

CS 585: INTRODUCTION TO NATURAL LANGUAGE PROCESSING

Guest Lecture: Gaja Jarosz

Gradient Phonotactics

1

□ How good are these as novel words of English?

□ stin

□ blick

□ mip

□ skell

□ blafe

□ bwip

□ shmy

□ smum

□ dlap

□ bzack

□ mrock

□ dmell

□ lnoot

□ mdap

Gradient Phonotactics

2

□ How good are these as novel words of English?

□ **stin**

□ **blick**

□ **mip**

□ **skell**

□ **blafe**

□ **bwip**

□ **shmy**

□ **smum**

□ **dlap**

□ **bzack**

□ **mrock**

□ **dmell**

□ **lnoot**

□ **mdap**

Gradient Phonotactics

3

- Knowledge of gradient phonotactics has been illustrated in a range of tasks
 - Lexical access
 - Faster recognition of more phonotactically probable words
 - Word learning
 - Faster acquisition of more phonotactically probable words
 - Word segmentation by children and adults
 - This is the main information we think infants use to segment speech
 - Example: Mattys & Jusczyk (2001)
 - ..beangaffehold.. [ng] and [fh] occur infrequently within words
 - ..fanggaffetine.. [ŋg] and [ft] occur frequently within words

What's Phonotactic Probability?

4

- What's the probability of?
 - blick vs. shmy vs. mrock
- Ideas?

N-grams!

5

- $\Pr(\#abcd\#)$
- $\Pr(a \mid b) = ?$
- Chain Rule
 - $\Pr(a,b) = \Pr(a \mid b)\Pr(b)$
 - $\Pr(a,b,c) = \Pr(a \mid b,c)\Pr(b \mid c)\Pr(c)$
 $= \Pr(c \mid a,b)\Pr(b \mid a)\Pr(a)$
- Apply chain rule to $\Pr(\#abcd\#)$?

N-grams!

6

- $\Pr(\#abcd\#)?$
 - $\Pr(a \mid \#) *$
 - $\Pr(b \mid \#,a) *$
 - $\Pr(c \mid \#,a,b) *$
 - $\Pr(d \mid \#,a,b,c) *$
 - $\Pr(\# \mid \#,a,b,c,d)$
- N-grams make an independence assumption that only a fixed amount of history matters...
 - Unigrams, bigrams, trigrams, etc.

N-grams

7

- We can't condition on every possible preceding sequence - there's too many!
 - Jack and Jill went up the hill
- N-gram models condition on N-1 previous symbols:
 - Unigrams: no history
 - $P(\text{Jack})P(\text{and})P(\text{Jill})P(\text{went})P(\text{up})P(\text{the})P(\text{hill})$
 - Bigrams: 1 previous symbol
 - $P(\text{Jack} \mid \#)P(\text{and} \mid \text{Jack})P(\text{Jill} \mid \text{and})P(\text{went} \mid \text{Jill})\dots$
 - Trigrams: 2 previous symbols
 - $P(\text{Jack} \mid \#\#)P(\text{and} \mid \# \text{ Jack})P(\text{Jill} \mid \text{Jack and})P(\text{went} \mid \text{and Jill})\dots$

What's Phonotactic Probability?

8

- What's the bigram probability of?
 - blick vs. shmy vs. mrock

What's Phonotactic Probability?

9

- What's the bigram probability of?
 - blick vs. shmy vs. mrock
- $P(b \mid \#)P(l \mid b)P(i \mid l)\dots$
 - How can we estimate these probabilities?

Baseline: Phoneme bi-grams

10

□ Example

biphone	stin	bleif
1	#s 0·118	#b 0·057
2	st 0·205	bl 0·106
3	ti 0·192	lei 0·042
4	in 0·108	eif 0·007
5	n# 0·151	f# 0·067
log (prob)	-4·12	-6·93

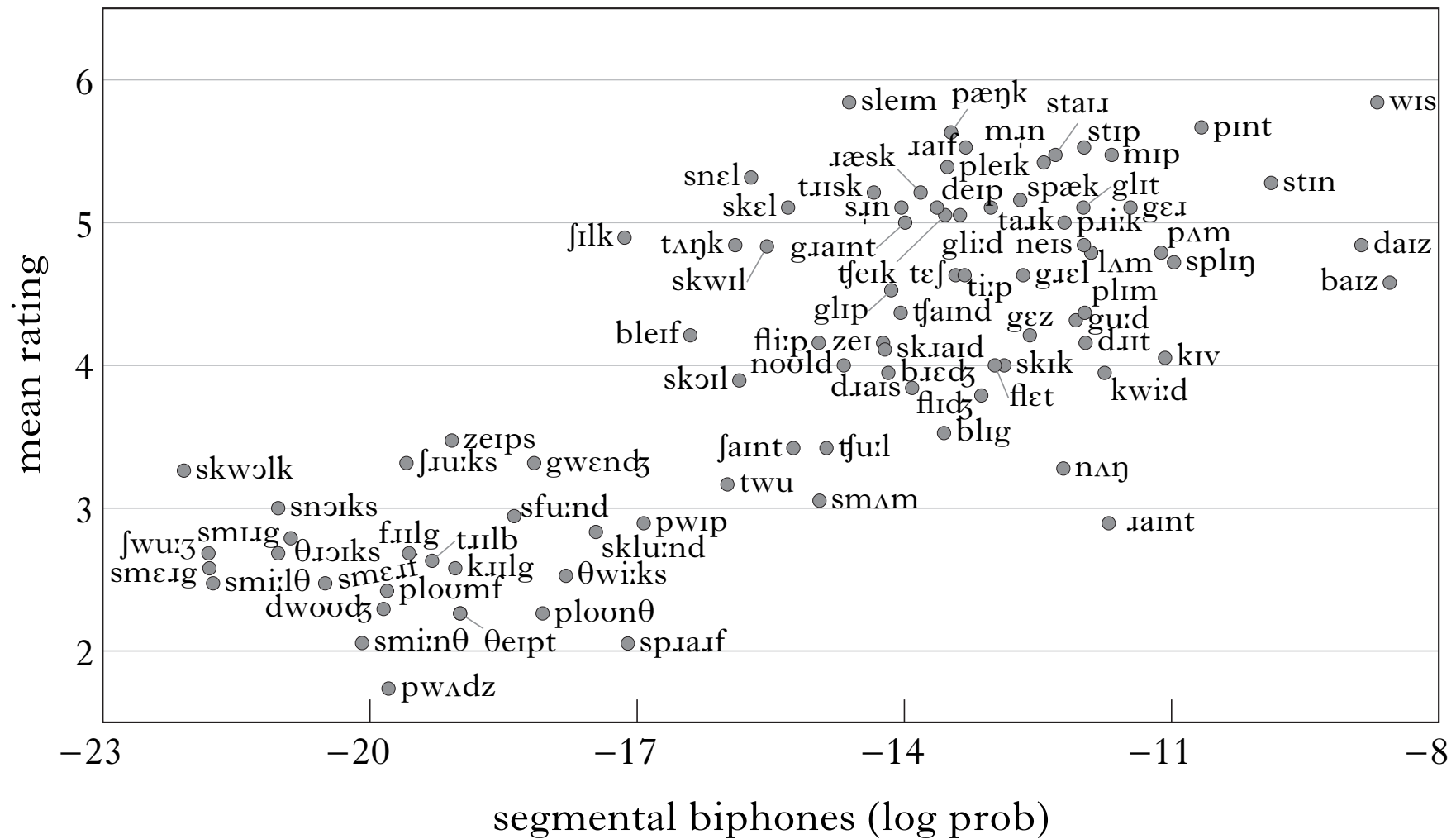
What's Phonotactic Probability?

11

- What's the bigram probability of?
 - ▣ blick vs. shmy vs. mrock
- If we're primarily interested in the initial consonant sequences, there are several ways we could use n-grams to model these sequences...
 - ▣ Pros and cons?

How do Bigrams do? (Albright 2009)

12



How do bigrams do?

13

- Overall correlations with human judgments?
 - Training on just onsets:
 - $r = 0.88$
 - Hayes & Wilson (2008)
 - Training on whole words:
 - $r = 0.78$ (with syllable boundaries), $= .50$ (without)
 - Daland et al. (2011)

Problem!

14

- People show gradient judgments of two kinds
- Attested onset clusters
 - stin > blin > bwin
- And unattested onset clusters
 - bnick > bzick > mbick
- High correlations due to differences across not within
- This latter category gives us crucial clues about how people *learn and represent* these patterns
- However,
 - Positional bigrams cannot model this. Why?
 - Whole-word bigrams do very poorly ($r = 0.22$). Why?

So what are people doing?

15

- Why do people have these judgments then?
 - bnick > bzick > mbick
- Notice: this is a problem with generalization
 - We want models to generalize correctly to new data
 - Cognitive scientists & Linguists are interested in modeling how humans generalize (correct = whatever humans do)
 - N-grams already generalize, just not in the right way!
- Ideas? Intuitions? Can we save this approach?
 - Hint: this is not a solved problem yet!

So what are people doing?

16

- Why do people have these judgments then?
 - bnick > bzick > mbick
- Two recent models
 - Albright 2009 featural bigrams
 - Hayes & Wilson 2008 Maximum Entropy model

Albright's Proposal

17

- bn is relatively good because it's similar to existing
 - bl
 - br
 - sn
 - ...
- bd is not as similar to existing sequences
 - bl is farther from bd than from bn
 - br is farther from bd than from bn
 - ...
- How to formalize *similar*?
 - Phonetic/phonological features!

Albright's Proposal

19

- Do bigrams over phonological *classes* not phonemes
- bn is b[+voi, +son, +cor, -strid, +nas]
 - bl is b[+voi, +son, +cor, -strid, ~~+nas~~]
 - We can group n and l into a small natural class b[n l r j]
 - Frequency of this specific combination is high
 - bl, br, bj are all attested
- bd is b[+voi, -son, +cor, -cont]
 - bl is b[+voi, ~~-son,~~ +cor, ~~-cont~~]
 - l and d have less in common, need larger class b[n l r j z n d]
 - Frequency of this combination is high, but not very related to bd

Albright's Proposal: Formally

20

- Phonotactic probability tradeoff between frequency and specificity

- $score(ab) = \max_{A,B \in Nat} \frac{count(AB)}{count(..)} \times P(a|A) \times P(b|B)$

- Where

- $P(a|A) = \frac{1}{|A|}$

- How does this deal with bd vs. bn vs. bl?

Hayes & Wilson (2008)

21

- Maximum Entropy Model of Phonotactics
 - closely related to regression models (in fact they are equivalent to multinomial logistic regression models)
 - You have some observation c to assign probability to
 - $$P(y = c) = \frac{e^{f_c(x)}}{Z} = \frac{e^{(w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)}}{Z}$$
 - You do so on the basis of various features $x_1 \dots x_n$
- For phonotactics, these features reflect phonological preferences and restrictions.
 - In Linguistics, features are called constraints
 - Constraint-based models of various kinds are widely used in theoretical phonology

Top MaxEnt Constraints (HW 2008)

22

Constraint	Weight	Comment	Examples
1. *[+ son, + dors]	5.64	*[ŋ]	*ŋ, *sŋ
2. *[+ cont, + voice, - ant]	3.28	*[ʒ]	*ʒ (see also #16)
3. * $\begin{bmatrix} \wedge - \text{voice} \\ + \text{ant} \\ + \text{strid} \end{bmatrix}$ [- approx]	5.91	Nasals and obstruents may only be preceded (within the onset) by [s].	*kt, *kk, *skt
4. *[] [+ cont]	5.17	Fricatives may not cluster with preceding C.	*sf, *sθ, *sh, *sfl
5. *[] [+ voice]	5.37	Voiced obstruents may not cluster with preceding C.	*sb, *sd, *sgr
6. *[+ son][]	6.66	Sonorants may only be onset-final.	*rt
7. *[- strid][+ cons]	4.40	Nonstrident coronals may not precede nonglides.	*dl, *tl, *θl
8. *[] [+ strid]	1.31	Stridents must be initial in a cluster.	*stʃ (see also #14, #22)
9. *[+ lab] $\begin{bmatrix} \wedge + \text{approx} \\ + \text{cor} \end{bmatrix}$	4.96	The only consonants that may follow labials are [l] and [r].	*pw vs. pl, pr
10. *[- ant] $\begin{bmatrix} \wedge + \text{approx} \\ - \text{ant} \end{bmatrix}$	4.84	Only [r] may follow nonanterior coronals.	*ʃl vs. ʃr
11. *[+ cont, + voice][]	4.84	Voiced fricatives must be final in an onset.	*vr, *vl vs. fr, fl
12. *[- cont, - ant][]	3.17	[tʃ] and [dʒ] must be final in an onset.	*tʃr, *dʒr vs. tr, dr (see also #22)

Results Comparison (Daland et al 2011)

23

model	syllabification				no syllabification			
	at-tested	mar-ginal	unat-tested	over-all	at-tested	mar-ginal	unat-tested	over-all
Albright	0·21	0·03	0·55	0·51	0·13	−0·07	0·18	0·26
bigram	0·19	0·16	0·22	0·78	0·23	0·01	−0·14	0·50
Coleman	0·35	0·31	−0·01	0·55	–	–	–	–
gnm.fig	0·07	0·25	−0·29	0·15	0·06	0·24	−0·32	0·08
gnm.oral	0·28	0·23	−0·28	0·22	0·26	0·23	−0·28	0·21
gnm.writ	0·17	0·24	−0·17	0·24	0·16	0·24	−0·30	0·15
gnm.lin	0·32	0·23	−0·22	0·31	0·30	0·22	−0·26	0·24
hw[100]	0·00	0·02	0·76	0·83	0·00	−0·31	0·79	0·68
hw[150]	0·00	0·06	0·69	0·82	0·00	0·04	0·67	0·75
hw[200]	−0·09	0·03	0·64	0·80	0·00	0·05	0·69	0·77
hw[250]	−0·09	0·13	0·64	0·84	0·00	0·00	0·70	0·80
hw[300]	−0·39	0·04	0·54	0·80	0·00	−0·02	0·70	0·81
hw[350]	−0·39	0·03	0·51	0·80	0·00	−0·10	0·67	0·81
hw[400]	−0·39	0·04	0·52	0·81	0·00	0·00	0·68	0·80
vl.uni	0·27	0·11	0·38	0·43	0·30	0·19	0·34	0·36
vl.bi	0·30	0·06	0·27	0·56	0·30	0·08	0·22	0·54

Results Discussion

24

□ Observations

- Bigram (.78) and Maxent (.83) do well overall
 - But it's easy to do well overall
- Bigrams (.22) do terribly on unattested
- MaxEnt does well on unattested (.76)
- No model does great on attested

□ Conclusions?

- Long way to go...
- But...
 - Features/similarity are important
 - Syllabification/higher level structure is important
- Ideas? Critiques? Questions?

Other Work

25

- One example of computational modeling to address cognitive questions
 - ▣ Main argument: human phonological generalization is feature-based and structure sensitive
- Modeling can also be used to address questions about learning *bias* to show how humans systematically ignore or depart from certain properties of the input

Sonority Sequencing Principle

26

- Sonority Sequencing Principle (SSP)
 - [lb]ack < [mb]ack < [bd]ack < [bn]ack < [bɹ]ack < [bj]ack
 - -2 -1 0 1 2 3
 - Typological Generalizations (Clements 1988, Selkirk 1984)
 - [bn]ack ⇒ [bl]ack
- Consistent findings of **Sonority Projection** in English
 - Preferences between unobserved clusters (nb vs db)
 - Production, perception, acceptability; aural, written
(Berent et al. 2007, Berent & Lennertz 2009, Berent et al. 2009, Davidson et al. 2004, Davidson 2006, Daland et al. 2011)
- Debate: Where do these preferences come from?
 - Daland et al examine SSP and argue models like MaxEnt can capture these effects...

How can this be? Statistics over classes!

27

	OO	ON	OL	OG
st	521	sn 109	pr 1046	kw 201
sp	313	sm 82	tr 515	sw 153
sk	278		kr 387	hw 111
			gr 331	tw 55
			br 319	dw 17
			fl 290	gw 11
			kl 285	θw 4
			fr 254	
			pl 238	
			bl 213	
			sl 213	
			dr 211	
			gl 131	
			θr 73	
			∫r 40	
	1112	191	4546	552
	(17.4%)	(3.0%)	(71.0%)	(8.6)%

- Plenty of SSP evidence

- LO none
- NO none
- OO only sO
- ON only sN
- OL 71%

- Higher level

- Son Rise 82.6%
- Son Plat. 17.4%
- Son Fall 0%

Polish Clusters

28

- My ongoing work
 - English is the wrong language to test!
 - Predictions based on input and based on bias are the same!
- Change tactic: look at a language that admits all these sequences!
 - What are predictions based on input?
 - What are predictions based on universal bias?
 - Build models for both and evaluate!
- Polish onsets
 - Whole Scale!
 - [wb]ack < [lb]ack < [mb]ack < [bd]ack < [bn]ack < [bɹ]ack < [bj]ack
 - -3 -2 -1 0 1 2 3

Polish Clusters: Examples

29

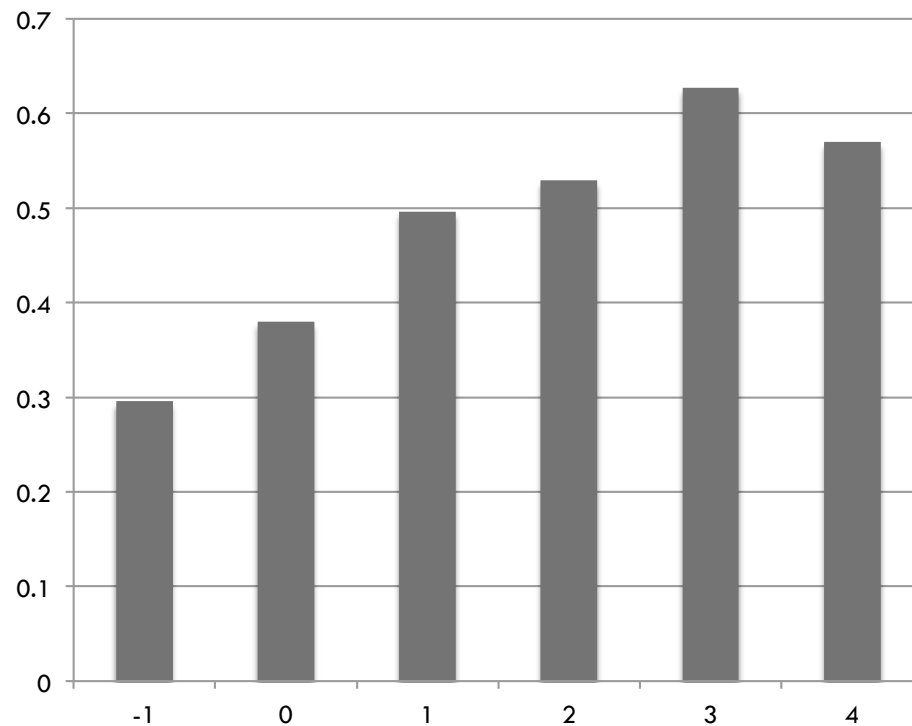
	P	F	N	L	G
P	[ptak] "bird"	[pʂijɛtɕ] "to come"	[dnɔ] "bottom"	[klutʃ] "key"	[gwɔva] "head"
	[ktɔ] "who"	[kɕɔwʂka] "book"	[dɲa] "day"	[drɔga] "road"	[bjawi] "white"
F	[ɕpi] "sleeps"	[xfila] "moment"	[ɕɲek] "snow"	[vlatɕ] "to pour"	[vwaɕɲe] "exactly"
	[stɔ] "one hundred"	[fʂistkɔ] "everything"	[smɔk] "dragon"	[fruvatɕ] "to fly"	[zwi] "bad"
N	[mdwɔ] "dull"	[mʂa] "mass"	[mɲe] "me (inst.)"	[mlɛkɔ] "milk"	[mjaw] "he had"
	[mgwa] "fog"	[mʂitsa] "mite"	[mnɔga] "multiple"	[mrufka] "ant"	[mjut] "honey"
L	[rtɛptɕ] "mercury"	[lvɨ] "lions"	[lɲani] "linen" (adj)		
	[rdza] "rust"	[rvatɕ] "tear"	[lnu] "linen" (gen)		
G	[wba] "head (gen)"	[wza] "tear"			
	[wkatɕ] "sob"				

What do kids learning Polish do?

30

SSP effects using scale: $P < F < N < L < G < V$

This is how accurate children are at producing onset clusters as a function of sonority rise degree. Strong relationship!



Polish Clusters: statistics

31

- Could kids get these preferences from input statistics?
 - Relative frequency very different from English
 - Half the input is Obstruent-Obstruent!
- I consider various ways of modeling the input
 - Segment bigrams
 - Class-based bigrams
 - MaxEnt
- None can capture the kids' preferences...

Sonority Profile	LO	OO	NN	ON	NL	OL	NG	OG
Token Frequency	0.04%	50.90%	0.50%	3.60%	0.20%	20.40%	3.00%	21.30%
Type Frequency	0.10%	47.70%	0.20%	6.70%	0.60%	19.60%	2.90%	22.30%

Thanks!

32

- Two brief examples of what we do
- Lots of other people in Linguistics interested in computational modeling (not just phonology)
- If you're interested in more...
 - Come talk to us!
 - Check out our classes!