Syntactic universals in the lab: New methods and approaches Jennifer Culbertson, University of Edinburgh

Advanced Core Training in Linguistics (ACTL), Summer 2015. University College London.

Lecture 5: Models, theory and explanation

Overview

FORMALIZATIONS OF UNIVERSAL 18, 20. Two cases of word order universals in which we have both typological and experimental data as well as well-specified theoretical models of the bias.

Universal 18: PCFG model¹

INFERRING A PCFG GIVEN SOME DATA.

• What is a PCFG? It's a CFG with probabilities attached to productions

Combining (by multiplication) a beta distribution with a binomial...

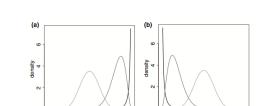
- Given a set of data, what probability should you infer?
- Binomial: (*c*=counts, *t*=total trials, *p*=probability)

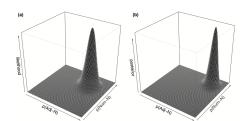
binomial(c|p) =
$$\binom{t}{c} p^c (1-p)^{t-c}$$

• From the experiment: given 12 instances of Adj-N and 28 instances of N-Num...?

binomial(12|0.3) = 0.14 binomial(28|0.7) = 0.14

What if you have a regularization bias...?





binomial(28|0.9) = 0.0003

BAYESIAN INFERENCE WITH PCFGs.

- We just did Bayesian inference!
- Bayes rule: contribution of input data and prior knowledge to learning

 $P(Grammar|Data) \propto P(Data|Grammar)P(Grammar)$

- *Grammar*: probabilistic re-write rules p(Adj-N), p(Num-N)
- Likelihood: binomial probability of training counts given grammar
- Prior:
 - Beta distribution as formalization of regularization bias (add counts)
 - Multinomial weights on grammar types as pattern bias
- What bias do learners have? Start with flat prior (no bias), fit to behavioral data
- Result:
 - Strong regularization bias; α , $\beta = (16.5, 0.001)$
 - Asymmetry among patterns; $\gamma = (0.62, 0.27, 0.0001, 0)$

¹ J. Culbertson, P. Smolensky, *Cognitive Science* **36**, 1468 (2012)

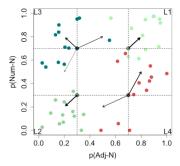


Figure 1: Basic result.

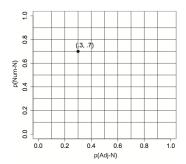


Figure 2: Grid approximation of space of PCFG grammars.

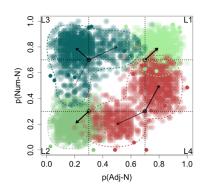
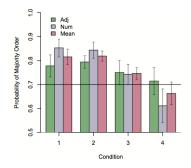


Figure 3: Predictive distribution of production grammars.

Universal 18: OT/PHG model²

BASIC RESULTS. Regularization and quantitative difference in use of majority order for numeral phrases across condition.



a. Num position in L3

| {N | Jum, N} | Num-L | HEAD-R | HEAD-L | | | | | | |
|-----------------------|---------|-------|--------|--------|--|--|--|--|--|--|
| a. E | Num-N | | * | | | | | | | |
| b. | N-Num | *! | | * | | | | | | |
| h Adi nosition in I ? | | | | | | | | | | |

| | {Adj | , N} | Num-L | HEAD-R | HEAD-L | | |
|----|------|-------|-------|--------|--------|--|--|
| a. | | Adj-N | | *! | | | |
| b. | 133 | N-Adj | | | * | | |

From OT to Probabilistic Harmonic Grammar (aka MaxEnt).

- OT constraints for Universal 18
 - HEAD-L, HEAD-R, NUM-L
 - Allows L1, L2, L3, prohibits L4 (no way to derive it without Nuм-R)
 - No way to predict harmonic > L₃
- Moving to PHG with bias
 - HEAD-L, HEAD-R, NUM-L, NUM-R
 - Allow all possible patterns
 - Prior penalty for use of specific → favor harmonic
 - Combine with prior penalty for use of Num-R \rightarrow disfavor (Adj-N, N-Num)

Num position in L3 ($Z \equiv e^{-0.85} + e^{-1.7}$)

| | | 1.7 | .85 | 0 | $H_G(x) = -\sum_k w_k C_k(x)$ | $P(x C) = e^{H_G(x)/7}$ | | | |
|----|--------------|-------|-------|--------|-------------------------------|---|-------------------------------------|--|--|
| | $\{Num, N\}$ | | Num-L | HEAD-R | HEAD-L | $HG(x) = -\sum_{k} w_k C_k(x)$ | $\Gamma(x G) = e^{-\phi \cdot x/Z}$ | | |
| a. | 13. | Num-N | | * | | $-[1.7 \cdot 0 + 0.85 \cdot 1 + 0 \cdot 0]$ | $e^{-0.85}/Z = 70\%$ | | |
| b. | | N-Num | *! | | * | $-[1.7 \cdot 1 + 0.85 \cdot 0 + 0 \cdot 1]$ | $e^{-1.7}/Z = 30\%$ | | |

- How to formalize regularization?
 - Bias over constraint weights (gaussian prior penalizing weights near zero)
 - Bias over resulting probabilities (beta distribution with asymmetric shape parameters, identical to PCFG model)
- Results:
 - Both models can capture average differences among conditions
 - Weight-based model captures asymmetry in numeral ordering better
 - Probability-based model is preferable on complexity grounds

OVERALL TAKE HOME MESSAGE

- The Bayesian view presented in both these treats underlying biases as soft
- All grammars (defined by the particular framework) are in principle part of the hypothesis space
- A priori biases make some less likely to be inferred
- Strong enough prior biases can be strong enough to prevent a particular grammar from being inferred (approximates hard/absolute constraints)

² J. Culbertson, P. Smolensky, C. Wilson, Topics in Cognitive Science 5, 392 (2013)

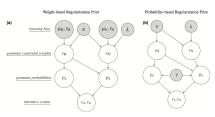


Figure 4: Alternative PHG models of regularization.

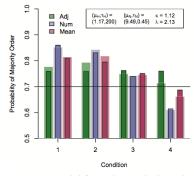


Figure 5: Model fit with weight-based prior.

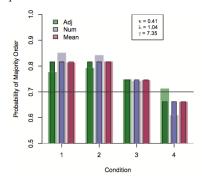


Figure 6: Model fit with probabilitybased prior.

Explanations and Universal 20

WHY? If U18 and U20 are the result of cognitive biases, what is the content of the biases? This brings us back to several of the larger issues we've been discussing:

- How to reconcile statistical generalizations with non-statistical theories?
- Are the biases part of the grammatical system or external to it?

THEORETICAL ACCOUNTS OF UNIVERSAL 20.

Original and Hawkins' reformulation³

Universal 20. When any or all of the items (demonstrative, numeral, and descriptive adjective) precede the noun, they are always found in that order. If they follow, the order is either the same or its exact opposite.

. For those that follow, no predictions are made, though the most frequent order is the mirror image of the order for preceding modifiers. In no case does the adjective precede the head when the demonstrative or numeral follows

Cinque⁴

Of the 24 mathematically possible orders of the four elements demonstrative, numeral, adjective, and noun, only 14 appear to be attested in the languages of the world. Some of these are unexpected under Greenberg's Universal 20. Here it is proposed that the actually attested orders, and none of the unattested ones, are derivable from a single, universal, order of Merge (Dem > Num > Adj > N) and from independent conditions on phrasal movement.

| | ✓ | a | Dem | Num | Α | N | MANY | 0 | m | Dem | Α | Num | N | zero |
|---|------------|---|-----|-----|-----|-----|------|---|---|-----|-----|-----|-----|------|
| | ✓ | b | Dem | Num | N | Α | many | ✓ | n | Dem | Α | N | Num | FEW |
| | ✓ | c | Dem | N | Num | Α | FEW | ✓ | 0 | Dem | N | Α | Num | many |
| 7 | ✓ | d | N | Dem | Num | Α | few | ✓ | p | N | Dem | Α | Num | FEW |
| ı | 0 | e | Num | Dem | Α | N | zero | 0 | q | Num | Α | Dem | N | zero |
| ı | 0 | f | Num | Dem | N | Α | zero | ✓ | r | Num | Α | N | Dem | FEW |
| ı | 0 | g | Num | N | Dem | Α | zero | ✓ | s | Num | N | Α | Dem | few |
| ı | 0 | h | N | Num | Dem | Α | zero | ✓ | t | N | Num | Α | Dem | few |
| ı | 0 | i | Α | Dem | Num | N | zero | 0 | u | Α | Num | Dem | N | zero |
| ı | 0 | j | A | Dem | N | Num | zero | 0 | v | A | Num | N | Dem | zero |
| ı | 1 | k | A | N | Dem | Num | FEW | ✓ | w | A | N | Num | Dem | FEW |
| ı | ✓ | 1 | N | Α | Dem | Num | few | ✓ | x | N | Α | Num | Dem | MANY |
| _ | _ | | | | | | | | | | | | | |

- Linear Correspondence Axiom plus universal merge order
- Movement of NP or XP containing it to c-commanding position
- Abels & Neeleman⁵
 - Any linearization of universal merge order plus leftward movement
 - Movement of sub-tree containing N to c-commanding position
- Steddy & Samek-Lodovici⁶: OT alignment constraints Dem-L, Num-L, Adj-L, N-L
- Issues⁷
 - Where to draw the line?
 - Do we actually want to predict that a type with frequency 1 is possible while an unattested type is impossible? Are we justified statistically?
 - Even worse, should we make this assumption at the cost of not explaining the huge frequency differences?
- Alternative explanation (with experimental support)
 - Universal semantic scope relations plus bias for isomorphic linearization
 - Bias against non-harmonic patterns
 - Bias against pre-nominal Adj (see also U18, where Adj-N, N-Num is WOW)

"It is a fact of logic that one cannot derive statistical predictions from a non-statistical theory. Therefore, providing an argument for any particular grammatical analysis on the basis of frequency information is an arduous task. The first step in constructing such an argument must consist in pairing the proposed analysis with a theory of markedness. The latter identifies grammatical structures as favored or disfavored, and thus generates statistical predictions. Unfortunately, little is known about this aspect of the language faculty. An independently motivated theory of markedness, which would allow us to test hypotheses about the grammar, is no more than a distant hope."

K. Abels, A. Neeleman, Ms., University of Tromsø and University College London

³ J. A. Hawkins, Word order universals (Academic Press, New York, 1983)

⁴ G. Cinque, Linguistic Inquiry 36, 315 (2005)

- ⁵ K. Abels, A. Neeleman, Syntax 15, 25
- ⁶ S. Steddy, V. Samek-Lodovici, Linguistic Inquiry 42, 445 (2011)

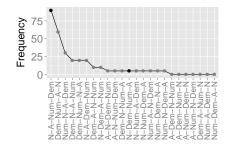


Figure 7: Distribution of 24 orders. 7 M. Cysouw, Linguistic Typology 14, 253 (2010)

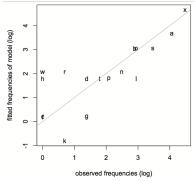


Figure 8: Fit of regression model with scope, harmony, N-Adj.