Structural Analysis of Word Collocation Networks

Shibamouli Lahiri University of North Texas

April 29, 2014

Shibamouli Lahiri Word Collocation Networks 1

・ロト ・回ト ・ヨト ・ヨト

Structuralist View	Speech View
Grammars POS tags Dependencies Co-reference resolution	n-gram language models PCFG language models Speech recognition Machine translation

<ロ> (四) (四) (三) (三) (三) (三)

Structuralist View	Speech View	Complex Networks View
Grammars	n-gram language models	Networks of words, sentences, documents,
POS tags	PCFG language models	Small world, scale-free networks
Dependencies	Speech recognition	Power-law degree distribution
Co-reference resolution	Machine translation	High clustering coefficient, low diameter

・ロン ・回 と ・ ヨン ・ ヨン

3

- Sentence networks
 - Summarization (LexRank, TextRank)
- Word networks
 - Keyword Extraction (TextRank, ExpandRank)
 - Authorship Attribution (Antiqueira et al., 2006)
 - Genre Identification (Stevanak et al., 2010)
 - Opinion Classification (Amancio et al., 2011)
 - Semantic Analysis (Biemann et al., 2012)

回 と く ヨ と く ヨ と

- Examples of word collocation networks
- Structural properties
- Three exploratory analyses

イロン イヨン イヨン イヨン

Э

The quick brown fox jumped over the lazy dog.



◆□ → ◆□ → ◆三 → ◆三 → ◆ ● ◆ ◆ ● ◆

6

Network Property	Notation
Number of vertices	V
Number of edges	<i>E</i>
Shrinkage exponent (Leskovec et al., 2007)	$\log_{ V } E $
Global clustering coefficient	Ċ
Small-worldliness (Walsh, 1999; Matsuo et al., 2001)	$\mu = (ar{ extsf{C}}/L)/(ar{ extsf{C}}_{ extsf{rand}}/L_{ extsf{rand}})$
Diameter (directed)	$d_{directed}$
Diameter (undirected)	$d_{undirected}$
Power-law exponent of degree distribution	α
Power-law exponent of in-degree distribution	α_{in}
Power-law exponent of out-degree distribution	α_{out}
p-value for $lpha$	p_{lpha}
p-value for α_{in}	$p_{\alpha_{\textit{in}}}$
p-value for α_{out}	$p_{lpha_{out}}$
Number of connected components	#CC
Number of strongly connected components	#SCC
Size of the largest connected component	-
Size of the largest strongly connected component	-

・ロト ・回 ト ・ヨト ・ヨト

7

- How do network properties vary across different genres of text?
- One work properties vary across different network types?
- **3** How do the properties **evolve** as a word network grows in size?

(1日) (日) (日)

Blog Authorship Corpus (Schler et al., 2006)

 19,320 blogs, 136.8 million words

Reuters-21578. Distribution 1.0

- 19,043 news articles, 2.6 million words
- NIPS Conference Papers Vols 0-12
 - 1,740 papers, 4.8 million words
- E-books from Project Gutenberg
 - 3,036 books, 210.9 million words

- 4 同 6 4 日 6 4 日 6

Exploratory Analysis 1: Across Genres



Shibamouli Lahiri Word Collocation Networks 10

イロン イロン イヨン イヨン 三日

Exploratory Analysis 2: Across Network Types



Shibamouli Lahiri Word Collocation Networks 11

< □ > < □ > < □ > < □ > < Ξ > < Ξ > = Ξ

Exploratory Analysis 3: Evolution of Properties



ヨト ヨ

Exploratory Analysis 3: Evolution of Properties



Shibamouli Lahiri Word Collocation Networks 13

< 🗗 🕨

12

표 🕨 🗉 표

- There are statistically significant variations between **genres** for distributions of network properties.
- There are also statistically significant variations between **network types** for distributions of network properties.
- As word networks grow in size, key structural properties **evolve** in *phases*, via spikes and drops.
- Networks **densify** as they grow, evidenced by the shrinkage in diameter.
- Networks from news articles have low small-worldliness (μ) and high α .

・ロン ・回 と ・ ヨ と ・ ヨ と

3

 https://drive.google.com/file/d/ OB2Mzhc7popBgODFKZVVnQTFMQkE/edit?usp=sharing

・ロン ・回 と ・ 回 と ・ 回 と

Э

• Visualization of collocation networks

- IBM Many Eyes project
 - (http://www.manyeyes.com/)

• Application of collocation networks

- Authorship Attribution (Lahiri and Mihalcea, 2013)
- Keyword Extraction (Lahiri et al., 2014)
- Gender and Genre Classification
- Stylometry

ヘロン 人間 とくほと くほとう

Questions?

Shibamouli Lahiri Word Collocation Networks 17

< □ > < □ > < □ > < □ > < □ > < □ > = □