

Grammatical inference and subregular phonology

Adam Jardine
Rutgers University

December 10, 2019 · Tel Aviv University

Review

Outline 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

Review of day 1

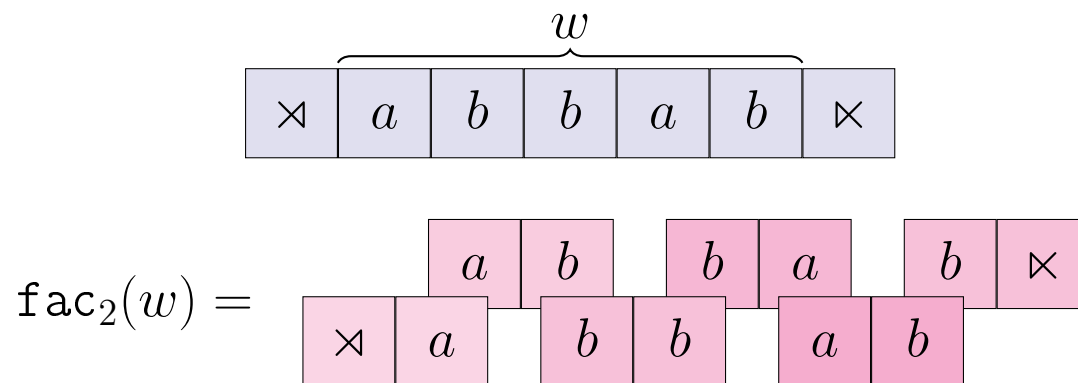
- Phonological patterns are governed by restrictive computational universals
- **Grammatical inference** connects these universals to solutions to the **learning problem**:

Problem

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

Review of day 1

- **Strictly local languages** are patterns computed solely by k -**factors** in a string

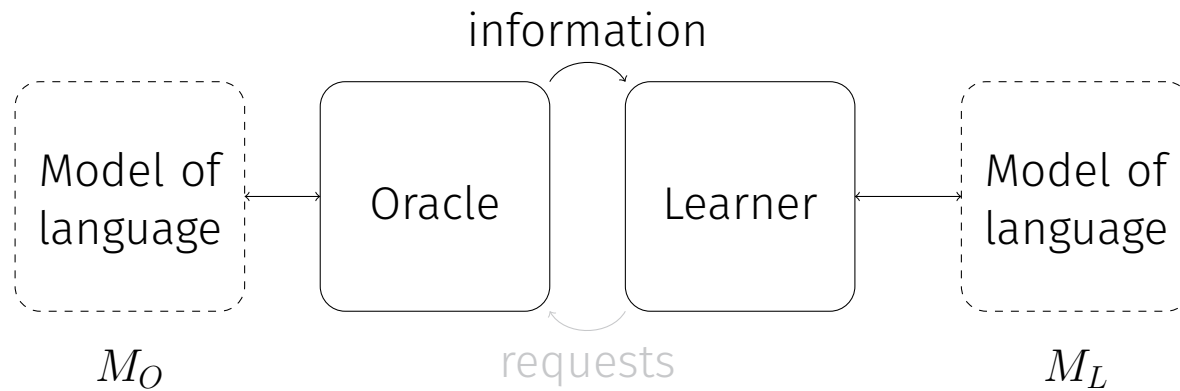


Today

- A provably correct method for learning SL_k languages
- The paradigm of **identification in the limit from positive data** (Gold, 1967; de la Higuera, 2010)
- Why learners target **classes** (not specific languages, or all possible languages)

Learning paradigm

Learning paradigm



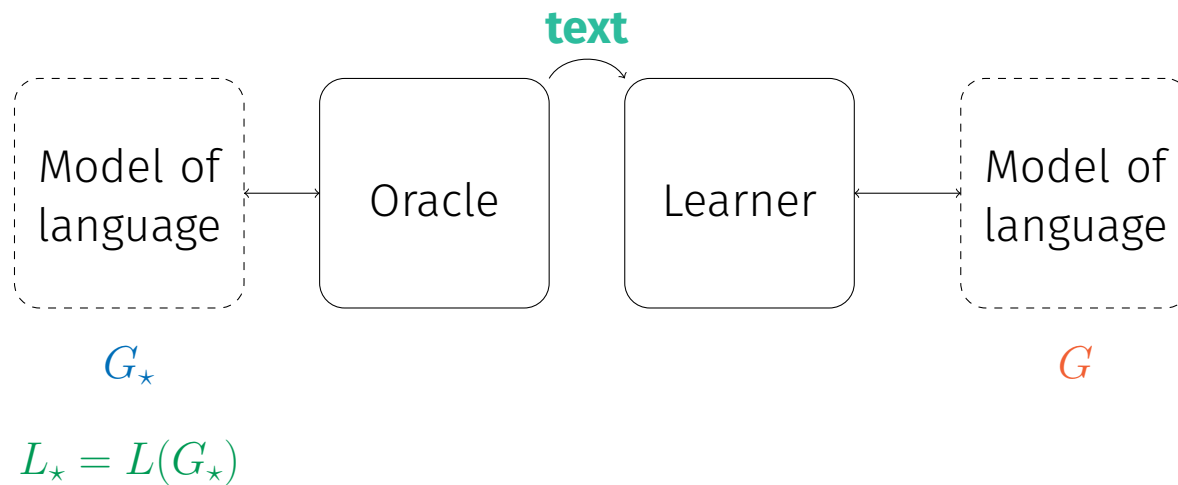
(from Heinz et al., 2016)

Problem

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

- This is **(exact) identification in the limit from positive data** (ILPD; Gold, 1967)

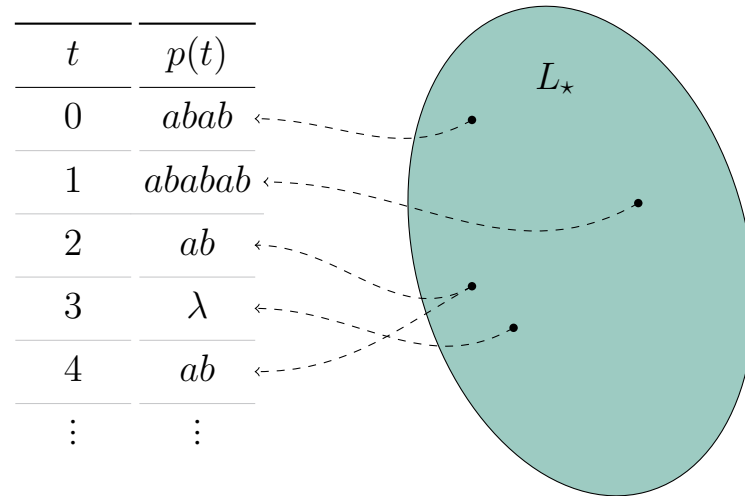
Identification in the limit from positive data (ILPD)



- A **text** of L_* is some sample of positive examples of L_*

Identification in the limit from positive data (ILPD)

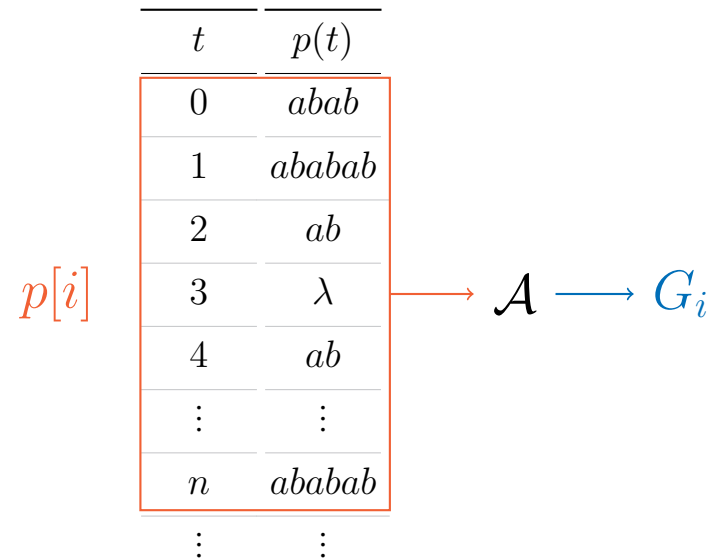
A **presentation** of L_* is a *sequence* p of examples drawn from L_*



(this is the ‘in the limit’ part)

Identification in the limit from positive data (ILPD)

A learner \mathcal{A} takes a finite sequence and outputs a grammar



Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_\star = \{ab, bab, aaa\}$

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_{\star} = \{ab, bab, aaa\}$

$$\frac{t \quad p(t) \quad G_t}{0 \quad bab}$$

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_\star = \{ab, bab, aaa\}$

$$\begin{array}{ccc} t & p(t) & G_t \\ \hline 0 & bab & \{bab\} \end{array}$$

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_\star = \{ab, bab, aaa\}$

t	$p(t)$	G_t
0	bab	$\{bab\}$
1	ab	

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_\star = \{ab, bab, aaa\}$

t	$p(t)$	G_t
0	bab	$\{bab\}$
1	ab	$\{ab, bab\}$

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_\star = \{ab, bab, aaa\}$

t	$p(t)$	G_t
0	bab	$\{bab\}$
1	ab	$\{ab, bab\}$
2	bab	

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_\star = \{ab, bab, aaa\}$

t	$p(t)$	G_t
0	bab	$\{bab\}$
1	ab	$\{ab, bab\}$
2	bab	$\{ab, bab\}$

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_\star = \{ab, bab, aaa\}$

t	$p(t)$	G_t
0	bab	$\{bab\}$
1	ab	$\{ab, bab\}$
2	bab	$\{ab, bab\}$
3	aaa	$\{ab, bab, aaa\}$

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_\star = \{ab, bab, aaa\}$

t	$p(t)$	G_t
0	bab	$\{bab\}$
1	ab	$\{ab, bab\}$
2	bab	$\{ab, bab\}$
3	aaa	$\{ab, bab, aaa\}$
4	ab	$\{ab, bab, aaa\}$

Identification in the limit from positive data (ILPD)

Let's take the learner \mathcal{A}_{Fin} :

$$\mathcal{A}_{\text{Fin}}(p[n]) = \{w \mid w = p(i) \text{ for some } i \leq n\}$$

Let's set $L_{\star} = \{ab, bab, aaa\}$

t	$p(t)$	G_t
0	bab	$\{bab\}$
1	ab	$\{ab, bab\}$
2	bab	$\{ab, bab\}$
3	aaa	$\{ab, bab, aaa\}$
4	ab	$\{ab, bab, aaa\}$
	...	
308	bab	$\{ab, bab, aaa\}$

Identification in the limit from positive data (ILPD)

A **converges** at point n if $G_m = G_n$ for any $m > n$

t	$p(t)$	G_t
0	abab	G_0
1	ababab	G_1
2	ab	G_2
\vdots	\vdots	\vdots
n	ababab	G_n
$n + 1$	abababab	G_n
\vdots	\vdots	\vdots
m	λ	G_n
\vdots	\vdots	\vdots

convergence

Identification in the limit from positive data (ILPD)

ILPD-learnability

A class \mathcal{C} is **ILPD-learnable** if there is some algorithm $\mathcal{A}_{\mathcal{C}}$ such that for *any* stringset $L \in \mathcal{C}$, given *any* positive presentation p of L , $\mathcal{A}_{\mathcal{C}}$ converges to a grammar G such that $L(G) = L$.

- How is ILPD learning an idealization?
- What are the advantages of using ILPD as a criterion for learning?

Learning strictly local languages

Learning SL languages

- Given any k , the class SL_k is IDLP-learnable.
- Using \mathcal{A}_{Fin} as an example, how might we learn a SL_k language?

Learning SL languages

$$G_{\star} = \{CC, C\bowtie\}$$

t	datum	hypothesis
0	V	
1	$CVCV$	
2	$CVVCVCV$	
3	$VCVCV$	
\vdots		

Learning SL languages

$$G_{\star} = \{CC, C\star\}$$

t	datum	hypothesis
0	V	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
1	$CVCV$	
2	$CVVCVVCV$	
3	$VCVVCV$	
	\vdots	

Learning SL languages

$$G_{\star} = \{CC, C\star\}$$

t	datum	hypothesis
0	V	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
1	$CVCV$	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
2	$CVVCVVCV$	
3	$VCVVCV$	
	\vdots	

Learning SL languages

$$G_{\star} = \{CC, C\star\}$$

t	datum	hypothesis
0	V	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
1	$CVCV$	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
2	$CVVCVVCV$	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
3	$VCVVCV$	
	\vdots	

Learning SL languages

$$G_{\star} = \{CC, C\star\}$$

t	datum	hypothesis
0	V	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
1	$CVCV$	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
2	$CVVCVVCV$	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
3	$VCVVCV$	$\{\star C, \star V, CC, CV, C\star, VC, VV, V\star\}$
	\vdots	

Learning SL languages

$$\mathcal{A}_{\text{SL}_k}(p[i]) = \text{fac}_k(\Sigma^*) - \text{fac}_k\{p(0), p(1), \dots, p(i)\}$$

Learning SL languages

$$\mathcal{A}_{\text{SL}_k}(p[i]) = \text{fac}_k(\Sigma^*) - \text{fac}_k\{p(0), p(1), \dots, p(i)\}$$

- The **characteristic sample** is ...

Learning SL languages

$$\mathcal{A}_{\text{SL}_k}(p[i]) = \text{fac}_k(\Sigma^*) - \text{fac}_k\{p(0), p(1), \dots, p(i)\}$$

- The **characteristic sample** is $\text{fac}_k(L_\star)$

Learning SL languages

$$\mathcal{A}_{\text{SL}_k}(p[i]) = \text{fac}_k(\Sigma^*) - \text{fac}_k\{p(0), p(1), \dots, p(i)\}$$

- The **characteristic sample** is $\text{fac}_k(L_\star)$
- The **time complexity** is **linear**—the time it takes to calculate is directly proportional to the size of the data sample.

Learning SL languages

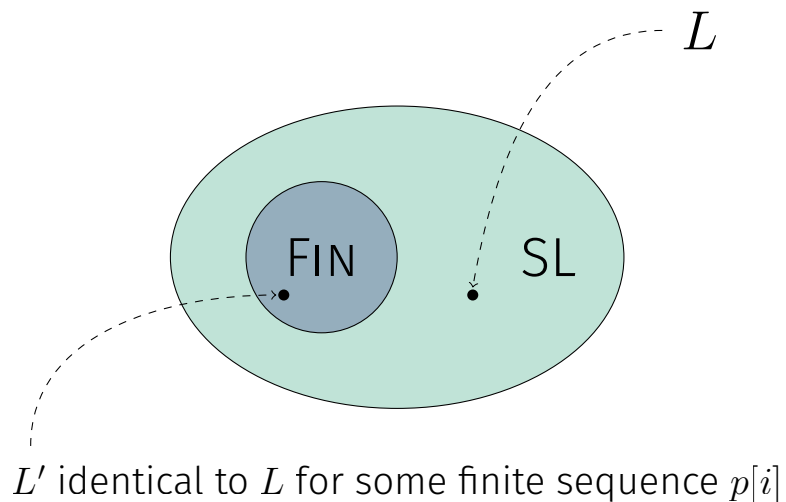
Let's learn Pintupi. Note that $k = 3$. What is the initial hypothesis?
At what point do we converge?

t	datum	hypothesis
0	σ	
1	$\sigma\sigma$	
2	$\sigma\sigma\sigma$	
3	$\sigma\sigma\sigma\sigma$	
4	$\sigma\sigma\sigma\sigma\sigma$	
5	$\sigma\sigma\sigma\sigma\sigma\sigma$	
\vdots		

The limits of SL learning

The limits of SL learning

- We must know k in advance



- Gold (1967): any class \mathcal{C} such that $\text{Fin} \subsetneq \mathcal{C}$ is not learnable from positive examples

The limits of SL learning

- Consider this pattern from Inseño Chumash:

ʃ-ʌpi-tʃ ^h ol-it	‘I have a stroke of good luck’
s-ʌpi-ts ^h ol-us	‘he has a stroke of good luck’
ʃ-ʌpi-tʃ ^h ol-uf-wʌʃ	‘he had a stroke of good luck’
hʌ-ʃxintila-wʌʃ	‘his former Indian name’
s-is-tisi-jep-us	‘they (two) show him’
k-ʃu-ʃojin	‘I darken it’

- What phonotactic constraints are active here?

The limits of SL learning

- Consider this pattern from Inseño Chumash:

ʃ-ʌpi-tʃ ^h ol-it	‘I have a stroke of good luck’
s-ʌpi-ts ^h ol-us	‘he has a stroke of good luck’
ʃ-ʌpi-tʃ ^h ol-uf-wʌʃ	‘he had a stroke of good luck’
hʌ-ʃxintila-wʌʃ	‘his former Indian name’
s-is-tisi-jep-us	‘they (two) show him’
k-ʃu-ʃojin	‘I darken it’

- What phonotactic constraints are active here?

*ʃ...s, *s...ʃ

The limits of SL learning

- Let's assume $L_{\star} = L_C$ for $\Sigma = \{s,o,t,f\}$ as given below

$$L_C = \{so, ss, \dots, sos, \int o \int, \int o \int o \int, sosos, \int o t o t o \int, s o t o t o s, \dots\}$$

t	datum	hypothesis
0	sos	$\{ss, so, sf, \dots, \int s, \int t, \int \int\}$
1	sotoss	$\{ss, so, sf, \dots, \int s, \int t, \int \int\}$
2	$\int o \int t o \int \int$	$\{ss, so, sf, \dots, \int s, \int t, \int \int\}$
\vdots		

The limits of SL learning

- Let's assume $L_* = L_C$ for $\Sigma = \{s,o,t,f\}$ as given below

$$L_C = \{so, ss, \dots, sos, \text{fof}, \text{fofof}, \text{sosos}, \text{fototof}, \text{sototos}, \dots\}$$

t	datum	hypothesis
0	sos	$\{ss, so, sf, \dots, fs, ft, ff\}$
1	sotoss	$\{ss, so, sf, \dots, fs, ft, ff\}$
2	fotofff	$\{ss, so, sf, \dots, fs, ft, ff\}$
\vdots		

- Learner will never see sf or fs , so in the limit $G = \{sf, fs\}$.

The limits of SL learning

$$L_C = \{so, ss, \dots, sos, foj, fofoj, sosos, fototoj, sototos, \dots\}$$

$$G = ? \{sf, js\}$$

- ✓ $sosos \in L_C$
- ✓ $sofs \notin L_C$
- ✗ $sofos \notin L_C$

The limits of SL learning

$$L_C = \{so, ss, \dots, sos, \text{fof}, \text{fofof}, \text{sosos}, \text{fototof}, \text{sototos}, \dots\}$$

$$G_{k=3} =? \{sof, ssf, stf, sff, \dots, fos, fss, fts, ffs\}$$

- ✓ $sosos \in L_C$
- ✓ $sofs \notin L_C$
- ✓ $sofos \notin L_C$
- ✗ $fotos \notin L_C$

The limits of SL learning

$$L_C = \{so, ss, \dots, sos, \text{fof}, \text{fofof}, \text{sosos}, \text{fototof}, \text{sototos}, \dots\}$$

- There is no k such that \mathcal{A}_{SL_k} learns a grammar for L_C
- This is because there is **no SL grammar** for L_C !

The limits of SL learning

- \mathcal{A}_{SL_k} **only learns SL_k languages**
- This is the advantage of studying learning with formal grammatical inference:
 - we what patterns it **can learn**,
 - what patterns it **cannot learn**,
 - on **exactly what data**

Review

- As a hypothesis of phonotactic learning, \mathcal{A}_{SL_k}
 - makes **restrictive** predictions about what patterns can and cannot be learned
 - suggests phonological learning is **modular** (Heinz, 2010)
 - directly connects computational typological generalizations with a theory of learning

Review

Problem

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

- We have formalized this problem as **identification in the limit from positive data**
- We have solved this problem for any SL_k class
- We'll find another solution with automata, and extend that to learn processes