

Structure matters: learning nonlocal phonotactics in a Strictly Piecewise phonotactic model

Huteng Dai

Rutgers University

Nonlocal phonotactics

- The speakers' knowledge of possible and impossible sound sequences:

legal *brick* [brɪk]

legal *blick* [blɪk]

illegal **bnick* [bnɪk]

(Chomsky & Halle, 1965; Gorman, 2013)

- Nonlocal phonotactics: the phonotactic knowledge of **nonadjacent** sound sequences at **arbitrary** distance.

Quechua nonlocal restrictions of laryngeal cooccurrence

Nonlocal stop-ejective and stop-aspirate pairs are illegal in Quechua.

- stop-ejective: *kut'u, *k'ut'u, *k^hut'u;
- stop-aspirate: *kut^hu, *k'ut^hu, *k^hut^hu;
- legal: k'utuj 'to cut', rit'i 'snow', jut^hu 'partridge'.

(Gouskova & Gallagher, 2020)

How do speakers learn a finite phonotactic grammar that distinguish legal and illegal words from an **infinite** set of possible sound sequences?

(Hayes & Wilson, 2008; Heinz, 2010)

Local n -grams and baseline Learner

- Local n -gram: contiguous sequence of n items;
- Previous works usually hypothesize **local** n -grams as the free parameters/constraints (grammar) of the phonotactic learner.

(Hayes & Wilson, 2008)

Local n -grams and baseline Learner

- Local n -gram: contiguous sequence of n items;
- Previous works usually hypothesize **local** n -grams as the free parameters/constraints (grammar) of the phonotactic learner.

(Hayes & Wilson, 2008)

- Imagine a learner only observed one word k'utuj:

n	observed local n -grams	unobserved local n -grams
2	k'u, ut, tu, uj	*uk'...
3	k'ut, utu, tuj	*tuk'...

- E.g. *tuk'u will be penalized by the bi-/trigram constraints (*uk', *tuk').

Challenge

- The parameter space will explode if the learner memorizes local n -grams to approximate nonlocal phonotactics;
- E.g. ***t**u**ʌ**k'u requires local 4-grams, ***t**u**p**u**k**'u requires local 5-grams, ...
(Hayes & Wilson, 2008; Gouskova & Gallagher, 2020)
- Any such approximation also completely misses the generalization of nonlocal interaction at arbitrary distance.

(Heinz, 2010)

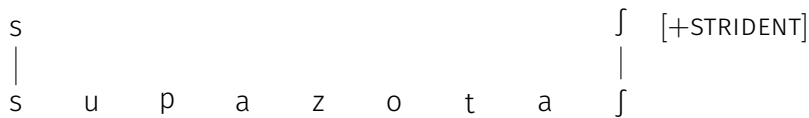


Figure 1: z blocks tier-based local bigram *s]

- Learning local n -gram on tiers/projections predicts unattested *blocking* effects.
Q & A (Heinz, 2010; Rose & Walker, 2004; Hansson, 2010)
- Tiers also leads to new problems e.g. learning tiers.
(Hayes & Wilson, 2008; Jardine & McMullin, 2017; Gouskova & Gallagher, 2020)

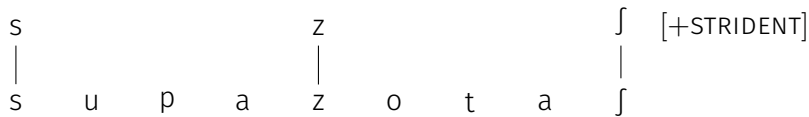


Figure 1: z blocks tier-based local bigram *sj

- Learning local n -gram on tiers/projections predicts unattested *blocking* effects.
Q & A (Heinz, 2010; Rose & Walker, 2004; Hansson, 2010)
- Tiers also leads to new problems e.g. learning tiers.
(Hayes & Wilson, 2008; Jardine & McMullin, 2017; Gouskova & Gallagher, 2020)

Goals and contributions

- Modeling and learning nonlocal phonotactics **without tiers**;
- Integrating Formal Language Theory and statistical learning to handle noisy corpus data, and to predict the gradient acceptability of nonce forms.

Strictly Piecewise phonotactic model

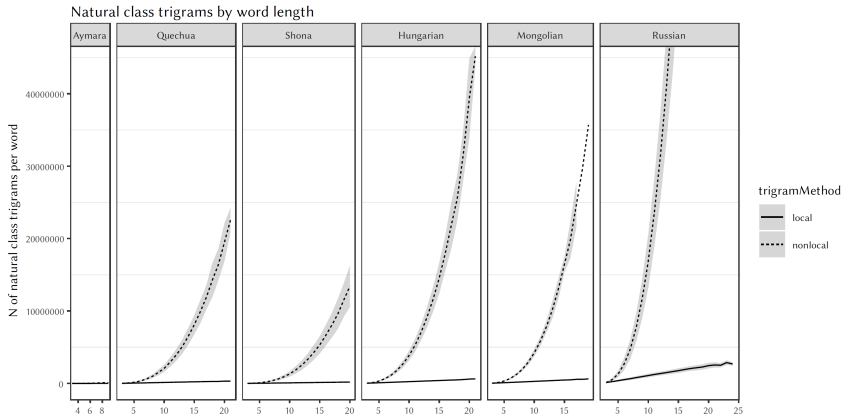
Subsequences

- Subsequences *aka.* **nonlocal** n -grams keep track of the **order** between symbols; e.g. if the learner observes $k'utuj$:

n	observed subsequences	unobserved subsequences
2	$k'...u, k'...t, k'...j, \dots$	$*t...k', \dots$
3	$k'...u...t, k'...u...j, \dots$	$*t...u...k', \dots$

- Strictly Piecewise (SP) grammar evaluates nonlocal n -grams; e.g. $*tuk'uj$, $*tu\wedge k'u$, $*tupuk'u$ are all penalized by nonlocal bigram $*t...k'$ (“ t precedes k' ”).
(Heinz & Rogers, 2010)

Problem of searching nonlocal n -grams

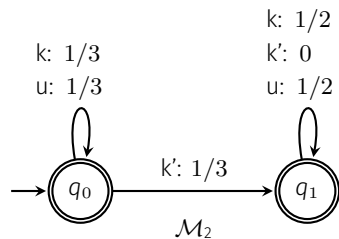
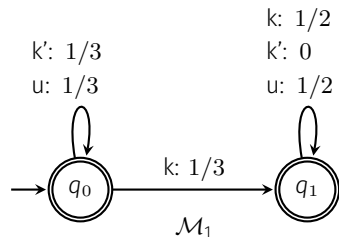


- “Devising a computationally efficient search...will require a sophisticated implementation that...is currently lacking.”

(Gouskova & Gallagher, 2020)

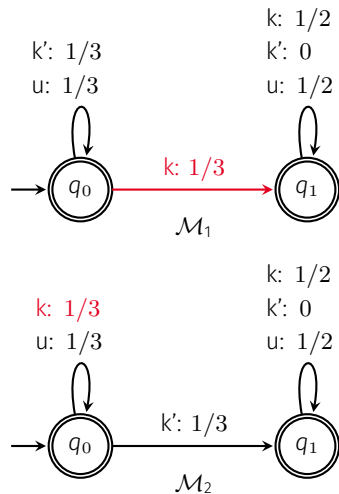
Solution: Parsing SP phonotactic model

- Strictly Piecewise grammar can be characterized by a set of Weighted Deterministic Finite-state Automata (WDFAs). (Shibata & Heinz, 2019)
- E.g. $\{\mathcal{M}_1, \mathcal{M}_2\}$ bans $\{*\textcolor{blue}{k}...\textcolor{red}{k}', *\textcolor{blue}{k}'...\textcolor{red}{k}'\}$ with a simplified alphabet $A = \{k, k', u\}$.



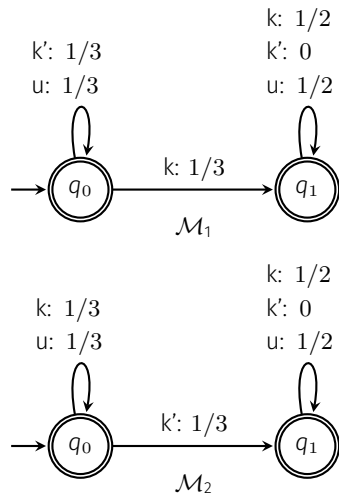
Parameters

- The parameters are transition weights $W(\mathcal{M}, q, \sigma)$ given the machine \mathcal{M} , state q , and segment σ .



Target symbol

- Each machine \mathcal{M} only checks if it has seen one specific **target symbol** σ ;
- No \Rightarrow stay in state q_0 ;
- Yes \Rightarrow go to state q_1 ;



Parsing: coemission probability

- Coemission probability synchronizes the parameters on different machines at the same time:

$$\text{Coemit}(\sigma_i) = \frac{\overbrace{\prod_{j=1}^K W(\mathcal{M}_j, q, \sigma_i)}^{\text{for one specific segment } \sigma_i}}{\underbrace{\sum_{\sigma' \in A} \prod_{j=1}^K W(\mathcal{M}_j, q, \sigma')}_{\text{normalizer}}}$$

(Shibata & Heinz, 2019)

$$\mathcal{M}_1: \quad q_0 \xrightarrow[1/3]{k} q_1 \xrightarrow[1/2]{u} q_1 \xrightarrow[0]{k'} q_1$$

$$\mathcal{M}_2: \quad q_0 \xrightarrow[1/3]{k} q_0 \xrightarrow[1/3]{u} q_0 \xrightarrow[1/3]{k'} q_1$$

$$\text{Coemit}(\sigma_i): \epsilon \xrightarrow[1/3]{k} \sigma_1 \xrightarrow[1/2]{u} \sigma_2 \xrightarrow[0]{k'} \sigma_3$$

$$\text{Time:} \quad t_0 \longrightarrow t_1 \longrightarrow t_2 \longrightarrow t_3$$

Parsing: graphical model over time

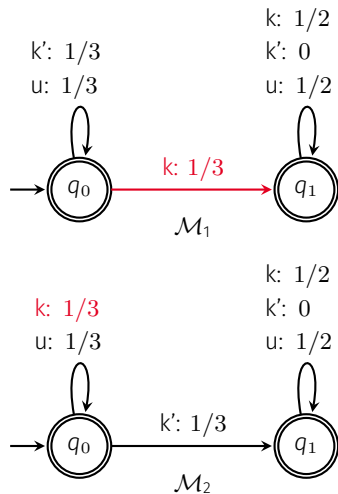
- E.g. The model reads *kuk':

$$\mathcal{M}_1: \quad q_0 \xrightarrow[k: 1/3]{k} q_1$$

$$\mathcal{M}_2: \quad q_0 \xrightarrow[k: 1/3]{k} q_0$$

$$\text{Coemit}(\sigma_i): \quad \epsilon \xrightarrow[k: 1/3]{k} \sigma_1$$

$$\text{Time:} \quad t_0 \longrightarrow t_1$$



Parsing: graphical model over time

$$\mathcal{M}_1: \quad q_0 \xrightarrow[\frac{1}{3}]{k} q_1 \xrightarrow[\frac{1}{2}]{u} q_1$$

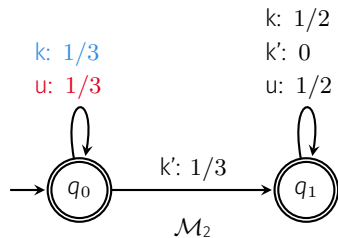
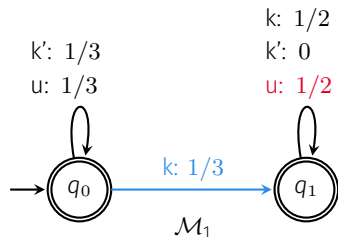
$$\mathcal{M}_2: \quad q_0 \xrightarrow[\frac{1}{3}]{k} q_0 \xrightarrow[\frac{1}{2}]{u} q_0$$

$$\text{Coemit}(\sigma_i): \quad \epsilon \xrightarrow[\frac{1}{3}]{k} \sigma_1 \xrightarrow[\frac{1}{2}]{u} \sigma_2$$

$$\text{Time:} \quad t_0 \xrightarrow{\quad} t_1 \xrightarrow{\quad} t_2$$

$\xrightarrow{\sigma:W}$ preceding transition;

$\xrightarrow{\sigma:W}$ current transition;



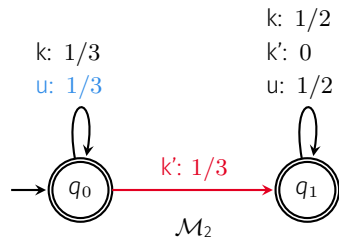
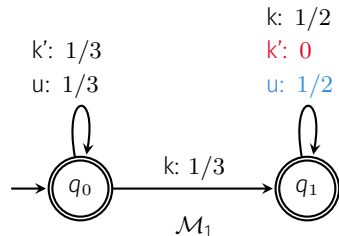
Parsing: graphical model over time

$$\mathcal{M}_1: \quad q_0 \xrightarrow[1/3]{k} q_1 \xrightarrow[1/2]{u} q_1 \xrightarrow[0]{k'} q_1$$

$$\mathcal{M}_2: \quad q_0 \xrightarrow[1/3]{k} q_0 \xrightarrow[1/3]{u} q_0 \xrightarrow[1/3]{k'} q_1$$

$$\text{Coemit}(\sigma_i): \quad \epsilon \xrightarrow[1/3]{k} \sigma_1 \xrightarrow[1/2]{u} \sigma_2 \xrightarrow[0]{k'} \sigma_3$$

$$\text{Time:} \quad t_0 \longrightarrow t_1 \xrightarrow{\text{blue}} t_2 \xrightarrow{\text{red}} t_3$$



$$\mathcal{M}_1: q_0 \xrightarrow[1/3]{k} q_1 \xrightarrow[1/2]{u} q_1 \xrightarrow[0]{k'} q_1$$

$$\mathcal{M}_2: q_0 \xrightarrow[1/3]{k} q_0 \xrightarrow[1/3]{u} q_0 \xrightarrow[1/3]{k'} q_1$$

$$\text{Coemit}(\sigma_i): \epsilon \xrightarrow[1/3]{k} \sigma_1 \xrightarrow[1/2]{u} \sigma_2 \xrightarrow[0]{k'} \sigma_3$$

- Word likelihood is the product of coemission probabilities of all the segments in a word:

$$\text{lhd}(w) = \text{lhd}(\sigma_1 \sigma_2 \dots \sigma_N) = \prod_{i=1}^N \text{Coemit}(\sigma_i)$$

- E.g. $\text{lhd}(kuk') = 1/3 \cdot 1/2 \cdot 0 = 0$, $\text{Coemit}(k') = 0$ given $k \Rightarrow *k\dots k'$ is penalized.

Learning

Learning problem

- Problem: to optimize parameters $\hat{W}(\mathcal{M}, q, \sigma)$ so that the generated distribution maximally approaches the target distribution \mathcal{D} .
- In practice, the parameters are optimized by minimizing the **negative log likelihood** (NLL) of a sample/wordlist S drawn from \mathcal{D} :

$$\hat{W}(\mathcal{M}, q, \sigma) = \arg \min_W - \sum_{w \in S} \log \text{lhd}(w).$$

↑

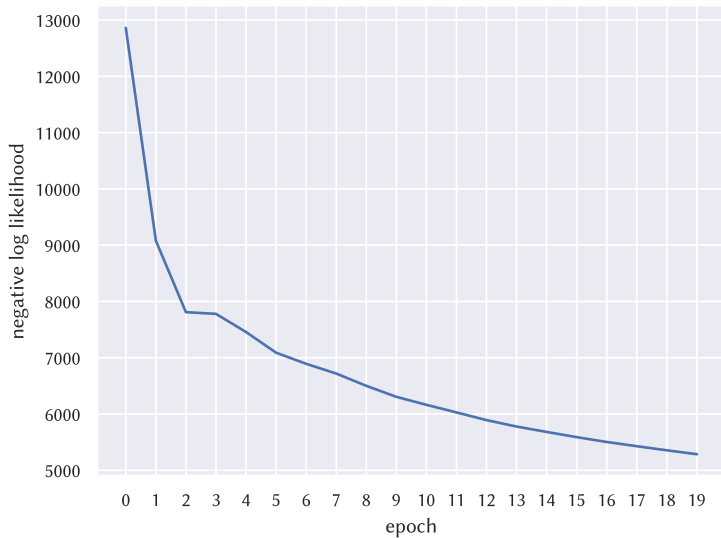
Maximum Likelihood Estimation \approx Maximum Entropy

- Training data: 10,848 unlabelled legal phonological words;
- Testing data: 24,352 nonce forms (C_1VC_2V and C_1VCC_2V) which were manually labelled as legal ($N = 18,502$), stop-aspirate ($N = 3,645$), stop-ejective ($N = 2,205$)¹. (Gouskova & Gallagher, 2020)

Training	Testing	
a h i n a \wedge a m a n t a q a	tʰ a tʰ a	illegal-aspirate
t' u k u tʃ i j a w a ŋ k i	tʰ a \wedge tʰ a	illegal-aspirate
qʰ e r k i ŋ tʃ o q a	tʰ a \wedge tʃ a	legal
...	...	

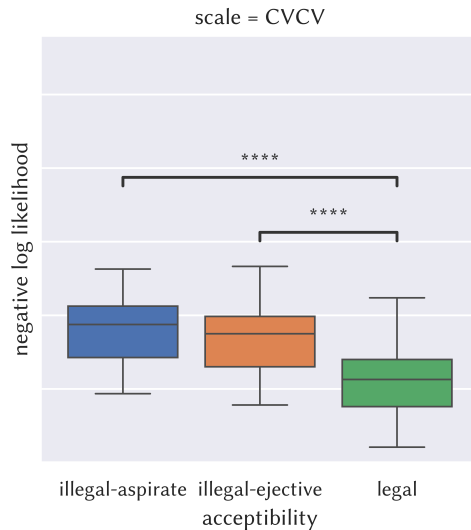
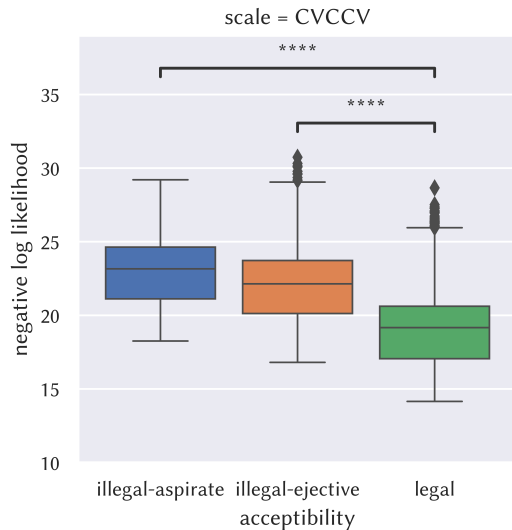
¹Thanks to Prof. Maria Gouskova and Prof. Gillian Gallagher for publishing their dataset
https://github.com/gouskova/inductive_projection_learner/tree/master/data/quechua!

Maximum Likelihood Estimation



- Can't test the accuracy since it's unsupervised learning with unlabelled data.
- Ask if the learned model distinguish the NLL of illegal words from legal words in testing data → Clustering + Mann-Whitney U test
- If the learning is successful, legal words should have lower NLL (higher likelihood).

Learning segment-based representation



Discussion and conclusion

- SP phonotactic model only keeps track of nonlocal n -grams, which guarantees the efficient learning of nonlocal phonotactics.
- The structure studied extensively in Formal Language Theory (FLT) is the conditions on the parameter space such as nonlocal n -grams.
(Heinz, 2018; Jardine & Heinz, 2016; Chandlee et al., 2019)

- There has been a gap between FLT and noisy corpus data; (Heinz & Rawski, in press; Gouskova & Gallagher, 2020)
- The computational learning theory grounded on FLT focuses on the theorem and proof of learnability instead of simulation;
- However, understanding the domain-specific, structural properties of small dataset can help us to handle large noisy dataset.
(Heinz, 2010; Jardine & Heinz, 2016; Jardine & McMullin, 2017)

Acknowledgement

I thank Jeff Heinz, Adam Jardine, Bruce Tesar, Adam McCollum, and Jon Rawski for their comments and insights. My special thanks are extended to Brian Pinsky, Liam Schramm, and Yu Cao for providing the valuable suggestions on the implemented Python code.



Ineseño Chumash nonlocal sibilant phonotactics

- In Ineseño Chumash, the co-occurrence of alveolar {s, z, t̪s, d̪z,...} and lamino-postalveolar {ʃ, ʒ, t̪ʃ, d̪ʒ,...} sibilants is illegal e.g. *ʃ...s, *s...ʃ.

(Applegate, 1972)

	3-grams	5-grams
(1) ʃapit̪ʰolit /s-api-t̪ʰo-it/ ‘I have a stroke of good luck’	ʃap	ʃapit̪ʰ
(2) ʃapit̪ʰoluʃwaf /ʃ-api-t̪ʰo-us-waf/ ‘He had a stroke of good luck’	api	apit̪ʰo
	pit̪ʰ	pit̪ʰol
(3) *sapit̪ʰolit, *ʃapit̪ʰoluswaf		...

Ineseño Chumash nonlocal sibilant phonotactics

- In Ineseño Chumash, the co-occurrence of alveolar {s, z, t̪s, d̪z,...} and lamino-postalveolar {ʃ, ʒ, t̪ʃ, d̪ʒ,...} sibilants is illegal e.g. *ʃ...s, *s...ʃ.

(Applegate, 1972)

		3-grams	5-grams
(4)	ʃapit̪ʰolit /s-api-t̪ʰo-it/ 'I have a stroke of good luck'	ʃap	ʃapit̪ʰ
(5)	ʃapit̪ʰoluʃwaf /ʃ-api-t̪ʰo-us-waf/ 'He had a stroke of good luck'	api	apit̪ʰo
		pit̪ʰ	pit̪ʰol
(6)	*sapit̪ʰolit, *ʃapit̪ʰoluswaf		...

- Trigrams won't work → difficult to choose current window n .

Ineseño Chumash and nonlocal n -grams

- (7) $\text{ʃapitʃ}^{\text{h}}\text{olit}$ $/\text{s-api-tʃ}^{\text{h}}\text{o-it}/$
‘I have a stroke of good luck’
- (8) $\text{ʃapitʃ}^{\text{h}}\text{oluʃwaʃ}$ $/\text{ʃ-api-tʃ}^{\text{h}}\text{o-us-waʃ}/$
‘He had a stroke of good luck’
- (9) $*\text{sapitʃ}^{\text{h}}\text{olit}, *ʃapitʃ^{\text{h}}\text{oluswaʃ}$

legal	illegal
ʃ...tʃ^{h}	$*\text{s...tʃ}^{\text{h}}$
$\text{tʃ}^{\text{h}}\text{...ʃ}$	$*\text{tʃ}^{\text{h}}\text{...s}$
ʃ...ʃ	$*\text{s...ʃ}$
...	...

Feature-based representation

- **Feature-based** model can be implemented by replacing the alphabet by a set of feature values $[\alpha F]$. For example, given the simple feature system below:

	<hr/>	
	F	G
<hr/>		
a	+	-
b	+	+
<hr/>		

Feature-based SP phonotactic model

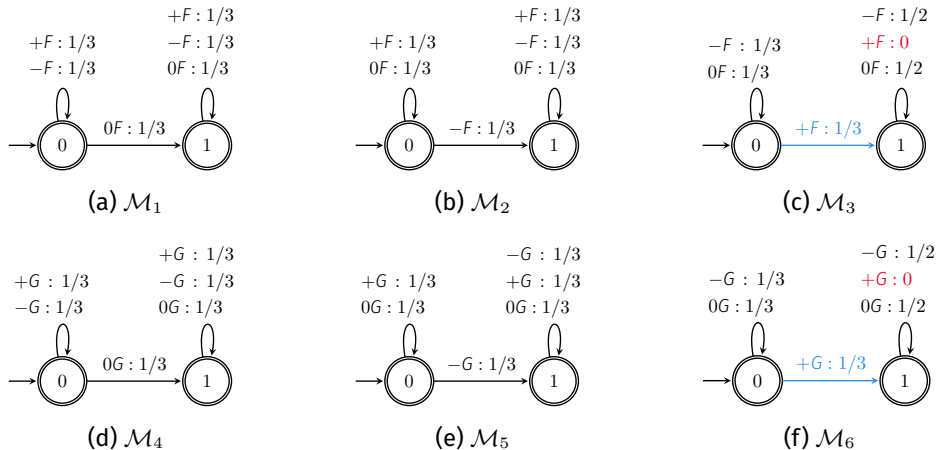
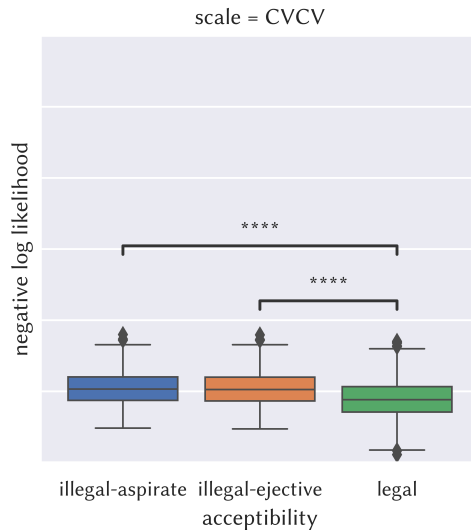
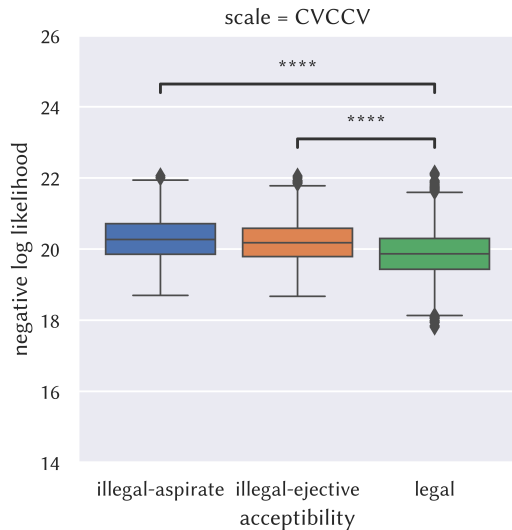


Figure 2: The feature-based SP phonotactic model which bans $*+F \dots +F$ and $*+G \dots +G$ with the simple feature system

Learning feature-based representation



Forward algorithm

$NLL \leftarrow 0;$

for *word* *in* S **do**

 state $\leftarrow 0$ in each automaton \mathcal{M}_j ;

for σ_i *in* *word* **do**

 Initialize a lookup dictionary D for $\prod_{j=1}^K T(\mathcal{M}_j, q, \sigma')$;

for \mathcal{M}_j *in* *automata* **do**

for σ' *in* *alphabet* **do**

 Update the lookup dictionary with σ' ;

 Update the state on \mathcal{M}_j ;

$NLL \leftarrow NLL - \log(\text{Coemit}(\sigma_i))$

Result: Negative log likelihood NLL of S

Reference

- Applegate, R. (1972). *Ineseño chumash grammar* (Unpublished doctoral dissertation). University of California, Berkeley.
- Chandlee, J., Eyraud, R., Heinz, J., Jardine, A., & Rawski, J. (2019). Learning with partially ordered representations. *arXiv preprint arXiv:1906.07886*.
- Chomsky, N., & Halle, M. (1965). Some controversial questions in phonological theory. *Journal of linguistics*, 1(2), 97–138.
- Gorman, K. (2013). *Generative phonotactics* (Unpublished doctoral dissertation). University of Pennsylvania.
- Gouskova, M., & Gallagher, G. (2020). Inducing nonlocal constraints from baseline phonotactics. *Natural Language & Linguistic Theory*, 1–40.
- Hansson, G. Ó. (2010). *Consonant harmony: Long-distance interactions in phonology* (Vol. 145). Univ of California Press.
- Hayes, B., & Wilson, C. (2008). A maximum entropy model of phonotactics and phonotactic learning. *Linguistic inquiry*, 39(3), 379–440.

- Heinz, J. (2010). Learning long-distance phonotactics. *Linguistic Inquiry*, 41(4), 623–661.
- Heinz, J. (2018). The computational nature of phonological generalizations. *Phonological Typology, Phonetics and Phonology*, 126–195.
- Heinz, J., & Rawski, J. (in press). History of phonology: Learnability. In E. Dresher & H. van der Hulst (Eds.), *Oxford handbook of the history of phonology* (chap. 32). Oxford University Press.
- Heinz, J., & Rogers, J. (2010). Estimating strictly piecewise distributions. In *Proceedings of the 48th annual meeting of the association for computational linguistics* (pp. 886–896).
- Jardine, A., & Heinz, J. (2016). Learning tier-based strictly 2-local languages. *Transactions of the Association for Computational Linguistics*, 4, 87–98.
- Jardine, A., & McMullin, K. (2017). Efficient learning of tier-based strictly k-local languages. In *International conference on language and automata theory and applications* (pp. 64–76).

- Rose, S., & Walker, R. (2004). A typology of consonant agreement as correspondence. *Language*, 475–531.
- Shibata, C., & Heinz, J. (2019). Maximum likelihood estimation of factored regular deterministic stochastic languages. In *Proceedings of the 16th meeting on the mathematics of language (mol 16)*.