# Constraint-driven analysis
# of formal languages

Dakotah Lambert

16 Feb 2024

- Classification
  - What relations are available?
  - How are they used?
- **M** operator :: multiple tiers
- Learning
  - What information is available?
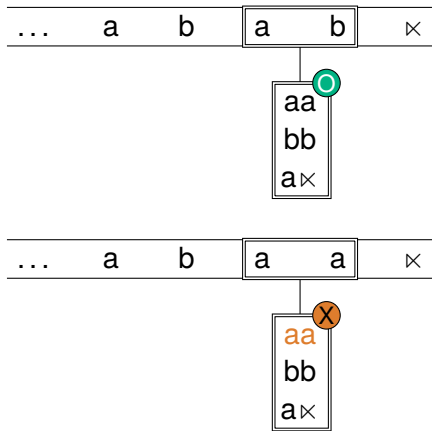  - How do we extrapolate from it?

Deriving a class hierarchy
as the linguists would

ab
bab
abab
. . .

(Asmat stress)

sik'is

ʃitʃedza

Navajo sibilant harmony

$aXb \in L$ and $cXd \in L \longrightarrow aXd \in L$
$X$ shared, length at least $k$

s(it)$^k$is and ʃ(it)$^k$iʃ,
but not s(it)$^k$iʃ

not strictly local

Acceptability based on set $S$ of factors.

Symmetric harmony: $\{s, \int\} \not\subseteq S$

Less than two 'b':
a, ab, abaa, aaaaab,
but not abab or abba

(basically every stress pattern ever)

$a^k ba^k$ and $a^k ba^k ba^k$ have same $k$-factor set

not locally testable

use first-order logic instead:
$$\neg(\exists x, y)[x \neq y \wedge b(x) \wedge b(y)]$$

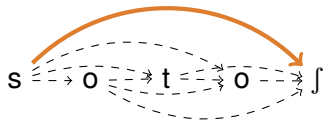FO: locally threshold testable
|
Prop: locally testable
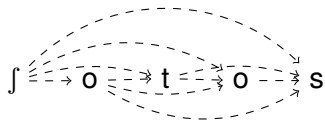|
CNL: strictly local

ʃotos but not sotoʃ

(attested in Sarcee)

# Asymmetric harmony via precedence

# Reanalyzing harmony: Tiers

$$s \xrightarrow{\quad\quad\quad\quad\quad} \int$$

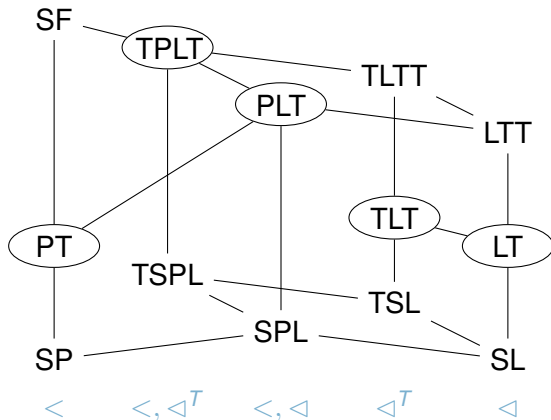o      t      o

other symbols → neutral

# Classes

factors "abc comes before def"
piecewise-locally testable = dot-depth one

propositional level subsumes LTT

# Varieties

### Variety
a class $\mathcal{V}$ where for each alphabet $\Sigma$,
if $L_1, L_2 \in \Sigma^*\mathcal{V}$:

- $\complement L_1 \in \Sigma^*\mathcal{V}$ and $L_1 \cup L_2 \in \Sigma^*\mathcal{V}$
  Boolean operations$^*$
- $\sigma^{-1}L_1 \in \Sigma^*\mathcal{V}$ and $L_1\sigma^{-1} \in \Sigma^*\mathcal{V}$
  Quotients
- $f : \Gamma \to \Sigma$ homomorphic, $f^{-1}(L_1) \in \Gamma^*\mathcal{V}$
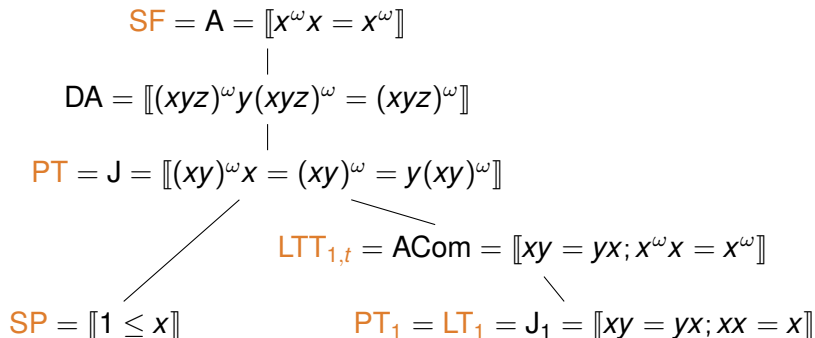  Inverse homomorphisms

$$\Sigma^* \mathcal{V} \sim \text{variety of monoids}$$
$$\Sigma^+ \mathcal{V} \sim \text{variety of semigroups}$$

collection closed under

- submonoid (subsemigroup)
- quotient
- finitary direct product

# Piecewise branch, expanded

$$\mathsf{SF} = \mathsf{A} = [\![x^\omega x = x^\omega]\!]$$
$$|$$
$$\mathsf{DA} = [\![(xyz)^\omega y(xyz)^\omega = (xyz)^\omega]\!]$$
$$|$$
$$\mathsf{PT} = \mathsf{J} = [\![(xy)^\omega x = (xy)^\omega = y(xy)^\omega]\!]$$

$$\mathsf{LTT}_{1,t} = \mathsf{ACom} = [\![xy = yx; x^\omega x = x^\omega]\!]$$

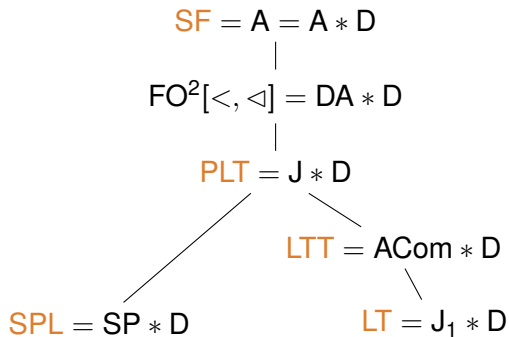$$\mathsf{SP} = [\![1 \le x]\!] \qquad\qquad \mathsf{PT}_1 = \mathsf{LT}_1 = \mathsf{J}_1 = [\![xy = yx; xx = x]\!]$$

convert *k*-factors to their own individual letters

$$\mathbf{V} \mapsto \mathbf{V} * \mathbf{D}$$

contains corresponding piecewise class

$$SF = A = A * D$$
$$|$$
$$FO^2[<, \lhd] = DA * D$$
$$|$$
$$PLT = J * D$$

$$LTT = ACom * D$$

$$SPL = SP * D$$

$$LT = J_1 * D$$

**MV** the variety of monoids
generated by $S^{\cdot}$ for $S \in$ **V**

linguistic "lift onto a tier" = algebraic "$S \mapsto S^{\cdot}$"

multiple tiers interacting (Boolean combinations):
converts $+$-variety $\mathcal{V}$ to $*$-variety $\mathcal{MV}$

What kinds of data do learners receive?

How do we extrapolate from that back to patterns?

# Limit-learnability with positive data

- Only valid words happen
- Every valid word will eventually happen
- Finite samples
- Incrementally: eventually hypothesis stops changing

# String extension learning

- Assume nothing is valid
- For each word, extract information
- Add that information into a "grammar"
- Information is never removed from the "grammar"
- How is this interpreted?

Information: "set of subsequences"
Insertion: set-union

Interpretation:
valid iff set of subsequences is subset of grammar

Information: "set of letters"
Insertion: element-insertion

Interpretation:
valid iff set of letters in grammar

Information: "thresholding multiset of letters"
Insertion: element-insertion

Interpretation:
valid iff multiset of letters in grammar

Information: "set of subsequences"
Insertion: element-insertion

Interpretation:
valid iff set of subsequences in grammar

choose *k*
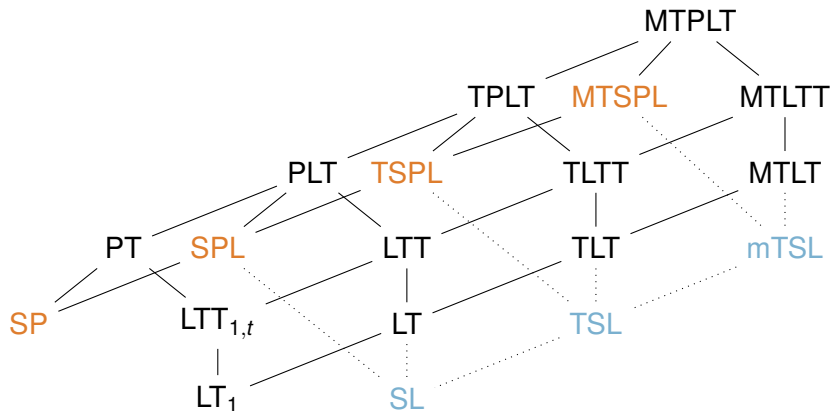apply factor–letter transformation
learn piecewise base class

choose *k*
for each subset of the alphabet:
apply erasing transformation then factor–letter transformation
information based on this collection

# The system

- Piecewise base class
- Close under inverse factor-collapse ("localize")
- Close under neutral-letter injection ("tierify")
- Close under Boolean operations ("multitierify")
- Result: a new Piecewise base class

$$\mathbf{A} = \mathbf{A} * \mathbf{D} = \mathbf{T}(\mathbf{A} * \mathbf{D}) = \mathbf{M}(\mathbf{A} * \mathbf{D})$$

and so much more

# Future directions

- Parameter-finding from machines
  known for some like D, Acom, . . .
- Decomposition of machines
- Same structural classes apply to functions
  inferring those? (SOSFIA/++)
- Other bases