

Local and global lexicon: a novel approach to quantifying persistence Jefrey Lijffijt

Department of Information and Computer Science Aalto University





Aalto University

UNIVERSITY OF HELSINKI

Introduction

- Development of data mining tools for study of language change
 - Joint work with Tanja Säily, Terttu Nevalainen

- <u>Today: Repetition</u>
 - General measure: Type/token ratio (over time)
 - Localization: Significance testing
 - Understanding: Unexpected repetition





Motivation

- <u>Persistence</u> of linguistic structures
 - Both conscious and subconscious (priming)
- Ample evidence for lexical and syntactic persistence (Bock 1986, Pickering & Branigan 1999), also from corpus linguistics (Gries 2005, Szmrecsanyi 2005, Dubey et al. 2008)
 - Always specific questions, never bottom up
- Frequency(word, text) = Author + Topic + Priming
 - Too hard → Deep understanding of word frequency needed first





Repetition at text level

- Type/token ratio (TTR) =
 # unique words / # words
- General measure for amount of repetition in a text

- British National Corpus (BNC-XML)
 - Plain words
 - Ignore punctuation, capitalization





Repetition depends on context



Sampling

- Solve text length bias by using sampling
 - For example 100 samples of 2,000 words







TTR over time

- Text as sequence of words
- Samples of 2000 words
 - Incremental sliding window





TTR over time



DATA

MINING

for changing modalities of communication

TOOLS



Significance testing

- Randomization approach
 - Produce 500 graphs based on random permutation
 - Equal frequency for every word
 - 1 sample = 1 graph
 - Dependency between windows
 - Compare equal ranks (Lowest vs. lowest points, highest vs. highest points)
 - Multiple hypothesis: Benjamini-Hochberg





Significant deviation

Type/token ratio over time, window = 2000



Significant deviation

Type/token ratio over time, window = 4000



What does this mean?

• Word frequency distribution inside a text changes over time!

- Can we explain the significant repetition?
- Most frequent words account for most repetition
 - This is also <u>expected</u> in randomized text
 - Significant drop \rightarrow unexpected repetition





Finding unexpected repetition

- Test probability for each word for each window
 - Likelihood function (Kleinberg 2004)
 - Likelihood ratio (Dunning 1993)
 - Chi-square inaccurate for very large or very small samples (Rayson et al. 2004)
- Likelihood ratio seems good choice
 - P(freq | prob-global) / P(freq | prob-local)





Summary of text EDN



Conclusion

- Significant repetition can be found using randomization approach
- Local explanation can be given using likelihood ratio
- We can construct a timeline of unexpectedly frequent words to summarize a text
- Open questions:
 - What is a good cut-off value for likelihood ratio
 - Can we detect topic shift in a text?



