

CORPUS LINGUISTICS - NGRAMS

NATURAL LANGUAGE PROCESSING - CS 322.00

Blake Howald
September 22, 2017

AGENDA

- Presentation – Lazar Zamurovic, Solomon Foster
- Logistics
- Questions
- Corpus Linguistics
- Ngrams (Part I)

CORPUS LINGUISTICS

- It's good to have data (generally) and lots of it.
- Can be use for the creation of a *Language Model (LM)* – Statistical model of word sequences (although can do at character level)
- Things to pay attention to when working with/ selecting corpora
 1. Representative of what you are looking to capture
 2. Robust enough to capture range of expected (and unexpected) variation
 3. 1. & 2. speak to the generalizeability of your model to unseen data
 4. Data "scrubbing" often required (garbage in – garbage out)
 5. Attention to inherent patterns of language and associated impacts on model (if any)

NGRAMS - APPROACH

*"just one more regret in a long list of **many**"*

The (conditional) probability of a word w given an associated history h .

$$P(w|h)$$

$$P(\text{many} \mid \text{just one more regret in a long list of})$$

$$C(\text{just one more regret in a long list of many}) / C(\text{just one more regret in a long list of})$$

NGRAMS – ASSUMPTIONS

- Markov assumption lessens the burden on collecting and maintaining a long history
- We can predict the probability of a future state without looking too far back in history

$P(\text{many} \mid \text{just one more regret in a long list of})$

$P(\text{many} \mid \text{of})$

$P(\text{many} \mid \text{of}) \approx P(\text{many} \mid \text{just one more regret in a long list of})$

NGRAMS – ASSUMPTIONS

$P(\text{many} \mid \text{of}) \approx P(\text{many} \mid \text{just one more regret in a long list of})$

$C(\text{of many})/C(\text{of})$

$P(\text{many} \mid \text{list of}) \approx P(\text{many} \mid \text{just one more regret in a long list of})$

$C(\text{list of many})/C(\text{list of})$

$P(\text{many} \mid \text{long list of}) \approx P(\text{many} \mid \text{just one more regret in a long list of})$

$C(\text{long list of many})/C(\text{long list of})$

Use of relative frequencies allows us to estimate not only the Maximum Likelihood Estimates (MLEs) of bigrams (or any ngram) but also the MLE of sequences of bigrams not explicitly seen in the corpus.

N(BI)GRAMS - CALCULATION

<s>Six sick bricks tick</s>

<s>Six sick chicks tock</s>

<s> six	$P(\text{six} \text{<s>})$	$C(\text{<s> six})/C(\text{<s>}) = 2/2 = 1$
six sick	$P(\text{sick} \text{six})$	$C(\text{six sick})/C(\text{six}) = 2/2 = 1$
sick bricks	$P(\text{bricks} \text{sick})$	$C(\text{sick bricks})/C(\text{sick}) = 1/2 = 0.5$
bricks tick	$P(\text{tick} \text{bricks})$	$C(\text{bricks tick})/C(\text{bricks}) = 1/1 = 1$
tick</s>	$P(\text{</s>} \text{tick})$	$C(\text{tick </s>})/C(\text{tick}) = 1/1 = 1$
sick chicks	$P(\text{chicks} \text{sick})$	$C(\text{sick chicks})/C(\text{sick}) = 1/2 = 0.5$
chicks tock	$P(\text{tock} \text{chicks})$	$C(\text{chicks tock})/C(\text{chicks}) = 1/1 = 1$
tock</s>	$P(\text{</s>} \text{tock})$	$C(\text{tock </s>})/C(\text{tock}) = 1/1 = 1$

NGRAMS - CALCULATION

Bigram	MLE	Bigram	MLE
have a	0.707	don't have	0.117
a box	0.589	I don't	0.117
box </s>	0.471	<s> I	0.589
I can't	0.471	can't be	0.471
to be	0.117	dumb </s>	0.117
so dumb	0.117	be so	0.333

- *I don't have a box*
- $P(I|\text{<s>}) \times P(\text{don't}|I) \times P(\text{have}|\text{don't}) \times P(a|\text{have}) \times P(\text{box}|a) \times P(\text{</s>}|\text{box})$
- $0.589 \times 0.117 \times 0.117 \times 0.707 \times 0.589 \times 0.471 = 0.001581$
- *I can't be so dumb*
- *We don't have a box*