

# Modeling morphological productivity with the Minimal Generalization Learner

Adam Albright (albright@mit.edu)

Massachusetts Institute of Technology

July 4, 2012 — DRAEM



# Outline

- 1 Model overview
  - Minimal generalization
  - Reliability
  - Confidence
- 2 A simple simulation
- 3 Extensions
  - Phonological rules
  - 'Impugnment'
  - Non-local and suprasegmental contexts

# Outline

## 1 Model overview

- Minimal generalization
- Reliability
- Confidence

## 2 A simple simulation

## 3 Extensions

- Phonological rules
- 'Impugnment'
- Non-local and suprasegmental contexts

# Outline

- 1 Model overview
  - Minimal generalization
  - Reliability
  - Confidence
- 2 A simple simulation
- 3 Extensions
  - Phonological rules
  - 'Impugnment'
  - Non-local and suprasegmental contexts

# Goal of the model

- Given: a set of morphologically related forms—e.g., English:

Present	Past		Present	Past	
mɪs	mɪst	' <i>miss(ed)</i> '	nɪd	nɪdəd	' <i>need(ed)</i> '
pɹɛs	pɹɛst	' <i>press(ed)</i> '	dʒʌmp	dʒʌmpt	' <i>jump(ed)</i> '
læf	læft	' <i>laugh(ed)</i> '	plæn	plænd	' <i>plan(ned)</i> '
hʌg	hʌgd	' <i>hug(ged)</i> '	sɪŋ	sæŋ	' <i>sing/sang</i> '
rʌb	rʌbd	' <i>rub(bed)</i> '	dɹɪŋk	dræŋk	' <i>drink/drank</i> '

- Goal: generalize to new items, such as novel verb [splɪŋk]<sub>Present</sub>

- Expected outputs

Very likely/acceptable: splɪŋkt  
Somewhat likely/acceptable: splʌŋk, splæŋk  
Quite unlikely/unacceptable: splɒt

- Not:

- \*splɪŋkd, \*splɪŋkəd (misapplication of phonology)
- \*splɒŋŋ (valid past, wrong context)
- \*splɪnd (idiosyncratic change: *make~made*)

# Desired information

## Some questions

- What string mappings relate the morphological categories? ( $\emptyset \rightarrow d$ ,  $i \rightarrow \text{æ}$ ,  $iŋk \rightarrow \text{ɔt}$ , etc.)
- Are some mappings phonological variants of others? ( $\emptyset \rightarrow d$ ,  $\emptyset \rightarrow t$ ,  $\emptyset \rightarrow \text{əd}$ )
- Are some mappings restricted (categorically, probabilistically) to specific phonological contexts?
- What is the relative strength of the mappings in various contexts?

# Induction strategy

- Parse pairs to discover changes and contexts—e.g.,
  - $\text{mis} \sim \text{mist} = \emptyset \rightarrow t / \text{mis} \_\_\_ \#$
  - $\text{d3}\lambda\text{mp} \sim \text{d3}\lambda\text{mpt} = \emptyset \rightarrow t / \text{d3}\lambda\text{mp} \_\_\_ \#$
- Compare contexts to discover broader contexts
  - $\text{mis}, \text{d3}\lambda\text{mp}: /[-\text{voi}] \_\_\_$
- Evaluate accuracy of resulting contexts

# Discovering changes

Pinker & Prince (1988), p. 130:

- 1 “Candidates for rules are hypothesized by comparing base and past tense versions of a word, and factoring apart the changing phonetic portion, which serves as the rule operation, from certain morphologically-relevant phonological components of the stem, which serve to define the class of stems over which the operation can apply.” (= the *context*)
- Pinker and Prince (1988), Ling and Marinov (1993), Albright and Hayes (2002)

# A simple parsing strategy

Albright and Hayes (2002, 2003)

- Edge-in alignment: start from both left and right, aligning until there's a mismatch
- Combine leftward and rightward alignments to align as many segments as possible

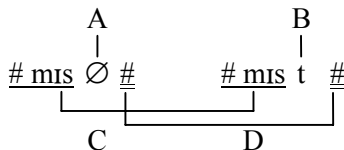
$$\begin{array}{ccc} m_1 & I_2 & s_3 \\ m_1 & I_2 & s_3 \end{array} \parallel \begin{array}{c} t \\ t \end{array} + \begin{array}{ccc} m & I & s \\ m & I & s \end{array} \parallel \begin{array}{c} t \\ t \end{array} \mapsto \begin{array}{ccc} m_1 & I_2 & s_3 \\ m_1 & I_2 & s_3 \end{array} \parallel \begin{array}{c} t \\ t \end{array}$$

- Medial changes: leftward and rightward scans are both able to align some segments; residue left in the middle

$$\begin{array}{ccc} s_1 & I & \eta \\ s_1 & \text{æ} & \eta \end{array} \parallel + \begin{array}{ccc} s & I & \eta_1 \\ s & \text{æ} & \eta_1 \end{array} \parallel \mapsto \begin{array}{ccc} s_1 & I & \eta_2 \\ s_1 & \text{æ} & \eta_2 \end{array} \parallel$$

# Factoring out changes and contexts

- Mismatched portions = change ( $A \rightarrow B$ )
- Aligned portions = context ( $C \text{ \_\_\_ } D$ )



- Resulting rule:  $\emptyset \rightarrow t / \# \text{mis} \text{ \_\_\_ } \#$

# Word-specific rules

Result: a set of word-specific rules

- a.  $\emptyset \rightarrow t / \#mɪs \_\_ \#$
- b.  $\emptyset \rightarrow t / \#prɛs \_\_ \#$
- c.  $\emptyset \rightarrow t / \#læf \_\_ \#$
- d.  $\emptyset \rightarrow d / \#hʌg \_\_ \#$
- e.  $\emptyset \rightarrow d / \#rʌb \_\_ \#$
- f.  $\emptyset \rightarrow əd / \#nɪd \_\_ \#$
- g.  $\emptyset \rightarrow t / \#dʒʌmp \_\_ \#$
- h.  $\emptyset \rightarrow d / \#plæn \_\_ \#$
- i.  $ɪ \rightarrow æ / \quad \quad \quad s \_\_ ɪ \#$
- j.  $ɪ \rightarrow æ / \quad \quad \quad dr \_\_ ɪ k \#$

# Comparing rules

- *miss* and *press* both take  $\emptyset \rightarrow t$ . Assume that this is because they have some crucial property in common
  - In both cases, the change is word-final
  - In both cases, the segment before the change is an [s]
  - In both cases, the segment before [s] is a non-low lax front V
  - etc...??? (sonorant before the vowel, monosyllabic, etc.)
- Albright and Hayes (2002): **Minimal generalization** approach
  - Conservative collapsing: new rule keeps *everything* that the pair has in common

# Minimal generalization

- Once again, some alignment scheme is needed. Assume *locality*

m	i	s	___	#
p	r	ε	___	#

Precludes many possible generalizations, such as:

- End in /s/
  - Vowel is front, lax
  - First segment is labial
- A pragmatic assumption: *myopic generalization*
    - Once a mismatch is encountered and featural generalization is needed, shared similarities farther away from the change location are not likely to be crucial
    - Convert more distant segments to free variables (X, Y)
    - We'll come back to this assumption later

# Minimal generalization

Comparing *miss* and *press*:  $\emptyset \rightarrow t / \dots$

	Residue	Shared Features	Shared Segments	Change Location	Shared Segments	Shared Features	Residue
A.	#m	I	S	—	#		
B.	#pr	ε	S	—	#		
C.	X	<div> <div>+syllabic</div> <div>—low</div> <div>—back</div> <div>—tense</div> <div>—round</div> </div>	S	—	#		

# Minimal generalization

Minimal generalization of features: retain all shared feature values

- Prevents generalization to unseen feature values
- Permits generalization to unseen feature *combinations*
- Example: comparing *b* and *n*

$$\begin{array}{c} \text{b} \\ \left[ \begin{array}{c} -\text{syllabic} \\ +\text{cons} \\ -\text{contin} \\ +\text{voice} \\ -\text{nasal} \\ +\text{labial} \\ -\text{coronal} \\ -\text{velar} \end{array} \right] \end{array} + \begin{array}{c} \text{n} \\ \left[ \begin{array}{c} -\text{syllabic} \\ +\text{cons} \\ -\text{contin} \\ +\text{voice} \\ +\text{nasal} \\ -\text{labial} \\ +\text{coronal} \\ -\text{velar} \end{array} \right] \end{array} \mapsto \begin{array}{c} \{ \text{b, m, d, n} \} \\ \left[ \begin{array}{c} -\text{syllabic} \\ +\text{cons} \\ -\text{contin} \\ +\text{voice} \\ \\ \\ -\text{velar} \end{array} \right] \end{array}$$

# Examples

Comparing *hug*, *plan*, and then *rub*:  $\emptyset \rightarrow d / \dots$

	Residue	Shared Features	Shared Segments	Change Location	Shared Segments	Shared Features	Residue
A.	#hʌ	g		—	#		
B.	#plæ	n		—	#		
C.	X	<div>             [             <ul style="list-style-type: none"> <li>—syllabic</li> <li>—continuant</li> <li>—labial</li> <li>—lateral</li> <li>+voice</li> <li>...</li> </ul> </div>		—	#		

And generalizing further with *rub*:

D.	#rʌ	b		—	#		
E.	X	<div>             [             <ul style="list-style-type: none"> <li>—syllabic</li> <li>—continuant</li> <li>—lateral</li> <li>+voice</li> <li>...</li> </ul> </div>		—	#		

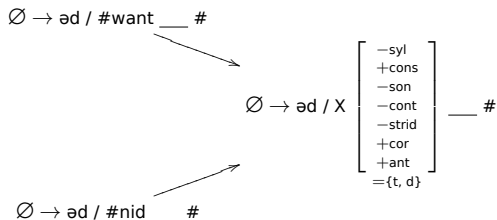
# Examples

Comparing *sing* and *drink*:  $\text{ɪ} \rightarrow \text{æ} / \dots$

	Residue	Shared Features	Shared Segments	Change Location	Shared Segments	Shared Features	Residue
A.	#	s		—	ŋ		#
B.	#d	r		—	ŋ		k#
C.	X	<div>                     −syllabic                      +coronal                      +continuant                 </div>		—	ŋ		Y

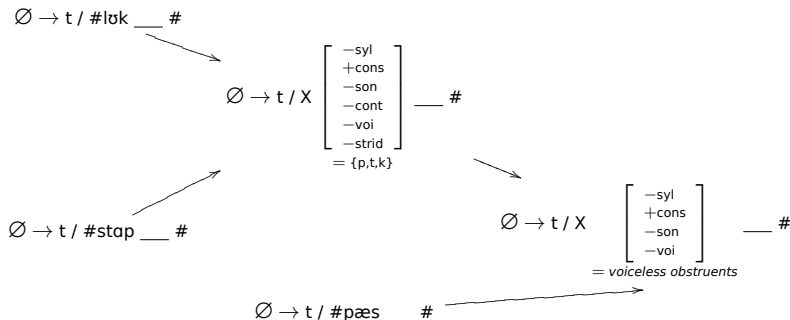
# What this model does well

Rapidly discovers most general environment for each pattern



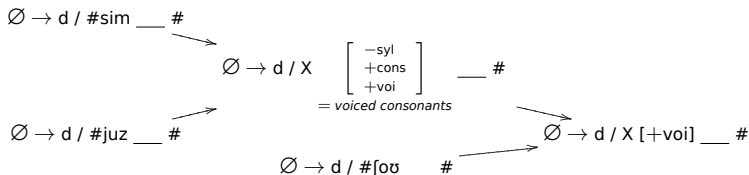
# What this model does well

Rapidly discovers most general environment for each pattern



# What this model does well

Rapidly discovers most general environment for each pattern

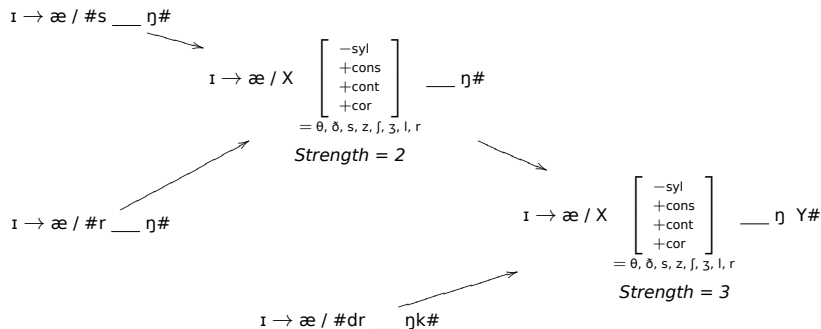


# Rapid generalization

- Generalization proceeds rapidly, given sufficiently diverse stems
- The pathways shown here can be found using verbs that are among the 75 most frequent verbs of English (according to CELEX) (of which the majority are actually irregular)

# Rule creation for irregulars

“Irregular” mappings are also compared and generalized



- Provides rules to generalize  $[\text{spl}\eta] \rightarrow [\text{spl}\text{æ}\eta]$  in addition to  $[\text{spl}\eta\text{d}]$

# Outline

## 1 Model overview

- Minimal generalization
- **Reliability**
- Confidence

## 2 A simple simulation

## 3 Extensions

- Phonological rules
- 'Impugnment'
- Non-local and suprasegmental contexts

# What determines the strength of a rule?

- Goal: suffixed output [splɪŋd] is more probable/acceptable than outputs like [splæŋ], [splʌŋ], etc.
- Irregular patterns (at least in English) tend to cover relatively few forms, which are similar to one another  $\approx$  share a set of phonological properties
- Result: rules that characterize them are more specific, and weaker strength
- These patterns are not productive, except possibly for inputs that fit a very particular phonological shape ( $\approx$  prototypical forms)

# What determines the strength of a rule?

Pinker and Prince (1988):

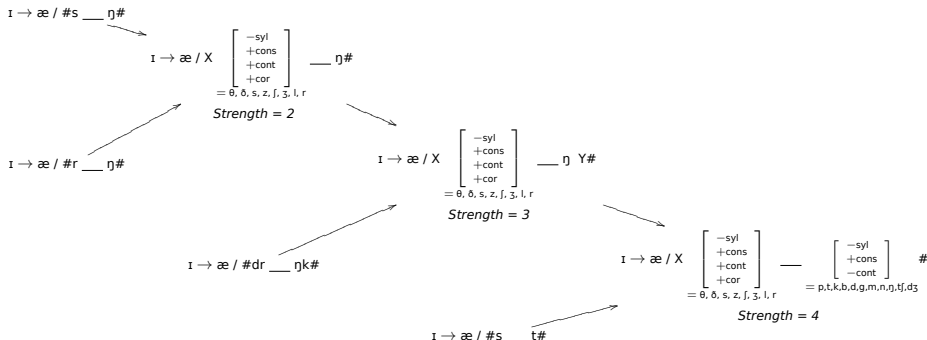
- “Rule candidates increase in strength each time they have been exemplified by an input pair”

Rule	Examples covered
1. $\emptyset \rightarrow t / \#mis \_ \#$	1
2. $\emptyset \rightarrow t / \#pres \_ \#$	1
3. $\emptyset \rightarrow t / X [+syl, -low, -bk, -tns, -rnd] s \_ \#$	2
4. $\emptyset \rightarrow d / \#hlg \_ \#$	1
5. $\emptyset \rightarrow d / \#plæn \_ \#$	1
6. $\emptyset \rightarrow d / \_ X [-syl, -cont, -lab, -lat, -del.rel.]$	2
7. $i \rightarrow æ / \#s \_ \eta \#$	1
8. $i \rightarrow æ / \#dr \_ \eta k \#$	1
9. $i \rightarrow æ / X [-syl, +cor, +cont] \_ \eta Y \#$	2

- p. 134: “[by various means], some regularities can be enshrined as permanent productive rules whereas others can be discarded or treated differently.”

# Simply counting examples is not enough

- Consider hypothesized rules after exposure to *sing*, *ring*, *drink*, *sit*



- “Prototypical” forms (*sing*, *ring*, *drink*) and “outlier” (*sit*) combine to yield more general rule
- More general rule encompasses more forms, so has higher strength

# Simply counting examples is not enough

- Prediction:  $\text{ɪ} \rightarrow \text{æ}$  change may apply equally well to any form within this broader context
  - Hypothetical *zick*  $\sim$  *zack* just as well supported as *spling*  $\sim$  *splang*

# Diagnosis

- *Strength*, measured purely in terms of number of examples covered, favors the broadest possible generalizations
- Under this view, “compatible with the data” = “encompasses all the examples”
- Broad generalizations are not always virtuous; sometimes, they incur numerous exceptions.
- A more intuitive notion of “compatible with the data” must incorporate not only generality, but also accuracy

# Rule reliability

**Reliability** of a hypothesized rule:

$$= \frac{\text{number of items that the rule works for ("hits")}}{\text{number of items that meet the rule's structural description (CAD) ("scope")}$$

- Also known as **accuracy**, or **coverage**

Rule	Hits	Exceptions
$I \rightarrow \text{æ} / X$ <div> <math>\begin{bmatrix} -\text{syl} \\ +\text{cons} \\ +\text{cont} \\ +\text{cor} \end{bmatrix}</math>  <math>= \theta, \delta, s, z, f, ʒ, l, r</math> </div> $\text{--- } \eta \text{ } Y\#$	$\text{ring, sing, drink}$  $= n$	$\text{wring, fling, ...}$  $= m$
$I \rightarrow \text{æ} / X$ <div> <math>\begin{bmatrix} -\text{syl} \\ +\text{cons} \\ +\text{cont} \\ +\text{cor} \end{bmatrix}</math>  <math>= \theta, \delta, s, z, f, ʒ, l, r</math> </div> $\text{--- } \begin{bmatrix} -\text{syl} \\ +\text{cons} \\ -\text{cont} \end{bmatrix} Y\#$ <div> <math>= p, t, k, b, d, g, m, n, \eta, tʃ, dʒ</math> </div>	$\text{ring, sing, drink, sat}$  $= n+1$	$\text{wring, fling, lick, rig, ship, rid ...}$  $= m+\text{lots}$

- Island of reliability:** context in which a change is especially likely
  - "especially likely" = more likely than in its most general context (Albright 2002)

# The intended effect

- Favoring rules with high reliability should let the model find specific contexts that *uniquely* characterize members of a particular class
- Overly broad contexts are discouraged, because of the cost of adding exceptions
  - Number of positive hits gained (numerator) must exceed number of exceptions added to scope (denominator)
  - For small, “irregular” class, optimal contexts will be narrow
  - For “regular” classes, there may be fewer distinguishing features that uniquely characterize most regular verbs, so no benefit to staying small
  - Model zooms to full generality, incurring some exceptions
- Model should settle on the right level of generality for each class

## How this plays out in practice

- Example: ran the model on the 200 most frequent verbs in CELEX
  - 134 regular (80 *-d*, 30 *-t*, 24 *-əd*); 66 irregular, of various types

## How this plays out in practice

- Example: ran the model on the 200 most frequent verbs in CELEX
  - 134 regular (80 -d, 30 -t, 24 -əd); 66 irregular, of various types
- Most reliable rules for a few of the major classes

Rule	Hits	Scope	Rel.	Examples	Exceptions	Excludes
$i \rightarrow \epsilon / \# \begin{bmatrix} -\text{syl} \\ +\text{cont} \end{bmatrix} \_\_\_\_ d\#$ $= \text{continuants}$	3	3	1.00	read, lead, feed	(none)	meet
Rule	Hits	Scope	Rel.	Examples	Exceptions	Excludes
$e \rightarrow o / \# \begin{bmatrix} -\text{syl} \\ +\text{voi} \\ +\text{lab} \end{bmatrix} \_\_\_\_ r$ $= b, m, v, r$	2	2	1.00	wear, bear	(none)	break
Rule	Hits	Scope	Rel.	Examples	Exceptions	Excludes
$\emptyset \rightarrow \emptyset / \# \begin{bmatrix} -\text{syl} \\ +\text{cons} \\ +\text{cont} \\ +\text{cor} \\ +\text{ant} \end{bmatrix} \epsilon t \_\_\_\_ \#$ $= \theta, \delta, l, s, z$	2	2	1.00	set, let	(none)	put, cut, get, sit, ...
Rule	Hits	Scope	Rel.	Examples	Exceptions	Excludes
$d \rightarrow t / X \begin{bmatrix} -\text{syl} \\ -\text{son} \\ -\text{voi} \end{bmatrix} \epsilon n \_\_\_\_ \#$ $= p, t, k, f, \theta, s, f, t_f$	2	3	0.67	send, spend	tend	build

# How this plays out in practice

- Example: ran the model on the 200 most frequent verbs in CELEX
  - 134 regular (80 *-d*, 30 *-t*, 24 *-əd*); 66 irregular, of various types
- Most reliable rules for a few of the major classes

Rule	Hits	Scope	Rel.	Examples	Exceptions	Excludes
$\emptyset \rightarrow t / X$ $\begin{bmatrix} -\text{syl} \\ -\text{son} \\ +\text{cont} \\ -\text{voi} \\ = f, \theta, s, j \end{bmatrix}$ $\_\_\_\#$	14	14	1.00	pass, produce, wish, laugh, ...	(none)	52 other [–voi]
$\emptyset \rightarrow t / X$ $\begin{bmatrix} +\text{syl} \\ -\text{low} \\ +\text{back} \end{bmatrix}$ $k \_\_\_\#$ $= u:, \sigma, \sigma\sigma, \sigma, \Lambda, \partial^*$	4	4	1.00	look, work, talk, walk	(none)	62 other [–voi]
Rule	Hits	Scope	Rel.	Examples	Exceptions	Excludes
$\emptyset \rightarrow \text{əd} / X$ $\begin{bmatrix} -\text{syl} \\ +\text{cons} \end{bmatrix}$ $t \_\_\_\#$ $= \text{consonants}$	10	10	1.00	want, start, inspect, suggest, ...	(none)	35 other <i>t, d</i>
$\emptyset \rightarrow \text{əd} / X$ $\begin{bmatrix} +\text{syl} \\ -\text{high} \\ -\text{back} \end{bmatrix}$ $d \_\_\_\#$ $= e, \varepsilon, \text{æ}, \sigma, i, \text{ai}$	4	4	1.00	provide, decide, add, avoid	(none)	41 other <i>t, d</i>

# How this plays out in practice

- Example: ran the model on the 200 most frequent verbs in CELEX
  - 134 regular (80 -d, 30 -t, 24 -əd); 66 irregular, of various types
- Most reliable rules for a few of the major classes

Rule	Hits	Scope	Rel.	Examples	Exceptions	Excludes
$\emptyset \rightarrow d / X \text{ ə} \_ \#$	11	11	1.00	remember, consider, offer, ...	(none)	69 other [+voi]
$\emptyset \rightarrow d / X \begin{bmatrix} +\text{syl} \\ -\text{low} \\ +\text{back} \end{bmatrix} \text{v} \_ \#$ $= u:, \sigma, o\sigma, \sigma, \Lambda, \text{ə}$	5	5	1.00	move, love, serve, prove, remove	(none)	75 other [+voi]
$\emptyset \rightarrow d / X \begin{bmatrix} +\text{syl} \\ -\text{high} \\ -\text{back} \end{bmatrix} \text{n} \_ \#$ $= e, \varepsilon, \text{æ}, \sigma, \text{ɪ}, \text{aɪ}$	5	5	1.00	remain, explain, plan, join, contain	(none)	75 other [+voi]

# The fate of the more general rules

The more general rules are not nearly so reliable

Rule	Hits	Scope	Rel.	Examples	Exceptions
$\emptyset \rightarrow d / X [+voi] \_\_ \#$	80	134	.60	seem, use, try, call, turn, ...	do, say, go, know, see, <b>need</b> , ...
$\emptyset \rightarrow t / X [-voi] \_\_ \#$	30	66	.46	look, ask, work, talk, help, ...	get, think, take, <b>want</b> , put, ...
$\emptyset \rightarrow \text{əd} / X$ <div style="display: inline-block; vertical-align: middle; margin-left: 10px;"> <math>\left[ \begin{array}{l} -\text{syl} \\ +\text{cons} \\ -\text{son} \\ -\text{cont} \\ -\text{strid} \\ +\text{cor} \\ +\text{ant} \\ =\{t, d\} \end{array} \right]</math> </div> $\_\_ \#$	24	45	.53	want, need, start, wait, expect, ...	get, find, put, sit, stand, ...

# An unintended consequence

- Focusing on reliability has something of the desired effect
  - Finds consists clusters of similar words
  - Identifies cases where it more distant words should be excluded
- However, it takes small-scale generalizations too seriously
- By favoring reliability instead of generality, we end up finding a long list of very accurate, but seemingly quite accidental contexts

# Why are the more general contexts so unreliable?

- Genuine exceptions: almost half of *t*, *d*-final verbs in this sample are irregular
  - Expect this proportion to go down in a larger sample
  - Larger sample will never make the problem disappear completely, however, since some subcontexts are 100% regular
    - E.g., verbs ending in voiceless fricatives, verbs ending in [əʊ], etc.
    - These will always be more reliable than the general context
  - In any system with irregularity, broad generalizations will necessarily involve some exceptions (decreased reliability)
- Inadequate characterization of rule interaction:
  - *t*, *d*-final verbs act as exceptions for more general *-d*, *-t* suffixes
  - More on this below...

# Outline

## 1 Model overview

- Minimal generalization
- Reliability
- Confidence

## 2 A simple simulation

## 3 Extensions

- Phonological rules
- 'Impugnment'
- Non-local and suprasegmental contexts

# Confidence

- Reliability is defined as a proportion
  - $2/2 = 7/7 = 100/100$
- Speaker intuitions: not all 'perfectly reliable' rules are equally productive
- 100/100 less likely to be coincidental than 2/2

# A comparison

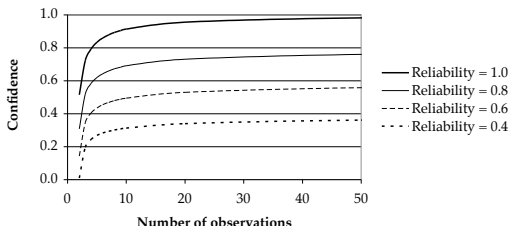
Reliability values from a larger set of English verbs:

- $\emptyset \rightarrow \text{əd} / \left[ \begin{array}{l} -\text{high} \\ -\text{low} \\ +\text{tense} \end{array} \right] \text{ \_\_\_\_\_\_ } (regular \text{ after } -ate/-ote) \quad 366/366$
- $\emptyset \rightarrow \text{t} / \left[ \begin{array}{l} -\text{sonorant} \\ +\text{continuant} \end{array} \right] \text{ \_\_\_\_\_\_ } (regular \text{ after voiceless frics.}) \quad 352/352$
- $\text{ɪ} \rightarrow \text{ʌ} / \left[ \begin{array}{l} +\text{coronal} \\ +\text{voiced} \\ +\text{anterior} \end{array} \right] \text{ \_\_\_\_\_\_ } \left[ \begin{array}{l} +\text{voiced} \\ +\text{dorsal} \end{array} \right] \quad 4/4$   
(*dig, cling, fling, sling*)
- $\emptyset \rightarrow \emptyset / \# \left[ \begin{array}{l} +\text{coronal} \\ +\text{continuant} \end{array} \right] \text{ɛt \_\_\_\_\_\_} \quad 2/2$   
(*let, set*)
- $\text{i:} \rightarrow \text{ɛ} / \# \left[ \begin{array}{l} +\text{sonorant} \\ -\text{nasal} \\ +\text{coronal} \end{array} \right] \text{ \_\_\_\_\_\_ } \text{d} \quad 2/2$   
(*read, lead*)

# Confidence

## Implementation: lower confidence intervals

- Although  $4/4 = 1$ , still likely that the next example will be an exception
  - Confidence interval: actual proportion might be somewhat lower or higher than observed
  - In this case, we can be 95% sure that the true value is between .875 and 1.125 (if values above 1 made sense here)
- Lower confidence limits:



# The effects of confidence

- Patterns involving small numbers of words may be weak, even if consistent
- Morphological categories that are rare may yield overall less confident rules

# Outline

- 1 Model overview
  - Minimal generalization
  - Reliability
  - Confidence
- 2 A simple simulation
- 3 Extensions
  - Phonological rules
  - 'Impugnment'
  - Non-local and suprasegmental contexts

# The training data

- Text file of training pairs and test items
  - Training pairs (and optionally, token frequency)
  - Test forms for 'wug testing'
- Phonological features
  - .fea file with feature specifications
  - "ASCII" column (numeric codes): values don't matter, should just be unique
  - Each segment must be a single character (ASCII, or Unicode)
  - Feature values: 0-1 (binary), 0-n (scalar)
  - Unspecified = -1

# Running the model

- Java archive: MinGenLearner.jar
  - From command line: `java -jar MinGenLearner.jar`
- Output files (tab-delimited text format)
  - `.out` file with list of changes
  - `.sum` file with wug test results
  - `.rules` file with list of rules discovered (if 'Save rules' option selected)

# Outline

- 1 Model overview
  - Minimal generalization
  - Reliability
  - Confidence
- 2 A simple simulation
- 3 Extensions
  - Phonological rules
  - 'Impugnment'
  - Non-local and suprasegmental contexts

# Outline

- 1 Model overview
  - Minimal generalization
  - Reliability
  - Confidence
- 2 A simple simulation
- 3 Extensions
  - Phonological rules
  - 'Impugnment'
  - Non-local and suprasegmental contexts

# Insufficient generalization without phonology

- Most general 'regular' rules in the sample simulation

Rule	Rel.	Conf.
$\emptyset \rightarrow d / X [+voi] \_\_\_ \#$	80/134	.568
$\emptyset \rightarrow t / X \begin{bmatrix} -syl \\ 0-2son \\ 0-1aper \end{bmatrix} \_\_\_ \#$ $= p,k,t,f,\theta,s,f,tj$ $b,d,g,v,\delta,z,\zeta,d\zeta,m,n,\eta$	31/142	.197
$\emptyset \rightarrow \text{əd} / X \begin{bmatrix} -syl \\ 0 son \\ 0 aper \\ +COR \\ +ant \end{bmatrix} \_\_\_ \#$ $= t,d$	24/45	.482

- Misses fact that t, əd could be derived by phonological rules
  - $d \rightarrow t / [-voi] \_\_\_$
  - $\emptyset \rightarrow \text{ə} / \{t,d\} \_\_\_ d$

# Unifying changes with phonology

Goal:

- Allow model to discover that adding *d* after other contexts would yield illegal sequences such as \*pd, \*dd
- Posit phonological rules that repair illegal sequences to generate the attested output

# Approach to discovering phonological rules

- Provide model with list of phonotactically (surface) illegal sequences
  - \*pd, \*kd, \*sd, \*td, \*dd, \*bt, \*gt, \*zt, \*tt, \*dt, ...
  - Listed in .ill file
- 'Doppelgänger' rules
  - When positing a rule  $A \rightarrow B / C \_\_\_ D$ , try out other changes  $A' \rightarrow B'$  in the same context  $C \_\_\_ D$
  - Original:  $\emptyset \rightarrow t / [-\text{voi}] \_\_\_ \#$
  - Doppelgänger  $\emptyset \rightarrow d / [-\text{voi}] \_\_\_ \#$
- Scan outputs of Doppelgänger rules for illegal sequences
- If difference between illegal output and attested output is a possible phonological mapping, posit a rule
  - Change/insertion/deletion of a single segment—e.g.,  $d \rightarrow t$
- Model options: 'Use Phonology', 'Use Doppelgängers'

# Some limitations

- Can't discover rule ordering
- Can't detect conflicting rules to see which is the productive one
  - $kd \rightarrow kt$ :  $l\text{ø}kd \rightarrow l\text{ø}kt$  'look'
  - $kd \rightarrow d$ :  $me\text{ɪ}kd \rightarrow me\text{ɪ}d$  'make'
- Assumes that rules hold in entire string (no NDEB or morphological conditioning)
  - $pd \rightarrow pt$ :  $\text{ʌ}pd\text{eɪ}t\text{əd} \rightarrow *ʌpt\text{eɪ}t\text{əd}$  'update'
  - Workaround: add word boundaries and augment list of illegal sequences to restrict to final position ( $*pd\#$ ,  $*dt\#$ , etc.)
- Also possible to pre-specify mappings by brute force in `.phon` file
  - End-run around limited discovery mechanism

# Results for sample file

Most general 'regular' rules, with phonology enabled

Rule	Rel.	Conf.
$\emptyset \rightarrow d / \text{X} \_\_\_ \#$	134/200	.647
$\emptyset \rightarrow t / \text{X} \left[ \begin{array}{l} -\text{syl} \\ 0-2\text{son} \\ 0-1\text{aper} \end{array} \right] \_\_\_ \#$ $= p, k, t, f, \theta, s, \int, \text{t}\int$ $b, d, g, v, \delta, z, \text{ʒ}, \text{dʒ}, m, n, \eta$	79/142	.528
$\emptyset \rightarrow \text{əd} / \text{X} \left[ \begin{array}{l} -\text{syl} \\ 0 \text{ son} \\ 0 \text{ aper} \\ +\text{COR} \\ +\text{ant} \end{array} \right] \_\_\_ \#$ $= t, d$	24/45	.482

- Model discovers more general  $\emptyset \rightarrow d$  rule
- Reliability of most general  $\emptyset \rightarrow t$  rule improves, too

# Outline

- 1 Model overview
  - Minimal generalization
  - Reliability
  - Confidence
- 2 A simple simulation
- 3 Extensions
  - Phonological rules
  - 'Impugnment'
  - Non-local and suprasegmental contexts

# Another problem

- Model outputs for [splɪŋ]

Output	Rule	Rel.	Conf.
splɪŋd	$\emptyset \rightarrow d / X \left[ \begin{array}{l} -\text{syl} \\ 1-2\text{son} \\ 0-1\text{aper} \end{array} \right] \_\_\#$ $= f, m, n, s, v, z, \delta, \eta, f, 3, \theta$	50/65	.729
splɪŋt	$\emptyset \rightarrow t / X \left[ \begin{array}{l} -\text{syl} \\ 0-2\text{son} \\ 0-1\text{aper} \end{array} \right] \_\_\#$ $= p, k, t, f, \theta, s, f, \eta$ $b, d, g, v, \delta, z, 3, d3, m, n, \eta$	79/142	.528
splæŋ	$i \rightarrow \text{æ} / X' \_\_\{m, n, \eta\}$	2/4	.310

- $\emptyset \rightarrow t$  rule covers voiceless segments + n (læn  $\rightarrow$  lænt 'learn')

## Diagnosis

- $\emptyset \rightarrow t / X \begin{bmatrix} \text{---syl} \\ 0\text{-2son} \\ 0\text{-1aper} \end{bmatrix}$  — # covers two distinct sets of words  
 $= p, k, t, f, \theta, s, f, \text{ʈ}$   
 $b, d, g, v, \delta, z, \text{ʒ}, \text{ɖ}, m, n, \eta$ 
  - Regular affixation after voiceless obstruents
  - Irregular forms *learnt*, *burnt*, etc.
- Comparison of [-voi] and *n* should reveal that reliability is very different in these two contexts

- $\emptyset \rightarrow t / \begin{bmatrix} -\text{syl} \\ 0-2\text{son} \\ 0-1\text{aper} \\ +\text{voi} \end{bmatrix} \text{ — } \# : 78/117 = .667$   
 $= p, k, t, f, \theta, s, \text{f}, \text{t}$

- $\emptyset \rightarrow t / \begin{bmatrix} -\text{syl} \\ 0\text{-2son} \\ 0\text{-1aper} \\ -\text{voi} \end{bmatrix}$  — #:  $1/25 = .04$  (*learnt*)  
 $= b, d, g, v, \delta, z, 3, d3, m, n, \eta$

# Impugnment

- When generalizing from rule  $R$  to superset  $R'$ 
  - Compare confidence of subset context ( $R$ ) and context outside subset ( $R' - R$ )
  - Be generous: **lower** confidence limit of subset, **upper** confidence limit of superset
- **Impugnment**: if superset confidence is lower, replace score of rule with (upper) confidence of superset region
- Result:  $\emptyset \rightarrow t$  rule is adjusted from .528 to .090

# Results of impugnment

Outputs for [splɪŋ] with impugnment:

Output	Rule	Rel.	Conf.
splɪŋd	$\emptyset \rightarrow d / X \left[ \begin{array}{l} -\text{syl} \\ 1\text{-2son} \\ 0\text{-1aper} \end{array} \right] \_\_\_ \#$ $= f, m, n, s, v, z, \delta, \eta, f, 3, \theta$	50/65	.700
splæŋ	$\text{ɪ} \rightarrow \text{æ} / X \text{ ' } \_\_\_ \{m, n, \eta\}$	2/4	.310
splɪŋt	$\emptyset \rightarrow t / X \left[ \begin{array}{l} -\text{syl} \\ 0\text{-2son} \\ 0\text{-1aper} \end{array} \right] \_\_\_ \#$ $= p, k, t, f, \theta, s, f, \eta$ $b, d, g, v, \delta, z, 3, d3, m, n, \eta$	79/142	.090

- Effect is stronger with larger training set

# Outline

- 1 Model overview
  - Minimal generalization
  - Reliability
  - Confidence
- 2 A simple simulation
- 3 Extensions
  - Phonological rules
  - 'Impugnment'
  - Non-local and suprasegmental contexts

# Some limitations of the model

- Input representations are 'flat'
  - No prosodic information (monosyllabic vs. polysyllabic)
  - No morphological information (inflection class, gender, etc.)
- Contexts are strictly local
  - No harmony, other non-local conditioning

# Simulating tiers in a flat representation

## Goal:

- Give the model access to prosodic size of the base
- E.g., monosyllabic vs. longer

## Approach

- Create dummy segments, specified for just for 'length' feature
  - E.g., M (monosyllabic) vs. P (polysyllabic)
  - Feature: [ $\pm$ polysyllabic]

Polysyllabic	
M	0
P	1

- Could likewise specify gender/inflection class, nearest vowel, etc.

# Augmenting inputs with dummy segments

Simple case: all changes are suffixal (or prefixal)

- Use 'empty space' on other side of affix to specify suprasegmental context

lekleb <sup>P</sup>	leklebi <sup>P</sup>	fegriv <sup>P</sup>	fegrivi <sup>P</sup>
globlef <sup>P</sup>	globlefi <sup>P</sup>	pom <sup>M</sup>	pomi <sup>M</sup>
nuv <sup>M</sup>	nuvi <sup>M</sup>	sivod <sup>P</sup>	sivodi <sup>P</sup>
gor <sup>M</sup>	gori <sup>M</sup>	sog <sup>M</sup>	sogi <sup>M</sup>
fabeg <sup>P</sup>	fabegi <sup>P</sup>	zibot <sup>P</sup>	ziboti <sup>P</sup>

- Resulting rule format
  - $\emptyset \rightarrow i$  / segmental context \_\_\_\_ prosodic context
- In practice, adding symbols to mark boundary can aid legibility
  - lekleb●<sup>P</sup> ~ leklebi○<sup>P</sup>, sog●<sup>M</sup> ~ sogi○<sup>M</sup>, ...
  - Change: ●→i○ (different input/output symbols ensure that they are not treated as part of the context)

# Discontinuous morphological contexts

- A particularly strong context for  $\tau \rightarrow \lambda$  (Bybee and Slobin 1982; Bybee and Moder 1983)
  - $sC(C)\_\_ \eta(k)$
- Model presented here cannot discover this, due to its locality restriction
  - Alignment outward from change can't skip segments, or find features of multiple segments
  - Also cannot posit rules with literal segment outside featurally underspecified segments (sCC)
- Bybee and Moder (1983): speakers do generalize to novel  $sC(C)$ ...items more than  $C(C)$

# A better alignment procedure

Albright and Hayes (2006): minimum string edit distance

- Aligns strings by finding optimal segmental correspondence between segments
- Align matching segments with one another
- Mismatches: find phonetically closest correspondent

	S		W	I	m	
	f		ɹ	ɪ	ŋ	k
⇒	[ -voi +strid ]		[ -syl 4+son ]	ɪ	[ +nas ]	k?
	s	p	ɹ	ɪ	ŋ	
⇒	[ -voi +strid ]	p	[ -syl 4+son ]	ɪ	[ +nas ]	k?

## A challenge...

- Discontinuous alignments retain much more information than the 'myopic' alignments of the local MGL
- Explosion in number of rules
- Albright and Hayes (2006) devise a procedure for pruning the grammar to retain the most reliable rules, using the Gradual Learning Algorithm (Boersma 1997)
- Possible to employ more sophisticated selection/weighting techniques (e.g., maximum entropy models)
- Interesting theoretical issue: what kinds of non-local contexts do human learners actually notice?
  - English v- initial verbs are all regular, but it appears that no effect of this is seen in wug verbs
- If you are interested in employing a non-local version of the model, let me know

# Links

- Model download: <http://www.mit.edu/~albright/mgl/>
- These slides: <http://www.mit.edu/~albright/mgl/MGL-Tutorial.pdf>
- Additional papers, data files, etc.:  
<http://www.linguistics.ucla.edu/people/hayes/learning/>

# References I

Albright, A. and B. Hayes (2002).

Modeling English past tense intuitions with minimal generalization.  
In *SIGPHON 6: Proceedings of the Sixth Meeting of the ACL Special Interest Group in Computational Phonology*, pp. 58–69. ACL.

Albright, A. and B. Hayes (2003).

Rules vs. analogy in English past tenses: A  
computational/experimental study.  
*Cognition* 90, 119–161.

Albright, A. and B. Hayes (2006).

Modeling productivity with the Gradual Learning Algorithm: The  
problem of accidentally exceptionless generalizations.  
In G. Fanselow, C. Féry, R. Vogel, and M. Schlesewsky (Eds.),  
*Gradience in Grammar: Generative Perspectives*, pp. 185–204.  
Oxford University Press.



## References II

Boersma, P. (1997).

How we learn variation, optionality, and probability.

*Proceedings of the Institute of Phonetic Sciences of the University of Amsterdam 21*, 43-58.

<http://fon.hum.uva.nl/paul/>.

Bybee, J. and C. Moder (1983).

Morphological classes as natural categories.

*Language 59*(2), 251-270.

Bybee, J. and D. Slobin (1982).

Rules and Schemas in the Development and Use of the English Past Tense.

*Language 58*(2), 265-289.

## References III

Ling, C. and M. Marinov (1993).

Answering the connectionist challenge: a symbolic model of learning the past tenses of English verbs.

*Cognition* 49, 235-290.

Pinker, S. and A. Prince (1988).

On language and connectionism: Analysis of a Parallel Distributed Processing model of language acquisition.

*Cognition* 28, 73-193.