

# Formal and Empirical Grammatical Inference

Jeffrey Heinz, Colin de la Higuera and Menno van Zaanen

`heinz@udel.edu, cdlh@univ-nantes.fr, mvzaanen@uvt.nl`

## Outline of the tutorial

- I. Formal GI and learning theory (de la Higuera)
- II. Empirical approaches to regular and subregular natural language classes (Heinz)
- III. Empirical approaches to nonregular natural language classes (van Zaanen)

## I Formal GI and learning theory

- What is grammatical inference?
- What does learning or having learnt imply?
- Reasons for considering formal learning
- Some criteria to study learning in a probabilistic and a non probabilistic setting

## A simple definition

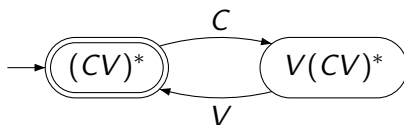
**Grammatical inference is about learning a grammar given information about a language**

## Vocabulary

- Learning = building, inferring
- Grammar= finite representation of a possibly infinite set of strings, or trees, or graphs
- Information=you can learn from text, from an informant, by actively querying
- Language= possibly infinite set of strings, or trees, or graphs

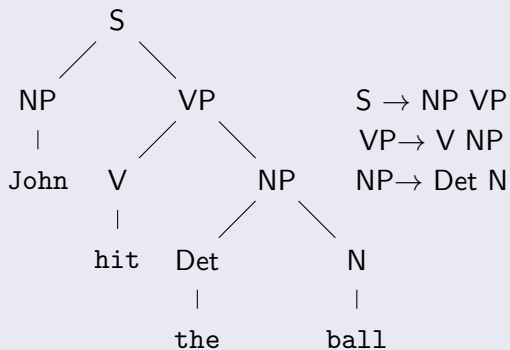
## A DFA (Ack: Jeffrey Heinz)

The  $(CV)^*$  language representing licit sequences of sounds in many languages in the world. Consonants and vowels must alternate; words must begin with C and must end with V. States show the regular expression indicating its “good tails”.



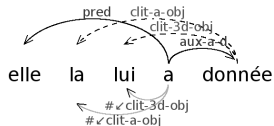
## A context free grammar and a parse tree

(de la Higuera 2010)



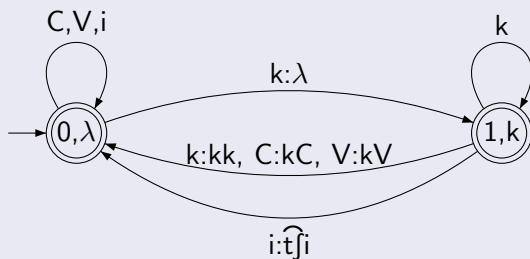
## A categorial dependency grammar (Béchet et al. 2011)

- $elle \mapsto [pred]$ ,
- $la \mapsto [\#(\surd \textit{clit} - a - \textit{obj})] \surd \textit{clit} - a - \textit{obj}$ ,
- $lui \mapsto [\#(\surd \textit{clit} - 3d - \textit{obj})] \surd \textit{clit} - 3d - \textit{obj}$ ,
- $a \mapsto [\#(\surd \textit{clit} - 3d - \textit{obj}) \setminus \#(\surd \textit{clit} - a - \textit{obj}) \setminus pred \setminus S / aux - a - d]$ ,
- $donnée \mapsto [aux - a - d] \surd \textit{clit} - 3d - \textit{obj} \surd \textit{clit} - a - \textit{obj}$



## A finite state transducer (Ack: Jeffrey Heinz)

A subsequential transducer illustrating a common phonological rule of palatalization ( $k \rightarrow \widehat{tj} / \text{---} i$ ). States are labelled with a number and then the output string given by the  $\sigma$  function for that state.



So for example:

$w$	$t(w)$
kata	kata
kita	$\widehat{t}jita$
tak	tak
taki	$tat\widehat{j}i$
...	

## Our definition

**Grammatical inference is about learning a grammar given information about a language**

## Questions

- Why **grammar** and not **language**?
- Why **a** and not **the**?

## Why not write “learn a language”?

Because you always learn a *representation* of a language

## Paradox

Take two learners learning a context-free language, one is learning a quadratic normal form and the other a Greibach normal form, they cannot agree that they have learnt the same thing (*undecidable question*).

Worth thinking about... is it a paradox? Do two English speakers agree they speak the same language?

## Our definition

**Grammatical inference is about learning a grammar given information about a language**

## How can *a* become *the*?

- Ask for the grammar to be the smallest, best (re a score). → Combinatorial characterisation
- The learning problem becomes an optimisation problem!
- Then we often have theorems saying that
  - If our algorithm does solve the optimisation problem, what we have learnt is correct
  - If we can prove that we can't solve the optimisation problem, then the class is not learnable

## Optimal with respect of some score

Score should take into account:

- Simplicity
- Coverage
- Usefulness

## What scores?

- Occam argument
- Compression argument
- Kolmogorov complexity
- MDL argument

## Moreover

- GI is not only about building a grammar from some data. It is concerned with saying something about:
  - the quality of the result,
  - the quality of the learning process,
  - the properties of the process.

## Naive example

- Suppose you are building a random number generator.
- How are you convinced that it works?
  - Because it follows sound principles as defined by number theory specialists?
  - Because you have tested and the number 772356191 has been produced?
  - Because you have proved that the series of numbers that will be produced is incompressible?
- Empirical approach
- Experimental approach
- Formal approach

## Empirical approach: using good (sound) ideas

- For example, genetic algorithms or neural networks
- Or some mathematical *principle* (Occam, Kolmogorov, MDL, ...)
- Can become a *principled* approach

## Alternative point of view

Empirical approach is about imitating what nature (or humans) do

## Experimental approach

- Benchmarks
- Competitions

- Necessary but not sufficient
- How do we know that all the cases are covered?
- How do we know that we dont have a hidden bias?

## Formal approach: showing that the algorithm has converged

- Is impossible:
  - Just one run
  - Can't prove that 23 is random
- But we can say something about the algorithm:
  - That in the near future, given some string, we can predict if this string belongs to the language or not;
  - Choose between defining clearly “near future” and accepting probable truths (or error bounds) or leaving it undefined and using identification.

What else would we like to say?

That if the solution we have returned is not good, then that is because the initial data was bad (insufficient, biased)

Idea:

Blame the data, not the algorithm

## Suppose we cannot say anything of the sort?

- Then that means that we may be terribly wrong even in a favourable setting
- Thus there is a *hidden bias*
- Hidden bias: the learning algorithm is supposed to be able to learn anything inside class  $\mathcal{L}_1$ , but can really only learn things inside class  $\mathcal{L}_2$ , with  $\mathcal{L}_2 \subset \mathcal{L}_1$

## Saying something about the process itself

- Key idea: if there is **something to learn** and the data is not corrupt, then, given enough time, we will learn it
- Replace the notion of *learning* by that of *identifying*

## In practise, does it make sense?

- No, because we never know if we are in the ideal conditions (something to learn + good data + enough of it)
- Yes, because at least we get to blame the data, not the algorithm

## Complexity issues

- Complexity theory should be used: the total or update runtime, the size of the data needed, the number of mind changes, the number and weight of errors. . .
- . . . should be measured and limited.

## A linguistic criterion

- One argument appealing to linguists (we hope) is that if the criteria are not met for some class of languages that a human is supposed to know how to learn, something is wrong somewhere
- (preposterously, the maths can't be wrong. . . )

## Non probabilistic settings

- Identification in the limit
- Resource bounded identification in the limit
- Active learning (query learning)

## Identification in the limit

- Information is presented to the learner who updates its hypothesis after inspecting each piece of data
- At some point, always, the learner will have found the correct concept and not change from it

(Gold 1967 & 1978)

## Example

Number	Presentation	Analysis of hypothesis	New hypothesis ( <i>regex</i> )
1	$a +$		$a$
2	$aaa +$	inconsistent	$a^*$
3	$aaaa -$	inconsistent	$a(aa)^*$
4	$aaaaaa -$	consistent	$a(aa)^*$
9234	$aaaaaaaa -$	consistent	$a(aa)^*$
45623416	$aaaaaaaa +$	consistent	$a(aa)^*$

## A presentation is

a function  $\phi : \mathbb{N} \rightarrow X$

- where  $X$  is some set,
- and such that  $\phi$  is associated to a language  $L$  through a function  $\text{YIELDS} : \text{YIELDS}(\phi) = L$
- If  $\phi(\mathbb{N}) = \psi(\mathbb{N})$  then  $\text{YIELDS}(\phi) = \text{YIELDS}(\psi)$

## *text* presentation

- A *text* presentation of a language  $L \subseteq \Sigma^*$  is a function  $\phi : \mathbb{N} \rightarrow \Sigma^*$  such that  $\phi(\mathbb{N}) = L$
- $\phi$  is an infinite succession of all the elements of  $L$
- (note : small technical difficulty with  $\emptyset$ )

## *informed* presentation

- An *informed* presentation (or an informant) of  $L \subseteq \Sigma^*$  is a function  $\phi : \mathbb{N} \rightarrow \Sigma^* \times \{-, +\}$  such that  $\phi(\mathbb{N}) = (L, +) \cup (\bar{L}, -)$
- $\phi$  is an infinite succession of all the elements of  $\Sigma^*$  labelled to indicate if they belong or not to  $L$

## Active presentation

- The learner interacts with the environment (modelled as an oracle) through queries
- A membership query
  - Learner presents string  $x$
  - Oracle answer yes or no
- A correction query (Becerra-Bonache et al. 2005 & 2008)
  - Learner presents string  $x$
  - Oracle answer yes or returns a close correction
- An equivalence query
  - Learner presents hypothesis  $H$
  - Oracle answer yes or returns a counter-example

### Example: presentations for $\{a^n b^n : n \in \mathbb{N}\}$

- Legal presentation from text:  $\lambda, a^2 b^2, a^7 b^7, \dots$
- Illegal presentation from text:  $ab, ab, ab, \dots$
- Legal presentation from informant :  $(\lambda, +), (abab, -), (a^2 b^2, +), (a^7 b^7, +), (aab, -), (abab, -), \dots$

## Example: presentation for Spanish

- Legal presentation from text: *En un lugar de la Mancha...*
- Illegal presentation from text: *Gooooool*
- Legal presentation from informant :  $(en,+)$ ,  $(whatever,-)$ ,  $(un,+)$ ,  $(lugar,+)$ ,  $(lugor,-)$ ,  $(xwszrrzt,-)$ ,

## What happens before convergence?

*On two occasions I have been asked [by members of Parliament], 'Pray, Mr. Babbage, if you put into the machine wrong figures, will the right answers come out?' I am not able rightly to apprehend the kind of confusion of ideas that could provoke such a question.*

**Charles Babbage**

## Further definitions

- Given a presentation  $\phi$ ,  $\phi_n$  is the set of the first  $n$  elements in  $\phi$ .
- A learning algorithm (learner) **A** is a function that takes as input a set  $\phi_n$  and returns a grammar of a language.
- Given a grammar  $G$ ,  $\mathbb{L}(G)$  is the language generated/recognised/ represented by  $G$ .

## Convergence to a hypothesis

- **A** converges to  $G$  with  $\phi$  if
  - $\forall n \in \mathbb{N} : \mathbf{A}(\phi_n)$  halts and gives an answer
  - $\exists n_0 \in \mathbb{N} : n \geq n_0 \implies \mathbf{A}(\phi_n) = G$
- If furthermore  $\mathbb{L}(G) = \text{YIELDS}(\phi)$  then we have identified.

## Identification in the limit

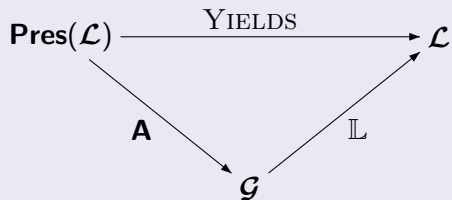


Figure: The learning setting.

from (de la Higuera 2010)

## Consistency and conservatism

- We say that the learner  $\mathbf{A}$  is *consistent* if  $\phi_n$  is consistent with  $\mathbf{A}(\phi_n) \forall n$
- A consistent learner is always consistent with the past

## Consistency and conservatism

- We say that the learner  $\mathbf{A}$  is *conservative* if whenever  $\phi(n+1)$  is consistent with  $\mathbf{A}(\phi_n)$ , we have  $\mathbf{A}(\phi_n) = \mathbf{A}(\phi_{n+1})$
- A conservative learner doesn't change his mind needlessly

## Learning from data

- A learner is order dependent if it learns something different depending on the order in which it receives the data.
- Usually an order independent learner is better.

## What about efficiency?

- We can try to bound
  - global time
  - update time
  - errors before converging (IPE)
  - mind changes (MC)
  - queries
  - good examples needed (characteristic samples)

(Pitt 1989, de la Higuera et al. 2008)

## Definition: polynomial number of implicit prediction errors

- Denote by  $G \not\equiv x$  if  $G$  is incorrect with respect to an element  $x$  of the presentation (i.e. the learner producing  $G$  has made an implicit prediction error).

$\mathcal{G}$  is polynomially identifiable in the limit from **Pres** if there exists an identification learner  $\mathbf{A}$  and a polynomial  $p()$  such that given any  $G$  in  $\mathcal{G}$ , and given any presentation  $\phi$  of  $\mathbb{L}(G)$ ,  
 $\#i : \mathbf{A}(\phi_i) \not\equiv \phi(i+1) \leq p(|G|)$ .

### Definition: polynomial characteristic sample

$\mathcal{G}$  has polynomial characteristic samples for identification learner  $\mathbf{A}$  if there exists a polynomial  $p(\cdot)$  such that: given any  $G$  in  $\mathcal{G}$ ,  $\exists Y$  correct sample for  $G$ , such that whenever  $Y \subset \phi_n$ ,  $\mathbf{A}(\phi_n) \equiv G$  and  $\|Y\| \leq p(\|G\|)$

- As soon as the CS is in the data, the result is correct;
- The CS is small.

## Polynomial queries

(Angluin 1987)

- Algorithm **A** learns with a polynomial number of queries if the number of queries made before halting with a correct grammar is polynomial in
  - the size of the target,
  - the size of the information received.

## Main negative results

- Cannot learn NFA, CFGs from an informant in most polynomial settings (Pitt 1989, de la Higuera 1997)
- Cannot learn DFA from text (Gold 1967)
- Cannot learn DFA from membership nor equivalence queries (Angluin 1981 & 1987).

## Main positive results

- Can learn DFA from an informant with polynomial resources (Oncina and García 1992);
- Can learn DFA from membership and equivalence queries (Angluin 1987).

## Probabilistic settings

- PAC learning (about learning yes-no machines with fixed but unknown distributions)
- Identification with probability 1 (about identifying distributions)
- PAC learning distributions (about approximately learning distributions)

## Learning a language from sampling

- We have a distribution over  $\Sigma^*$
- We sample twice:
  - once to learn,
  - once to see how well we have learned
- The PAC setting: Les Valiant, Turing award 2010



## PAC-learning

(Valiant 1984, Pitt 1989)

- $\mathcal{L}$  a class of languages
- $\mathcal{G}$  a class of grammars
- $\epsilon > 0$  and  $\delta > 0$
- $m$  a maximal length over the strings
- $n$  a maximal size of machines

$H$  is  $\epsilon$ -AC (approximately correct)\*  
if  
 $Pr_D[H(x) \neq G(x)] < \epsilon$

## Polynomial PAC learning

There is a polynomial  $p(\cdot, \cdot, \cdot, \cdot)$  such that

- in order to learn  $\epsilon$ -AC machines of size at most  $n$  with error at most  $\delta$  we require at most  $p(m, n, \frac{1}{\delta}, \frac{1}{\delta})$  data and time;
- we want the errors to be less than  $\epsilon$  and bad luck to be less than  $\delta$ .

## (French radio)

- Unless there is a surprise there should be no surprise
- French radio, (after the last primary elections, on 3rd of June 2008)
- First surprise is  $\delta$ , second surprise is  $\epsilon$

## Results (Kearns and Valiant 1989, Kearns and Vazirani 1994)

- Using cryptographic assumptions, we cannot PAC-learn DFA
- Cannot PAC-learn NFA, CFGs with membership queries either
- Learning can be seen as finding the encryption function from examples (Kearns & Vazirani)

## Alternatively

- Instead of learning classifiers in a probabilistic world, learn directly the distributions!
- Learn probabilistic finite automata (deterministic or not)

## No error (Angluin 1988)

- This calls for identification in the limit with probability 1
- Means that the probability of not converging is 0
- Goal is to identify the structure and the probabilities
- Mainly a (nice) theoretic setting

## Results

- If probabilities are computable, we can learn with probability 1 finite state automata (Carrasco and Oncina, 1994)
- But not with bounded (polynomial) resources (de la Higuera and Oncina, 2004)

## With error

- PAC definition applies
- But error should be measured by a distance between the *target distribution* and the *hypothesis*
- How do we measure the distance:  $L_1$ ,  $L_2$ ,  $L_\infty$ , Kullback-Leibler?

## Results

- Too easy to learn with  $L_\infty$
- Too hard to learn with  $L_1$
- Both results hold for the same algorithm! (de la Higuera and Oncina, 2004)
- Nice algorithms for biased classes of distributions

## Open problems

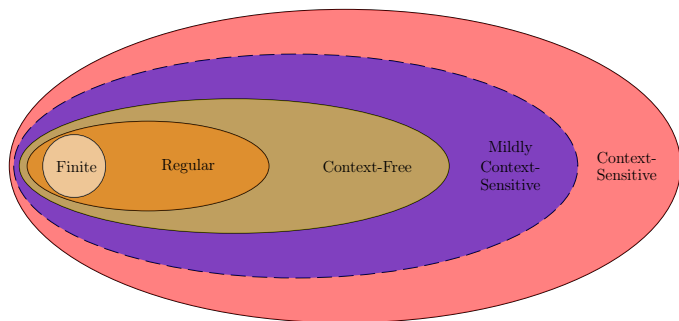
We conclude this section on “what is language learning about” with some open questions:

- What is a good definition of polynomial identification?
- How do we deal with shifting targets? (robustness issues)
- Alternative views on learnability?
- Is being learnable a good indicator of being linguistically reasonable?
- Can we learn transducers? Probabilistic transducers?

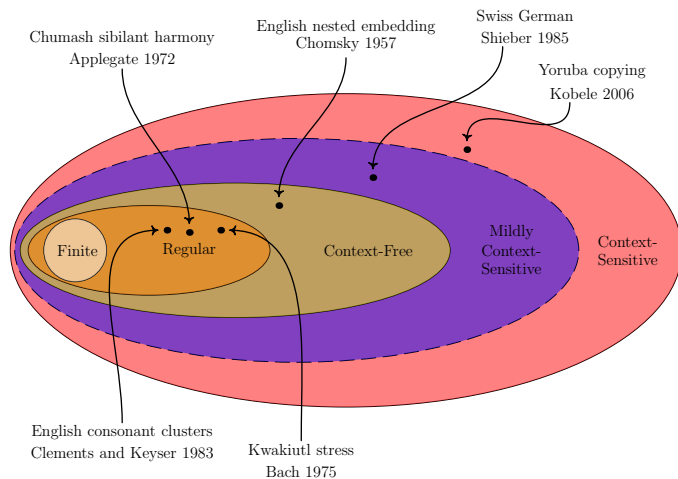
## II. GI of Regular Patterns

- Why regular?
- What are the general GI strategies?
- What are the main results?
- The main techniques?
- The main lessons?

# Logically Possible Computable Patterns

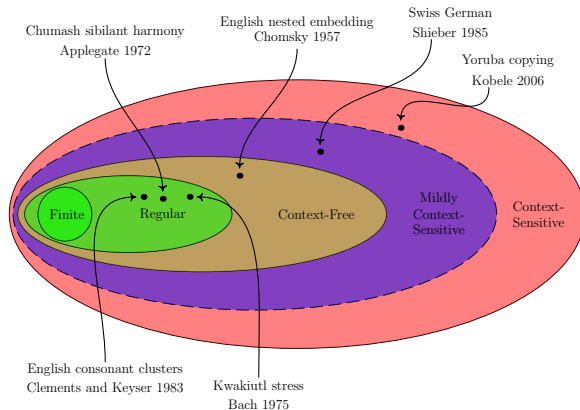


# Logically Possible Computable Patterns



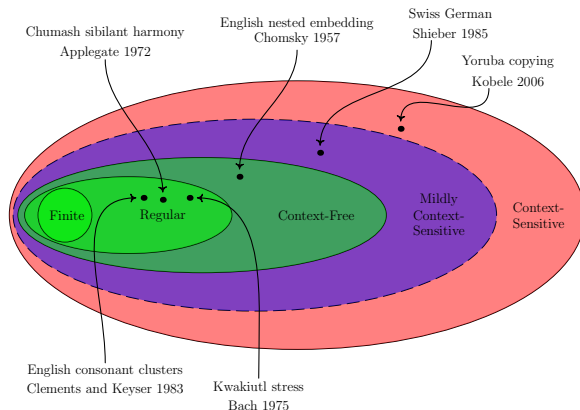
# GI Strategies

#1. Define “learning” so that large regions can be learned



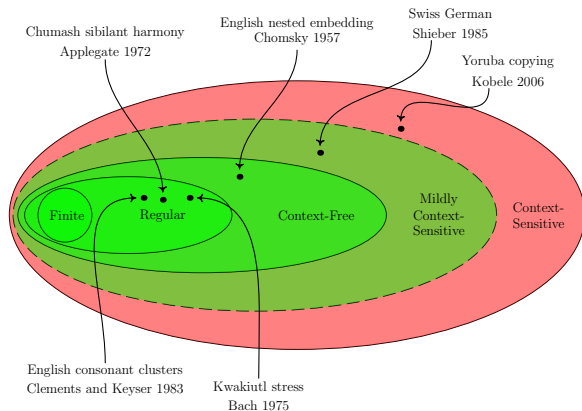
# GI Strategies

#1. Define “learning” so that large regions can be learned



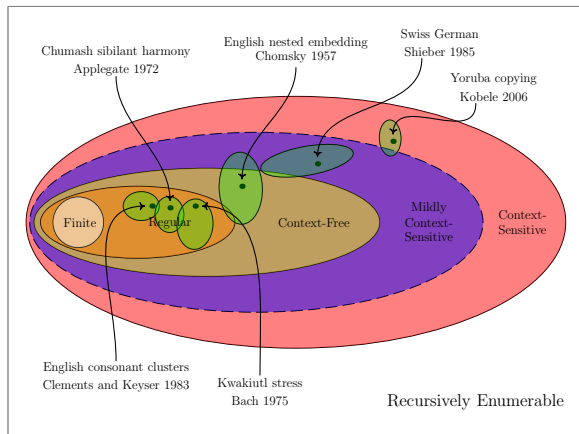
# GI Strategies

#1. Define “learning” so that large regions can be learned



# GI Strategies

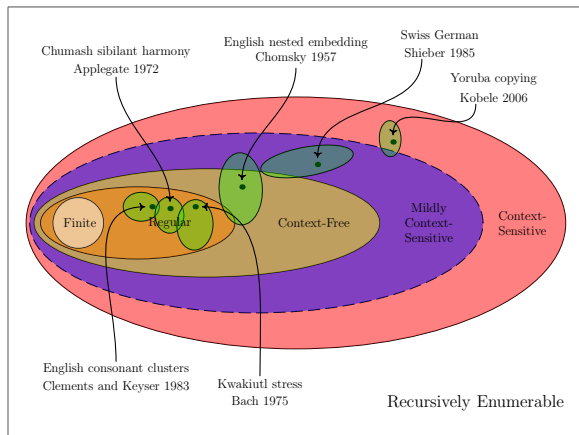
## #2. Target cross-cutting classes



# GI Strategies

## #2. Target cross-cutting classes

(instructor's bias)



# Common Theme

- 1 Different learning frameworks may *better characterize* the data presentations learners actually get (strategy #1).
- 2 Classes of formal languages may exist which *better characterize* the patterns we are interested in (strategy #2).
- 3 Hard problems are *easier to solve* with better characterizations because the instance space of the problem is smaller.

## Why Begin with Regular?

Insights obtained here can be (and have been) applied fruitfully to nonregular classes.

- Angluin 1982 showed a subclass of regular languages (the reversible languages) was identifiable in the limit from positive data by an incremental learner.
- Yokomori's (2004) Very Simple Languages are a subclass of the context-free languages, but draws on ideas from the reversible languages.
- Similarly, Clark and Eryaud's (2007) substitutable languages (also subclass of context-free) are also based on insights from this paper.

# Section Outline

- 1 Targets of Learning
- 2 Learning Frameworks
- 3 State-merging
- 4 Results for learning regular languages, relations, and distributions

# Targets of Learning: Regular Languages

Multiple grammars (i.e. representations) for regular languages:

- 1 Regular expressions
- 2 Generalized regular expressions
- 3 Finite state acceptors
- 4 Words which satisfy formulae in monadic second order logic
- 5 Right or left branching rewrite rules
- 6 ...

# Targets of Learning: Regular Relations

Multiple grammars (i.e. representations) for regular relations:

- Regular expressions (for relations, e.g. Beesley and Karttunen 2003)
- Generalized regular expressions (for relations)
- Finite state transducers
- ...

# Targets of Learning: Regular distributions

Multiple grammars (i.e. representations) for distributions over regular sets and relations:

- Weighted finite state automata
- Hidden Markov Models
- Weighted right or left branching rewrite rules
- ...

# This tutorial: Finite State Automata

## Acceptors and subsequential transducers admit canonical forms

- 1 The smallest deterministic acceptor, syntactic monoids, ...
- 2 Canonical forms relate to algebraic properties (Nerode equivalence relation, i.e. states represent sets of “good tails”)
- 3 In contrast, canonical regular expressions have yet to be determined. For example, there are no canonical (e.g. shortest) regular expressions for regular languages.

# Learning Frameworks: Main Choices

Success required on which input data streams?

All possible vs. some restricted set  
i.e. “distribution-free” vs. “non distribution-free”

What kind of samples?

Positive data vs. positive and negative data

- Other choices (e.g. query learning) are not discussed here.

# Learning Frameworks: Main Results

“Distribution-free” w/ positive and negative data

- 1 Recursive languages are identifiable in the limit (Gold 1967)
- 2 Non-enumerative algorithms for regular languages:
  - 1 Gold (1978)
  - 2 RPNI (Oncina and García 1992)

# Learning Frameworks: Main Results

## “Distribution-free” with positive data only

- 1 No superfinite class (including regular, cf, etc.) is identifiable in the limit (Gold 1967)
- 2 Not even the finite class is PAC-learnable (Blumer et al. 1989)
- 3 No superfinite class is identifiable in the limit with probability  $p$  ( $p > 2/3$ ) (Pitt 1985, Wiehagen et al. 1986, Angluin 1988)

# Learning Frameworks: Main Results

## “Distribution-free” with positive data only

- 1 No superfinite class (including regular, cf, etc.) is identifiable in the limit (Gold 1967)
- 2 Not even the finite class is PAC-learnable (Blumer et al. 1989)
- 3 No superfinite class is identifiable in the limit with probability  $p$  ( $p > 2/3$ ) (Pitt 1985, Wiehagen et al. 1986, Angluin 1988)
- 4 But many **subregular** classes are learnable in this difficult setting.

# Learning Frameworks: Main Results

“Distribution-free” with positive data only: learnable subregular classes

- 1 reversible languages (Angluin 1982)
- 2 strictly local languages (Garcia et al. 1990)
- 3 locally testable and piecewise testable (Garcia and Ruiz 2004)
- 4 left-to-right and right-to-left iterative languages (Heinz 2008)
- 5 strictly piecewise languages (Heinz 2010)
- 6 ...
- 7 subsequential functions (Oncina et al. 1993)
- 8 ...

# Learning Frameworks: Main Results

## “Non distribution-free” w/ positive data only

- 1 The class of r.e. **languages** are identifiable in the limit from computable classes of r.e. texts (Gold 1967)
- 2 The class of r.e. **distributions** are identifiable from “approximately computable” sequences (Angluin 1988, Chater and Vitanyí 2007)
- 3 The class of **distributions** describable with Probabilistic Deterministic FSAs (PDFAs) is learnable with probability one (de la Higuera and Thollard 2000)
- 4 The class of **distributions** describable with PDFAs is learnable in a modified PAC setting (Clark and Thollard, 2004)

# Learning regular languages: Key technique

## State-merging

- Angluin 1982 (reversible languages)
- Muggleton 1990 (contextual languages)
- Garcia et al. 1990 (strictly local languages)
- Oncina et al. 1993 (subsequential functions)
- Clark and Thollard 2004 (PDFA distributions)
- ...

# Other techniques

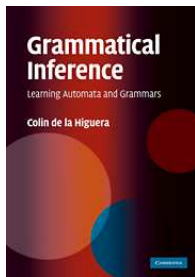
## Lattice-climbing

- Heinz 2010 (strictly local languages, strictly piecewise languages, many others)
- Kasprizk and Kötzing 2010 (function-distinguishable lanaguages, pattern languages, many others)

## State-splitting

- Tellier (2008)

Only so much can be covered. . .



It's impossible to be fair to all those who have contributed and to cover all the variants, even all the algorithms in a short tutorial. That's why there are books!

# Overview of State-merging

- 1 Builds a FSA representation of the input
- 2 Generalize by merging states

# Illustrative Example: Stress pattern of Pintupi

a.	páŋa	'earth'	acute acute
b.	tʰútaya	'many'	acute acute acute
c.	málawàna	'through from behind'	acute acute grave acute
d.	púlĩŋkàlatʰu	'we (sat) on the hill'	acute acute grave acute acute
e.	tʰámulĩmpatʰũŋku	'our relation'	acute acute grave acute grave acute
f.	tʰĩrĩŋulàmpatʰu	'the fire for our benefit flared up'	acute acute grave acute grave acute acute
g.	kúranʰũlũlĩmpatʰũŋa	'the first one who is our relation'	acute acute grave acute grave acute grave acute
h.	yúmaũŋkamàratʰũŋaka	'because of mother-in-law'	acute acute grave acute grave acute grave acute acute

Generalization (Hayes (1995:62) citing Hansen and Hansen (1969:163)):

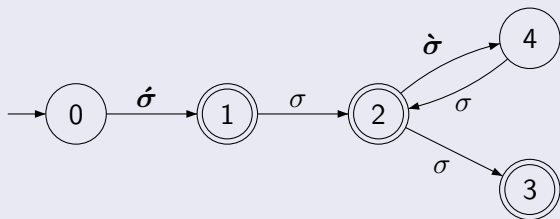
- Primary stress falls on the initial syllable
- Secondary stress falls on alternating nonfinal syllables

# Illustrative Example: Stress pattern of Pintupi

Generalization (Hayes (1995:62) citing Hansen and Hansen (1969:163)):

- Primary stress falls on the initial syllable
- Secondary stress falls on alternating nonfinal syllables

Minimal deterministic FSA for Pintupi Stress



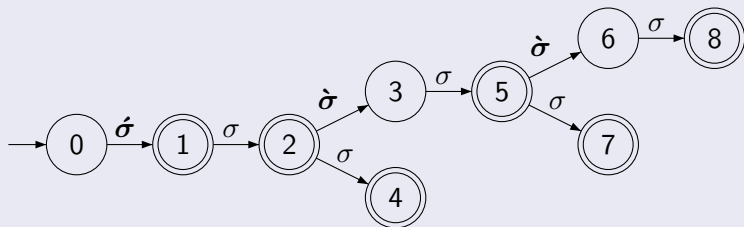
# Structured representations of Input

- 1 Each word its own FSA (Nondeterministic)
- 2 Prefix Trees (deterministic)
- 3 Suffix Trees (reverse deterministic)

# Examples of Prefix and Suffix Trees

$$S = \left\{ \begin{array}{ll} \acute{\sigma} & \acute{\sigma} \sigma \\ \acute{\sigma} \sigma \sigma & \acute{\sigma} \sigma \grave{\sigma} \sigma \\ \acute{\sigma} \sigma \grave{\sigma} \sigma \sigma & \acute{\sigma} \sigma \grave{\sigma} \sigma \grave{\sigma} \sigma \end{array} \right\}$$

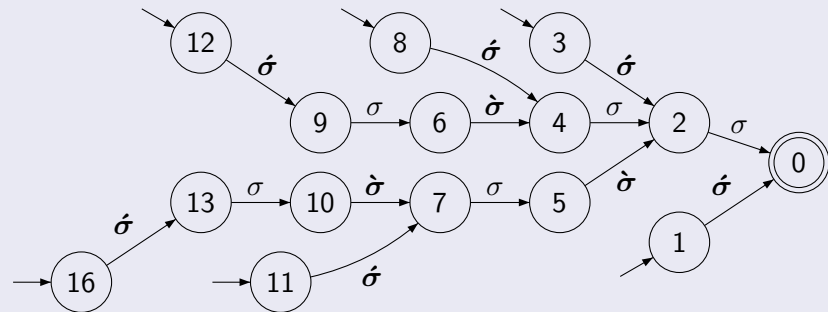
PT(S)



# Examples of Prefix and Suffix Trees

$$S = \left\{ \begin{array}{ll} \acute{\sigma} & \acute{\sigma} \sigma \\ \acute{\sigma} \sigma \sigma & \acute{\sigma} \sigma \grave{\sigma} \sigma \\ \acute{\sigma} \sigma \grave{\sigma} \sigma \sigma & \acute{\sigma} \sigma \grave{\sigma} \sigma \grave{\sigma} \sigma \end{array} \right\}$$

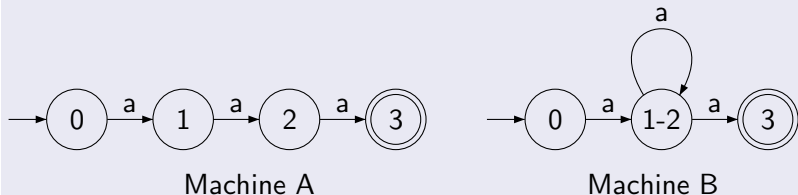
ST(S)



# State-merging Informally

Eliminate redundant environments by state-merging.

- States are identified as equivalent and then *merged*.
- **All** transitions are preserved.
- This is one way in which generalizations may occur—because the post-merged machine accepts everything the pre-merged machine accepts, possibly more.

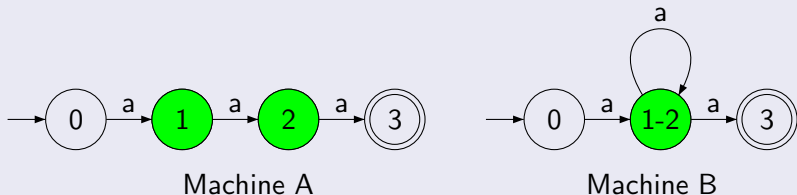


- The merged machine may not be deterministic.

# State-merging Informally

Eliminate redundant environments by state-merging.

- States are identified as equivalent and then *merged*.
- **All** transitions are preserved.
- This is one way in which generalizations may occur—because the post-merged machine accepts everything the pre-merged machine accepts, possibly more.



- The merged machine may not be deterministic.

# State-merging Formally

## Definition

Given an acceptor  $A = (Q, I, F, \delta)$  and a partition  $\pi$  of its states *state-merging* returns the acceptor  $A/\pi = (Q', I', F', \delta')$ :

- 1  $Q' = \pi$  (the states are the *blocks* of  $\pi$ )
- 2  $I' = \{B \in \pi : I \cap B \neq \emptyset\}$
- 3  $F' = \{B \in \pi : F \cap B \neq \emptyset\}$
- 4 For all  $B \in \pi$  and  $a \in \Sigma$ ,  
 $\delta'(B, a) = \{B' \in \pi : \exists q \in B, q' \in B' \text{ such that } q' \in \delta(q, a)\}$

# Theorem

## Theorem

*Given any regular language  $L$ , let  $A(L)$  denote the minimal deterministic acceptor recognizing  $L$ . There exists a finite sample  $S \subseteq L$  and a partition  $\pi$  over  $PT(S)$  such that  $PT(S)/\pi = A(L)$ .*

## Notes

- The finite sample need only exercise every transition in  $A(L)$ .
- What is  $\pi$ ?

# Theorem

## Theorem

*Given any regular language  $L$ , let  $A(L)$  denote the minimal deterministic acceptor recognizing  $L$ . There exists a finite sample  $S \subseteq L$  and a partition  $\pi$  over  $PT(S)$  such that  $PT(S)/\pi = A(L)$ .*

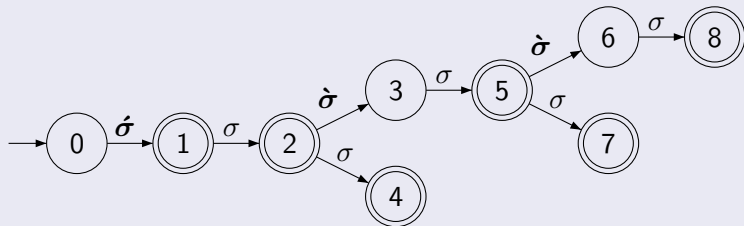
## Notes

- The finite sample need only exercise every transition in  $A(L)$ .
- What is  $\pi$ ?

# Illustrative Example

Let's merge states with the same incoming paths of length 2!

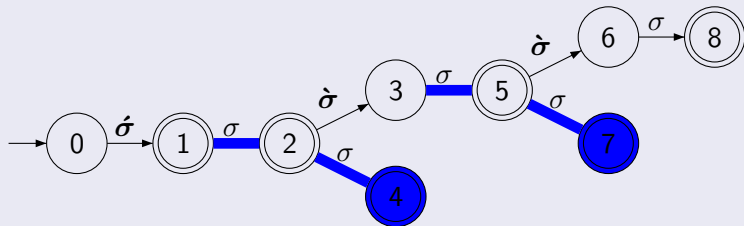
PT(S)



# Illustrative Example

Let's merge states with the same incoming paths of length 2!

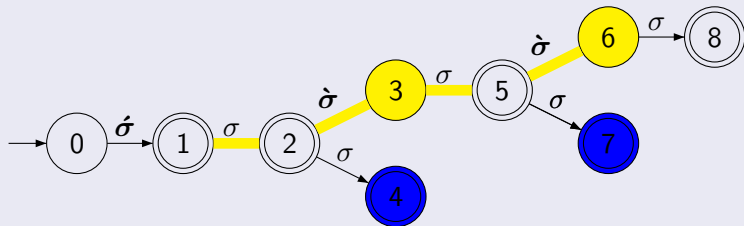
PT(S)



# Illustrative Example

Let's merge states with the same incoming paths of length 2!

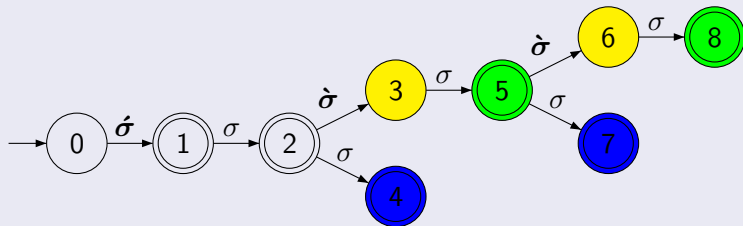
PT(S)



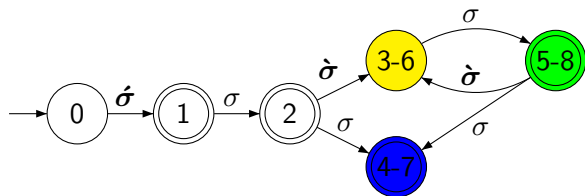
# Illustrative Example

Let's merge states with the same incoming paths of length 2!

PT(S)



## Result of State Merging



This acceptor is not the canonical acceptor we saw earlier but it recognizes the same language.

Generalization (Hayes (1995:62) citing Hansen and Hansen (1969:163)):

- Primary stress falls on the initial syllable
- Secondary stress falls on alternating nonfinal syllables

# Summary of Algorithm

- 1 States in the prefix tree are merged if they have the same  $k$ -length suffix.

$$u \sim v \stackrel{\text{def}}{\iff} \exists x, y, w \text{ such that } |w| = k, u = xw, v = yw$$

- 2 The algorithm then is simply:

$$G = PT(S)/\pi_{\sim}$$

- 3 This algorithm provably identifies in the limit from positive data the Strictly  $(k + 1)$ -Local class of languages (Garcia et al. 1990).

## Back to the Illustrative Example

### Results for stress patterns more generally

- Out of 109 distinct stress patterns in the world's languages (encoded as FSAs), this state-merging strategy works for only 44 of them
- If we merge states with the same paths up to length 5(!), only 81 are learned.
- This is the case even permitting very generous input samples.

## Back to the Illustrative Example

### Results for stress patterns more generally

- Out of 109 distinct stress patterns in the world's languages (encoded as FSAs), this state-merging strategy works for only 44 of them
- If we merge states with the same paths up to length 5(!), only 81 are learned.
- This is the case even permitting very generous input samples.

In other words, 44 attested stress patterns are Strictly 3-Local and 81 are Strictly 6-Local. 28 are not Strictly 6-Local. In fact those 28 are not Strictly  $k$ -Local for any  $k$  (Edlefsen et al. 2008).

## Other ways to merge states

If the current structure is “ill-formed” then merge states to eliminate source of ill-formedness

### State equivalence relations

- 1 merge state with same incoming paths of length  $k$  (Garcia et. al 1990)
- 2 recursively eliminate reverse non-determinism (Angluin 1982)
- 3 merge states with same “contexts” (Muggleton 1990, Clark and Eryaud 2007)
- 4 merge final states (Heinz 2008)
- 5 merge states with same “neighborhood” (Heinz 2009)
- 6 ...
- 7 merge states to maximize posterior probability (for HMMs, Stolcke 1994)
- 8 ...

## Other ways to merge states

Merge states in a particular order unless “ill-formedness” arises

Merge unless something tells us not to

- 1 unless “onward subsequentiality” is lost (for transducers, Oncina et al. 1993)
- 2 unless they are “ $\mu$ -distinguishable” (Clark and Thollard 2004)
- 3 ...

# State-merging as inference rules

Strictly  $k$ -Local languages (Garcia et al. 1990

)

merge states with same incoming paths of length  $k$

$$\forall u, v, w \in \Sigma^* : uv, wv \in \text{Prefix}(L), |v| = k$$

$\Downarrow$

$$\text{Tails}_L(uv) = \text{Tails}_L(wv)$$

# State-merging as inference rules

## 0-Reversible languages (Angluin 1982)

recursively eliminate reverse non-determinism

$$\forall u, v, w, y \in \Sigma^* : uv, wv, uy \in L \Rightarrow wy \in L$$

## State-merging summary

- 1 Distinctions maintained in the prefix tree are lost by state merging, which results in generalizations.
- 2 The choice of partition corresponds to the generalization strategy (i.e. which distinctions will be maintained and which will be lost)

Gleitman (1990:12):

*The trouble is that an observer who notices everything can learn nothing for there is no end of categories known and constructible to describe a situation [emphasis in original].*

# Results for regular languages

## Distribution-free with positive data

### Identification in the limit from positive data

- 1 strictly  $k$ -local languages (each state corresponds to suffixes of up to length  $k$ ) (Garcia et al. 1990)
- 2 reversible languages (acceptors are both forward and reverse  $k$ -deterministic for some  $k$ ) (Angluin 1982)
- 3  $k$ -contextual languages (Muggleton 1990)
- 4 ...

# Regular relations

## Regular relations in CL

- 1 transliteration
- 2 translation
- 3 ...
- 4 anything with finite state transducers

# OSTIA (Oncina et al. 1993)

distribution-free with positive data

## OSTIA

- 1 identifies *subsequential functions* in the limit from positive data.
- 2 Merges states greedily unless subsequentiality is violated
- 3 If the function is partial, exactness is guaranteed only where the function is defined.

# OSTIA (Oncina et al. 1993)

## Subsequential relations

- 1 are a subclass of the regular relations, recognizing functions.
- 2 are those which are recognized by subsequential transducers, which are deterministic on the input and which have an “output” string associated with every state.
- 3 have a canonical form.
- 4 have been generalized to permit up to  $p$  outputs for each input (Mohri 1997).

# OSTIA for learning phonological rules

## Gildea and Jurafsky 1996

- 1 Show that OSTIA doesn't learn the English tapping rule or German word-final devoicing rule from data present in adapted dictionaries of English or German
- 2 Applied additional phonologically motivated heuristics to improve state-merging choices and obtained significantly better results.

# OSTIA for learning phonological rules

## Gildea and Jurafsky 1996

- 1 Show that OSTIA doesn't learn the English tapping rule or German word-final devoicing rule from data present in adapted dictionaries of English or German
- 2 Applied additional phonologically motivated heuristics to improve state-merging choices and obtained significantly better results.

What about well-defined subclasses of subsequential relations?

# Weighted finite-state automata

non-distribution-free with positive data

## The problem

Given a finite multiset of words drawn independently from the target distribution, what grammar accurately describes the distribution?

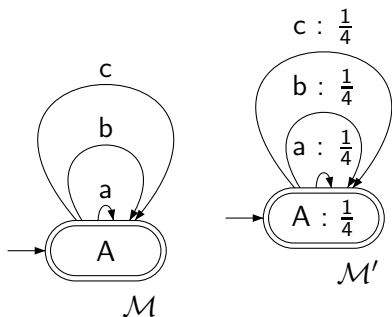
# Weighted finite-state automata

non-distribution-free with positive data

## Theorem

*The class of distributions describable with Non-deterministic Probabilistic Finite-State Automata (NPFA) exactly matches the class of distributions describable with Hidden Markov Models (Vidal et al. 2005).*

# Maximum Likelihood Estimation



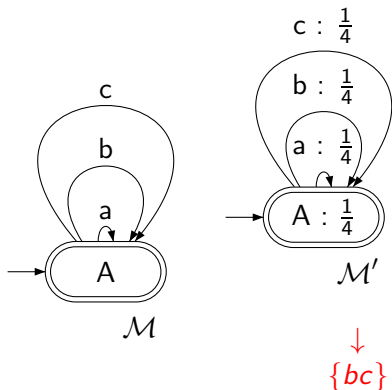
$\mathcal{M}$  represents a family of distributions with 4 parameters.  $\mathcal{M}'$  represents a particular distribution in this family.

## Theorem

For a sample  $S$  and deterministic finite-state acceptor  $\mathcal{M}$ , **counting the parse of  $S$  through  $\mathcal{M}$  and normalizing at each state optimizes the maximum-likelihood estimate.**

(Vidal et. al 2005, de la Higuera 2010)

# Maximum Likelihood Estimation



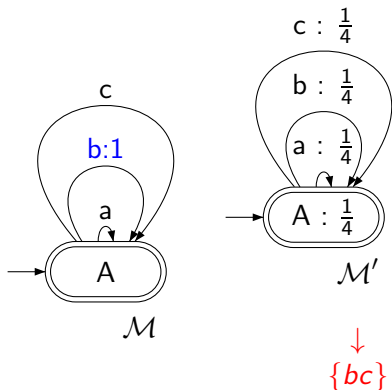
$\mathcal{M}$  represents a family of distributions with 4 parameters.  $\mathcal{M}'$  represents a particular distribution in this family.

## Theorem

For a sample  $S$  and deterministic finite-state acceptor  $\mathcal{M}$ , **counting the parse of  $S$  through  $\mathcal{M}$  and normalizing at each state optimizes the maximum-likelihood estimate.**

(Vidal et. al 2005, de la Higuera 2010)

# Maximum Likelihood Estimation



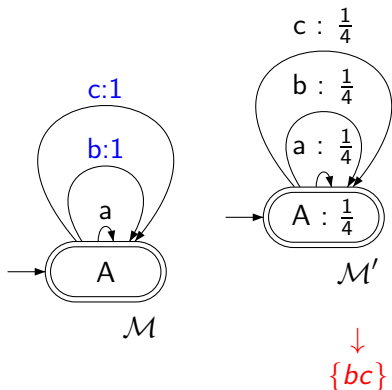
$\mathcal{M}$  represents a family of distributions with 4 parameters.  $\mathcal{M}'$  represents a particular distribution in this family.

## Theorem

For a sample  $S$  and deterministic finite-state acceptor  $\mathcal{M}$ , **counting the parse of  $S$  through  $\mathcal{M}$  and normalizing at each state optimizes the maximum-likelihood estimate.**

(Vidal et. al 2005, de la Higuera 2010)

# Maximum Likelihood Estimation



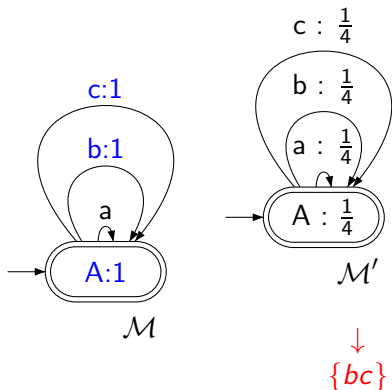
$\mathcal{M}$  represents a family of distributions with 4 parameters.  $\mathcal{M}'$  represents a particular distribution in this family.

## Theorem

For a sample  $S$  and deterministic finite-state acceptor  $\mathcal{M}$ , **counting the parse of  $S$  through  $\mathcal{M}$  and normalizing at each state optimizes the maximum-likelihood estimate.**

(Vidal et. al 2005, de la Higuera 2010)

# Maximum Likelihood Estimation



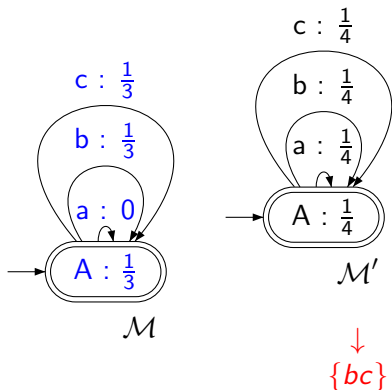
$\mathcal{M}$  represents a family of distributions with 4 parameters.  
 $\mathcal{M}'$  represents a particular distribution in this family.

## Theorem

For a sample  $S$  and deterministic finite-state acceptor  $\mathcal{M}$ , **counting the parse of  $S$  through  $\mathcal{M}$  and normalizing at each state optimizes the maximum-likelihood estimate.**

(Vidal et. al 2005, de la Higuera 2010)

# Maximum Likelihood Estimation



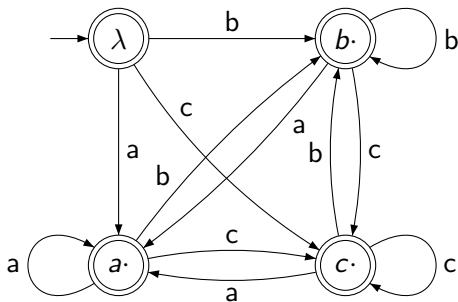
$\mathcal{M}$  represents a family of distributions with 4 parameters.  $\mathcal{M}'$  represents a particular distribution in this family.

## Theorem

For a sample  $S$  and deterministic finite-state acceptor  $\mathcal{M}$ , **counting the parse of  $S$  through  $\mathcal{M}$  and normalizing at each state optimizes the maximum-likelihood estimate.**

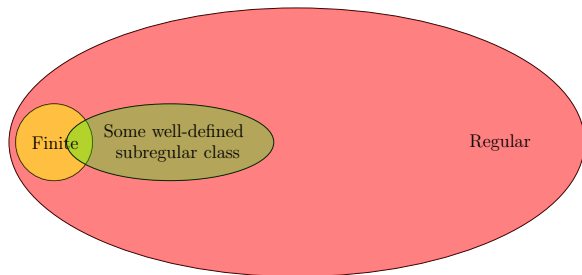
(Vidal et. al 2005, de la Higuera 2010)

# Strictly 2-Local Distributions are bigram models



**Figure:** The structure of a bigram model. The 16 parameters of this model are given by associating probabilities to each transition and to “ending” at each state.

# Subregular distributions



- 1 When the structure of a Deterministic FSA is known in advance, MLE is easy to do.
- 2 The DFA represents a subregular class of distributions.

# Strictly Piecewise Distributions

- 1  $N$ -gram models can't describe long-distance dependencies.

## Long-distance dependencies in phonology

- 1 **Consonantal harmony**

(Jensen 1974, Odden 1994, Hansson 2001, Rose and Walker 2004, and many others)

- 2 **Vowel harmony**

(Ringen 1988, Baković 2000, and many others)

# Sibilant Harmony example from Samala (Ineseño Chumash)

[ʃtojonowonowaʃ] 'it stood upright' (Applegate 1972:72)

cf. \*[stojonowonowaʃ] and

cf. \*[ʃtojonowonowas]

**Hypothesis:** \*[stojonowonowaʃ] and \*[ʃtojonowonowas] are ill-formed because the *discontiguous subsequences* *sf* and *ʃs* are ill-formed.

# Strictly Piecewise languages

Rogers et al. 2010

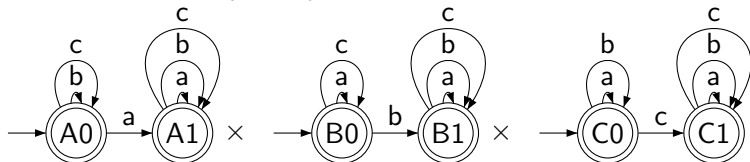
- 1 solely make distinctions on the basis of potentially discontinuous subsequences up to some length  $k$
- 2 are mathematically natural. They have several characterizations in terms of formal language theory, automata theory, logic, model theory, and the
- 3 algebraic theory of automata (Fu et al. 2011)

# Strictly Piecewise Distributions

## Heinz and Rogers 2010

- 1 are defined in terms of the factored automata-theoretic representations (Rogers et al. 2010)
- 2 along with the co-emission probability as the product (Vidal et al. 2005)
- 3 Estimation over the factors permits learnability of the patterns like the ones in Samala.

Example with  $\Sigma = \{a, b, c\}$  and  $k = 2$ .



# SP2 learning results for Chumash

- Training corpus 4800 words from a dictionary of Samala

$P(x   y <)$		x			
		s	$\widehat{ts}$	$\int$	$\widehat{tj}$
y	s	0.0325	0.0051	0.0013	0.0002
	$\widehat{ts}$	0.0212	0.0114	0.0008	0.
	$\int$	0.0011	0.	0.067	0.0359
	$\widehat{tj}$	0.0006	0.	0.0458	0.0314

**Table:** SP2 probabilities of sibilant occurring sometime after another one (collapsing laryngeal distinctions)

# Learning larger classes of regular distributions

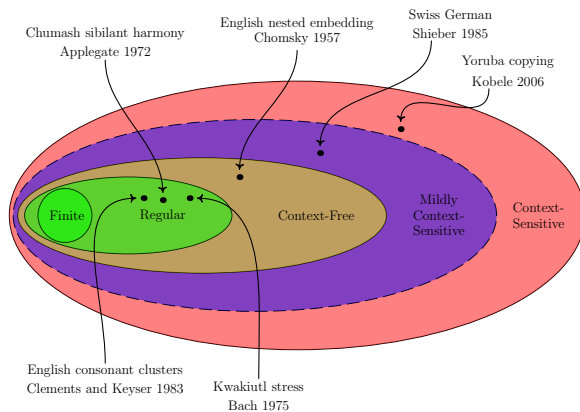
More non-distribution-free with positive data

## The class of distributions describable with PDFAs

- 1 are identifiable in the limit with probability one (de la Higuera and Thollard 2000).
  - 2 are learnable in modified-PAC setting (Clark and Thollard 2004).
  - 3 The algorithms presented employ state-merging methods.
- 
- 1 This is a (much!) larger class than that which is describable with n-gram distributions or with SP distributions.
  - 2 To my knowledge these approaches have not been applied to tasks in CL.

# Summary

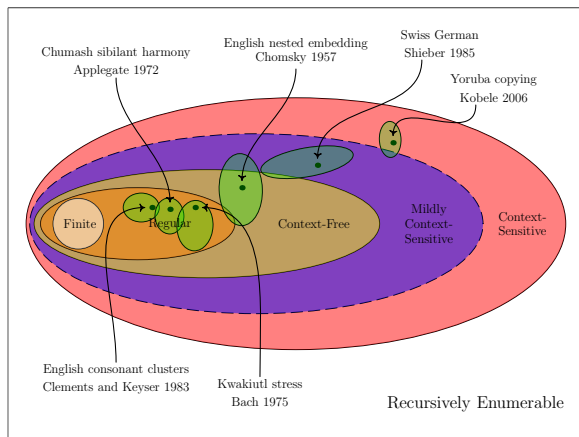
#1. Define “learning” so that large regions can be learned



Oncina et al. 1993, de la Higuera and Thollard 2000, Clark and Thollard 2004, ...

# Summary

## #2. Target cross-cutting classes



Angluin 1982, Muggleton 1990, Garcia et al. 1990, Heinz 2010, ...

# Have we put the cart before the horse?



- 1 So far we have discussed algorithms that learn various classes of languages.
- 2 But shouldn't we *first* know *which* classes are relevant for our goals?
- 3 E.g. for phonology, while “being regular” may be a *necessary* property of phonological patterns, it certainly is not *sufficient*.

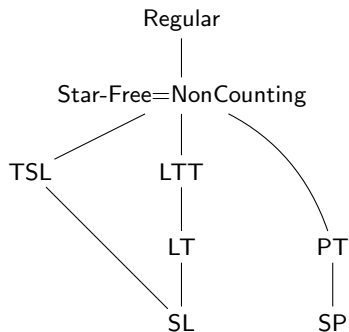
# Have we put the cart before the horse?

## Research strategy

Patterns  $\Rightarrow$  Characterizations  $\Rightarrow$  Learning algorithms

- 1 Identify the range and kind of patterns (linguistics).
- 2 Characterize the range and kind of patterns (computational linguistics).
- 3 Create learning algorithms for these classes, prove their success in a variety of settings, and otherwise demonstrate their success (grammatical inference, formal learning theory, computational linguistics)

# Subregular classes of regular sets



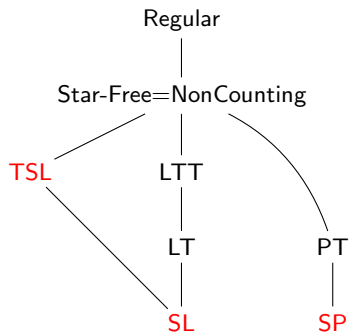
Proper inclusion relationships among subregular language classes.

TSL Tier-based Strictly Local  
 LTT Locally Threshold Testable  
 LT Locally Testable

PT Piecewise Testable  
 SL Strictly Local  
 SP Strictly Piecewise

(McNaughton and Papert 1971, Simon 1975, Rogers and Pullum 2007, in press, Rogers et al. 2010, Heinz et al. 2011)

# Subregular classes of regular sets



Proper inclusion relationships among subregular language classes.

*instructor's hunch for phonology*

TSL Tier-based Strictly Local  
 LTT Locally Threshold Testable  
 LT Locally Testable

PT Piecewise Testable  
 SL Strictly Local  
 SP Strictly Piecewise

(McNaughton and Papert 1971, Simon 1975, Rogers and Pullum 2007, in press, Rogers et al. 2010, Heinz et al. 2011)

# Conclusion to section 2 : part 1

## State-merging

- 1** is a well-studied strategy for inferring automata, including acceptors, transducers, and weighted acceptors and transducers.
- 2** has yielded theoretical results in many learning frameworks including both distribution-free and non-distribution-free learning frameworks.

## Conclusion to section 2 : part 2

- 1 Many subclasses of regular languages are learnable even in the hardest learning settings.
- 2 Recent advances yield algorithms for large classes (probabilistic DFAs)
- 3 Computational linguists can explore which are relevant to natural language and consequently which are useful for NLP!
- 4 There is a rich literature in GI which speaks to these classes, and how such patterns in these classes can be learned.

## Overview

- Empirical grammatical inference
- Family of languages
- Information contained in input
- Overview of systems
- Evaluation issues
- From empirical to formal GI

## Introduction

- Language learning
  - Starting from family of languages
  - Given set of samples
  - Identify language that is used to generate samples
- Formal grammatical inference
  - Identify family of languages that can be learned efficiently
  - Under certain restrictions
- Empirical grammatical inference
  - Exact underlying family of languages is unknown
  - Target language is approximation

## Empirical GI

- Try to identify language given samples
  - E.g. sentences (syntax), words (morphology), ...
- Underlying language class is unknown
  - For algorithm we still need to make a choice
- If identification is impossible, provide approximation
  - Evaluation of empirical GI is different from formal GI

## Family of languages

- What is the underlying family of languages?
- Choice has impact on learning algorithm
- Many possibilities
  - Use simple, fixed structures ( $n$ -grams)
    - Find probabilities
  - Extract structure from treebanks
    - Slightly more flexible structure
    - Find probabilities
  - Learn structure
    - Flexible structure
    - Find probabilities

## $N$ -grams

- 1 Starting from a plain text or collection of texts (corpus)
- 2 Extract all substrings of length  $n$  ( $n$ -grams)
- 3 Count occurrences of  $n$ -grams in texts
- 4 Assign probabilities to each  $n$ -gram based on counts

## Issues

- Unseen  $n$ -grams
  - Back-off: use  $n$ -grams with smaller  $n$
  - Smoothing: adjust probabilities for unseen  $n$ -grams

## Using $n$ -gram models

- How likely is the sentence 'John likes Mary' ?
  - Unigram language model
    - $P(\text{John likes Mary}) \approx P(\text{John})P(\text{likes})P(\text{Mary})$
  - Bigram language model
    - $P(\text{John likes Mary}) \approx P(\text{John}|\langle s \rangle)P(\text{likes}|\text{John})P(\text{Mary}|\text{likes})$
  - Trigram language model
    - $P(\text{John likes Mary}) \approx P(\text{John}|\langle s \rangle \langle s \rangle)P(\text{likes}|\langle s \rangle \text{John})P(\text{Mary}|\text{John likes})$
  - N-gram language model
    - $P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-N+1}^{k-1})$
- $N$ -grams provide a probability for each sequence
  - Probability describes how well sequence fits language

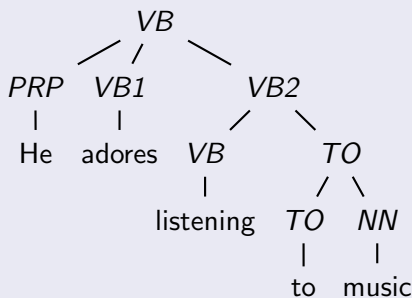
## Extract structure from treebanks

- 1 Starting from a treebank (sentences with structure)
- 2 Extract grammar rules that are used to create tree structures
  - For instance, context-free grammars (Charniak 1993)
  - or sub-trees (Data-Oriented Parsing) (Bod 1998)
- 3 Count occurrences of grammar rules in treebank
- 4 Assign probabilities to grammar rules based on counts

## Issues

- Over-generalization, “incorrect” probabilities
  - Add information on applicability of grammar rules (Johnson 1998)
  - Reestimate probabilities (EM) (Dempster et al 1977, Lari and Young 1990)

## Extract structure from tree



*VB* → *PRP VB1 VB2*

*PRP* → He

*VB1* → adores

*VB2* → *VB TO*

*VB* → listening

*TO* → *TO NN*

*TO* → to

*NN* → music

- Extract counts from treebank → probabilities
- Reestimate probabilities
  - Improve fit of grammar and sentences

## Learn structure

- 1 Starting from a corpus
- 2 Identify regularities that may serve as grammar rules
- 3 Output:
  - Structure assigned to sentences → extract grammar
  - Extracted grammar rules (and probabilities) → parse

## Issues

- Learning system has to deal with both
  - flexibility in structure
  - probabilities of structure

## Summarizing fixed versus flexible structure

- Fixed versus flexible is really a sliding scale
- Language modelling using  $n$ -grams
  - Structure is very simple and very rigid
  - Requires plain sequences as input
  - Corresponds to  $k$ -testable languages (García 1990)
- Language modelling using extracted grammar rules
  - Structure is more flexible, but restricted by treebank
  - Requires structured sequences as input
  - Corresponds to e.g. (limited) context-free languages
- “*Learning structure*”
  - Structure is flexible, restricted by learning algorithm
  - Requires plain sequences as input
  - Corresponds to e.g. context-free languages

## Empirical grammatical inference

- Choices:
  - What type of grammar are we learning?
    - Regular language
    - $K$ -testable language ( $n$ -grams)
    - Context-free language
    - ...
  - What kind of input do we require?
    - Sequence of words (sentence)
    - Sequence of part-of-speech tags
    - (Partial) tree structures
    - ...
  - What kind of output do we want?
    - Structured version of input
    - Explicit grammar
    - Binary or  $n$ -ary (context-free rules)
    - ...

## Overview of systems

- EMILE
- Alignment-Based Learning (ABL)
- ADIOS
- CCM+DMV
- U-DOP
- ...

## Underlying approach

- Given a collection of plain sentences
- On what basis are we going to assign structure?
- Should structure be linguistically motivated?
  - or similar to what linguists would assign?
- Perhaps we can use tests for constituency to find structure

## Substitutability

### Elements of the same type are substitutable

Test for constituency (Harris, 1951)

*What is (a family fare)<sub>NP</sub>*

Replace noun phrase with another noun phrase

## Substitutability

### Elements of the same type are substitutable

Test for constituency (Harris, 1951)

*What is (a family fare)<sub>NP</sub>*

Replace noun phrase with another noun phrase

*What is (the payload of an African Swallow)<sub>NP</sub>*

## Substitutability

### Elements of the same type are substitutable

Test for constituency (Harris, 1951)

*What is (a family fare)<sub>NP</sub>*

Replace noun phrase with another noun phrase

*What is (the payload of an African Swallow)<sub>NP</sub>*

### Learning by reversing test

What is a family fare

What is the payload of an African Swallow

## Substitutability

### Elements of the same type are substitutable

Test for constituency (Harris, 1951)

*What is (a family fare)<sub>NP</sub>*

Replace noun phrase with another noun phrase

*What is (the payload of an African Swallow)<sub>NP</sub>*

### Learning by reversing test

What is (a family fare)<sub>X</sub>

What is (the payload of an African Swallow)<sub>X</sub>

## EMILE

- Learns context-free grammars
- Using plain sentences
- Originally used to show formal learnability of (a form of) Categorical Grammars in a PAC learning setting
- (Adriaans 1992, Adriaans and Vervoort 2002, Vervoort 2000)

## Approach

- 1 Starting from simple sentences
  - identify recurring substrings
- 2 Store recurring substrings and contexts
- 3 Introduce grammar rules when there is enough evidence
  - Practical implementation allows for several constraints
    - Context length, subsequence length, . . .

## Example matrix

John walks

Mary walks

John sees Mary

	(.) walks	John (.)	(.) sees Mary	...	<i>contexts</i>
John	x			x	...
walks		x			...
Mary	x				...
sees					...
⋮	⋮	⋮	⋮	⋮	⋮
<i>terms</i>					

## Learn grammar rules

- Terms that share (approximately) same context are clustered
  - “John” and “Mary” are grouped together
- Occurrences of terms in cluster are replaced by new symbol
- Modified sequences may again contain terms/contexts
- Terms may consist of multiple words

## Example

John walks  $\Rightarrow X$  walks

Mary walks  $\Rightarrow X$  walks

John sees Mary  $\Rightarrow X$  sees  $X$

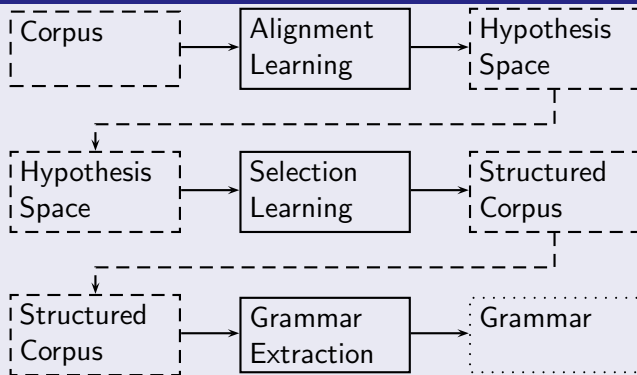
Mary slaps John  $\Rightarrow X$  slaps  $X$

“sees” and “slaps” now also share the same context

## Alignment-Based Learning (ABL)

- Based on substitutability test
- Using plain sentences
- Similar to EMILE, but
  - Clustered terms are not explicitly replaced by symbol
  - Terms and contexts are *always* separated
  - All terms are considered (and only selected afterwards)
- Output is structured version of input or grammar
- (van Zaanen 2000a, b, 2002)

## Alignment-Based Learning (ABL)



## Alignment-Based Learning (ABL)

- Alignment learning
  - Align pairs of sentences
  - Unequal parts of sentences are stored as hypotheses
- (Clustering)
  - Group hypotheses in same context together
- Selection learning
  - Remove overlapping hypotheses

## Alignment learning

- Align pairs of sentences
  - using edit distance (Wagner and Fischer 1974)
  - or suffixtrees (Geertzen and van Zaanen 2004, Ukkonen 1995)
- Unequal parts of sentences are stored as hypotheses
- Align all sentences in a corpus to all others

## Example

I need ( $x_1$  a dinner during the flight) $x_1$

I need ( $x_1$  to return on tuesday) $x_1$

he wants to return on wednesday

## Alignment learning

- Align pairs of sentences
  - using edit distance (Wagner and Fischer 1974)
  - or suffixtrees (Geertzen and van Zaanen 2004, Ukkonen 1995)
- Unequal parts of sentences are stored as hypotheses
- Align all sentences in a corpus to all others

## Example

$(Y_1 I \text{ need } (X_1 \text{ a dinner during the flight}) X_1) Y_1$

$I \text{ need } (X_1 \text{ to return on tuesday}) X_1$

$(Y_1 \text{ he wants to return on wednesday}) Y_1$

## Alignment learning

- Align pairs of sentences
  - using edit distance (Wagner and Fischer 1974)
  - or suffixtrees (Geertzen and van Zaanen 2004, Ukkonen 1995)
- Unequal parts of sentences are stored as hypotheses
- Align all sentences in a corpus to all others

## Example

$(Y_1 \text{I need } (X_1 \text{a dinner during the flight})_{X_1})_{Y_1}$   
 $(Z_1 \text{I need})_{Z_1} (X_1 \underline{\text{to return on}} (Z_2 \text{tuesday})_{Z_2})_{X_1}$   
 $(Y_1 (Z_1 \text{he wants})_{Z_1} \underline{\text{to return on}} (Z_2 \text{wednesday})_{Z_2})_{Y_1}$

## Selection Learning

- Alignment learning can generate overlapping brackets
- Underlying grammar is considered context-free
- Structure describes parse according to underlying grammar
- “Wrong” brackets have to be removed
  - Based on e.g. chronological order or statistics

## Example

from  $(y_1 \text{ Tilburg } (x_2 \text{ to})_{y_1} \text{ Portland})_{x_2}$

from  $(x_1 \text{ Portland } (y_2 \text{ to})_{x_1} \text{ Tilburg})_{y_2}$

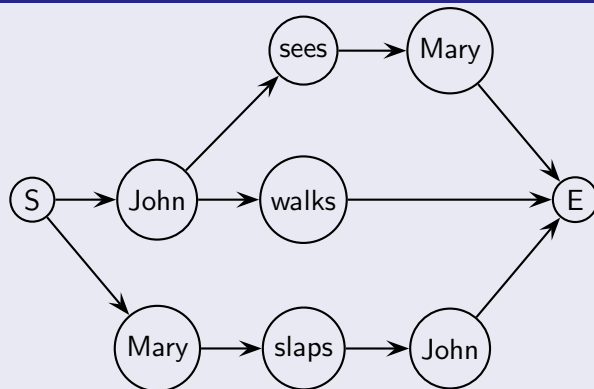
## ADIOS

- Automatic Distillation of Structure (ADIOS) (Solan 2005)

### Idea

- 1 Represent language as a graph
- 2 Compress graph
- 3 As long as possible, find significant patterns in paths
  - Using substitutability and *significance tests*
- 4 (Recursion may be added as a post-processing step)

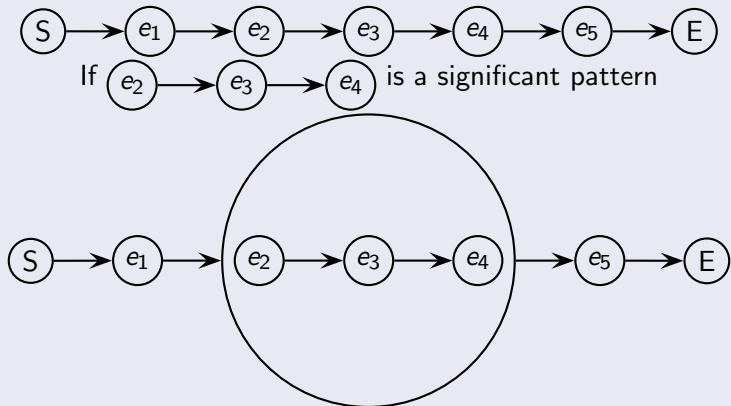
## Graph



## Phases

- 1 Initialization
  - Load all sentences (as paths) in the graph
- 2 Pattern distillation
  - Find sub-paths
    - shared by significant number of partially-aligned paths
    - using motif-extraction (MEX) algorithm
- 3 Generalization
  - Group all nodes that occur in same pattern together
  - Cluster words/substrings similarly to EMILE
- 4 Repeat 2 and 3 until no new patterns are found

## Graph



## MEX

- Compute probabilities depending on in-/out-degree of nodes

$$P_R(e_1; e_2) = \frac{\# \text{ paths from } e_1 \text{ to } e_2}{\# \text{ paths to } e_1}$$

$$P_R(e_1; e_3) = \frac{\# \text{ paths from } e_1 \text{ to } e_3}{\# \text{ paths to } e_1}$$

$$D_R(e_1; e_3) = \frac{P_R(e_1; e_4)}{P_R(e_1; e_3)}$$

- $P_R$  describes path to the right  
similarly  $P_L$  describes path to the left
- Significance is computed based on  $D_R$  and  $D_L$  wrt parameter
  - Informally: find significant changes in number of paths
- Pick most significant pattern

## Constituent-Context Model (CCM)

- Consider all possible binary tree structures on POS sequences
- Define a probability distribution over the possible bracketings
  - A bracketing is a particular structure on a sequence

$$P(s, B) = P_{\text{bin}}(B)P(s|B)$$

$$P(s|B) = \prod_{i,j:i \leq j} P_{\text{span}}(s_{ij}|B_{ij})P_{\text{ctx}}(s_{i-1}, s_j|B_{ij})$$

- Run (iterative) Expectation-Maximization (EM) algorithm
  - to maximize likelihood  $\prod_{s \in S} P(s)$
- (Klein 2002)

## Dependency Model with Valence (DMV)

- DMV aims to learn dependency relations in contrast to CCM which learns context-free grammar rules
- Dependency parse links words in a head-dependent relation
- Model describes likelihood of
  - left dependencies
  - right dependencies
  - stop condition (no more dependencies)
- Again, iterative EM is used to maximize likelihood of corpus

## CCM+DMV

- CCM and DMV can be combined
- Both models have different view on structure
- Results of combined system are better than either systems
  - Strengths of both systems are combined
- (Klein 2004)

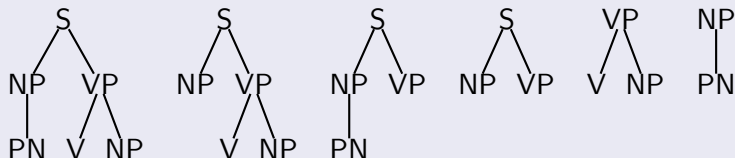
## U-DOP

- Similar to CCM in that it
  - finds probability distribution over “all” structures
  - uses POS sequences
- U-DOP uses Data-Oriented Parsing (DOP) as formalism
  - Extends probabilistic model of context-free grammars
- Requires practical implementation choices
  - Random sampling due to huge size of search space
- (Bod 2006a, b)

## Procedure

- 1 Generate all possible binary trees on example sentences
- 2 Extract all *subtrees*
- 3 Estimate probabilities on subtrees using EM

## Subtrees

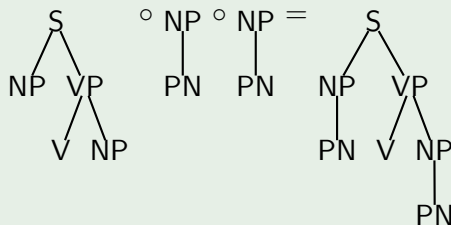


- Remove either all or no elements on a level
  - Leads to *many* subtrees
- Each subtree receives a probability
  - Longer distance dependencies may be modeled

## Parsing

- Subtrees can be recombined into a larger tree
  - Similar to context-free grammar rules
- Same parse may be created using different derivations
  - Statistical model has to take this into account

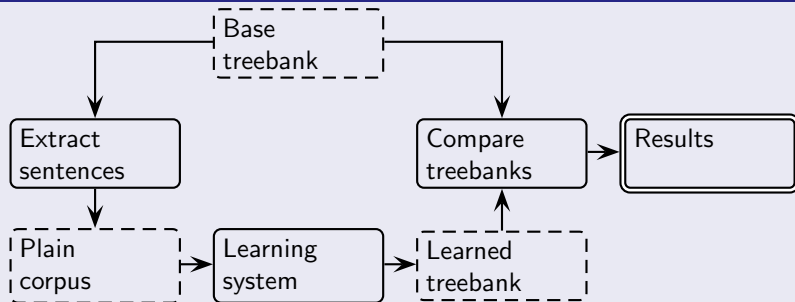
## Example



## Underlying idea

- U-DOP works because span of subtrees reoccur in a corpus
  - Likelihood of “useful” spans increase
  - Hence, likelihood of contexts (also subtrees) increase
- Essentially, U-DOP uses implied substitutability
  - while system leans heavily on probabilities

## Evaluation



- Recall (completeness)
- Precision (correctness)
- F-Score (combination of Precision and Recall)
- (van Zaanen and Adriaans 2001)

## Evaluation settings

- Air Travel Information System (ATIS)
- Taken from Penn Treebank II
- 568 English sentences

## Example

list the flights from baltimore to seattle that stop in minneapolis  
does this flight serve dinner  
the flight should arrive at eleven a.m. tomorrow  
what airline is this

## Results on ATIS

Precision	47.01
Recall	44.94
F-Score	44.60

## Results on ATIS

	Micro	Macro
Precision	47.01	46.18
Recall	44.94	50.98
F-Score	44.60	47.10

## Explanation

**Micro** Count constituents, weighted average per sentence

**Macro** Count constituents and average per sentence

## Results on ATIS

	Micro	Macro	Macro <sup>2</sup>
Precision	47.01	46.18	46.18
Recall	44.94	50.98	50.98
F-Score	44.60	47.10	48.46

## Explanation

**Micro** Count constituents, weighted average per sentence

**Macro** Count constituents and average per sentence

**Macro<sup>2</sup>** Compute Macro Precision/Recall, average at end

## Results on ATIS

Micro Precision	47.01
Micro Recall	44.94
Micro F-Score	44.60
Macro Precision	46.18
Macro Recall	50.98
Macro F-Score	47.10
Macro <sup>2</sup> F-Score	48.46

## Results on ATIS

		remove sentence
Micro Precision	47.01	47.67
Micro Recall	44.94	45.30
Micro F-Score	44.60	45.09
Macro Precision	46.18	47.66
Macro Recall	50.98	52.96
Macro F-Score	47.10	48.62
Macro <sup>2</sup> F-Score	48.46	50.17

## Example

(bla bla bla)→bla bla bla

## Results on ATIS

		remove sentence	remove empty
Micro Precision	47.01	47.67	77.10
Micro Recall	44.94	45.30	44.95
Micro F-Score	44.60	45.09	55.31
Macro Precision	46.18	47.66	77.08
Macro Recall	50.98	52.96	51.07
Macro F-Score	47.10	48.62	60.00
Macro <sup>2</sup> F-Score	48.46	50.17	61.43

## Example

bla () bla → bla bla

## Results on ATIS

		remove sentence	remove empty	remove both
Micro Precision	47.01	47.67	77.10	79.07
Micro Recall	44.94	45.30	44.95	45.29
Micro F-Score	44.60	45.09	55.31	56.13
Macro Precision	46.18	47.66	77.08	81.18
Macro Recall	50.98	52.96	51.07	52.80
Macro F-Score	47.10	48.62	60.00	62.47
Macro <sup>2</sup> F-Score	48.46	50.17	61.43	63.99

## Example

(bla () bla)→bla bla

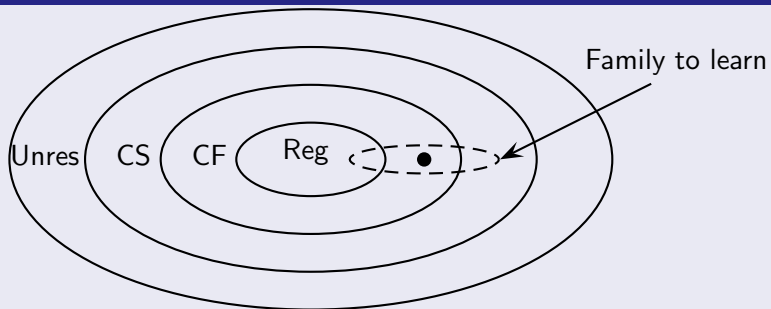
## Evaluation insights

- No standard evaluation exists but de facto evaluation datasets arise
  - ATIS (van Zaanen and Adriaans 2001)
  - WSJ10, WSJ40 (WSJ with sentence length limitations)
  - NEGRA10 (German)
  - CTB10 (Chinese)
- Systems have different input/output
- Evaluation settings influence results
  - Different metrics (micro/macro/macro<sup>2</sup>)
  - Included constituents (sentence/empty)
- Formal grammatical inference does not have this problem
  - Evaluation performed through formal proofs

## Context-sensitive grammars

- Learning context-free grammars is hard
- Is learning context-sensitive grammars impossible?
  - That depends
  - To what degree is the grammar context-sensitive?
- We may not need “full” context-sensitiveness
  - Grammar rules:  $\alpha A \beta \rightarrow \alpha \gamma \beta$
- Mildly context-sensitive grammars may be enough for NL (Huybrechts 1984, Shieber 1985)
- Perhaps the full power of context-freeness is not needed

## Family of languages



## Learning context-sensitive languages

- Open research area
- Some work has already been done
  - Augmented Regular Expressions (Alquézar 1997)
  - Variants of substitutability (Yoshinaka 2009)
  - Distributional Lattice Grammars (Clark 2010)

## Relationship between empirical and formal GI

- Is there a relationship between empirical GI and formal GI?
- Example: consider the case of substitutability
- There are situations in which substitutability breaks:
  - John eats meat
  - John eats much
- This suggests that learning based on substitutability learns a different family of languages (not CFG)
- Non-terminally separated (NTS) languages
  - Subclass of deterministic context-free grammars

## Learning NTS grammars

- Grammar  $G = \langle \Sigma, V, P, S \rangle$  is NTS
  - $\Sigma$  is vocabulary
  - $V$  is set of non-terminals
  - $P$  is set of production rules
  - $S \in V$  is the start symbol
- Additional restriction:
  - If  $N \in V$
  - $N \xRightarrow{*} \alpha\beta\gamma$
  - $M \xRightarrow{*} \beta$
  - then  $N \xRightarrow{*} \alpha M \gamma$
- In other words:  
non-terminals correspond exactly with substitutability
- (Clark and Eyraud 2005, Clark 2006, Clark and Eyraud 2007)

## Learning NTS grammars

- It can be shown that NTS grammars are
  - identifiable in the limit
  - PAC learnable
- Unfortunately, natural language is not an NTS language
- Ultimate goal:
  - Find family of languages that fits natural language
  - and is learnable in the right learning setting

## Formal GI and empirical GI

- Relation between formal GI and empirical GI
  - Formal GI can show learnability
    - Under certain conditions
  - Empirical GI tries to learn structure from real data
    - Practically shows possibilities and limitations
- Ultimate aim: Find family of languages that is
  - learnable under different conditions
  - fits natural languages

## Conclusions

- 1 There are new strong positive results in the recent past for learning (subclasses of) DFA, PFA, transducers, CFGs, and MCSGs.
- 2 The use of GI techniques both in computational linguistics is taking place and the future is bright!

# International Colloquium of Grammatical Inference



September 2012  
Washington, D.C. area

- Talks, tutorials, competitions, ...
- Please keep an eye out for announcements.