

Department of Informatics, University of Rijeka Radmile Matejčić 2, 51000 Rijeka, Croatia Tel.: + 385 51 584 700 Fax: + 385 51 584 749 www.langnet.uniri.hr

### **Multilaverd Language Networks**

### Sanda Martinčić-Ipšić

#### smarti@uniri.hr

#### LangNet team

Ana Meštrović Slobodan Beliga Tajana Ban Kirigin Tanja Miličić

### Language

main tool of communication

- reflects our history and culture
- evolving in parallel with our society
- can be seen as a complex adaptive system
- written (as well as spoken) language can be modeled via complex networks
- the lingual units (words) are represented by vertices and their linguistic interactions by links
- allows systematic quantitative analyses



### Language networks



model the various language subsystems (levels)

- examine unique **function** through complex networks
- examine various linguistic units
- deepening the understanding of conceptual similarities. differences and universalities in natural languages
- cognitive representation of the language in the human brain

#### establish a bridge:

 linguistics, complex networks science, computer science, and natural language processing

### Language networks - levels

• various language subsystems - represented as complex networks

-syllabic

-graphemic

• sub-word level:

-phonetic (phonology)

-morphology (morphosyntactic)

- vertices/nodes linguistic units
- edges/links model their relationships

#### • word level:

- -co-occurance
- -syntax semantics
- pragmatics

- present: focus on isolated linguistic subsystems
- -lacking to explain (or even explore) the mechanism of their mutual interaction, interplay or inheritance

### Outline

 $\mathbb{N}$ 



### Word-level networks

#### co-occurrence networks

- directed or undirected [ITIS 2013a]
- weighted or unweighted [ITIS 2013a] • stopwords preserved [ITIS 2013a, MIPRO2014a]
- not lemmatized
- in the full variety of flective word forms

• size of the co-occurrence window: 2 [ITIS 2013a]

• within boundaries: words and sentences [ITIS 2013a, CompleNet 2014]

#### • sensitive to used corpus





### **Syntax Datasets**

- Croatian & English Dependency Treebanks parsed syntax tree
- sentences HR: 3.465 EN: 3.829
- tokens HR: 88.045 tokens EN: 94.084 tokens S





Ń



### Subword-level networks

### • syllables network [MIPRO2013]

- syllables that co-occur in the same word -also syllables across words - toward speech
- Croatian has two possible syllabifications
- phonological and phonetic
- -phonological syllabification: our algorithm [Speech, 2016]
- -phonetic syllabification: our grapheme-to-phoneme method

graphemes that co-occur in the same word



#### • graphemes network



### Multilayer Language Network Experiment

- datasets: Croatian and English Dependency Treebanks
- •10 networks (5 HR+ 5 EN): directed and weighted
- not lemmatized, stopwords included
- word level: sentence boundaries
- co-occurrance window size 2
- shuffle
- syntax

#### subword level: words boundaries

- syllables from words in original sentences
- graphemes from words in original sentences

- Multilayer Language Network M is a quintuple M = (V<sub>M</sub>, E<sub>M</sub>, V, L, C) • V is a non-empty set of nodes;
- C is a nonempty set of perspective elements;
- L is a set of perspects L, where {L<sub>0</sub>, L<sub>1</sub>; L<sub>2</sub>} is a partition of C.
- L<sub>0</sub> language perspect, L<sub>1</sub> hierarchy perspect and L<sub>2</sub> construction perspect;
- For perspect L<sub>1</sub> = {g<sub>1</sub>,..., g<sub>k</sub>} sequence of its elements g<sub>1</sub>,..., g<sub>k</sub> is the subsequence of the following sequence - hierarchy:
- -discourse, sentence; phrase; syntagm; word; morphem; syllable; phoneme; grapheme
- An element of the set  $L_0 \times L_1 \times L_2$  is called a **layer**;
- $V_M \subseteq V \times L_0 \times L_1 \times L_2$  is the set whose elements are called **MLN**nodes;
- $E_M \subseteq V_M \times V_M$  is the set of edges.

## **Multilayer Network definition**





### **Results: Croatian Dataset**

	W				
	Co-occurr		Syntax		
Number of nodes (N)	23359	23359	23359	2634	34
Number of edges (K)	71860	86214	70155	18849	491
Number of components	2	2	2	17	1
Average path length (L)	4.01	3.74	1.81	1.86	1,58
Diameter (D)	16	17	12	8	3
Average clustering coefficient(C)	0.17	0.19	0.12	0.26	0.64
Transitivity	0.004	0.013	0.003	0.120	0.522
Density	0.00013	0.00016	0.00013	0.00272	0.43761

12

Langnet

15

• avg. path length - degree of separation between linguistic units

diameter – maximal separation

• density - probability of connecting 2 units

• transitivity - realized number of triangls (among possible ones)

## **Results: English Dataset**

	W				
Number of nodes (N)	10930	10930	10930	2599	26
Number of edges (K)	50299	58920	52221	6053	333
Number of components	3	1	3	54	1
Average path length (L)	3.47	3.45	1.96	1.88	1.51
Average clustering coefficient(C)	0.286	0.295	0.153	0.057	0.0838
Transitivity	0.009	0.016	0.014	0.020	0.654
Density	0.00042	0.00049	0.00044	0.0009	0.5123

13

avg. path length - degree of separation between linguistic units diameter – maximal separation density – probability of connecting 2 units transitivity – realized number of triangls (among possible ones)



### **Complex network measures**

### selectivity

- average weight distribution on the links of the single node
- the average strength of the node

$$s_i^{in/out} = \sum_j w_{ji/ij}$$
  $e_i^{in/out} = \frac{s_i^{in/out}}{k_i^{in/out}}$ 





3

### Word-level layers overlap

between two network layers α and α'

۲

S

21

• Jaccard overlap:



#### • Preserved weighted ratio:



#### Preserved weighted overlap:

 normalized preserved weighted ratio by the number of intersected links



Results:	Word	lavers	over	lan
nesuits.	vvoru	layers	over	ap

		CROATIAN			ENGLISH	
	CO-SIN	CO-SHU	SIN-SHU	CO-SIN	CO-SHU	SIN-SHU
Jaccard	16.72%	5.47%	4.81%	13.44%	6.31%	5.34%
W	18.96%	6.43%	5.63%	13.58%	6.28%	4.82%
WO	90.60%	76.6%	74.6%	90.00%	74.72%	73.81%



10

22



Pearson correlations of motif's motif frequencies and normalized triad



### Recapitulation

• co-occurance: traditional. not sufficient

- •shuffled: reveals interesting behavior, boundary layer
- syntax: more credible for linguistic insights
- syllables: like syntax
- syllables like morphological root
- morphological networks should be constructed
- graphemes: completely different (complex network??)

### **MLN Model**

- language networks can be viewed through different perspects:
- different levels (e.g. word-level, subword-level),
- different construction rules (e.g. co-occurrence, shuffle),
- different languages (e.g Croatian, English)
- there is a need for a general network model that can capture all language aspects in one single framework
- we propose an application of general multilayer networks model introduced by Kivela et al. 2014 to the multilayer language networks (MLN)



### **Open questions?**

test on **other** data sets and languages
comments from linguists

### language network model

explanation of M:N relationships between layers
quantification of language emergent properties

 S. Martinčić-ipšić, D. Margan, A. Meštrović. Multilayer Network of Language: a Unified Framework for Structural Analysis of Linguistic Subsystems, Physica A: Statistical Mechanics and its Applications, 457, 117-128, 2016,

 T. Ban Kirigin, A. Meštrović, S. Martinčić-Ipšić. Towards a Formal Model of Language Networks, ICIST 2015, Communications in Computer and Information Science, Springer, Vol. 538, 469–479, 2015.

24



Department of Informatics, University of Rijeka Radmile Matejčić 2, 51000 Rijeka, Croatia Tel.: + 385 51 584 700 Fax: + 385 51 584 749 www.langnet.uniri.hr

### **Application: Keyword Extraction**



### Introduction 1/2



#### **Keyword extraction**

automatically identify a **set of terms** that best describe the document

• identify and rank the most representative features of the source applications in text:

• summarization, indexing, labeling, categorization, clustering

#### keyword extraction is traditionally supervised

based on statistical methods

• learning from hand-annotated data sets

### Introduction 2/2

#### network-based approach

- or graph, since the number of words in isolated documents is limited
- the source (document, text) is modelled in a network
   words are nodes and their co-occurrence is represented with links
- the keyword extraction can exploit: network, subnetwork or node level measures:
- -coreness, clustering coefficient
- PageRank motivated ranking score or HITS motivated hub and authority score
   communities
- -strength, centrality
- degree, betweenness, closeness and eigenvector centrality
- **selectivity** the average strength of the node

3. Complex network analysis 1/2



#### selectivity

- average weight distribution on the links of the single node

#### generalized selectivity



adapted from [Opsahl et al. (2010)] [Beliga, Meštrović, Martinčić-Ipšić (2016)]



 $\alpha \in [0,1]$  - prefers nodes with higher degree  $\alpha \in [1, +\infty]$  - prefers nodes with lower degrees  $\alpha = 1$  - node strength

### Data

# Langnet

- collections with manually annotated keywords by human experts
- HINA Croatian News Agency

- Wikipedia: technical reports covering different aspects of CS

HINA	WIKI-20
8 human experts	15 teams (2 undergrad. stud.)
60 texts	20 texts
Croatian	English
10 keywords on AVG per human	5,7 keywords on AVG per team
human consistency≈ 39.5%	team consistency $\approx 30.5\%$

### **HINA Dataset**



- learning: 960 documents
- testing: 60 documents keywords manually annotated 8 experts
- per document:
- 60 to 800 tokens (318 on average)

### - 9 to 42 keywords (24 on average)

#### inter-annotator agreement

- average IIC score 39.5% ( 31% 44%)
- no predefined set of keywords list annotators could make up their own
   in some cases annotated with keywords, which were not present in the original article (out-of-vocabulary words - 57%).

30

#### preprocessing:

- parsing only text and title (excluding annotations)
- cleaning diacritics and symbols (w inst. of vv, ! inst. of I, etc.)
- -lematization and NSW types preserved (numbers, acronyms, abbreviations, etc.)

### WIKI-20 dataset

#### 20 technical reports covering different aspects of CS

- 20 documents: keywords manually annotated
- annotators: 15 teams of 2 senior CS undergraduate students independetly
- per document:
- -5 terms to each report
- -controled vocabulary: Wikipedia articles
- per document 5.7 keyphrases on average

#### inter-annotator agreement

-average IIC score 30.5% (21.4% - 37.1%)

#### network construction:

- 20 individual networks per document
- 1 integral network for whole collection





#### directed and weighted co-occurrence networks:

-60 + 20 individual networks: network per document +

-1+1 integral network : from all documents (collection extraction) node: each word

**links:** two words are linked if they are adjacent in the sentence **weights:** proportional to the overall co-occurrence frequencies of the corresponding word pairs

Python + NetworkX

### Comparison of centrality measures

$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$								
DATA	MEASURE	TOP	TOP 5			TOP 10		
SET	WILASURE	Ravg	Pavg	F1 <sub>avg</sub>	Ravg	Pavg	F1 <sub>avg</sub>	
	Closeness: cc <sub>i</sub>	4.3	39.3	7.2	8.4	42.9	12.8	
₹	Betweenness: bci	3.6	40.8	6.5	9.2	52.8	13.9	
Ę	in/out-degree: $dc_i^{in/out}$	23.0	38.4	17.0	62.1	25.1	26.6	
	in/out-selectivity: $e_i^{in/out}$	14.0	35.9	19.3	16.8	38.0	22.1	
0	Closeness: cc <sub>i</sub>	1.6	42.5	3.1	5.3	63.3	9.7	
-2	Betweenness: bci	0.3	10.0	0.5	2.6	55.0	4.9	
WIKI-	in/out-degree: $dc_i^{in/out}$	0.1	5.0	0.3	2.5	57.5	4.8	
	in/out-selectivity: $e_i^{in/out}$	2.3	17.8	3.9	8.6	18.2	10.8	

Selectivity

differentiate between two types of nodes (words)
high strength and high degree values > low selectivity
-closed-class words: stop-words, conjunctions, prepositions

high strength and low degree > high selectivity -open-class words: nouns, adjectives, verbs and -words that are part of collocations, keyphrases, names, etc.

#### selectivity: + +

efficiently detect semantically rich open-class words extract better keyword candidates

# gnet

34

31

### SBKE method



32

 out-selectivity CANDIDATE EXPANSION EXTRACTION K2E: two word-tuples - SET K3E: three words-tuples - SET3 max out-selectivit SETI keyword predecessor 2 candidate keyword candidate SET max in-selectivi K3E: three words-tunies - SET3 K2E: two word-turdes - SET2 CANDIDATE EXPANSION EXTRACTION

in-degree: to be, and, in, on, which, for, but, this, self, of betweenness: to be, and, in, on, self, this, which, for, <u>Croatian</u>, but

in/out selectivity: Bratislava, area, Tuesday, <u>inland</u>, revolution, verification, decade, Balkan, freedom, <u>Universe</u> 33

### **SBKE method**





SBKE method can be easily adjusted by incorporating new measures

### •generalized selectivity measure: $ge_i^{\alpha \ in/out}$

- α adjust the relationship between the node degree and strength provide different set of keywords
- Find the  $\alpha$  parameter which best fits the experimental settings according to  $\mathit{IIC}_{0}$  scores



The performance of the first step (SET1) of the SBKE method in terms of  $IIC_0$  scores for the TOP 5 and TOP 10 keywords measured for generalized selectivity  $ge_1^{e_1(n)out}$  with different values of parameter  $\alpha$  (plotted as lines) and selectivity ( $e^{in/out}$ ) values (plotted as dots on the y-axis) for HINA and WIK1-20 datasets





Evaluation - IIC

#### Consistency between any two annotators

- different evaluation approaches compare the obtained results against the gold standard

### - in case of multiple annotations, gold standard is not clear

#### -Inter-indexer consistency (Rolling, 1981) human and machine, human and human, or machine and machine

numan and machine, numan and numan, or machine and machine
 equivalent to the harmonic mean of R and P - F1 score

$$IIC = \frac{2c}{a+b}$$

 where a and b are the number of terms assigned by each annotator, and c is the number of terms they have in common

$$IIC_{overall} = \frac{\sum_{i=1}^{N_{doc}} \sum_{j=1}^{N_{an}} IIC_{ij}}{N_{doc} N_{an}}$$

Langnet

37

### **HINA Results**

A	
Langnet	

Annotator	IIC <sub>o</sub>			SBKE		lic
1	31.9%				е	geα, α=5.5
2	33.5%		S	SET1	20.5%	19.6%
3	39.3%		Р	SET2	22.1%	21.1%
4	40.2%	P	SET3	22.4%	21.9%	
5	41.4%				е	ge <sup>α</sup> , α=1.0
6	41.6%		0	SET1	22.5%	26.1%
7	43.9%		DP 1	SET2	22.7%	19.8%
8	44.4%		P	SET3	23.0%	19.7%
AVERAGE	39.5%					

Annotator	IIC <sub>o</sub>		SBKE	IIC <sub>o</sub>	
1	21.4%			e	ge <sup>α</sup> , α=2.5
2	24.1%		SET1	3.5%	6.1%
3	26.2%	TOP	SET2	7.5%	10.8%
4	28.7%	Ĕ	SET3	10.3%	12.0%
5	30.2%	0	SET1	6.9%	13.6%
6	30.8%	DP 1	SET2	11.4%	16.3%
7	31.0%	μ¥	SET3	12.5%	17.2%
8	31.2%				
9	31.6%				
10	31.6%				
11	31.6%				
12	32.4%				
13	33.8%				
14	35.5%				
15	37.1%				
AVERAGE	30.5%				

42

langnet

45

### **Document vs. Collection-Oriented Extraction Results**

	60 indiv	vidual net	tworks	integral network				
піла	R <sub>avg</sub>	Pavg	F1 <sub>avg</sub>	R	Р	F1		
SET1	18.90	39.35	23.60	30.71	35.80	33.06		
SET2	19.74	39.15	24.76	33.46	33.97	33.71		
SET3	29.55	22.96	22.71	60.47	19.89	28.89		

WIKI- 20	20 n (e	individu etworks ei <sup>in/out</sup> >1	ial 5 )	inte (	gral network e <sub>i</sub> <sup>in/out</sup> >1)		integral network ( e <sub>i</sub> <sup>in/out</sup> >2)		vork !)
	Ravg	Pavg	F1 <sub>avg</sub>	R	Р	F1	R	Р	F1
SET1	59.06	13.40	21.48	76.17	19.08	30.51	32.71	30.30	31.46
SET2	60.12	12.33	20.46	76.64	18.87	30.29	36.45	31.97	34.06
SET3	62.17	12.05	20.19	77.50	15.75	26.18	36.68	32.04	34.21

43





Comparison of the performance of the SBKE method in terms of F1 and  $IIC_0$  score with other approaches

			WIKI-20		
Method	Approach	F1	Method	Approach	
tf-idf (Ahel et al., 2009)	unsupervised	13.2%	tf-idf (Medelyan, 2009)	unsupervised	
MDL+POS (Ahel et al., 2009)	supervised	17.2%	SBKE, α=2.5	graph-based	
SBKE - SET2	graph-based	24.8%	KEA++ (Medelyan, 2009)	supervised	
			Humans	gold standard	
Method	Approach	IIC <sub>o</sub>	Maui (Wang et al., 2014)	supervised	
SBKE - SET1, α=1.0	graph-based	26.1%	GA (Joorabchi et al., 2013)	supervised	
Humans	gold standard	39.5%	Maui+ (Wang et al., 2014)	supervised	
*for TOP 10 KE task			*for TOP 5 KE task		

### **SBKE Conclusion**

- Selectivity-Based Keyword Extraction SBKE method purely from statistical and structural information
- the source text is reflected into the structure of the network low precision vs. high recall
- beside keywords, personal names and entities not marked as keywords - comparable with existing supervised and unsupervised method
- better to longer than to shorter texts collection extraction task
- Keyword annotation is highly subjective task even human experts have difficulties to agree upon keyphrases human annotators F2 score of inter-annotator agreement is 46% (SBKE 42%)
- Other graph-based (centrality) approaches
- similar results but they incorporate linguistic knowledge in a form of different syntactic filters
- POS tagging, stop-words filtering, noun-phrase parsing, etc. - generally more demanding to implement



 S. Beliga, A. Meštrović, S. Martinčić-Ipšić." Selectivity-Based Keyword Extraction Method", <u>International Journal on Semantic Web and Information Systems</u>, vol. 12, No. 3, pp. 1-26, 2016, doi: 10.4018/IJSWIS.2016070101.

S. Beliga, A. Meštrović, S. Martinčić-Ipšić. "An Overview of Graph-Based Keyword Extraction Methods and Approaches". *Journal of Information and Organizational* <u>Sciences</u>, vol. 39, No 1, pages 1-20, 2015.

• S. Beliga, A. Meštrović, S. Martinčić-Ipšić. "Toward Selectivity-Based Keyword Extraction for Croatian News". CEUR Proceedings of the Workshop on Surfacing the Deep and the Social Web (SDSW 2014), Vol. 1310, pp. 1-8, Riva del Garda, Trentino, Italy, 2014.

### **Other applications**

- Text genres differentiation
- · Can we determine text genre using only network structure properties? . Legislation or literature? Blogs vs. literature?
- Authorship attribution and language detection •
- Text classification
- dimensionality reduction in BoW model using structural properties of text
- Wikipedia
- knowledge extraction and linking concepts .

#### Twitter

- polarity (positive/negative) detection prediction of missing links
- Social networks
- •
- Coautorship networks analyis .
- scientometrics



50







#### universal language model

- include many languages
- add perspectives /layers: semantics, cognitive representation

#### • text quality evaluation:

- derive an assessment model for the evaluation of the quality of texts from complex networks parameters
- creativity?
- keywords extraction, summarization

### **Generalized selectivity**

Additional observations

- selectivity takes into account weights which represent the frequencies of word bigrams
- captures the importance of the node strength which is crutial for weighted networks
- degree centrality also yields optimistic results
- combine selectivity + degree = generalized selectivity



 $\alpha \in [1, +\infty]$  – prefers nodes with  $\alpha = 1$ - node strength

51