mathcal{L}(\hat{w}) = \{w \in \Sigma^n \mid \forall i \in \{0, \dots, n-1\} \quad w[i] \in \hat{w}[i]\}.$$

Example 1

$$\text{Let } \hat{w} = \begin{Bmatrix} a \\ b \end{Bmatrix} \cdot \begin{Bmatrix} b \\ c \end{Bmatrix} \cdot \begin{Bmatrix} b \\ c \end{Bmatrix} \cdot \{c\} \cdot \begin{Bmatrix} a \\ c \end{Bmatrix}.$$

$$\mathcal{L}(\hat{w}) = \{abbca, abbcc, abcca, abccc, acbca, acbcc, accca, acccc, \\ bbbca, bbbcc, bbcca, bbccc, bcbca, bcbcc, bccca, bcccc\}$$

Definition (Strong period)

A degenerate string \hat{w} has strong period p if there exists a string $w \in \mathcal{L}(\hat{w})$ with period p .

- To model one specific string, whose letters are not precisely known.
- $P^s(\hat{w})$ denotes the set of strong periods of \hat{w} .

Example 2

$$\text{Let } \hat{w} = \begin{Bmatrix} a \\ b \end{Bmatrix} \cdot \begin{Bmatrix} b \\ c \end{Bmatrix} \cdot \begin{Bmatrix} b \\ c \end{Bmatrix} \cdot \{c\} \cdot \begin{Bmatrix} a \\ c \end{Bmatrix}.$$

$accca \in \mathcal{L}(\hat{w})$, $accca$ has periods 0, 4. $P^s(\hat{w}) = \{0, 4\}$.

Definition (Weak period)

A degenerate string $\hat{w} = \hat{w}[0 \dots n-1]$ has weak period $p \in \{0, \dots, n-1\}$ if and only if $\hat{w}[0 \dots n-p-1]$ matches $\hat{w}[p \dots n-1]$

- Weak periods can model variations in a set of related strings.
- $P^w(\hat{w})$ denotes the set of weak periods of \hat{w} .
- Running example: $\hat{w} = \begin{Bmatrix} a \\ b \end{Bmatrix} \cdot \begin{Bmatrix} b \\ c \end{Bmatrix} \cdot \begin{Bmatrix} b \\ c \end{Bmatrix} \cdot \{c\} \cdot \begin{Bmatrix} a \\ c \end{Bmatrix}$.

Example

	0	1	2	3	4	5	6	7	8
\hat{w}	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-	-	-	-
\hat{w}	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-	-	-	-
	-	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-	-	-
	-	-	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-	-
	-	-	-	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$	-
	-	-	-	-	$\begin{Bmatrix} a \\ b \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\begin{Bmatrix} b \\ c \end{Bmatrix}$	$\{c\}$	$\begin{Bmatrix} a \\ c \end{Bmatrix}$

Definition (Weak period)

A degenerate string $\hat{w} = \hat{w}[0 \dots n-1]$ has weak period $p \in \{0, 1, \dots, n-1\}$ if and only if $\hat{w}[0 \dots n-p-1]$ matches $\hat{w}[p \dots n-1]$

- Weak periods can model variations in a set of related strings.
- $P^w(\hat{w})$ denotes the set of weak periods of \hat{w} .

Example 3

$$\hat{w} = \begin{Bmatrix} a \\ b \end{Bmatrix} \cdot \begin{Bmatrix} b \\ c \end{Bmatrix} \cdot \begin{Bmatrix} b \\ c \end{Bmatrix} \cdot \{c\} \cdot \begin{Bmatrix} a \\ c \end{Bmatrix}.$$

$$\text{Then } P^w(\hat{w}) = \{0, 1, 2, 4\},$$

Definition (Medium period)

A degenerate string $\hat{w} = \hat{w}[0n-1]$ has medium period $p \in \{0, 1, \dots, n-1\}$ if and only if for any $0 \leq i, j \leq n-1$ such that $i \equiv j \pmod{p}$ we have $\hat{w}[i] \cap \hat{w}[j] \neq \emptyset$.

- $P^m(\hat{w})$ denotes the set of medium periods of \hat{w} .
- $P^w(\hat{w}) = \{0, 1, 2, 4\} \implies P^m(\hat{w}) = \{0, 2, 4\}$
- For each degenerate string, $P^s \subseteq P^m \subseteq P^w$.
- The sets of weak, medium and strong period sets of all degenerate strings of length n are denoted by Ω_n^w , Ω_n^m , and Ω_n^s .

Table of Contents

- 1 Degenerate strings, three types of periodicity
- 2 Characterization of periodicities
- 3 Counting the number of period sets

Theorem (Characterization of PS)

Let $P^s \subseteq P^m \subseteq P^w \subseteq \{0, \dots, n-1\}$. Then P^w , P^m and P^s are respectively the weak, medium, and strong period sets of some non-hollow degenerate string \hat{w} of length n *if and only if*

- (i) $0 \in P^s$,
- (ii) $\forall p \in P^w$ with $p \geq n/2 \implies p \in P^s$,
- (iii) $p \in P^m$ iff $\forall k \in \mathbb{N}$ with $kp \in \{0, \dots, n-1\} \implies kp \in P^w$
- (iv) $p \in P^s$ iff $\forall k \in \mathbb{N}$ with $kp \in \{0, \dots, n-1\} \implies kp \in P^s$.

Proof of characterization of periodicities, part 1

Given a non-hollow degenerate string \hat{w} .

Proof: Necessity (\Rightarrow).

- (i) Every classical string has period 0, so take any string from $\mathcal{L}(\hat{w})$.
- (ii) $\hat{w} = \hat{w}[0 \dots n-p-1] \hat{w}[n-p \dots p-1] \hat{w}[p \dots n-1]$ with $\hat{w}[0 \dots n-p-1] \cap \hat{w}[p \dots n-1] \neq \emptyset$
 $\implies \exists$ string $w = w[0 \dots n-p-1] w[n-p \dots p-1] w[p \dots n-1]$:
 - $w[0 \dots n-p-1] = w[p \dots n-1] \in \hat{w}[0 \dots n-p-1] \cap \hat{w}[p \dots n-1]$
 - and $w[n-p \dots p-1] \in \hat{w}[n-p \dots p-1]$.
- (iii) This is the definition of medium periods.
- (iv) If a classical string $w \in \mathcal{L}(\hat{w})$ has period p , then it also has period kp for $k \in \mathbb{N}$, giving \hat{w} strong period kp as well.



Proof of characterization of periodicities, part 2

Proof: Sufficiency (\Leftarrow).

Sufficiency (\Leftarrow): Proven by constructing \hat{w} given the period sets:

$$\hat{w}[i] = \begin{cases} \{a, b\} & \text{if } i = 0 \\ \{a, c\} & \text{if } i \in P^s \setminus \{0\} \\ \{b, c\} & \text{if } i \in P^w \setminus P^s \\ \{c\} & \text{otherwise} \end{cases}$$

- $\hat{w}[p] \neq c \iff p \in (P^s \setminus \{0\}) \cup (P^w \setminus P^s) \iff p \in P^w \setminus \{0\}$.
- For every $p \in P^m$, every multiple is a weak period with symbol $\{b, c\}$.
- For $p \in P^s \setminus \{0\}$, there exists a string $w \in \mathcal{L}(\hat{w})$ such that:

$$w[i] = \begin{cases} a & \text{if } p \mid i \\ c & \text{otherwise.} \end{cases}$$

Proof.

...hence \hat{w} has strong period p . If $p \notin P^s$, then either

- $p \notin P^w$, or
- $p \in P^w \setminus P^s \iff \exists k \in \mathbb{N}$ such that $kp \in [\frac{n}{2}, n-1]$, $kp \in P^s \iff$

$$\hat{w}[0] \cap \hat{w}[p] \cap \hat{w}[kp] = \{a, b\} \cap \{b, c\} \cap \{a, c\} = \emptyset.$$



Corollary

Because the constructed string \hat{w} has alphabet size 3, period sets in general are independent of their alphabet Σ when $|\Sigma| \geq 3$.

Table of Contents

- 1 Degenerate strings, three types of periodicity
- 2 Characterization of periodicities
- 3 Counting the number of period sets**

Primitive sets and relation to period sets

- By the characterization theorem, we notice that P^m, P^s are multiplicative subsets of $[0, \dots, n-1]$ and hence $\Omega_n^m = \Omega_n^s$.
- For $P \subseteq [0, n-1]$ we write $\langle P \rangle = \{kp \in [0, n-1] \mid p \in P, k \in \mathbb{N}\}$. We say that P is a generator of $\langle P \rangle$.

Definition (Primitive set)

A set P of integers is primitive if it does not contain a pair $i \neq j$ such that i divides j or j divides i .

Lemma

For any $P \subseteq [0, n-1]$, there exists a unique P_{prim} such that P_{prim} is a primitive set and the *minimum generator* of P . Hence, $P \mapsto P_{\text{prim}}$ is a *one-to-one mapping* between primitive subsets and multiplicative subsets of $[0, n-1]$.

Define $Q(n)$ as the number of primitive sets with greatest element at most n .

Counting the number of period sets

Theorem (Counting and convergence)

- 1 The number of weak period sets is: $|\Omega_n^w| = 2^{n-1}$.
- 2 For medium/strong period sets: for any $\varepsilon > 0$, we have

$$|\Omega_n^m| = |\Omega_n^s| = \alpha^{n(1+O(\exp((-1+\varepsilon)\sqrt{\log n \log \log n})))}.$$

Proof.

- 1 The set of weak period sets “corresponds to” the power set of $\{1, \dots, n-1\}$.
- 2 We applied the knowledge from number theory and $|\Omega_n^m| = |\Omega_n^s| = Q(n-1)$.



Conclusion

- We provided three notions of periodicity (weak, medium, and strong) for degenerate strings.
- We characterized the period sets for each, exhibiting necessary and sufficient conditions.
- We counted the number of period sets for strings of length n for each type and studied its convergence using recent results from number theory.
- We investigated the structure of the families of period sets. i.e. Proved that all the types of period sets form lattices under set intersection and set union.
- We computed how many degenerate strings share a given period set (i.e. its population) using graph theory.

- 1 Investigate the combinatorial part of periodicity of languages,.i.e, any set of strings.
- 2 Study the algorithmic aspect of periodicity of languages,e.g.,Improve the complexity to determine period sets of degenerate string.
- 3 Apply our notions of periodicity of degenerate string to pattern matching, or string-graph matching problem.




Funding and acknowledgements

*This work is part of a project that has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreements No 872539 and No 956229, from the Netherlands Organisation for Scientific Research (NWO) through Gravitation-grant NETWORKS-024.002.003 and from the Constance van Eeden PhD Fellowship. Moreover, we would like to thank Solon P. Pissis for his helpful advice and suggestions.



Thank you for your attention!

Questions?

-  Guibas, L. and Odlyzko, A. (1981).
Periods in strings.
[Journal of Combinatorial Theory, Series A, 30:19–43.](#)
-  Rivals, E. and Rahmann, S. (2003).
Combinatorics of periods in strings.
[Journal of Combinatorial Theory, Series A, 104\(1\):95–113.](#)
-  Rivals, E., Sweering, M., and Wang, P. (2022).
Convergence of the number of period sets in strings.