# Asymptotic Approximation in Formal Languages

**17th Workshop Computational Logic and Applications @the Jagiellonian University in Kraków** 

### Ryoma Sin'ya (Akita University), 2023 Dec 14.



Akita University



### CIAA 2024 in Akita (Septermber 3 to 6) **28th International Conference on Implementation and Application of Automata**

- Topics: algorithms on automata, automata and logic, bioinformatics, complexity of automata and the world-wide web.
- The proceedings will be published by in the series of LNCS
  - Tentative submission deadline: April 14, 2024

operations, compilers, computer-aided verification, concurrency, data structure design for automata, data and image compression, design and architecture of automata software, digital libraries, DNA/molecular/membrane computing, document engineering, editors, environments, experimental studies and practical experience, implementation of verification methods and model checking, industrial applications, natural language and speech processing, networking, new algorithms for manipulating automata, object-oriented modeling, pattern-matching, pushdown automata and context-free grammars, quantum computing, structured and semi-structured documents, symbolic manipulation environments for automata, transducers and multi-tape automata, techniques for graphical display of automata, VLSI, viruses and related phenomena,

### About Akita

- In Tohoku (northern Japan) region
  - The latitude of Akita City is 39.72°N
- Reachable from Tokyo, Nagoya, Osaka:
  - By plane (1~1.5 hours)
  - By bullet train (3.5~6 hours)
- Lot of snow in winter: There is a ski resort 15 minutes by car from Akita university







# Snowboarding in Akita (Febrary 16, 2023)





### Outline

- 1. Background I: density and measurability
- 2. Background II: known properties
- 3. Complexity results
- 4. Conclusion

### nd measurability roperties

(7 min.) (5 min.) (8 min.) (3 min.)

### Outline

- 1. <u>Background I: density and measurability</u>
- 2. Background II: known properties
- 3. Complexity results
- 4. Conclusion

(7 min.) (5 min.) (8 min.) (3 min.)

The density of a language L over A is defined as

$$\delta_A(L) = \lim_{n \to \infty} \frac{1}{n} \sum_{i=0}^{n-1} \frac{\#(L \cap A^i)}{\#(A^i)}$$

Here #(X) denotes the cardinality of X.

 $\delta_A(L)$  can be regarded as the (average) **probability** that a randomly chosen word is in *L*.

The density of a language L over A is defined as

$$\delta_A(L) = \lim_{n \to \infty} \frac{1}{n} \sum_{i=0}^{n-1} \frac{\#(L \cap A^i)}{\#(A^i)}$$

Theorem [Berstel 1973]:

Every unambiguous context-free language do have an algebraic density.

Theorem [Kemp 1980]: There is a context-free language whose density is transcendental.

# Example 1: $\delta_A((AA)^*) = \frac{1}{2}$ .

Example 2:  $\delta_A(A^*wA^*) = 1$  for any w.

Example 3:  $L_1 = \{ w \in A^* \mid 3^n \le |w| < 3^{n+1} \text{ for some even } n \}$ does **not** have a density.

Every regular language do have a rational density, and It is computable.

The density of a language L over A is defined as

$$\delta_A(L) = \lim_{n \to \infty} \frac{1}{n} \sum_{i=0}^{n-1} \frac{\#(L \cap A^i)}{\#(A^i)}$$

Example 3:  $L_{\perp} = \{ w \in A^* \mid 3^n \le |w| < 3^{n+1} \text{ for some even } n \}$ does **not** have a density.

Problem:

Does every context-free language *L* have a density? ( $\delta_A(L)$  converges?)

Theorem [Nakamura 2019]:

it is undecidable whether a given context-free language L satisfies  $\delta_A(L) = 1$ .

# **Juages** Example 1: $\delta_A((AA)^*) = \frac{1}{2}$ .

Example 2:  $\delta_A(A^*wA^*) = 1$  for any w.

The density of a language L over A is defined as

$$\delta_A(L) = \lim_{n \to \infty} \frac{1}{n} \sum_{i=0}^{n-1} \frac{\#(L \cap A^i)}{\#(A^i)}$$

Example 3:  $L_{\perp} = \{ w \in A^* \mid 3^n \le |w| < 3^{n+1} \text{ for some even } n \}$ does **not** have a density.

Theorem [Kozik, CLA'05]: it is undecidable whether  $\lim_{n \to \infty} \frac{\#(L \cap A^n)}{\#(A^n)}$ 

converges or not for a given context-free language L.

# **Juages** Example 1: $\delta_A((AA)^*) = \frac{1}{2}$ .

Example 2:  $\delta_A(A^*wA^*) = 1$  for any w.

### $\mathscr{C}$ -measurability [S., SOFSEM'21] (cf. [Buck, 1946]) $A^*$



*L* is said to be  $\mathscr{C}$ -measurable if there exists an *infinite sequence of pairs of languages*  $(M_n, K_n)_{n \in \mathbb{N}}$  in  $\mathscr{C}$  such that  $M_n \subseteq L \subseteq K_n$  and  $\lim_{n \to \infty} \delta_A(K_n \setminus M_n) = 0$ .

### Example of a regular measurable language

Theorem:

$$\mathsf{B} = \{ w \in A \mid |w|_a = |w|_b \} \text{ over } A$$

the # of occurrences of a in w

Proof: Let  $L_k = \{w \in A^* \mid |w|_a = |w|_a = |w|_a$ 

Then, for each  $k \ge 1$ ,  $B \subseteq L_k$  as Thus the infinite sequence (Ø, l

 $A = \{a, b\}$  is regular measurable.

$$w|_{k} \mod k$$
 for each  $k \ge 1$ .

nd 
$$\delta_A(L_k) = \frac{1}{k} \to 0$$
 (if  $k \to \infty$ ).  
 $L_k)_{k \ge 1}$  converges to B.

### $\mathscr{C}$ -measurability [S., SOFSEM'21] (cf. [Buck, 1946]) $A^*$



*L* is said to be  $\mathscr{C}$ -measurable if there exists an *infinite sequence of pairs of languages*  $(M_n, K_n)_{n \in \mathbb{N}}$  in  $\mathscr{C}$  such that  $M_n \subseteq L \subseteq K_n$  and  $\lim_{n \to \infty} \delta_A(K_n \setminus M_n) = 0$ .

### **Original motivation of mesurability**

A non-empty word w is said to be *primitive* if it can not be represented as a power of shorter words, i.e., w = u<sup>n</sup> ⇒ u = w (and n = 1).
Q denotes the set of all primitive words over {a, b}.

### Example : $ababa \in Q$

- Primitive worsd conjecture [Dömösi-Horvath-Ito 1991]: Q is *not* context-free.
- My naive idea: while every context-free language is regular measurable, Q is regular immeasurable.

$$ababab = (ab)^3 \notin \mathbb{Q}$$

### Summary of [S., SOFSEM'21] **Regular measurable languages** A (simple) deterministic CFL $\mathbf{B} = \{ w \in A \mid |w|_a = |w|_b \}$ Many complex context-free languages. $= \{ w \in \{a, b\}^* \mid |w|_a > |w|_b \}$ There are **uncountably many** regular measurable languages. The set of all primitive words



### Outline

- 1. Background I: density and measurability
- 2. <u>Background II: known properties</u>
- 3. Complexity results
- 4. Conclusion

(7 min.)(5 min.) (8 min.) (3 min.)

# Some properties of *C*-measurability [S. DLT'21] Notation: $\mathcal{M}_A(\mathscr{C}) = \{L \subseteq A^* \mid L \text{ is } \mathscr{C}\text{-measurable}\}$

- $\mathcal{M}_A(\mathscr{C})$  can be defined as the Carathéodory extension of  $\mathscr{C}$ , a standard notion from measure theory.
- $\mathscr{C}$  is closed under these operations.

 "Is a given CFG generates a regular measurable languages?" is undecidable.

•  $\mathcal{M}_A(\mathscr{C})$  is closed under Boolean operations and left-and-right quotients if

# Some properties of *C*-measurability [S. DLT'21] Notation: $\mathcal{M}_A(\mathscr{C}) = \{L \subseteq A^* \mid L \text{ is } \mathscr{C}\text{-measurable}\}$

### Q: How about the decidability of $\mathscr{C}$ -measurability for some subclass *C* of regular languages?

 "Is a given CFG generates a regular measurable languages?" is undecidable.

• PT-measurability for DFAs is decidable in linear time [SYN 2022], where PT is the class of all *piecewise testable* languages.

Definition: *L* is *piecewise testable* [Simon 1972] if it can be represented as a finite Boolean combination of languages of the form  $L_{w} = A^{*}a_{1}A^{*}a_{2}...A^{*}a_{n}A^{*}$  where  $w = a_{1}a_{2}...a_{n}A^{*}a_{n}A^{*}$ .

- PT-measurability for DFAs is decidable in linear time [SYN 2022], where PT is the class of all *piecewise testable* languages.
- AT-measurability for DFAs is **coNP-complete** [SYN 2022], where AT is the class of all *alphabet testable* languages.

Definition: L is alphabet testable if it can be represented as a finite Boolean combination of languages of the form  $A^*aA^*$  (where  $a \in A$ ).

- PT-measurability for DFAs is decidable in linear time [SYN 2022], where PT is the class of all *piecewise testable* languages.
- AT-measurability for DFAs is coNP-complete [SYN 2022], where AT is the class of all *alphabet testable* languages.
- where GD is the class of all generalised definite languages. Boolean combination of languages of the form  $uA^*, A^*v$ .

• The decidability of GD-measurability for DFAs is **decidable** [S. CIAA'23],

Definition: L is generalised definitie if it can be represented as a finite

### $\mathcal{M}_A(\mathscr{C}) = \{L \subseteq A^* \mid L \text{ is } \mathscr{C}\text{-measurable}\}$ Notation:

- PT-measurability for DFAs is decidable in linear time [SYN 2022], where PT is the class of all *piecewise testable* languages.
- AT-measurability for DFAs is coNP-complete [SYN 2022], where AT is the class of all *alphabet testable* languages.
- The decidability of GD-measurability for DFAs is decidable [S. CIAA'23], where GD is the class of all generalised definite languages.
- Hierarchy is strict [S. DLT'22]:  $\mathcal{M}_A(AT) \subsetneq \mathcal{M}_A(PT) \subsetneq \mathcal{M}_A(GD)$ .

### Outline

- 1. Background I: density and measurability
- 2. Background II: known properties
- 3. <u>Complexity results</u>
- 4. Conclusion

### nd measurability roperties

(7 min.) (5 min.) (8 min.) (3 min.)

### Characterisation of AT- and PT-meaurability

Theorem [S. DLT'22]:

A language *L* is AT-measurable if and only if *L* or  $\overline{L}$  contains the language  $\bigcap_{a \in A} A^* a A^*$ .

A language L is PT-measurable if and only if L or  $\overline{L}$  contains the language  $L_w$  for some  $w \in A^*$ , where  $L_w$  is the set of all words containing w as a subsequence (scattered subword).

Corollary:

If a given language L is regular, then AT-measurability is decidable (because the inclusion problem is decidable for regular languages).



### Characterisation of AT- and PT-meaurability

Theorem [S. DLT'22]:

A language *L* is AT-measurable if and only if *L* or  $\overline{L}$  contains the language  $\bigcap_{a \in A} A^* a A^*$ .

A language L is PT-measurable if and only if L or  $\overline{L}$  contains the language  $L_w$  for some  $w \in A^*$ , where  $L_w$  is the set of all words containing w as a subsequence (scattered subword).

Theorem [SYN 2022]:

The AT-measurability for DFAs is **coNP-complete**, while the PT-measurability for DFAs is decidable in **linear time**.



### Definite languages [Brzozowski 1962] [Ginzburg 1966] Notation: $\mathscr{B}(\mathscr{C})$ denotes the (finite) Boolean closure of $\mathscr{C}$ . Reverse definite: Definite: $\mathsf{D} = \mathscr{B}\{A^*w \mid w \in A^*\} \qquad \mathsf{RD} = \mathscr{B}\{wA^* \mid w \in A^*\}$ Generalised definite: $GD = \mathscr{B}\{ uA^*v \mid u, v \in A^* \}$





(1) *S* is strongly connected:  $\forall p, q \in S \exists w \in A^* p \cdot w = q$ 

Sink components can be considered as a **minimal** SCC with respect to the reachability relation.



### **Characterisation of RD-measurable REGs**

Theorem: Let  $\mathscr{A} = (Q, \cdot, q_0, F)$  be a n (1)  $L(\mathscr{A})$  is RD-measurable. (2) Every sink component of



### RD-immeasurable

Theorem: Let  $\mathscr{A} = (Q, \cdot, q_0, F)$  be a minimal deterministic automaton. TFAE:

(2) Every sink component of  $\mathscr{A}$  is a singleton (contains only one state).



**RD**-measurable

### **Characterisation of RD-measurable REGs**

Theorem: Let  $\mathscr{A} = (Q, \cdot, q_0, F)$  be a n (1)  $L(\mathscr{A})$  is RD-measurable. (2) Every sink component of

Corollary: RD-measurability for minimal DFAs are decidable in linear time.

Corollary: D-measurability for DFAs are decidable.

Theorem: Let  $\mathscr{A} = (Q, \cdot, q_0, F)$  be a minimal deterministic automaton. TFAE:

(2) Every sink component of  $\mathscr{A}$  is a singleton (contains only one state).

### Characterisation of GD-measurable REGs

Theorem: Let  $\mathscr{A} = (Q, \cdot, q_0, F)$  be a deterministic automaton, and  $Q_1, \ldots, Q_k$  be its all sink components. Define  $P_i = \{ w \in A^* \mid q_0 \cdot w \in Q_i \},\$ 

- $S_{i} = \{ w \in A^{*} \mid Q_{i} \cdot w \subseteq F \} \text{ and } S_{i}' = \{ w \in A^{*} \mid Q_{i} \cdot w \cap F = \emptyset \}.$ Let  $M = \bigcup P_{i}S_{i}$  and  $M' = \bigcup P_{i}S_{i}'.$ i=1 i=1Then  $L(\mathscr{A})$  is GD-measurable if and only if  $\delta_A(M) + \delta_A(M') = 1$ .

### **Characterisation of GD-measurable REGs**

**Intuition**: *M* is a largest subset of  $L(\mathscr{A})$  that can be represented as a (possibly infinite) union of languages of the form  $uA^*w$ .



This condition means  $\delta_A(L(\mathscr{A}))$ 

### $L(\mathscr{A})$ That is, M is a largest GD-measurable subset of $L(\mathscr{A})$ . es Also, M' is a largest GD-measurable subset of $\overline{L(\mathscr{A})}$ .

and 
$$M' = \bigcup_{i=1}^{k} P_i S'_i$$
.  
ole if and only if  $\delta_A(M) + \delta_A(M') = 1$ .  
 $\delta = \delta_A(M)$  and  $\delta_A(\overline{L(\mathscr{A})}) = \delta_A(M')$ .





## **Characterisation of GD-measurable REGs**

Theorem: Let  $\mathscr{A} = (Q, \cdot, q_0, F)$  be a deterministic automaton, and  $Q_1, \ldots, Q_k$  be its all sink components. Define  $P_i = \{ w \in A^* \mid q_0 \cdot w \in Q_i \},\$ i = 1

Corollary: GD-measurability for DFAs is **decidable** (in PSPACE). and M, M' are regular by the construction)

- $S_i = \{ w \in A^* \mid Q_i \cdot w \subseteq F \} \text{ and } S'_i = \{ w \in A^* \mid Q_i \cdot w \cap F = \emptyset \}.$ Let  $M = \bigcup P_i S_i$  and  $M' = \bigcup P_i S'_i.$ i=1
- Then  $L(\mathscr{A})$  is GD-measurable if and only if  $\delta_A(M) + \delta_A(M') = 1$ .
- (because the density of a regular language is computable

### Outline

- 1. Background I: density and measurability
- 2. Background II: known properties
- 3. Complexity results
- 4. <u>Conclusion</u>

### nd measurability roperties

(7 min.) (5 min.) (8 min.) (3 min.)

### Summary

while AT-measurability for DFA is coNP-complete.

GD-measurability is decidable for DFAs (in PSPACE).

# PT-measurability and RD-measurability for DFAs is decidable in linear time,

### Progress: we (S., Y. Nakamura and Y. Yamaguchi) found that it is in PTIME.



# **Open problem**

equivalent or not?

Definition: L is star-free if and only if it can be represented as a finite Boolean combination and concatenation finite languages.

[S. DLT'22]:  $\mathcal{M}_A(AT) \subsetneq \mathcal{M}_A(P)$ 

How much GD-measurability and SF are weaker than regular measurability?

- Is the measuring power of GD and SF (the class of all star-free languages)

$$\mathsf{T}) \subsetneq \mathscr{M}_A(\mathsf{GD}) \subseteq \mathscr{M}_A(\mathsf{SF}).$$

Is this inclusion strict?

## **Application?**

The decidable characterisation of GD-measurability gives us the following **approximation scheme**:

Input : an automaton  $\mathscr{A}$  and an admissible error ratio  $\epsilon > 0$ . Output: an automaton  $\mathscr{B}$  (if exists) such that (1)  $L(\mathscr{A}) \subseteq L(\mathscr{B})$ , (2)  $L(\mathscr{B})$  is generalised definite, and (3)  $|\delta_A(L(\mathscr{A})) - \delta_A(L(\mathscr{A}))| \le \epsilon$ .

Can we apply this scheme to, say, obtain an efficient regular expression matching algorithm? (or other decision problems, e.g., inclusion ?)



