Grammatical inference and subregular phonology

Adam Jardine Rutgers University

December 9, 2019 · Tel Aviv University

Overview

"[V]arious formal and substantive universals are intrinsic properties of the language-acquisition system, these providing a schema that is applied to data and that determines in a highly restricted way the general form and, in part, even the substantive features of the grammar that may emerge upon presentation of appropriate data."

Chomsky, Aspects

"[I]f an algorithm performs well on a certain class of problems then it necessarily pays for that with degraded performance on the set of all remaining problems." Wolpert and Macready (1997), NFL Thms.

- Phonological patterns are governed by restrictive computational universals
- Formal language theory gives us tools to discover and state these universals
- **Grammatical inference** allows us to develop and study learning procedures that derive from these universals
- The result is algorithms...
 - that directly connect linguistic universals with learning
 - whose behavior in the general case is well-understood
 - that make typological and psycholinguistic predictions

Rough breakdown of course

- **Day 1:** Learning, languages, and grammars
- **Day 2:** Learning strictly local grammars
- **Day 3:** Automata and input strictly local functions
- **Day 4:** Learning functions and stochastic patterns, other open questions

By the end of this course, you should be able to engage with the literature, and start your own research project!

• Collaborators/Mentors:





Jeff Heinz Jim Rogers Rémi Eyraud Jane Chandlee Kevin McMullin (Stony Brook) (Earlham) (Marseilles) (Haverford) (Ottowa)

...at Rutgers:

Eileen Blum Chris Oakden Nate Koser

Dine Mamadou

Wenyue Hua Huteng Dai

- What do we mean when we say a child/animal/machine has 'learned' something?
- What do we mean when we say a child has learned their language?

- What do we mean when we say a child/animal/machine has 'learned' something?
- What do we mean when we say a child has learned their language?

- What is the nature of the sample?
- When is learning successful?

Grammatical inference

• Formal GI studies solutions to specific learning problems

Grammatical inference

Problem

Given a **positive** sample of a language, return a grammar that describes that language **exactly**

Languages and grammars

- Two kinds of phonological patterns:
 - Well-formedness (phonotactics)
 ex. *NC
 - Transformations (processes) ex. $/NC/ \rightarrow [NC]$

 Well-formedness patterns are sets ex. *N^C

well-formed:	{an, anda, amba, lalalalanda, blik, ffffff,}
ill-formed:	{anta, ampa, lalalalaŋka,}

• Processes are **relations** $/NC/ \rightarrow [NC]$

{(an, an), (anda, anda), (anta, anda), (lalalalampa, lalalalamba),...}

\cdot This is true regardless of how we describe them

 $C \rightarrow [+voice] / N _ \approx *N_{\circ} \gg ID[\pm voice]$

• We're going to first focus on sets as **formal languages**, and then move on to **(functional) relations**.

- An **alphabet** Σ is a finite set of symbols

 $\{0, 1\}$ $\{a, b, c\}$ $\{a, b, c, ..., \mathfrak{B}, \beta, \beta, ..., z\}$ $\{N, V, ADJ, ..., C\}$

- A string w over Σ is some sequence $\sigma_1 \sigma_2 \dots \sigma_n$ of symbols in Σ .
- Σ^* is all strings over Σ

$$\Sigma = \{a, b, c\}$$
$$\Sigma^* =$$

- A string w over Σ is some sequence $\sigma_1 \sigma_2 \dots \sigma_n$ of symbols in Σ .
- Σ^* is all strings over Σ

 $\Sigma = \{a, b, c\}$

$$\Sigma^* = \{ \begin{array}{l} \lambda, a, b, c, aa, ab, ac, \\ ba, bb, bc, ca, cb, cc, \\ aaa, aab, aac, \dots, \\ abbaaacccbabacb, \dots \end{array} \}$$

- A (formal) language some subset $L \subseteq \Sigma^*$
- Some formal languages for $\Sigma = \{a, b, c\}$:

-
$$(ab)^n = \{\lambda, ab, abab, ababab, ...\}$$

-
$$a^n b^n = \{\lambda, ab, aabb, aaabbb, aaaabbbb, ...\}$$

- ...

- Equivalently, a formal language maps strings in Σ^* to \top or \bot

Formal language classes

all possible languages

• How would you compute the *NC language?¹

 $\{ an, and a, amba, lalalaland a, blik, ffffff, ... <math display="inline">\}$

 $^{1}\Sigma = \{a, b, c, ..., a, \beta, o, ..., z\}$

- How would you compute the *NC language?¹
 { an, anda, amba, lalalalanda, blik, ffffff, ... }
- \cdot Make sure the string doesn't contain NC sequences!

 $\{anta, ampa, lalalalanka, ...\}$

 $^{1}\Sigma = \{a, b, c, ..., a, \beta, z, ..., z\}$

• u is a **substring** of w iff $w = v_1 u v_2$

a	b	b	a	b		
$\underbrace{v_1}$	<u>u</u>		$\underbrace{v_2}$			
a	h	h	a	h		

• u is a k-factor of w iff it is a substring of $\rtimes w \ltimes$ of size k

• A SL $_k$ grammar is a set of forbidden k-factors

 $G = \{bb, aa\}$

• L(G) is the set of strings $w \in \Sigma^*$ such that $w \models G$

 $G = \{bb, aa\}$ $w \quad w \models G?$ $w \qquad w \models G?$ λ abbbaa \bot \bot a \bot baaaa \bot aa. . . \top ababab \bot abba \bot aaa \perp aabbaba \top aba \bot • • •

- A language is **strictly local** iff it can be described by a SL_k grammar for some k
- Let's do some examples...

- A good many (but not all!) phonotactics are SL (Heinz, 2010)
- Long-distance phonotactics can be captured with two similar classes:
 - Strictly piecewise (SP) languages (Heinz, 2010)
 - Tier-based strictly local (TSL) languages

(Heinz et al., 2011; McMullin, 2016)

• For a general, formal review see Rogers et al. (2013)

Review

Problem

Given a **positive** sample of a language, return a grammar that describes that language **exactly**

- We're going to learn how SL languages have a solution to this problem
- We're going to learn other language classes that have a similar solution