

# APPROXIMATIONS OF PREDICTIVE ENTROPY CORRELATE WITH READING TIMES

---

Marten van Schijndel   William Schuler

July 29, 2017

Department of Linguistics, The Ohio State University

Angele et al. (2015)

A | child | XXXXXXX | the | fish

Angele et al. (2015)

A	child <sup>*</sup>	XXXXXXX	the	fish
A	child	annoyed <sup>*</sup>	XXX	fish

Angele et al. (2015)

A	child <sup>*</sup>	XXXXXXX	the	fish
A	child	annoyed <sup>*</sup>	XXX	fish
A	child	annoyed	the <sup>*</sup>	XXXX

Angele et al. (2015)

A	child <sup>*</sup>	XXXXXXX	the	fish
A	child	annoyed <sup>*</sup>	XXX	fish
A	child	annoyed	the <sup>*</sup>	XXXX

Lexical frequency of the upcoming masked word affects processing

Angele et al. (2015)

A	child <sup>*</sup>	XXXXXXX	the	fish
A	child	annoyed <sup>*</sup>	XXX	fish
A	child	annoyed	the <sup>*</sup>	XXXX

Lexical frequency of the upcoming masked word affects processing

Hypothesis: Effect is due to uncertainty over continuations

Angele et al. (2015)

A	child <sup>*</sup>	XXXXXXX	the	fish
A	child	annoyed <sup>*</sup>	XXX	fish
A	child	annoyed	the <sup>*</sup>	XXXX

Lexical frequency of the upcoming masked word affects processing

Hypothesis: Effect is due to uncertainty over continuations

Problem: Uncertainty is expensive to calculate

Shannon (1948)

$$H(X) \stackrel{\text{def}}{=} - \sum_{x \in X} P(x) \log P(x) \quad (1)$$



Shannon (1948)

$$H(X) \stackrel{\text{def}}{=} - \sum_{x \in X} P(x) \log P(x) \quad (1)$$

Roark et al. (2009) distinguishes two kinds of entropy  
(over words and preterminals)

$$\text{Lex}H(w_{1..i-1}) \stackrel{\text{def}}{=} - \sum_{w_i \in V} P_G(w_i | w_{1..i-1}) \log P_G(w_i | w_{1..i-1}) \quad (2)$$

$$\text{Syn}H(w_{1..i-1}) \stackrel{\text{def}}{=} - \sum_{p_i \in G} P_G(p_i | w_{1..i-1}) \log P_G(p_i | w_{1..i-1}) \quad (3)$$

Roark et al. (2009) showed

- $SynH$  predicts self-paced reading times
- $LexH$  is not predictive of SPR times

Roark et al. (2009) showed

- $SynH$  predicts self-paced reading times
- $LexH$  is not predictive of SPR times  
(No Angele et al., 2015, effect)

Roark et al. (2009) showed

- $SynH$  predicts self-paced reading times
- $LexH$  is not predictive of SPR times  
(No Angele et al., 2015, effect)

But

- Small training corpus ( $V$  is poor)
- Small test corpus:  
~ 200 sentences, ~ 4000 words, 23 subjects

Natural Stories self-paced reading corpus (Futrell et al., in prep)

- 181 subjects
- 10 narrative texts
- 485 sentences (10256 words)
- Each text followed by 6 comprehension questions
- Events removed if  $<100$  ms or  $>3000$  ms

Parsed using Roark (2001) parser

Fitted with *lmer*

# SPACES WERE MASKED

-----

A -----

- child -----



----- annoyed -----

----- the -----

----- fish.

# SYNTACTIC ENTROPY PREDICTS RTs

Predictor	$\hat{\beta}$	$\hat{\sigma}$
<b>Syntactic <math>H</math></b>	<b>4.53*</b>	<b>0.54</b>
Lexical $H$	-1.05	0.41

Replication of Roark et al. (2009)

# SYNTACTIC ENTROPY PREDICTS RTs

Predictor	$\hat{\beta}$	$\hat{\sigma}$
<b>Syntactic <math>H</math></b>	<b>4.53*</b>	<b>0.54</b>
Lexical $H$	-1.05	0.41

Replication of Roark et al. (2009)

But Angele et al. (2015) found a *lexical* frequency effect

## CAN WE MAKE LEXH MORE TRACTABLE?

$$S_G(w_i, w_{1..i-1}) \stackrel{\text{def}}{=} -\log P_G(w_i | w_{1..i-1}) \quad (4)$$

$$\text{Lex}H_G(w_{1..i-1}) \stackrel{\text{def}}{=} \sum_{w_i \in V} -P_G(w_i | w_{1..i-1}) \log P_G(w_i | w_{1..i-1}) \quad (5)$$

$$= \sum_{w_i \in V} P_G(w_i | w_{1..i-1}) S_G(w_i, w_{1..i-1}) \quad (6)$$

$$= E[S_G(w_i, w_{1..i-1})] \quad (7)$$

## CAN WE MAKE LEXH MORE TRACTABLE?

$$S_G(w_i, w_{1..i-1}) \stackrel{\text{def}}{=} -\log P_G(w_i | w_{1..i-1}) \quad (4)$$

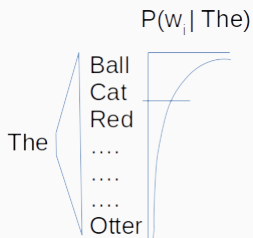
$$\text{Lex}H_G(w_{1..i-1}) \stackrel{\text{def}}{=} \sum_{w_i \in V} -P_G(w_i | w_{1..i-1}) \log P_G(w_i | w_{1..i-1}) \quad (5)$$

$$= \sum_{w_i \in V} P_G(w_i | w_{1..i-1}) S_G(w_i, w_{1..i-1}) \quad (6)$$

$$= E[S_G(w_i, w_{1..i-1})] \quad (7)$$

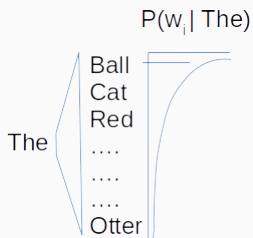
We can use a corpus instead of explicitly computing the expectation

# ENTROPY GIVES MEAN SURPRISAL

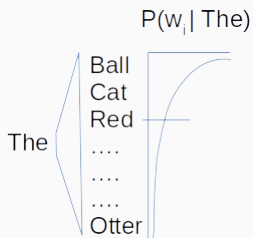




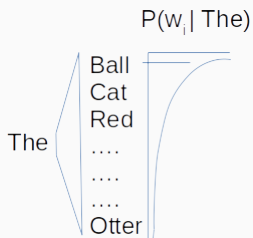
# SURPRISAL APPROXIMATES ENTROPY IN THE AGGREGATE



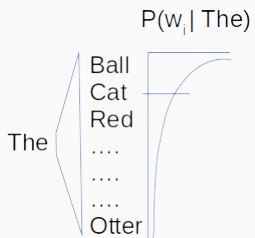
# SURPRISAL APPROXIMATES ENTROPY IN THE AGGREGATE



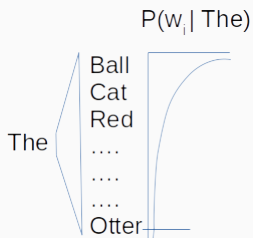
# SURPRISAL APPROXIMATES ENTROPY IN THE AGGREGATE



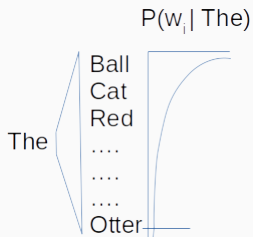
# SURPRISAL APPROXIMATES ENTROPY IN THE AGGREGATE



# SURPRISAL APPROXIMATES ENTROPY IN THE AGGREGATE

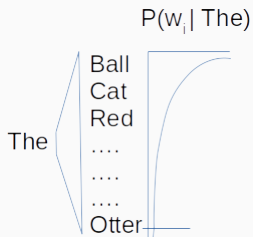


# SURPRISAL APPROXIMATES ENTROPY IN THE AGGREGATE



Ex: The boy annoyed the fish.

# SURPRISAL APPROXIMATES ENTROPY IN THE AGGREGATE



We can treat large corpora as our samplers.

We can try:

- Future Roark surprisal  
(same distribution as SynH)



We can try:

- Future Roark surprisal  
(same distribution as SynH)
- Future 5-gram Surprisal  
(similar to what Angele et al., observed)

We can try:

- Future Roark surprisal  
(same distribution as SynH)
- Future 5-gram Surprisal  
(similar to what Angele et al., observed)
- Future categorial grammar surprisal  
(tests how specific syntactic prediction is)

# UNCERTAINTY OVER BOTH WORDS AND SYNTAX

Predictor	$\hat{\beta}$	$\hat{\sigma}$
<b>Syntactic <math>H</math></b>	<b>4.62*</b>	<b>0.53</b>
Future Roark Surprisal	0.33	0.40
<b>Future <math>N</math>-gram Surprisal</b>	<b>4.05*</b>	<b>0.58</b>
<b>Future Categorical Grammar Surprisal</b>	<b>4.10*</b>	<b>0.74</b>

# WHY DOES THIS PRE-SLOWING OCCUR?

- Better encoding of  $w_i$  to help with  $w_{i+1}$

# WHY DOES THIS PRE-SLOWING OCCUR?

- Better encoding of  $w_i$  to help with  $w_{i+1}$
- A kind of Uniform Information Density (UID; Jaeger, 2010)
  - Optimizes per-millisecond informativity

# CONCLUSIONS

- Uncertainty about upcoming words slows processing

- Uncertainty about upcoming words slows processing
- That influence can be detected prior to any expectation violation



- Uncertainty about upcoming words slows processing
- That influence can be detected prior to any expectation violation
- Future surprisal can efficiently approximate that uncertainty

- Uncertainty about upcoming words slows processing
- That influence can be detected prior to any expectation violation
- Future surprisal can efficiently approximate that uncertainty
- Syntactic uncertainty is fine-grained

Thanks to:

- The reviewers for their very helpful comments
- National Science Foundation (DGE-1343012)