

**Prof. Dr. Andreas Hotho,**  
**M.Sc. Janna Omeliyanenko**  
Lecture Chair X for Data Science, Universität Würzburg

## 12. Exercise for “Sprachverarbeitung und Text Mining”

11.02.2022

### 1 Knowledge Questions

1. Describe the features and structure of WordNet in your own words.

WordNet is a lexical database for the English language. WordNet groups English words according to their meaning into synonym groups called synsets. WordNet provides definitions and usage examples for synsets and defines a set of semantic relationships between synsets.

Describe and give an example for WordNets synsets.

A synset is a set of all synonyms that represent a particular meaning or concept.  
Example: *Dog*

- Synset1 (First meaning): dog, domestic dog, Canis familiaris
  - Definition: A member of the genus Canis (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds)
  - Example: "the dog barked all night"
- Synset2 (Second meaning): andiron, firedog, dog, dog-iron
  - Definition: metal supports for logs in a fireplace
  - Example: "the andirons were too hot to touch"

2. What are potential problems of standard similarity metrics that are based on path length?

These similarity metrics do not take into account information about:

- Depth of a concept in the hierarchy
- network density
- Meaning of the word

Name one possible way to fix these problems.

Use of hierarchy information (e.g. from WordNet) and corpus statistics (e.g. from Brown Corpus) → Information Content

3. Define the term *Information Content* in your own words.

The Information Content of an event indicates how "surprising" the occurrence of the event is, or how much information (in bits) is gained through the occurrence of the event.

4. Explain why the similarity metrics from the lecture that use Information Content use the Information Content of the lowest common subsumer of the words being compared in the concept hierarchy.

The lower a concept is in the concept hierarchy, the more specific and informative it is (higher Information Content). The words that share a common subsumer low in the hierarchy are more similar to each other than those that share only more abstract subsumers.

## 2 Word Similarity

### 2.1 Similarity measures

1. For the following ontology, compute the similarities of the following concepts:

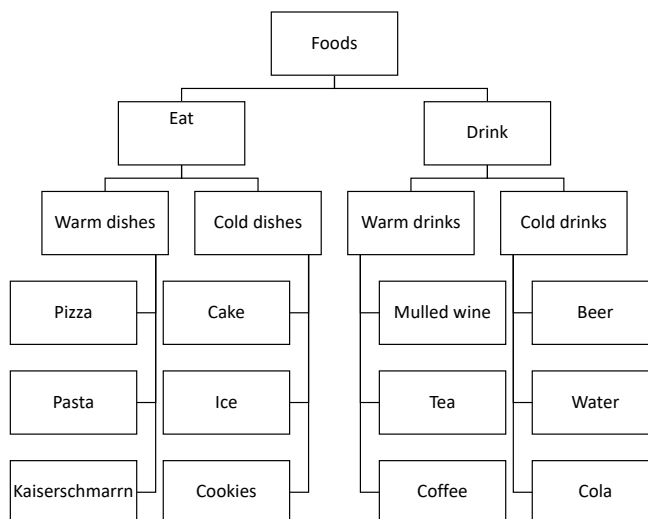
- $\text{sim}(\text{Pizza}, \text{Pasta})$
- $\text{sim}(\text{Pizza}, \text{Beer})$
- $\text{sim}(\text{Cola}, \text{Coffee})$
- $\text{sim}(\text{Cookies}, \text{Beer})$

Use the measures introduced in the lecture for this purpose:

- a)  $\text{sim}_{\text{path}}(c_1, c_2)$   
b)  $\text{sim}_{\text{jiang-conrath}}(c_1, c_2)$

Use the following text for your calculations, if needed:

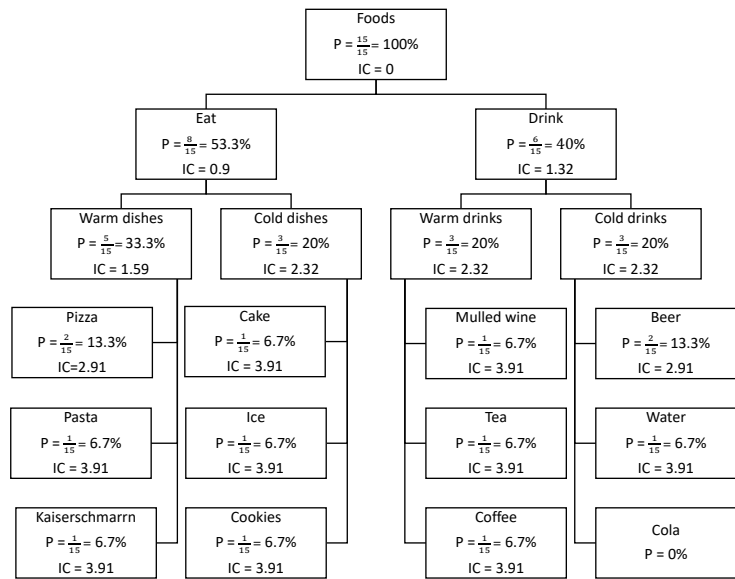
My favorite foods are pizza and pasta. I have either beer or tea with it. Warm dishes are part of every meal for me. With cake I like to have coffee and cookies. The water is already frozen to ice when I have mulled wine with Kaiserschmarrn at the Christmas market. I would rather have pizza and beer, though.



We define the probability of a word from the taxonomy as:

$$P(c) = \frac{\sum_{w \in \text{words}(c)} \text{count}(w)}{N}$$

where  $\text{words}(c)$  contains all sub-concepts of the taxonomy from concept  $c$  including concept  $c$  itself.  $N$  is the number of words in the text that occur as a concept in the taxonomy.



### Path-Similarity

- $\text{sim}(\text{Pizza}, \text{Pasta}) = \frac{1}{\text{pathlen}(c_1, c_2)} = \frac{1}{3}$
- $\text{sim}(\text{Pizza}, \text{Beer}) = \frac{1}{7}$
- $\text{sim}(\text{Cola}, \text{Coffee}) = \frac{1}{5}$
- $\text{sim}(\text{Cookies}, \text{Bier}) = \frac{1}{7}$

### Jiang-Conrath

- $\text{sim}(\text{Pizza}, \text{Pasta}) = \frac{1}{2 \log_2 P(\text{LCS}(c_1, c_2)) - \log_2 P(c_1) - \log_2 P(c_2)} = 0.274$

- $sim(\text{Pizza}, \text{Beer}) = 0.172$
- $sim(\text{Cola}, \text{Coffee}) = \text{Error}$
- $sim(\text{Cookies}, \text{Beer}) = 0.147$

2. Another similarity metric is the weighted path length (wpath) measure, which is defined as follows:

$$sim_{wpath}(c_1, c_2) = \frac{1}{1 + pathlen'(c_1, c_2) * k^{IC(LCS(c_1, c_2))}}$$

$pathlen'(c_1, c_2)$  = Number of edges in graph between nodes  $c_1$  and  $c_2$

- Describe the intuition behind wpath in your own words. What types of information that we introduced previously are used in wpath? What is the role of wpath's parameter  $k$ ?

wpath combines information about shortest path length and the Information Content of the deepest common subsumer of two concepts. The lower a node is in the hierarchy, the more specific the concept it describes (and the higher the Information Content).

Two concepts are more similar if they can be described by a more specific super-concept. Additionally, the shorter the shortest path between two concepts, the more similar they are.

The parameter  $k$  specifies how much the Information Content contributes to the weighting of the path length. If  $k = 1$ , wpath is identical to the path metric.

- Calculate the wpath similarity for the concept pairs from subtask 1.

### Wpath-Similarity

- $sim_{wpath}(\text{Pizza}, \text{Pasta}) = \frac{1}{1 + pathlen'(c_1, c_2) * 0.8^{1.59}} = \frac{1}{1 + 2 * 0.8^{1.59}} = 0,416$
- $sim_{wpath}(\text{Pizza}, \text{Beer}) = 0.143$
- $sim_{wpath}(\text{Cola}, \text{Coffee}) = 0.251$
- $sim_{wpath}(\text{Cookies}, \text{Beer}) = 0.143$

### 3 String Kernel

Calculate the String-Kernel of the words *Text* and *Test*. A completely simplified term depending on  $\lambda$  is sufficient as a result.

Text	Test	○
T	T	$\lambda^2$
e	e	$\lambda^2$
x	s	0
t	t	$\lambda^2$
Te	Te	$\lambda^4$
ex	es	0
xt	st	0
Tx	Ts	0
et	et	$\lambda^6$
Tt	Tt	$\lambda^8$
Tex	Tes	0
ext	est	0
Tet	Tet	$\lambda^8$
Txt	Tst	0

$$\text{Kernel}(\text{Text}, \text{Test}) = 2\lambda^8 + \lambda^6 + \lambda^4 + 3\lambda^2$$

### 4 Minimum Edit Distance

Determine the minimum edit distance between the following words. We define the costs as triples  $w = (w_1, w_2, w_3)$ , where  $w_1$  is the cost for Insertion,  $w_2$  the cost for Deletion, and  $w_3$  the cost for Replace.

- Initial-state: datamining Goal-state: textmining with  $w = (1, 1, 2)$
- Initial-state: gehend Goal-state: verstehen with  $w = (3, 7, 9)$

I = Insert, D = Delete, K = Keep (replace for  $X[i] = Y[j]$ ),  
R = Replace (replace for  $X[i] \neq Y[j]$ )

		t	e	x	t	m	i	n	i	n	g
	0	1	2	3	4	5	6	7	8	9	10
<b>d</b>	<sup>D</sup> 1	2	3	4	5	6	7	8	9	10	11
<b>a</b>	<sup>D</sup> 2	3	4	5	6	7	8	9	10	11	12
<b>t</b>	3	<sup>K</sup> 2	<sup>I</sup> 3	<sup>I</sup> 4	5	6	7	8	9	10	11
<b>a</b>	4	3	4	5	<sup>R</sup> 6	7	8	9	10	11	12
<b>m</b>	5	4	5	6	7	<sup>K</sup> 6	7	8	9	10	11
<b>i</b>	6	5	6	7	8	7	<sup>K</sup> 6	7	8	9	10
<b>n</b>	7	6	7	8	9	8	7	<sup>K</sup> 6	7	8	9
<b>i</b>	8	7	8	9	10	9	8	7	<sup>K</sup> 6	7	8
<b>n</b>	9	8	9	10	11	10	9	8	7	<sup>K</sup> 6	7
<b>g</b>	10	9	10	11	12	11	10	9	8	7	<sup>K</sup> 6
		v	e	r	s	t	e	h	e	n	
	0	3	6	9	12	15	18	21	24	27	
<b>g</b>	7	<sup>R</sup> 9	12	15	18	21	24	27	30	33	
<b>e</b>	14	16	<sup>K</sup> 9	<sup>I</sup> 12	<sup>I</sup> 15	<sup>I</sup> 18	<sup>I</sup> 21	24	27	30	
<b>h</b>	21	23	16	18	21	24	27	<sup>K</sup> 21	24	27	
<b>e</b>	28	30	23	25	27	30	24	27	<sup>K</sup> 21	24	
<b>n</b>	35	37	30	32	34	36	31	33	28	<sup>K</sup> 21	
<b>d</b>	42	44	37	39	41	43	38	40	35	<sup>D</sup> 28	