



Text Analytics

Pilsung Kang

School of Industrial Management Engineering

Korea University

Text Analytics

- What we have done so far...

Cornell University Library

We gratefully acknowledge support from the Simons Foundation and member institutions

arXiv.org > search

Search or Article ID All papers

(Help | Advanced search)

arXiv.org Search Results

Back to Search form | Next 25 results

The URL for this search is http://arxiv.org:443/find/all/1/all:+EXACT+text_mining/0/1/0/all/0/1

Showing results 1 through 25 (of 168 total) for all:"text mining"

1. arXiv:1703.05692 [pdf, other]
OncoScore: a novel, Internet-based tool to assess the oncogenic potential of genes
Rocco Piazza, Daniele Ramazzotti, Roberta Spinelli, Alessandra Pirola, Luca De Sano, Pierangelo Ferrari, Vera Magistroni, Nicoletta Cordani, Nitesh Sharma, Carlo Gambacorti-Passerni
Subjects: Genomics (q-bio.GN); Quantitative Methods (q-bio.QM)

2. arXiv:1703.04213 [pdf, other]
MetaPAD: Meta Pattern Discovery from Massive Text Corpora
Meng Jiang, Jingbo Shang, Taylor Cassidy, Xiang Ren, Lance M. Kaplan, Timothy P. Hanratty, Jiawei Han
Comments: 9 pages
Subjects: Computation and Language (cs.CL)

3. arXiv:1703.02819 [pdf, other]
Introduction to Formal Concept Analysis and Its Applications in Information Retrieval and Related Fields
Dmitry I. Ignatov
Journal-ref: RussIR 2014, Nizhniy Novgorod, Russia, CCIS vol. 505, Springer 42-141
Subjects: Information Retrieval (cs.IR); Artificial Intelligence (cs.AI); Computation and Language (cs.CL); Discrete Mathematics (cs.DM); Machine Learning (stat.ML)

4. arXiv:1702.07117 [pdf, other]
LTSG: Latent Topical Skip-Gram for Mutually Learning Topic Model and Vector Representations
Jianan Law, Hankui Zhuo, Junhua He, Erhu Rong (Dept. of Computer Science, Sun Yat-Sen University, GuangZhou, China.)
Subjects: Computation and Language (cs.CL)

5. arXiv:1702.03519 [pdf, ps, other]
A Technical Report: Entity Extraction using Both Character-based and Token-based Similarity
Zeyi Wen, Dong Deng, Rui Zhang, Kotagiri Ramamohanarao
Comments: 12 pages, 6 figures, technical report
Subjects: Databases (cs.DB)



The complicated, evolving landscape of cancer

Mining textual patterns in news, tweets, papers, and many more... This paper is a tutorial on Formal Concept Analysis (FCA) and its applications. FCA is an applied branch of Lattice Theory, a mathematical discipline which enables formalisation of concepts as basic units of human thinking and analysing data in the object-attribute form. Originated in early 80s, during the last three decades, it became a popular human-centred tool for knowledge representation and data analysis with numerous applications. Since the tutorial was specially prepared for RuSSIR 2014, the covered FCA topics include Information Retrieval with a focus on visualisation aspects, Machine Learning, Data Mining and Knowledge Discovery, Text Mining and several others.

pattern quality assessment function, which avoids costly dependency parsing and generates high-quality patterns; (2) it identifies and groups synonymous meta patterns from multiple facets---their types, contexts, and extractions; and (3) it examines type distributions of entities in the instances extracted by each group of patterns, and looks for appropriate type levels to make discovered patterns precise. Experiments demonstrate that our proposed framework discovers high-quality typed textual patterns efficiently from different genres of massive corpora and facilitates information extraction.

Text Analytics

- Preprocessing with some NLP techniques

The complicated, evolving landscape of cancer

Mining textual patterns in news, tweets, papers, and

This paper is a tutorial on Formal Concept Analysis (FCA) and its applications. FCA is an applied branch of Lattice Theory, a mathematical discipline which enables formalisation of concepts as basic units of human thinking and analysing data in the object-attribute form. Originated in early 80s, during the last three decades, it became a popular human-centred tool for knowledge representation and data analysis with numerous applications. Since the tutorial was specially prepared for RuSSIR 2014, the covered FCA topics include Information Retrieval with a focus on visualisation aspects, Machine Learning, Data Mining and Knowledge Discovery, Text Mining and several others.

pattern quality assessment function, which avoids costly dependency parsing and generates high-quality patterns; (2) it identifies and groups synonymous meta patterns from multiple facets---their types, contexts, and extractions; and (3) it examines type distributions of entities in the instances extracted by each group of patterns, and looks for appropriate type levels to make discovered patterns precise. Experiments demonstrate that our proposed framework discovers high-quality typed textual patterns efficiently from different genres of massive corpora and facilitates information extraction.

the complic evolv landscap of cancer mutat pose a

form mine textual pattern in news tweet paper and mani

list oth this paper is a tutori on formal concept analysi fca and
priori tex it applic fca is an appli branch of lattic theori a
mathemat disciplin which enabl formalis of concept as
basic unit of human think and analys data in the
objectattribut form origin in earli s dure the last three
decad it becam a popular humancentr tool for
knowledg represent and data analysi with numer applic
sinc the tutori was special prepar for russir the cover
fca topic includ inform retriev with a focus on visualis
aspect machin learn data mine and knowledg discoveri
text mine and sever other

Text Analytics

- What we will do...

- ✓ Transform unstructured data into structured data

the complic evolv landscap of cancer mutat pose a
form mine textual pattern in news tweet paper and mani
list oth
prior tex
to a dep
too the
base on
curv and
curat da
oncosco
oncosco
priorit o

this paper is a tutori on formal concept analysi fca and
it applic fca is an appli branch of lattic theori a
mathemat disciplin which enabl formalis of concept as
basic unit of human think and analys data in the
objectattribut form origin in earli s dure the last three
decad it becam a popular humancentr tool for
knowledg represent and data analysi with numer applic
sinc the tutori was special prepar for russir the cover
fca topic includ inform retriev with a focus on visualis
aspect machin learn data mine and knowledg discoveri
text mine and sever other



	Var 1	Var 2	Var P
Doc 1					
Doc 2					
Doc 3					
...					
...					
...					
Doc D					

Text Analytics

- Install and load necessary packages

```
install.packages("textstem")
install.packages("tm")
install.packages("wordcloud")
install.packages("plyr")
install.packages("igraph")
install.packages("treemap")

library(textstem)
library(tm)
library(wordcloud)
library(plyr)
library(igraph)
library(dplyr)
library(treemap)
```

Text Analytics

- Load two datasets
 - ✓ Dataset 1: The abstracts of 500 arXiv papers with the search term "natural language processing"
 - ✓ Dataset 2: The abstracts of 500 arXiv papers with the search term "image processing"

```
load("arxiv_NLP.RData")
load("arxiv_Image_Processing.RData")

nlp_papers <- paste(arxiv_nlp$title, arxiv_nlp$abstract)
nlp_papers[1]

image_papers <- paste(arxiv_image_processing$title,
arxiv_image_processing$abstract)
image_papers[1]

# Remove invalid characters
nlp_papers <- sapply(nlp_papers,function(row) iconv(row, "latin1", "ASCII",
sub=""))
image_papers <- sapply(image_papers,function(row) iconv(row, "latin1",
"ASCII", sub=""))
```

Text Analytics I: Tokenization

- Goals of lexical analysis
 - ✓ Convert a sequence of characters into a sequence of **tokens**, i.e., meaningful character strings.
 - In natural language processing, **morpheme** is a basic unit
 - In text mining, **word** is commonly used as a basic unit for analysis
- Process of lexical analysis
 - ✓ Tokenizing
 - ✓ Part-of-Speech (POS) tagging
 - ✓ Additional analysis: named entity recognition (NER), noun phrase recognition, sentence split, chunking, etc.

Text Analytics I: Tokenization

- Text is split into basic units called Tokens

- ✓ word tokens, number tokens, space tokens, ...

Diagram illustrating the tokenization process:

```

> crude[[1]]
<<PlainTextDocument (metadata: 15)>>
Diamond Shamrock Corp said that
effective today it had cut its contract prices for crude oil by
1.50 dlr$ a barrel.
The reduction brings its posted price for West Texas
Intermediate to 16.00 dlr$ a barrel, the company said.
"The price reduction today was made in the light of falling
oil product prices and a weak crude oil market," a company
spokeswoman said.
Diamond is the latest in a line of U.S. oil companies that
have cut its contract, or posted, prices over the last two days
citing weak oil markets.
Reuter
  
```

The text above is shown with two parallel tokenization results below it, separated by a horizontal line.

	MC_tokenizer(crude[[1]])	scan_tokenizer(crude[[1]])		
1	"diamond"	"Diamond"		
2	"shamrock"	"Shamrock"		
3	"Corp"	"Corp"		
4	"said"	"said"		
5	"that"	"that"		
6	"effective"	"effective"		
7	"today"	"today"		
8	"it"	"it"		
9	"had"	"had"		
10	"cut"	"cut"		
11	"its"	"its"		
12	"contract"	"contract"		
13	"prices"	"prices"		
14	"for"	"for"		
15	"crude"	"crude"		
16	""	"1.50"		
17	"dlrs"	"dlrs"		
18	"a"	"by"		
19	"barrel"	"reduction"		
20	""	"brings"		
21	"dlrs"	"for"		
22	"a"	"West"		
23	"barrel"	"Texas"		
24	""	"dlrs"		
25	"reduction"	"a"		
26	"brings"	"\The"		
27	"its"	"copany"		
28	"posted"	"said."		
29	"west"	"made"		
30	"Texas"	"in"		
31	"Intermediate"	"the"		
32	""	"light"		
33	""	"of"		
34	"dlrs"	"oil"		
35	"a"	"product"		
36	"barrel"	"prices"		
37	"copy"	"and"		
38	"said"	"crude"		
39	""	"oil"		
40	""	"market,\\"		
41	"The"	"company"	"diamond"	"is"
42	"price"	"spokeswoman"	"said."	"line"
43	"for"	"latest"	"in"	"a"
44	""	"in"	"a"	"that"
45	"to"	"U.S."	"oil"	"contract,"
46	""	"oil"	"companies"	"or"
47	""	"cut"	"its"	"the"
48	"the"	"over"	"contract,"	"last"
49	"the"	"days"	"over"	"the"
50	"weak"	"citing"	"the"	"last"
51	"spokeswoman"	"days"	"citing"	"weak"
52	"company"	"markets."	"Reuter"	"oil"
53	"spokeswoman"			
54	"latest"			
55	"in"			
56	"of"			
57	"U"			
58	"companies"			
59	"that"			
60	"contract"			
61	"prices"			
62	"over"			
63	"days"			
64	"citing"			
65	""			
66	"markets."			
67	"Reuter"			
68				
69				
70				
71				
72				
73				
74				
75				
76				
77				
78				
79				
80				
81				
82				
83				
84				
85				
86				
87				
88				
89				
90				
91				
92				
93				
94				
95				
96				
97				
98				
99				
100				
101				
102				
103				
104				
105				
106				
107				
108				
109				
110				
111				
112				
113				
114				
115				
116				
117				
118				
119				
120				
121				
122				
123				
124				
125				
126				
127				
128				
129				
130				
131				
132				
133				
134				
135				
136				
137				
138				
139				
140				
141				
142				
143				
144				
145				
146				
147				
148				
149				
150				
151				
152				
153				
154				
155				
156				
157				
158				
159				
160				
161				
162				
163				
164				
165				
166				
167				
168				
169				
170				
171				
172				
173				
174				
175				
176				
177				
178				
179				
180				
181				
182				
183				
184				
185				
186				
187				
188				
189				
190				
191				
192				
193				
194				
195				
196				
197				
198				
199				
200				
201				
202				
203				
204				
205				
206				
207				
208				
209				
210				
211				
212				
213				
214				
215				
216				
217				
218				
219				
220				
221				
222				
223				
224				
225				
226				
227				
228				
229				
230				
231				
232				
233				
234				
235				
236				
237				
238				
239				
240				
241				
242				
243				
244				
245				
246				
247				
248				
249				
250				
251				
252				
253				
254				
255				
256				
257				
258				
259				
260				
261				
262				
263				
264				
265				
266				
267				
268				
269				
270				
271				
272				
273				
274				
275				
276				
277				
278				
279				
280				
281				
282				
283				
284				
285				
286				
287				
288				
289				
290				
291				
292				
293				
294				
295				
296				
297				
298				
299				
300				
301				
302				
303				
304				
305				
306				
307				
308				
309				
310				
311				
312				
313				
314				
315				
316				
317				
318				
319				
320				
321				
322				
323				
324				
325				
326				
327				
328				
329				
330				
331				
332				
333				
334				
335				
336				
337				
338				
339				
340				
341				
342				
343				
344				
345				
346				
347				
348				
349				
350				
351				
352				
353				
354				
355				
356				
357				
358				
359				
360				
361				
362				
363				
364				
365				
366				
367				
368				
369				
370				
371				
372				
373				
374				
375				
376				
377				
378				
379				
380				
381				
382				
383				
384				
385				
386				
387				
388				
389				
390				
391				
392				
393				
394				
395				
396				
397				
398				
399				
400				
401				
402				
403				
404				
405				
406				
407				
408				
409				
410				
411				
412				

Text Analytics I: Tokenization

- Even tokenization can be difficult

✓ Is John's sick one token or two?

- If one → problems in parsing (where is the verb?)
- If two → what do we do with John's house?

✓ What to do with hyphens?

- database vs. data-base vs. data base

✓ What to do with “C++”, “A/C”, “:-)”, “...”, “ㅋㅋㅋㅋㅋㅋㅋㅋ”?

✓ Some languages do not use whitespace (e.g., Chinese)

2013年5月，习主席在视察成都战区时，郑重提出在适当时候召开全军政治工作会议，并明确提出到古田召开这次会议，以更好弘扬我党我军的光荣传统和优良作风。6月，总政治部向中央军委提交《关于筹备召开全军政治工作会议的请示》，提出要通过召开会议形成一个指导性文件。习主席随即批示同意，明确要求这个文件要充分体现深厚的历史积淀和政治意蕴，能够管一个时期，起到历史性作用。

- Consistent tokenization is important for all later processing steps.

Text Analytics 2: Morphological Analysis

- Morphological Variants: Stemming and Lemmatization

Morphological Variants

Words are changed through a morphological process called *inflection*:

- typically indicates changes in case, gender, number, tense, etc.
- example *car* → *cars*, *give* → *gives*, *gave*, *given*

Goal: “normalize” words

Stemming and Lemmatization

Two main approaches to normalization:

Stemming reduce words to a *base form*

Lemmatization reduce words to their *lemma*

Main difference: stemming just finds **any** base form, which doesn't even need to be a word in the language! Lemmatization find the actual *root* of a word, but requires morphological analysis.

Text Analytics 2: Morphological Analysis

- Stemming

Stemming

Commonly used in Information Retrieval:

- Can be achieved with rule-based algorithms, usually based on suffix-stripping
- Standard algorithm for English: the *Porter* stemmer
- Advantages: simple & fast
- Disadvantages:
 - Rules are language-dependent
 - Can create words that do not exist in the language, e.g., *computers* → *comput*
 - Often reduces different words to the same stem, e.g., *army, arm* → *arm*
 - *stocks, stockings* → *stock*
- Stemming for German: German stemmer in the full-text search engine *Lucene*, *Snowball* stemmer with German rule file

Text Analytics 2: Morphological Analysis

- Lemmatization

Lemmatization

Lemmatization is the process of deriving the base form, or *lemma*, of a word from one of its inflected forms. This requires a morphological analysis, which in turn typically requires a *lexicon*.

- Advantages:
 - identifies the *lemma* (root form), which is an actual word
 - less errors than in stemming
- Disadvantages:
 - more complex than stemming, slower
 - requires additional language-dependent resources
- While stemming is good enough for Information Retrieval, Text Mining often requires lemmatization
 - Semantics is more important (we need to distinguish an *army* and an *arm!*)
 - Errors in low-level components can multiply when running downstream

Text Analytics 2: Morphological Analysis

- Stemming vs. Lemmatization

Word	Stemming	Lemmatization
Love	Lov	Love
Loves	Lov	Love
Loved	Lov	Love
Loving	Lov	Love
Innovation	Innovat	Innovation
Innovations	Innovat	Innovation
Innovate	Innovat	Innovate
Innovates	Innovat	Innovate
Innovative	Innovat	Innovative

Text Analytics 2: Morphological Analysis

- Lemmatization

```
nlp_lemma <- NULL  
image_lemma <- NULL  
  
# Lemmatization  
for (i in 1:length(nlp_papers)){  
  cat("Lemmatizing", i, "-th paper. \n")  
  nlp_lemma[i] <- lemmatize_strings(nlp_papers[i])  
}  
  
for (j in 1:length(image_papers)){  
  cat("Lemmatizing", j, "-th paper. \n")  
  image_lemma[j] <- lemmatize_strings(image_papers[j])  
}
```

Text Analytics 2: Morphological Analysis

- Before lemmatization

```
> nlp_papers[1]
```

```
[1] "Towards non-toxic landscapes: Automatic toxic comment detection using DNN The spectacular expansion of the Internet led to the development of a new research problem in the natural language processing field: automatic toxic comment detection, since many countries prohibit hate speech in public media. There is no clear and formal definition of hate, offensive, toxic and abusive speeches. In this article, we put all these terms under the \"umbrella\" of toxic speech. The contribution of this paper is the design of binary classification and regression-based approaches aiming to predict whether a comment is toxic or not. We compare different unsupervised word representations and different DNN classifiers. Moreover, we study the robustness of the proposed approaches to adversarial attacks by adding one (healthy or toxic) word. We evaluate the proposed methodology on the English Wikipedia Detox corpus. Our experiments show that using BERT fine-tuning outperforms feature-based BERT, Mikolov's word embedding or fastText representations with different DNN classifiers."
```

- After lemmatization

```
> nlp_lemma[1]
```

```
[1] "Towards non - toxic landscape: Automatic toxic comment detection use DNN The spectacular expansion of the Internet lead to the development of a new research problem in the natural language process field: automatic toxic comment detection, since many country prohibit hate speech in public medium. There be no clear and formal definition of hate, offensive, toxic and abusive speech. In this articl e, we put all this term under the \" umbrella \" of toxic speech. The contribution of this paper be t he design of binary classification and regression - base approach aim to predict whether a comment be toxic or not. We compare different unsupervised word representation and different DNN classifier. Mo reover, we study the robustness of the propose approach to adversarial attack by add one ( healthy or toxic ) word. We evaluate the propose methodology on the English Wikipedia Detox corpus. Our experim ent show that use BERT fine - tune outperform feature - base BERT, Mikolov's word embed or fastText r epresentation with different DNN classifier."
```

Text Representation: Bag-of-Words

- Document Representation

- ✓ How to represent a document in a structured way?
- ✓ How to convert a unstructured text into a matrix form to apply those machine learning algorithms based on a vector space?

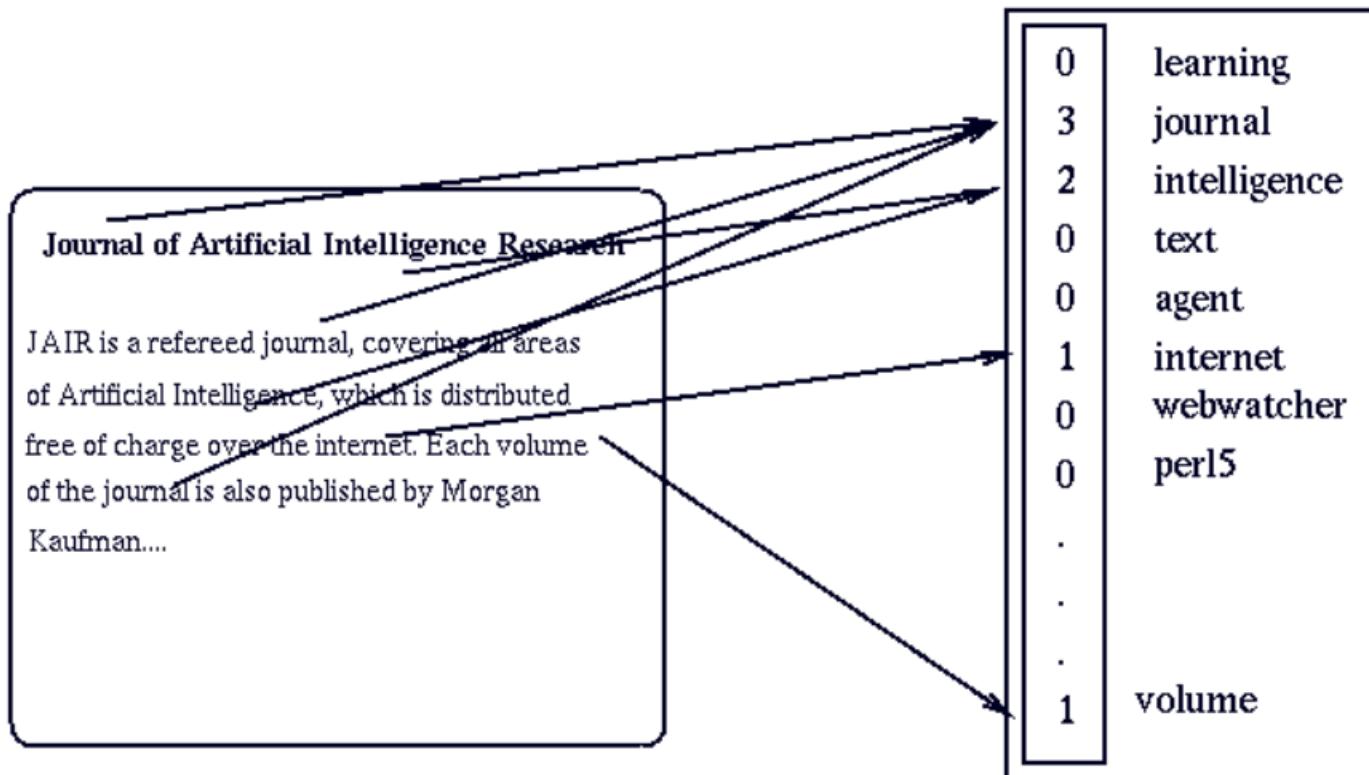
Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach the brain from the eyes. For a long time it was believed that the retinal image was sent directly from the eye to the cerebral cortex. In 1960, David Hubel and Torsten Wiesel, working at the Massachusetts Institute of Technology, demonstrated that the visual system is much more complex than previously thought. They found that the visual cortex contains many layers of nerve cells, each with a specific function. The visual information is processed through a series of stages, starting with the retina and ending in the cerebral cortex. The process involves several types of nerve cells, including bipolar cells, ganglion cells, and interneurons. The visual system is also influenced by other sensory inputs, such as touch and sound, which can affect the way we perceive visual information.

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% jump in exports compared with a 18% rise in imports. The figure is likely to annoy the US, which has been pressuring China to allow its currency, the yuan, to appreciate. The US has threatened to impose trade sanctions if China does not allow the yuan to rise. The Chinese government has agreed to allow the yuan to appreciate, but it has not yet done so. The Chinese government has also announced that it will not allow the yuan to rise further in value.

Text Representation: Bag-of-Words

- Bag-of-words

- ✓ Simplifying representation method for documents where a text is represented in a vector of an unordered collection of words



Text Representation: Bag-of-Words

- **Bag-of-words: Term-Document Matrix**

- ✓ Simplifying representation method for documents where a text is represented in a vector of an unordered collection of words

S₁: John likes to watch movies. Mary likes too.

S₂: John also likes to watch football game.

Binary representation

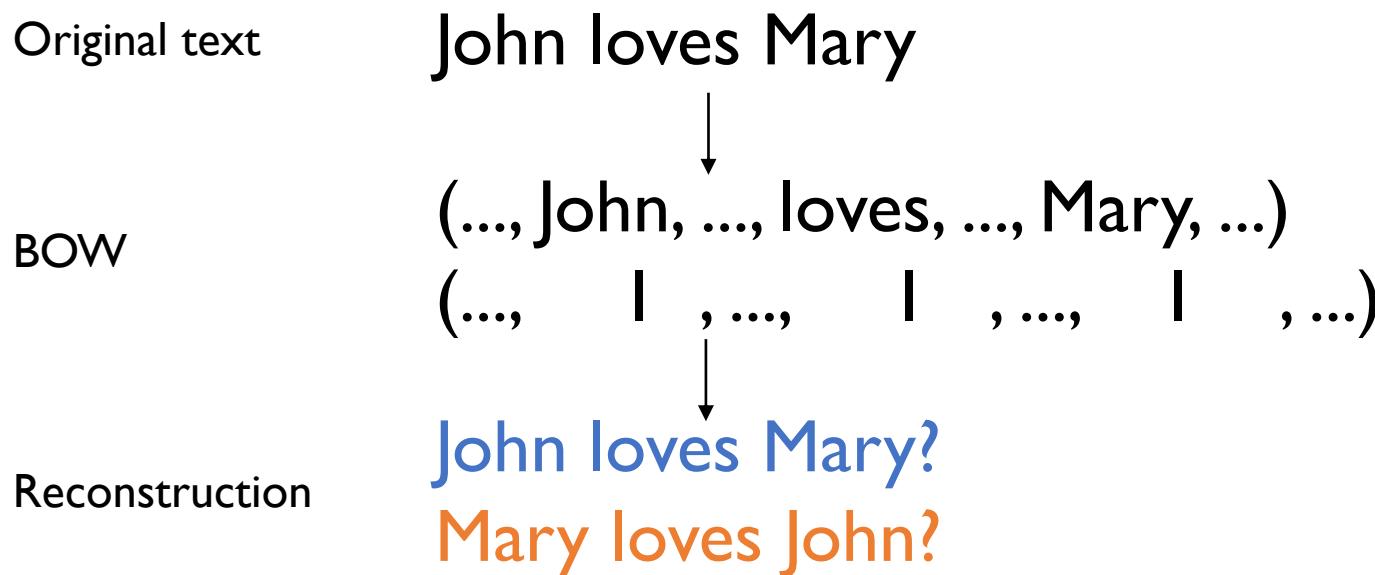
Word	S ₁	S ₂
John	1	1
Likes	1	1
To	1	1
Watch	1	1
Movies	1	0
Also	0	1
Football	0	1
Games	0	1
Mary	1	0
too	1	0

Frequency representation

Word	S ₁	S ₂
John	1	1
Likes	2	1
To	1	1
Watch	1	1
Movies	1	0
Also	0	1
Football	0	1
Games	0	1
Mary	1	0
too	1	0

Text Representation: Bag-of-Words

- Bag of words Representation in a Vector Space
 - ✓ The contents can be inferred from the frequency of words
 - ✓ Vector representation **does not consider the ordering of words** in a document
 - Visual words = independent features
 - John is quicker than Mary = Mary is quicker than John in BOW representation
 - ✓ We cannot reconstruct the original text based on the term-document matrix



Text Preprocessing

- Remove unnecessary information
 - ✓ They vs. they: different words in many systems
 - lower case is commonly used
 - ✓ Punctuation
 - Punctuations do not contain significant information → Remove them!
 - ✓ Numbers
 - Numbers are not critical in some domains but critical in other domains
 - Removing numbers should be carefully determined based on the domain for which a collection of text is about to be analyzed

Text Preprocessing

- Remove unnecessary information: to lower case

```
# Data preprocessing # 1: to lower case
nlp_corpus <- tm_map(nlp_corpus, content_transformer(tolower))
nlp_corpus[[1]][1]

image_corpus <- tm_map(image_corpus, content_transformer(tolower))
image_corpus[[1]][1]
```

✓ Before

```
> nlp_papers[1]
[1] "Towards non-toxic landscapes: Automatic toxic comment detection using DNN The spectacular expansion of the Internet led to the development of a new research problem in the natural language processing field: automatic toxic comment detection, since many countries prohibit hate speech in public media. There is no clear and formal definition of hate, offensive, toxic and abusive speeches. In this article, we put all these terms under the \"umbrella\" of toxic speech. The contribution of this paper is the design of binary classification and regression-based approaches aiming to predict whether a comment is toxic or not. We compare different unsupervised word representations and different DNN classifiers. Moreover, we study the robustness of the proposed approaches to adversarial attacks by adding one (healthy or toxic) word. We evaluate the proposed methodology on the English Wikipedia Detox corpus. Our experiments show that using BERT fine-tuning outperforms feature-based BERT, Mikolov's word embedding or fastText representations with different DNN classifiers."
```

Text Preprocessing

- Remove unnecessary information: to lower case

```
# Data preprocessing # 1: to lower case
nlp_corpus <- tm_map(nlp_corpus, content_transformer(tolower))
nlp_corpus[[1]][1]

image_corpus <- tm_map(image_corpus, content_transformer(tolower))
image_corpus[[1]][1]
```

✓ After

```
> nlp_corpus[[1]][1]
$content
[1] "towards non - toxic landscape: automatic toxic comment detection use dnn the spectacular expansion of the internet lead to the development of a new research problem in the natural language process field: automatic toxic comment detection, since many country prohibit hate speech in public medium. there be no clear and formal definition of hate, offensive, toxic and abusive speech. in this article, we put all this term under the \" umbrella \" of toxic speech. the contribution of this paper be the design of binary classification and regression - base approach aim to predict whether a comment be toxic or not. we compare different unsupervised word representation and different dnn classifier. moreover, we study the robustness of the propose approach to adversarial attack by add one ( healthy or toxic ) word. we evaluate the propose methodology on the english wikipedia detox corpus. our experiment show that use bert fine - tune outperform feature - base bert, mikolov's word embed or fasttext representation with different dnn classifier."
```

Text Preprocessing

- Remove unnecessary information: remove punctuations

```
# 2: remove puntuations
nlp_corpus <- tm_map(nlp_corpus, content_transformer(removePunctuation))
nlp_corpus[[1]][1]

image_corpus <- tm_map(image_corpus, content_transformer(removePunctuation))
image_corpus[[1]][1]
```

```
> nlp_corpus[[1]][1]
$content
[1] "towards non toxic landscape automatic toxic comment detection use dnn the spectacular expansion
of the internet lead to the development of a new research problem in the natural language process field
automatic toxic comment detection since many country prohibit hate speech in public medium there
be no clear and formal definition of hate offensive toxic and abusive speech in this article we put
all this term under the umbrella of toxic speech the contribution of this paper be the design of b
inary classification and regression base approach aim to predict whether a comment be toxic or not w
e compare different unsupervised word representation and different dnn classifier moreover we study t
he robustness of the propose approach to adversarial attack by add one healthy or toxic word we eva
luate the propose methodology on the english wikipedia detox corpus our experiment show that use bert
fine tune outperform feature base bert mikolovs word embed or fasttext representation with differe
nt dnn classifier"
```

Text Preprocessing

- Remove unnecessary information: remove numbers

```
# 3. remove numbers
nlp_corpus <- tm_map(nlp_corpus, content_transformer(removeNumbers))
nlp_corpus[[5]][1]

image_corpus <- tm_map(image_corpus, content_transformer(removeNumbers))
image_corpus[[1]][1]
```

✓ Before

```
> nlp_corpus[[5]][1]
$content
[1] "a subword level language model for bangla language language model be at the core of natural lang
uage process the ability to represent natural language give rise to its application in numerous nlp t
ask include text classification summarization and translation research in this area be very limit in
bangla due to the scarcity of resource except for some count base model and very recent neural lang
uage model be propose which be all base on word and limit in practical task due to their high perplex
ity this paper attempt to approach this issue of perplexity and propose a subword level neural langua
ge model with the awd lstm architecture and various other technique suitable for train in bangla lan
guage the model be train on a corpus of bangla newspaper article of a appreciable size consist of muc
h than 285 million word token the performance comparison with various other model depict the signific
ant reduction in perplexity the propose model provide reach as low as 39 84 in just 20 epoch"
```

Text Preprocessing

- Remove unnecessary information: remove numbers

```
# 3. remove numbers
nlp_corpus <- tm_map(nlp_corpus, content_transformer(removeNumbers))
nlp_corpus[[5]][1]

image_corpus <- tm_map(image_corpus, content_transformer(removeNumbers))
image_corpus[[1]][1]
```

✓ After

```
> nlp_corpus[[5]][1]
$content
[1] "a subword level language model for bangla language language model be at the core of natural lang
uage process the ability to represent natural language give rise to its application in numerous nlp t
ask include text classification summarization and translation research in this area be very limit in
bangla due to the scarcity of resource except for some count base model and very recent neural lang
uage model be propose which be all base on word and limit in practical task due to their high perplex
ity this paper attempt to approach this issue of perplexity and propose a subword level neural langua
ge model with the awd lstm architecture and various other technique suitable for train in bangla lan
guage the model be train on a corpus of bangla newspaper article of a appreciable size consist of muc
h than million word token the performance comparison with various other model depict the significant
reduction in perplexity the propose model provide reach as low as in just epoch"
```

Text Preprocessing: Stopwords

- What are stop words?

✓ Words that **do not carry any information**

- Mainly functional role
- Usually remove them to help the machine learning algorithms to perform better

✓ Natural language dependent

- English: a, about, above, across, after, again, against, all, also, etc.
- 한국어: ...습니다, ...로서(써), ...를 등

[Original text]

Information Systems Asia Web - provides research, IS-related commercial materials, interaction, **and even** research sponsorship **by** interested corporations **with a focus on** Asia Pacific region.

[After removing stop words]

Information Systems Asia Web provides research IS-related commercial materials interaction research sponsorship interested corporations focus Asia Pacific region

Text Preprocessing: Stopwords

- Example I: SMART stop words list

✓ SMART: System for the Mechanical Analysis and Retrieval of Text

- A total of 571 stop words

[1]	"a"	"a's"	"able"	"about"	"above"	"according"	"accordingly"	"across"	"actually"	"after"	"afterwards"
[12]	"again"	"against"	"ain't"	"all"	"allow"	"allows"	"almost"	"alone"	"along"	"already"	"also"
[23]	"although"	"always"	"am"	"among"	"amongst"	"an"	"and"	"another"	"any"	"anybody"	"anyhow"
[34]	"anyone"	"anything"	"anyway"	"anyways"	"anywhere"	"apart"	"ask"	"appear"	"appropriate"	"are"	"aren't"
[45]	"around"	"as"	"aside"	"ask"	"asking"	"associated"	"at"	"available"	"away"	"awfully"	"b"
[56]	"be"	"became"	"because"	"become"	"becomes"	"becoming"	"been"	"before"	"beforehand"	"behind"	"being"
[67]	"believe"	"below"	"beside"	"besides"	"best"	"better"	"between"	"beyond"	"both"	"brief"	"but"
[78]	"by"	"c"	"c'mon"	"c's"	"came"	"can"	"can't"	"cannot"	"cant"	"cause"	"causes"
[89]	"certain"	"certainly"	"changes"	"clearly"	"co"	"com"	"come"	"comes"	"concerning"	"consequently"	"consider"
[100]	"considering"	"contain"	"containing"	"contains"	"corresponding"	"could"	"couldn't"	"course"	"currently"	"d"	"definitely"
[111]	"described"	"despite"	"did"	"didn't"	"different"	"do"	"does"	"doesn't"	"doing"	"don't"	"done"
[122]	"down"	"downwards"	"during"	"e"	"each"	"edu"	"eg"	"eight"	"either"	"else"	"elsewhere"
[133]	"enough"	"entirely"	"especially"	"et"	"etc"	"even"	"ever"	"every"	"everybody"	"everyone"	"everything"
[144]	"everywhere"	"ex"	"exactly"	"example"	"except"	"f"	"far"	"few"	"fifth"	"first"	"five"
[155]	"followed"	"following"	"follows"	"for"	"former"	"formerly"	"forth"	"four"	"from"	"further"	"furthermore"
[166]	"g"	"get"	"gets"	"getting"	"given"	"gives"	"go"	"goes"	"going"	"gone"	"got"
[177]	"gotten"	"greetings"	"h"	"had"	"hadn't"	"happens"	"hardly"	"has"	"haven't"	"hereafter"	"heresy"
[188]	"having"	"he"	"he's"	"hello"	"help"	"hence"	"her"	"here"	"here's"	"hitherto"	"how"
[199]	"herein"	"hereupon"	"hers"	"herself"	"hi"	"him"	"himself"	"his"	"hopefully"	"ignored"	"immediate"
[210]	"however"	"however"	"i"	"i'd"	"i'll"	"i'm"	"i've"	"ie"	"if"	"instead"	"into"
[221]	"in"	"inasmuch"	"inc"	"indeed"	"indicate"	"indicated"	"indicates"	"inner"	"itself"	"j"	"just"
[232]	"inward"	"is"	"isn't"	"it"	"it'd"	"it'll"	"it's"	"its"	"last"	"lately"	"later"
[243]	"K"	"keep"	"keeps"	"kept"	"know"	"knows"	"known"	""	"liked"	"likely"	"little"
[254]	"latter"	"latterly"	"least"	"less"	"lest"	"let"	"let's"	"like"	"maybe"	"mean"	"myself"
[265]	"look"	"looking"	"looks"	"ltd"	"m"	"mainly"	"many"	"may"	"must"	"my"	"never"
[276]	"meanwhile"	"merely"	"might"	"more"	"moreover"	"most"	"mostly"	"much"	"needs"	"neither"	"normally"
[287]	"n"	"name"	"namely"	"nd"	"near"	"nearly"	"necessary"	"need"	"noone"	"nor"	"oh"
[298]	"nevertheless"	"new"	"next"	"nine"	"no"	"nobody"	"non"	"none"	"off"	"often"	"other"
[309]	"not"	"nothing"	"novel"	"now"	"nowhere"	"o"	"obviously"	"of"	"onto"	"or"	"own"
[320]	"ok"	"okay"	"old"	"on"	"once"	"one"	"ones"	"only"	"over"	"overall"	"probably"
[331]	"others"	"otherwise"	"ought"	"our"	"ours"	"ourselves"	"out"	"outside"	"possible"	"presumably"	"reasonably"
[342]	"p"	"particular"	"particularly"	"per"	"perhaps"	"placed"	"please"	"plus"	"re"	"really"	"say"
[353]	"provides"	"q"	"que"	"quite"	"qv"	"r"	"rather"	"rd"	"same"	"seem"	"seen"
[364]	"regarding"	"regardless"	"regards"	"relatively"	"respectively"	"right"	"s"	"said"	"seemed"	"seems"	"seen"
[375]	"saying"	"says"	"second"	"secondly"	"see"	"seeing"	"seven"	"several"	"shall"	"she"	"should"
[386]	"self"	"selves"	"sensible"	"sent"	"serious"	"seriously"	"several"	"someone"	"something"	"sometimes"	"sometimes"
[397]	"shouldn't"	"since"	"six"	"so"	"some"	"somebody"	"somehow"	"someone"	"sub"	"such"	"sup"
[408]	"somewhat"	"somewhere"	"soon"	"sorry"	"specified"	"specify"	"specifying"	"still"	"than"	"thank"	"thanks"
[419]	"sure"	"t"	"t's"	"take"	"taken"	"tell"	"tends"	"thi"	"themselves"	"then"	"thence"
[430]	"thanx"	"that"	"that's"	"thats"	"the"	"their"	"theirs"	"them"	"these"	"they"	"they'd"
[441]	"there"	"there's"	"thereafter"	"thereby"	"therefore"	"therein"	"theres"	"thereupon"	"those"	"though"	"three"
[452]	"they'll"	"they're"	"they've"	"think"	"third"	"this"	"thorough"	"thoroughly"	"toward"	"towards"	"tried"
[463]	"through"	"throughout"	"thru"	"thus"	"to"	"together"	"too"	"took"	"under"	"unfortunately"	"unless"
[474]	"tries"	"truly"	"try"	"trying"	"twice"	"two"	"u"	"un"	"uses"	"using"	"want"
[485]	"unlikely"	"until"	"unto"	"up"	"upon"	"us"	"use"	"used"	"w"	"welcome"	"well"
[496]	"usually"	"uucp"	"v"	"value"	"various"	"very"	"via"	"viz"	"vs"	"where"	"where's"
[507]	"wants"	"was"	"wasn't"	"way"	"we"	"we'd"	"we'll"	"we're"	"we've"	"whence"	"who"
[518]	"went"	"were"	"weren't"	"what"	"what's"	"whatever"	"when"	"whence"	"whenever"	"while"	"whither"
[529]	"whereafter"	"whereas"	"whereby"	"wherein"	"whereupon"	"wherever"	"whether"	"which"	"will"	"with"	"within"
[540]	"who's"	"whoever"	"whole"	"whom"	"whose"	"why"	"would"	"willing"	"wish"	"yet"	"you"
[551]	"without"	"won't"	"wonder"	"would"	"wouldn't"	"x"	"y"	"yes"	"z"	"zero"	
[562]	"you'd"	"you're"	"you've"	"your"	"yours"	"yourself"	"yourselves"				

Text Preprocessing: Stopwords

- Example 2: MySQL Stop words list

✓ <http://dev.mysql.com/doc/refman/5.1/en/fulltext-stopwords.html>

- A total of 543 stop words

a's	able	about	above	according	her	here	here's	hereafter	hereby	serious	seriously	seven	several	shall
accordingly	across	actually	after	afterwards	herein	hereupon	hers	herself	hi	she	should	shouldn't	since	six
again	against	ain't	all	allow	him	himself	his	hither	hopefully	so	some	somebody	somewhat	somewhere
allows	almost	alone	along	already	how	howbeit	however	i'd	i'll	something	sometime	sometimes	somewhat	somewhere
also	although	always	am	among	i'm	i've	ie	if	ignored	soon	sorry	specified	specify	specifying
amongst	an	and	another	any	immediate	in	inasmuch	inc	indeed	still	sub	such	sup	sure
anybody	anyhow	anyone	anything	anyway	indicate	indicated	indicates	inner	insofar	t's	take	taken	tell	tends
anyways	anywhere	apart	appear	appreciate	instead	into	inward	is	isn't	th	than	thank	thanks	thanx
appropriate	are	aren't	around	as	it	it'd	it'll	it's	its	that	that's	thats	the	their
aside	ask	asking	associated	at	itself	just	keep	keeps	kept	theirs	them	themselves	then	thence
available	away	awfully	be	became	know	known	knows	last	lately	there	there's	thereafter	thereby	therefore
because	become	becomes	becoming	been	later	latter	latterly	least	less	therein	theres	thereupon	these	they
before	beforehand	behind	being	believe	lest	let	let's	like	liked	they'd	they'll	they're	they've	think
below	beside	besides	best	better	likely	little	look	looking	looks	third	this	thorough	thoroughly	those
between	beyond	both	brief	but	ltd	mainly	many	may	maybe	though	three	through	throughout	thru
by	c'mon	c's	came	can	me	mean	meanwhile	merely	might	thus	to	together	too	took
can't	cannot	cant	cause	causes	more	moreover	most	mostly	much	toward	towards	tried	tries	truly
certain	certainly	changes	clearly	co	must	my	myself	name	namely	try	trying	twice	two	un
com	come	comes	concerning	consequently	nd	near	nearly	necessary	need	under	unfortunately	unless	unlikely	until
consider	considering	contain	containing	contains	needs	neither	never	nevertheless	new	unto	up	upon	us	use
corresponding	could	couldn't	course	currently	next	nine	no	nobody	non	used	useful	uses	using	usually
definitely	described	despite	did	didn't	none	noone	nor	normally	not	value	various	very	via	viz
different	do	does	doesn't	doing	nothing	novel	now	nowhere	obviously	vs	want	wants	was	wasn't
don't	done	down	downwards	during	of	off	often	oh	ok	way	we	we'd	we'll	we're
each	edu	eg	eight	either	okay	old	on	once	one	we've	welcome	well	went	were
else	elsewhere	enough	entirely	especially	ones	only	onto	or	other	weren't	what	what's	whatever	when
et	etc	even	ever	every	others	otherwise	ought	our	ours	whence	whenever	where	where's	whereafter
everybody	everyone	everything	everywhere	ex	ourselves	out	outside	over	overall	whereas	whereby	wherein	whereupon	wherever
exactly	example	except	far	few	own	particular	particularly	per	perhaps	whether	which	while	whither	who
fifth	first	five	followed	following	placed	please	plus	possible	presumably	who's	whoever	whole	whom	whose
follows	for	former	formerly	forth	probably	provides	que	quite	qv	why	will	willing	wish	with
four	from	further	furthermore	get	rather	rd	re	really	reasonably	within	without	won't	wonder	would
gets	getting	given	gives	go	regarding	regardless	regards	relatively	respectively	wouldn't	yes	yet	you	you'd
goes	going	gone	got	gotten	right	said	same	saw	say	you'll	you're	you've	your	yours
greetings	had	hadn't	happens	hardly	saying	says	second	secondly	see	yourself	yourselves	zero		
has	hasn't	have	haven't	having	seeing	seem	seemed	seeming	seems					
he	he's	hello	help	hence	seen	self	selves	sensible	sent					

Text Preprocessing: Stopwords

- Example 3: Stop words list in Korean

✓ <http://www.ranks.nl/stopwords/korean>

- A total of 677 stop words

아	어찌거든	하기보다는	만나 아니다 다시 말하지면	끼닭으로	할 생각이다 즐음하여	본대로	얼마간	너	혼자
휴	그위에	차라리	만이 아니다 바꿔 말하면	이유만으로	하려고다 다른	본대로	약간	너희	자기
아이구	개다가	하는 편이 낫다	만은 아니다 즉	이로 인하여	다른 방면으로	자	당신	자기집	
아이쿠	절에서 보아	흐흔	막하고 구체적으로	그래서	이리하여	이로	좀	어찌	자신
아이고	비추어 보아	놀라다	관계없이 말하자면	이 때문에	그리하여	이쪽	조금	설마	우에 종합한것과
어	고려하면	상대적으로 말하	그치지 않다 시작하여	그러므로	그렇게 함으 습니까	여기	다수	차단리	같이
나	하게된다	자면	그러나 시초에	그런 까닭에	했어요	이것	몇	활자언정	총적으로 보면
우리	일것이다	마치	그런데 이상	알 수 있다	하지만	말할것도 없고	이번	활자지도	총적으로 말하면
저희	비교적	아니라면	하지만 허	결론을 낼 수 있	일때	무를쓰고	얼마	지만	활망정
따라	좀	헛	든간에 학	앞에서	할때	개의치않고	지만	활자연경	총적으로
의해	보다더	그렇지 않으면	둔하지 않다 허걱	하는것은 못하다	하는것은 못하다	이려한	한들며	대로 하다	대로 하다
을	비하면	그렇지 않다면	둔하지 않다 허걱	으로 인하여	중에서	하는것이 낫다	또한	구토하다	으로서
를	시킨다	안 그러면	둔하지 않다 허걱	보는데서	이와 같은	이와 같은	그러나	게우다	찰
에	하게하다	아니었다면	둔하지 않다 허걱	매	요만큼	요만큼	그렇지만	토하다	그만이다
의	할만하다	하든지	둔하지 않다 허걱	매번	요만한 것	요만한 것	마스겁다	할 뜻이다	
가	의해서	아니면	둔하지 않다 허걱	로써	얼마 안 되는 것	얼마 안 되는 것	하지만	활자연경	
으로	연이어서	이라면	둔하지 않다 허걱	모	이만큼	이만큼	이외에도	열사광	총팅
로	이어서	좋아	둔하지 않다 허걱	이제	이 정도의	이 정도의	대해 말하자 웨	활자연경	활자연경
에게	잇따라	알았어	둔하지 않다 허걱	로써	이렇게 많은 것	이렇게 많은 것	면	쳇	쾅쾅
뿐이다	뒤따라	하는것도	둔하지 않다 허걱	반드시	이와 같다	이와 같다	뿐이다	의거하여	뚱뚱
의거하여	뒤이어	그만이다	둔하지 않다 허걱	로써	갖고말하자면	갖고말하자면	다음에	근거하여	봐
근거하여	결국	어쩔수 없다	둔하지 않다 허걱	반드시	이때	이때	반대로	의해	보란
입각하여	의지하여	하나	둔하지 않다 허걱	로써	이렇게 많다	이렇게 많다	반대로	의해	아이야
기준으로	기대여	일	둔하지 않다 허걱	반드시	이제	이제	반대로	말하	아니
예하연	통하여	일반적으로	둔하지 않다 허걱	로써	이제는	이제는	반대로	힐입어	힐입어
예를 들면	자마자	듣타	둔하지 않다 허걱	반드시	다르게	다르게	이와 반대로 그	와아	
예를 들자면	더더욱	문운	둔하지 않다 허걱	반드시	다르게	다르게	바꾸어서 말 다음	등	
저	불구하고	이용하여	둔하지 않다 허걱	반드시	다르게	다르게	바꾸어서 한 두번째로	아이	
소인	오자마자	여서	둔하지 않다 허걱	반드시	다르게	다르게	다면	점나	
소생	얼마든지	이용해보면	둔하지 않다 허걱	반드시	다르게	다르게	기타		
저희	이렇게되면	제외하고	둔하지 않다 허걱	반드시	다르게	다르게	만약		
지말고	이렇게되면	예하면	둔하지 않다 허걱	반드시	다르게	다르게	첫번째로		
하지마	즉시	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	그렇지않으		
하지마라	비로	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	나머지는		
다른	근거로	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	그중에서		
물론	하기에	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	까악		
또한	하자마자	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	견지에서		
그리고	방에 안된다	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	록		
비길수	않아	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	형식으로 쓰여		
해서는	원된	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	딱		
다	요컨대	다시말하면	둔하지 않다 허걱	반드시	다르게	다르게	입장에서		

Text Preprocessing

- Remove unnecessary information: remove stopwords
 - ✓ You can add your own stopwords to an existing stopwords list

```
# 4. remove stopwords (SMART stopwords list)
arxiv_stopwords <- c("process", "paper", "model", "learn", "propose", "task",
                     "method", "show", "approach", "result", "system")
myStopwords <- c(stopwords("SMART"), arxiv_stopwords)

nlp_corpus <- tm_map(nlp_corpus, removeWords, myStopwords)
nlp_corpus[[1]][1]

image_corpus <- tm_map(image_corpus, removeWords, myStopwords)
image_corpus[[1]][1]

> nlp_corpus[[1]][1]
$content
[1] " toxic landscape automatic toxic comment detection dnn spectacular expansion internet lead
     development research problem natural language process field automatic toxic comment detection
     country prohibit hate speech public medium clear formal definition hate offensive toxic abusive
     speech article put term umbrella toxic speech contribution design binary classification
     regression base approach aim predict comment toxic compare unsupervised word representation
     dnn classifier study robustness propose approach adversarial attack add healthy
     toxic word evaluate propose methodology english wikipedia detox corpus experiment show bert
     fine tune outperform feature base bert mikolovs word embed fasttext representation dnn classifier"
```

Text Preprocessing

- Bag-of-Words Representation: Construct a Term-Document Matrix

```
# Construct Term-Document Matrix
nlpTDM <- TermDocumentMatrix(nlp_corpus, control = list(minWordLength = 2))
imageTDM <- TermDocumentMatrix(image_corpus, control = list(minWordLength = 2))

# Check the Term-Document Matrix
as.matrix(nlpTDM)[11:30,11:30]
as.matrix(imageTDM)[11:30,11:30]
```

```
> as.matrix(nlpTDM)[11:30,11:30]
    Docs
Terms   11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30
binary  0  0  0  0  0  0  0  0  0  0  0  0  0  2  0  0  0  0  0  0
classification 0  0  0  1  1  0  0  2  0  0  0  0  0  1  0  0  0  0  6  0
classifier  0  0  0  0  0  0  0  2  0  0  0  0  0  0  0  0  0  0  1  0
clear     0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
comment   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
compare   1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  1  0
contribution 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
corpus    0  0  0  0  0  0  0  0  0  1  1  0  0  0  0  1  0  0  0  0
country   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
definition 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
design    0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
detection  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0
detox     0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
development 0  0  0  0  0  1  0  0  1  0  0  0  0  0  0  0  0  0  0  0
dnn       0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
embed     0  0  0  0  0  0  0  0  3  0  0  0  0  0  0  2  0  9  0
english   0  0  0  0  0  1  0  0  0  0  0  0  0  0  1  0  0  0  0  0
evaluate  0  1  0  0  0  0  1  1  0  0  0  2  0  0  1  0  0  0  0  0
expansion 0  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0
experiment 0  0  0  0  0  0  0  0  0  0  0  1  0  1  0  0  0  0  1
```

```
> as.matrix(imageTDM)[11:30,11:30]
    Docs
Terms   11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30
datum   0  2  1  1  0  2  1  0  0  1  0  1  0  0  0  0  0  1  0  0
design   0  1  1  3  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
development 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
device   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
domain   0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  3  0  0  0
ease     0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
easily   0  0  0  1  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
edge     0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
effective 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
environment 0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  1
fog      0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
framework 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  1  1
gateway  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
grow    0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
healthcare 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
heterogeneous 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
igatelink 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
image    1  1  9  4  1  1  2  1  2  7  5  2  1  10  4  9  2  1  1  12
important 0  0  0  0  0  0  0  0  0  0  0  0  0  0  1  0  0  0  1  2
increasingly 0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
```

Text Analytics

- Highly Frequent Tokens

```
# High frequent words
nlpTDM_matrix <- as.matrix(nlpTDM)
nlp_word_count <- rowSums(nlpTDM_matrix)
nlp_word_count <- sort(nlp_word_count, decreasing = TRUE)
nlp_word_count[1:30]
```

```
> nlp_word_count[1:30]
```

language	natural	train	datum	word	base
1035	641	449	428	392	384
network	neural	text	deep	performance	representation
359	321	304	274	272	253
dataset	good	nlp	information	state	problem
250	237	234	223	219	200
work	machine	art	domain	large	classification
191	190	183	177	174	170
attention	embeddings	feature	analysis	achieve	technique
170	162	161	161	158	156

Text Analytics

- Highly Frequent Tokens

```
# High frequent words
imageTDM_matrix <- as.matrix(imageTDM)
image_word_count <- rowSums(imageTDM_matrix)
image_word_count <- sort(image_word_count, decreasing = TRUE)
image_word_count[1:30]
```

image	datum	network	base	deep	algorithm	problem	
1927	490	485	462	366	311	299	
good	application	neural	high	train	time	dataset	
251	248	247	226	224	219	218	
feature	performance	segmentation	technique	present	detection	provide	
214	210	210	207	200	191	168	
analysis	work	study	state	quality	field	function	
163	161	157	151	150	143	140	
convolutional	information						
140	140						

Text Analytics

- Highly Frequent Tokens

```
nlp_freq_words <- findFreqTerms(nlpTDM, lowfreq=100)
image_freq_words <- findFreqTerms(imageTDM, lowfreq=100)

intersect(nlp_freq_words, image_freq_words)
setdiff(nlp_freq_words, image_freq_words)
setdiff(image_freq_words, nlp_freq_words)
```

✓ Words appearing more than 100 times in the both datasets

```
> intersect(nlp_freq_words, image_freq_words)
[1] "base"           "classification" "compare"      "feature"       "problem"
[6] "study"          "improve"        "technique"    "accuracy"     "algorithm"
[11] "information"   "performance"   "provide"     "dataset"      "domain"
[16] "application"   "architecture"  "high"        "neural"       "train"
[21] "analysis"       "apply"         "art"         "datum"        "deep"
[26] "image"          "large"         "network"     "state"        "work"
[31] "good"           "present"       "time"        "achieve"     "demonstrate"
[36] "recognition"   "computer"     "vision"      "machine"     "set"
[41] "challenge"
```

Text Analytics

- Highly Frequent Tokens

```
nlp_freq_words <- findFreqTerms(nlpTDM, lowfreq=100)
image_freq_words <- findFreqTerms(imageTDM, lowfreq=100)

intersect(nlp_freq_words, image_freq_words)
setdiff(nlp_freq_words, image_freq_words)
setdiff(image_freq_words, nlp_freq_words)
```

- ✓ Words appearing more than 100 times in the “nlp” dataset but not in the “image” dataset

```
> setdiff(nlp_freq_words, image_freq_words)
[1] "bert"           "corpus"         "embed"          "evaluate"        "experiment"
[6] "language"       "natural"        "representation" "research"        "speech"
[11] "word"          "generate"      "nlp"            "understand"     "test"
[16] "text"          "exist"         "level"          "recent"         "label"
[21] "pre"           "make"          "embeddings"     "knowledge"      "human"
[26] "semantic"       "sequence"      "transformer"    "context"        "sentence"
[31] "question"      "graph"         "attention"     "answer"
```

Text Analytics

- Highly Frequent Tokens

```
nlp_freq_words <- findFreqTerms(nlpTDM, lowfreq=100)
image_freq_words <- findFreqTerms(imageTDM, lowfreq=100)

intersect(nlp_freq_words, image_freq_words)
setdiff(nlp_freq_words, image_freq_words)
setdiff(image_freq_words, nlp_freq_words)
```

- ✓ Words appearing more than 100 times in the “image” dataset but not in the “nlp” dataset

```
> setdiff(image_freq_words, nlp_freq_words)
[1] "design"          "framework"       "function"        "quality"         "medical"
[6] "convolutional"   "include"         "matrix"          "noise"           "point"
[11] "parameter"       "sample"          "signal"          "structure"       "real"
[16] "low"              "detection"       "object"          "segmentation"    "field"
[21] "space"
```

Text Analytics

- Associated words

```
# Which words are associated with "image"?
findAssocs(nlpTDM, 'image', 0.2)
findAssocs(imageTDM, 'image', 0.2)
```

```
> findAssocs(nlpTDM, 'image', 0.2)
```

```
$image
```

caption	correlate	characterisation	flickr	observers
0.68	0.53	0.52	0.52	0.52
unexpectedly	facial	subjective	visual	imitation
0.52	0.49	0.41	0.38	0.36
description	fidelity	fifteen	fool	perceptual
0.33	0.32	0.32	0.32	0.32
reliably	substitute	criterion	adjective	proportional
0.32	0.32	0.29	0.29	0.28
vision	gans	retrieve	metric	judgment
0.28	0.28	0.27	0.27	0.27
perturb	generate	reconstruction	coco	dnnns
0.25	0.24	0.24	0.23	0.23
endeavor	localize	equivalent	sagan	dissimilarity
0.23	0.23	0.23	0.23	0.23
neuroscore	psychoperceptual	induction	gan	creation
0.23	0.23	0.22	0.22	0.22
visually	nonetheless	feature	norm	generative
0.22	0.22	0.21	0.21	0.20

Text Analytics

- Associated words

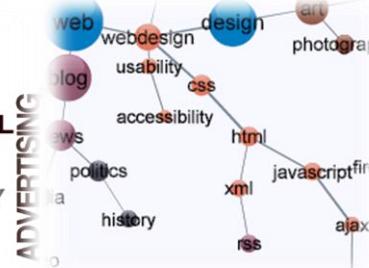
```
# Which words are associated with "image"?
findAssocs(nlpTDM, 'image', 0.2)
findAssocs(imageTDM, 'image', 0.2)
```

```
> findAssocs(imageTDM, 'image', 0.2)
```

\$image	medical	quality	dataset	registration	manual	treatment
	0.29	0.27	0.24	0.24	0.22	0.22
option	enhancement		category	curate	diva	heterogeneity
0.21	0.20		0.20	0.20	0.20	0.20
lisa	longitudinal	multiplicity	parametrically		convey	naked
0.20	0.20		0.20	0.20	0.20	0.20
personnel	translate					
	0.20	0.20				

Text Summarization: Wordcloud

- Wordcloud
 - ✓ A tool used to visually show the popularity of words in a collection of documents and how often they have been used
 - ✓ Conceptually resemble histograms, but can represent more items
 - The way it works
 - ✓ The more a word is presented, the **larger** it will appear within the cloud
 - ✓ Words that are similar appear **next to each other** in the cloud
 - Various algorithms/designs exist



<http://www.edudemic.com/9-word-cloud-generators-that-arent-wordle/>

Text Summarization: Wordcloud

- Creation of wordcloud
 - ✓ The font size of a word in a wordcloud is determined by its **incidence**.
 - ✓ For smaller frequencies, one can specify font size directly, from one to whatever the maximum font size.
 - ✓ For larger values, a scaling should be made
 - ✓ An example of font size computation

$$\text{font size } (w(i)) = \frac{\text{freq}(w(i)) - \text{min. freq}(w, D)}{\text{max. freq}(w, D) - \text{min. freq}(w, D)} + \text{constant}$$

Text Summarization: Wordcloud

- Word cloud for “nlp” dataset

```
# Word cloud pal <- brewer.pal(12, "Paired")
wordcloud(names(nlp_word_count)[1:200],
          nlp_word_count[1:200],
          min.freq=10,
          scale = c(3, 1),
          rot.per = 0.1,
          col=pal,
          random.order=F)
```

- ✓ Draw a wordcloud with top 200 words
- ✓ Each word must appear at least 10 times
- ✓ scale option controls the size of words in the wordclouds
- ✓ rot.per controls the ratio of words that are rotated

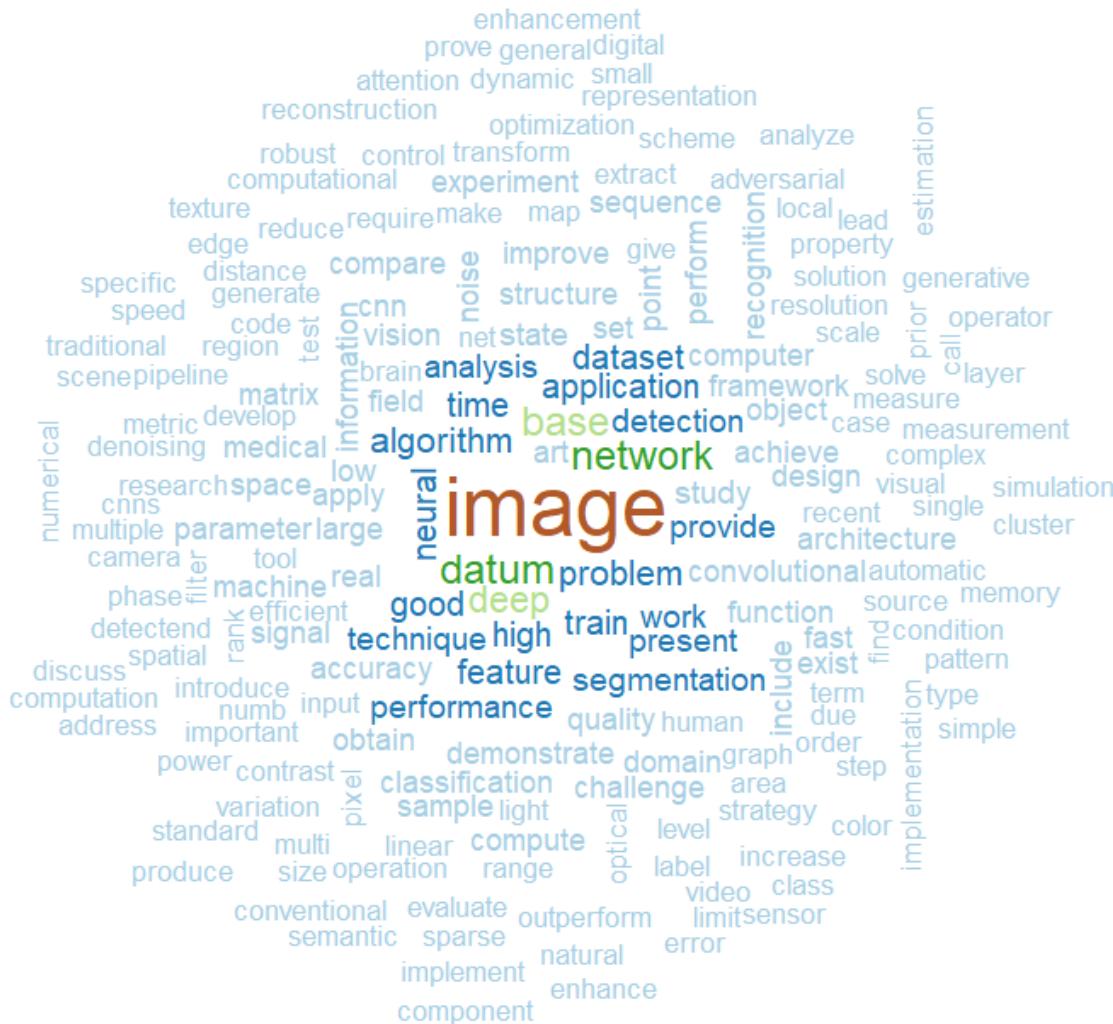
Text Summarization: Wordcloud

- Word cloud for “nlp” dataset



Text Summarization: Wordcloud

- Word cloud for “image” dataset



Text Summarization: Treemap

- Treemap

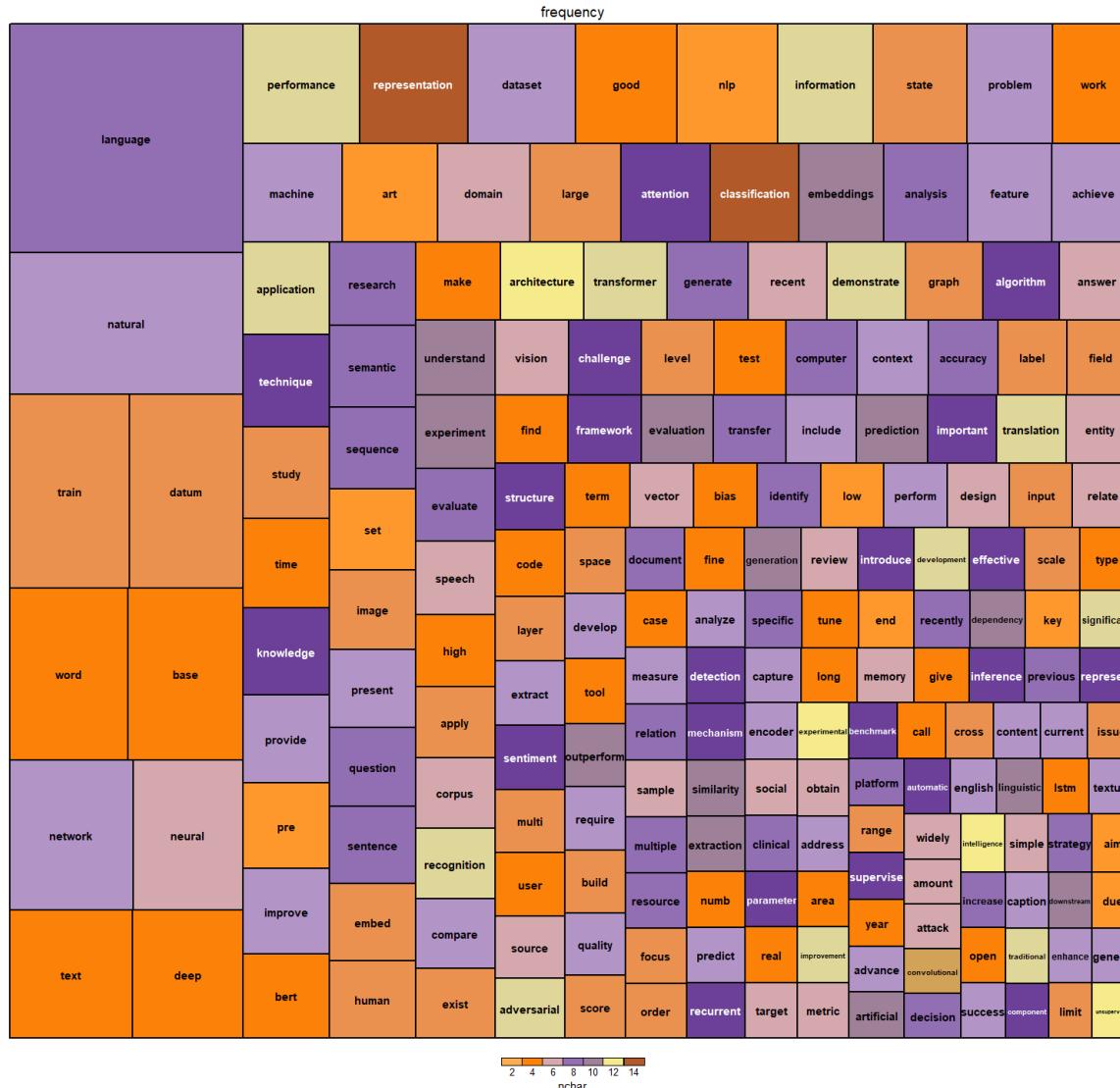
```
# Treemap
nlp_treemap_data <- data.frame(names(nlp_word_count), nlp_word_count,
nchar(names(nlp_word_count)))
names(nlp_treemap_data) <- c("word", "frequency", "nchar")
str(nlp_treemap_data)

treemap(nlp_treemap_data[1:200,],
       index = "word",
       vSize = "frequency",
       vColor = "nchar",
       type = "value",
       palette = brewer.pal(12, "Paired"))
```

- ✓ Size of rectangle: frequency of words
- ✓ Color: number of characters in each word

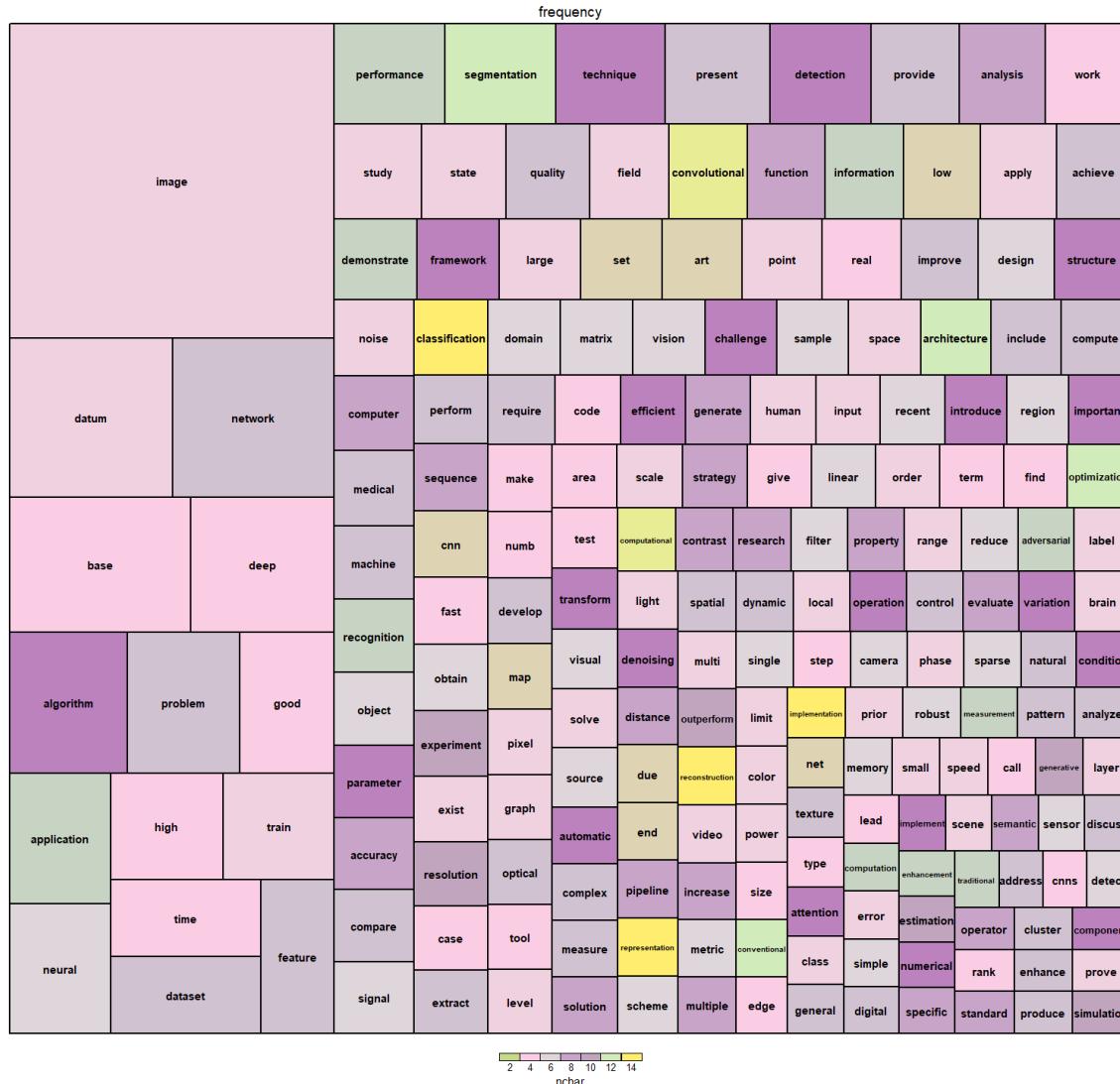
Text Summarization: Treemap

- Treemap for “nlp” dataset



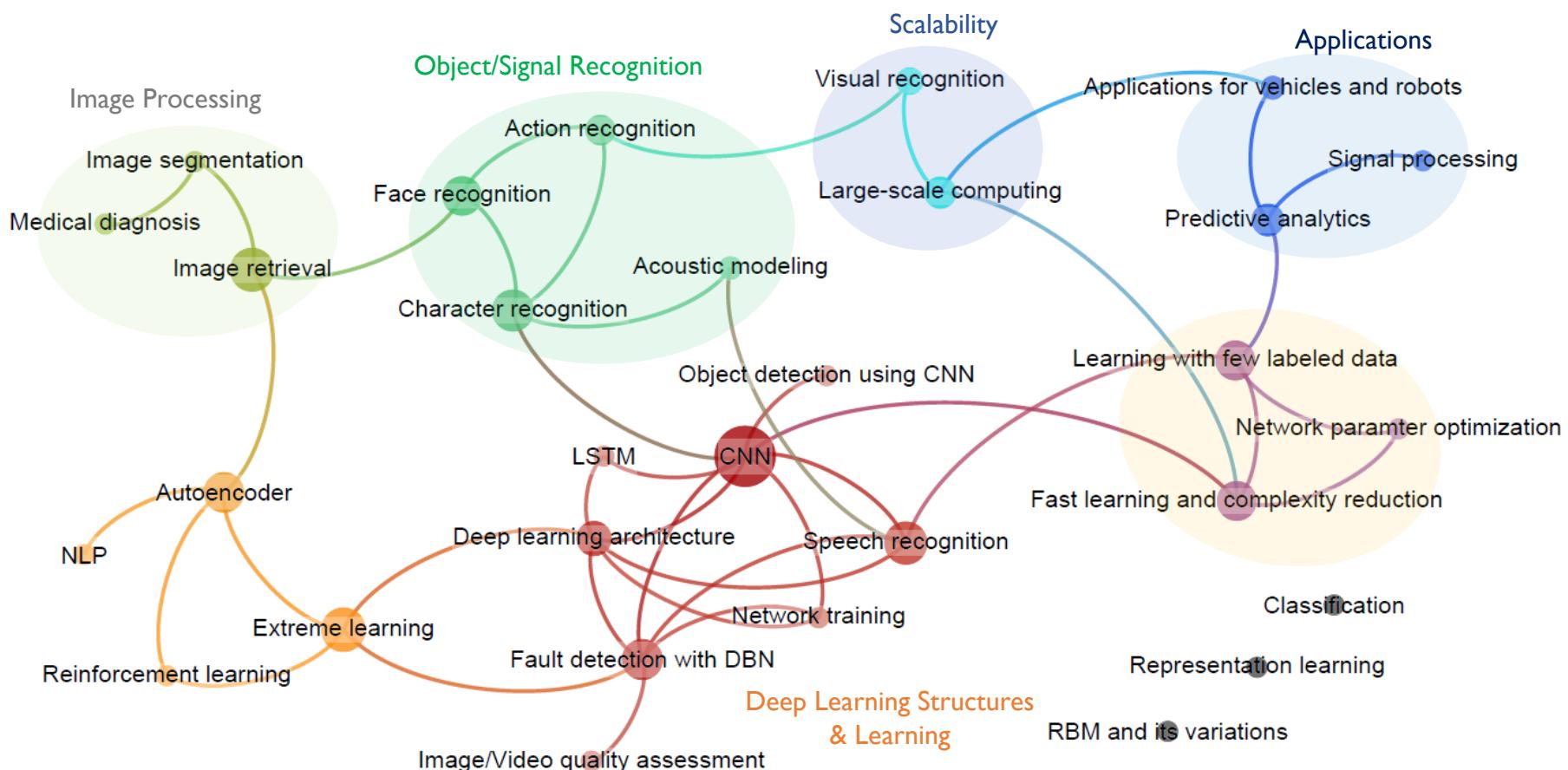
Text Summarization: Treemap

- Treemap for “image” dataset



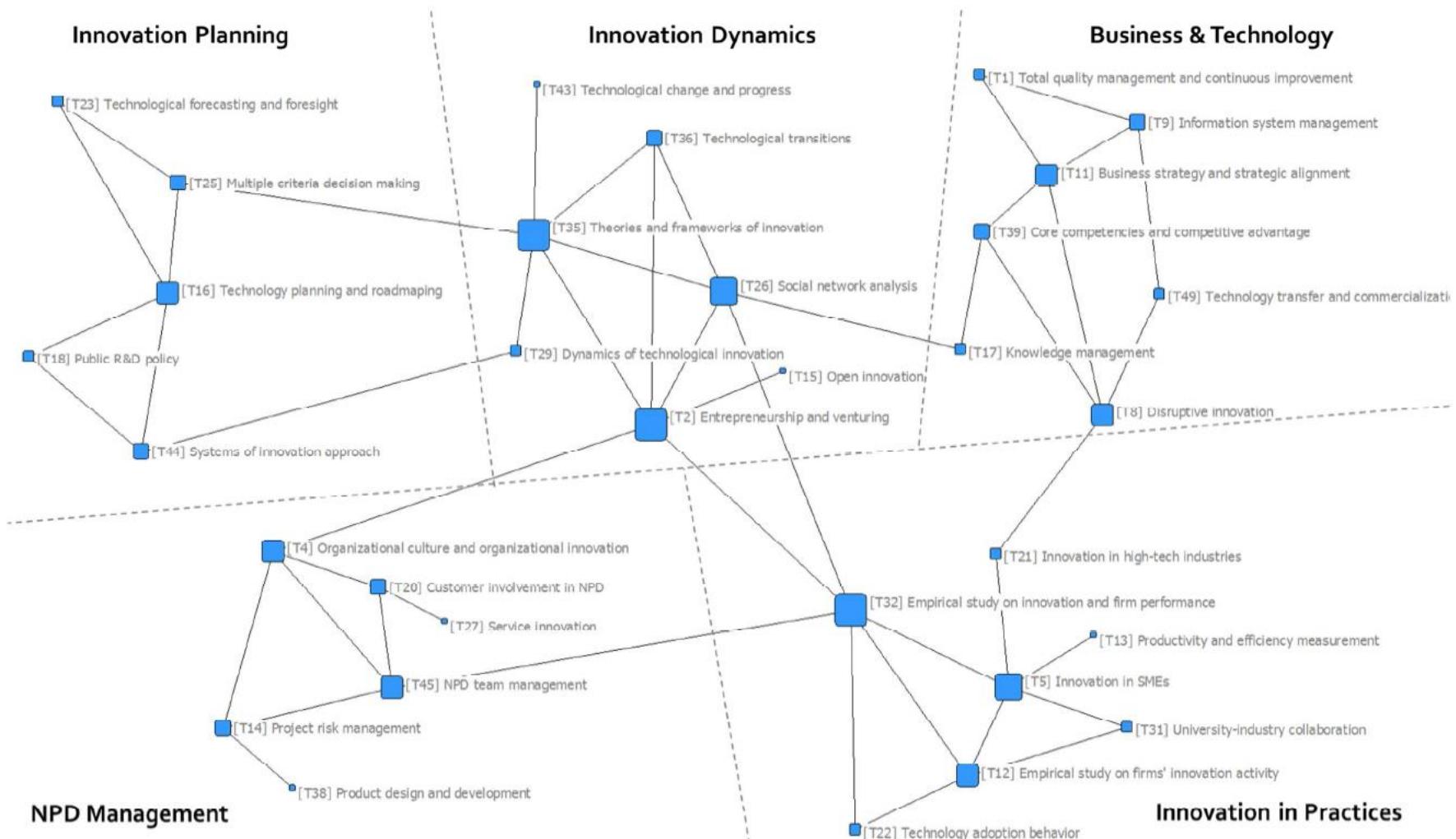
Text Summarization: Word Network

- Network of Deep Learning Topics Until Feb. 2016



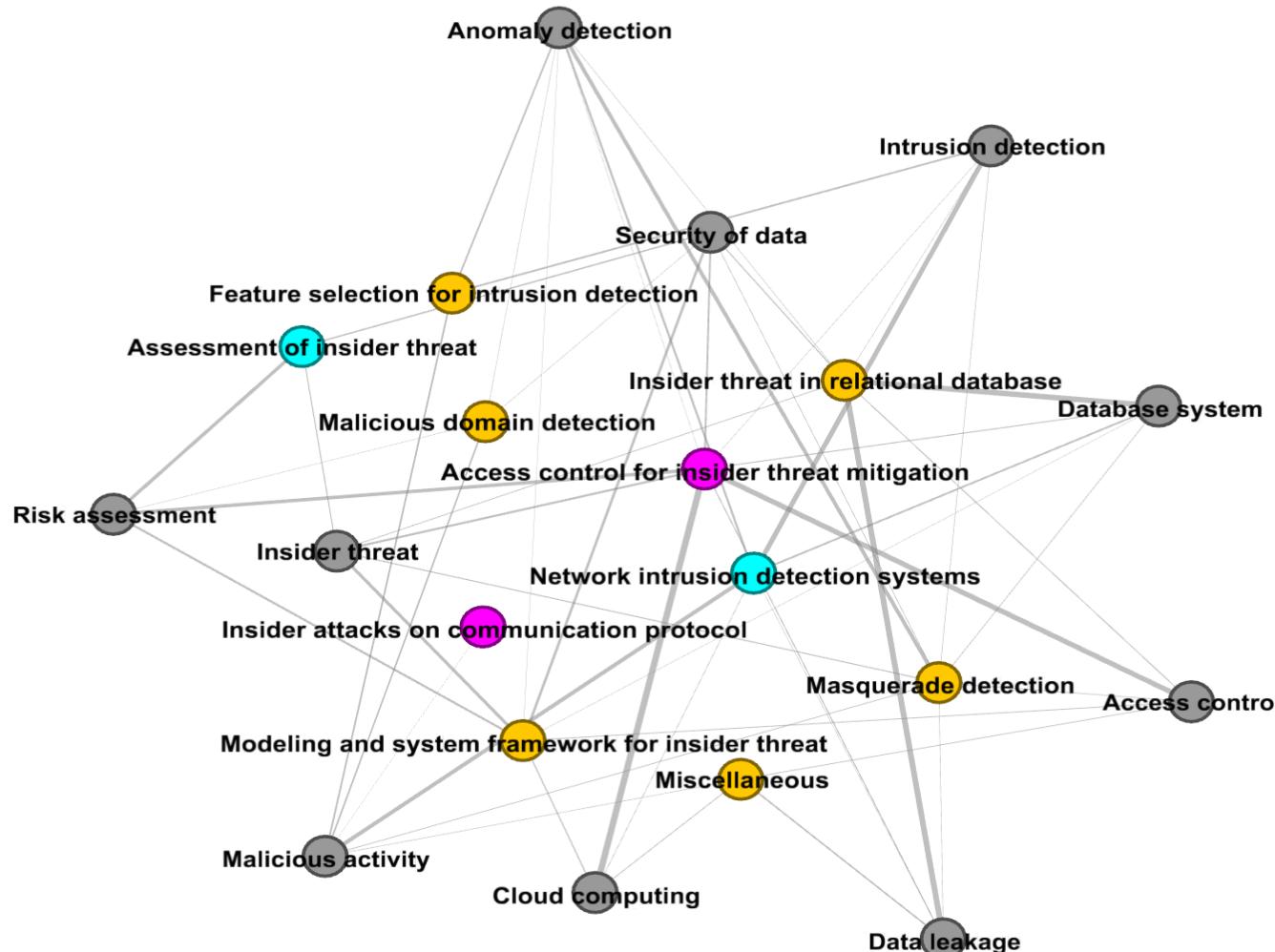
Text Summarization: Word Network

- Network of research topics for “technology and innovation management”



Text Summarization: Word Network

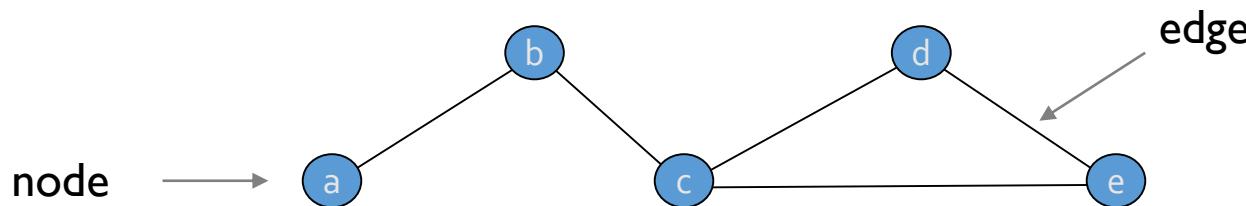
- Network of research topics for “Insider threats”



Text Summarization: Word Network

- What is a Network?

- ✓ A network is a combined set of nodes connected by edges
- ✓ Network ≡ Graph

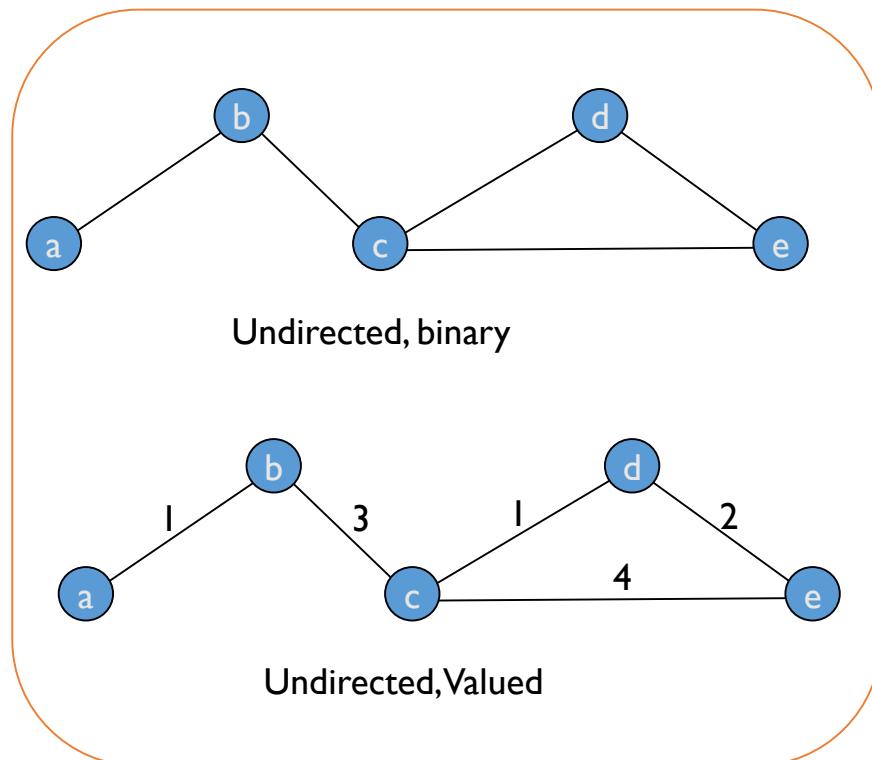


Points	Lines	Domain
Vertices	Edges, arcs	Math
Nodes	Links	Computer Science
Sites	Bonds	Physics
Actors	Ties, relations	Sociology
Keyword	Co-occurrence	Text Mining

Text Summarization: Word Network

- Type of Connections

- ✓ A relation can be binary or valued, directed or undirected

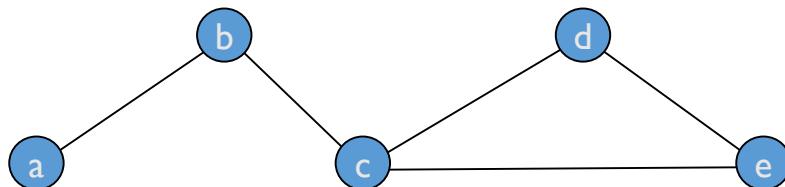


“Undirected networks are commonly used for Text Mining”

Text Summarization: Word Network

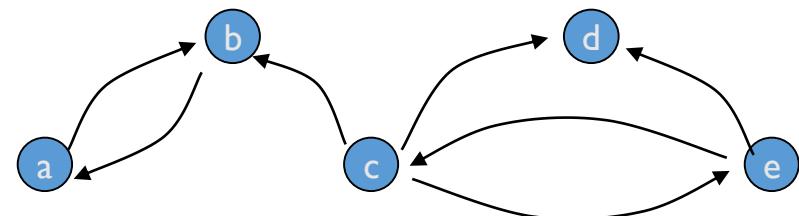
- Basic Data Structure

- ✓ From pictures to matrices



Undirected, binary

	a	b	c	d	e
a		I			
b	I			I	
c		I		I	I
d			I		I
e			I	I	



Directed, binary

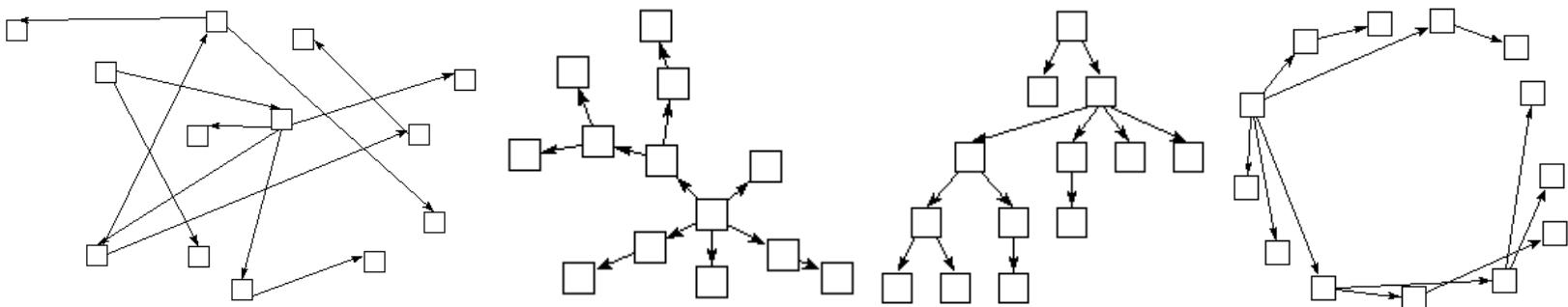
	a	b	c	d	e
a		I			
b	I				
c		I		I	I
d					
e			I	I	

Text Summarization: Word Network

- Graphical Representation

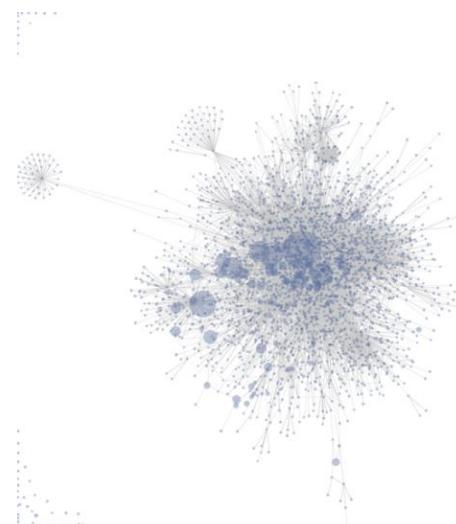
- ✓ No standard way to draw a sociogram

- Each of these are equal



- ✓ Force-directed layout (Fruchterman and Reingold, 1991)

- Distribute the vertices evenly in the frame
 - Minimize edge crossings
 - Make edge lengths uniform
 - Reflect inherent symmetry
 - Conform to the frame



Text Summarization: Word Network

- Measuring Network: Network Level

- ✓ Size
 - Number of nodes

- ✓ Density
 - Number of ties that are present

- ✓ Out-degree (directed network)
 - Sum of connections from an actor to other

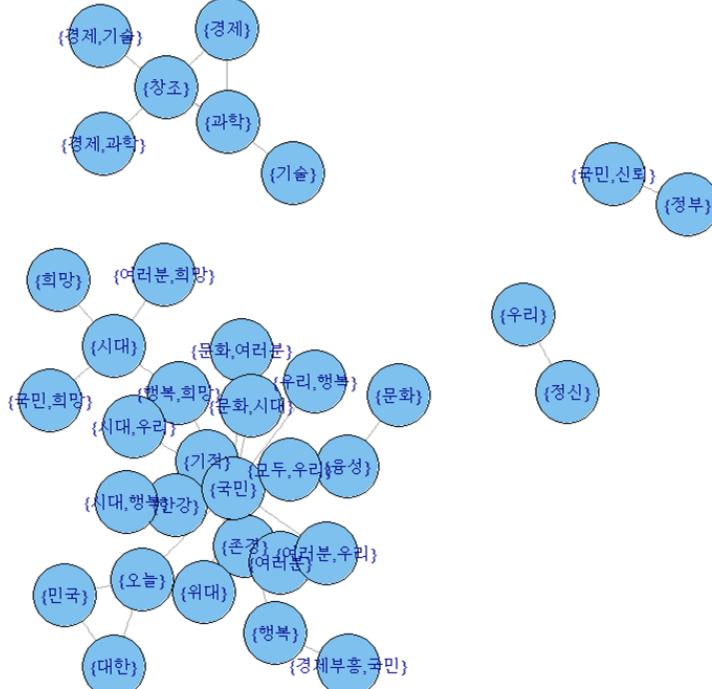
- ✓ In-degree (undirected network)
 - Sum of connections to an actor

Text Summarization: Word Network

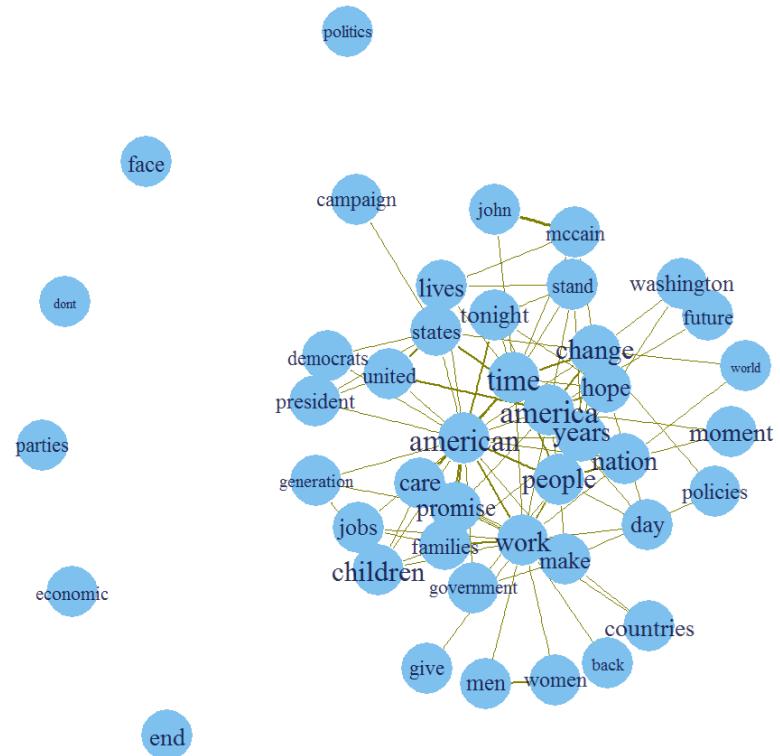
- Community Structure

- ✓ Nodes of a network can be grouped into (potentially overlapping) sets of nodes such that each set of nodes is densely connected internally.

Four non-overlapping communities

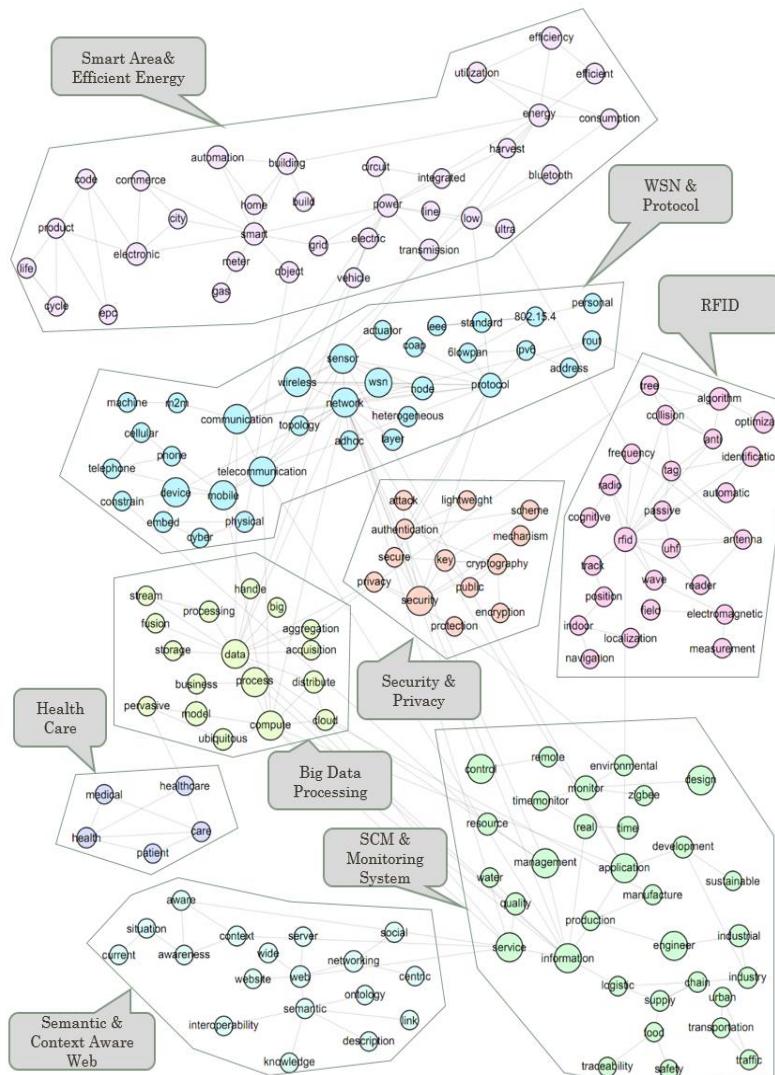


Difficult to identify the community structure



Text Summarization: Word Network

- Keyword network of Research Papers for “Internet of Things”



Text Summarization: Word Network

- Conduct some processing to draw word networks
 - ✓ Convert the frequency of each TDM to binary representation
 - ✓ Remove loops, delete edges whose edge weights are smaller than 50, etc.

```
# Construct a word network for the nlp papers
# Change it to a Boolean matrix
nlpTDM_matrix[nlpTDM_matrix >= 1] <- 1
nlp_adj_mat <- nlpTDM_matrix %*% t(nlpTDM_matrix)

# inspect terms numbered 1 to 10
nlp_adj_mat[1:10,1:10]

# Build a graph from the above matrix
nlp_graph <- graph.adjacency(nlp_adj_mat, weighted=T, mode = "undirected")
# remove loops
nlp_graph <- simplify(nlp_graph)
# set labels and degrees of vertices
V(nlp_graph)$label <- V(nlp_graph)$name
nlp_graph <- delete.edges(nlp_graph, which(E(nlp_graph)$weight <= 50))
nlp_graph <- delete.vertices(nlp_graph, which(degree(nlp_graph) == 0))
V(nlp_graph)$degree <- degree(nlp_graph)
```

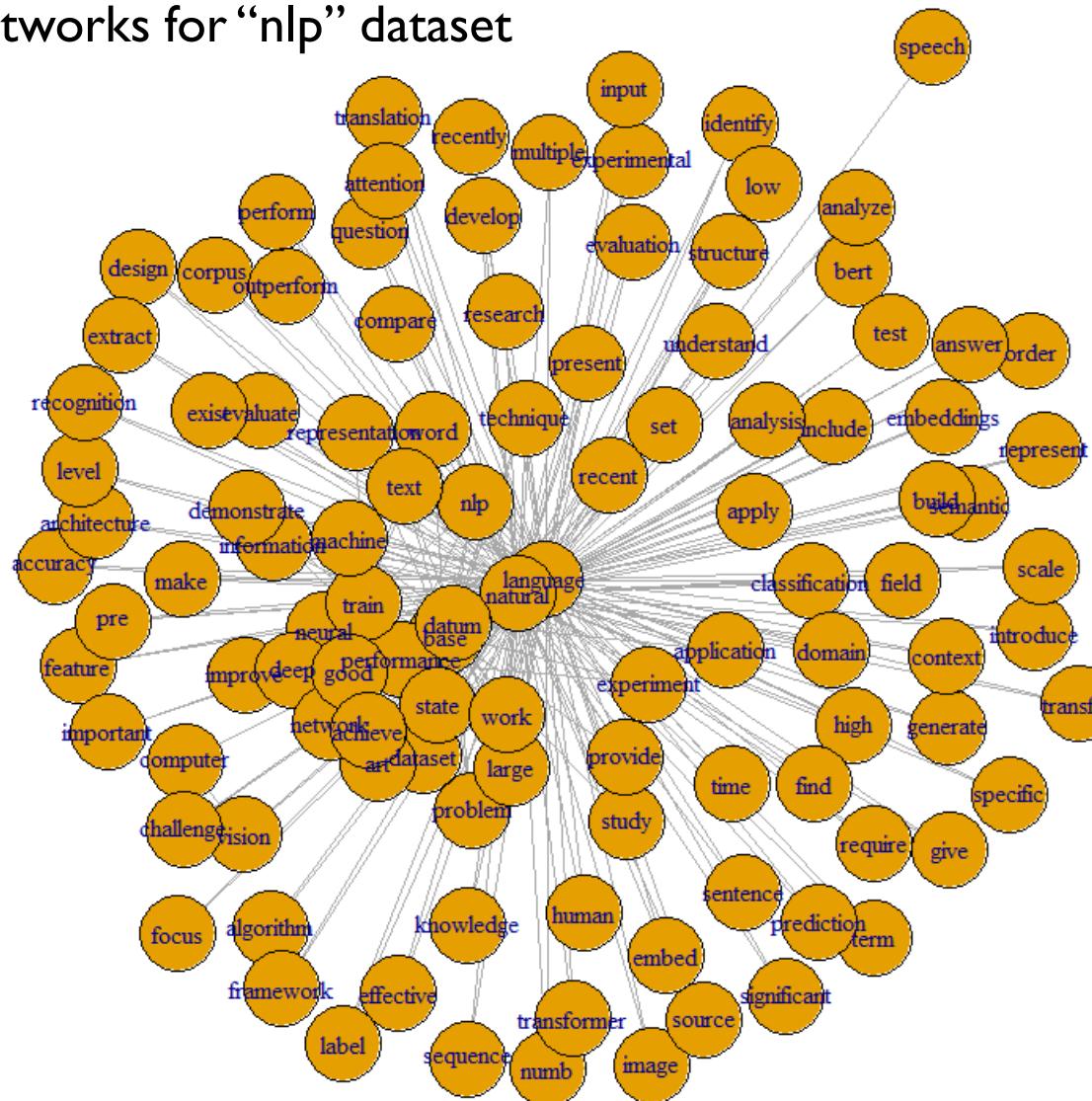
Text Summarization: Word Network

- Plot some networks for “nlp” dataset

```
# set seed to make the layout reproducible
set.seed(3952)
plot(nlp_graph, layout=layout.fruchterman.reingold)
plot(nlp_graph, layout=layout.kamada.kawai, vertex.size = 10,
      vertex.color = 8, vertex.label.cex = 1)
```

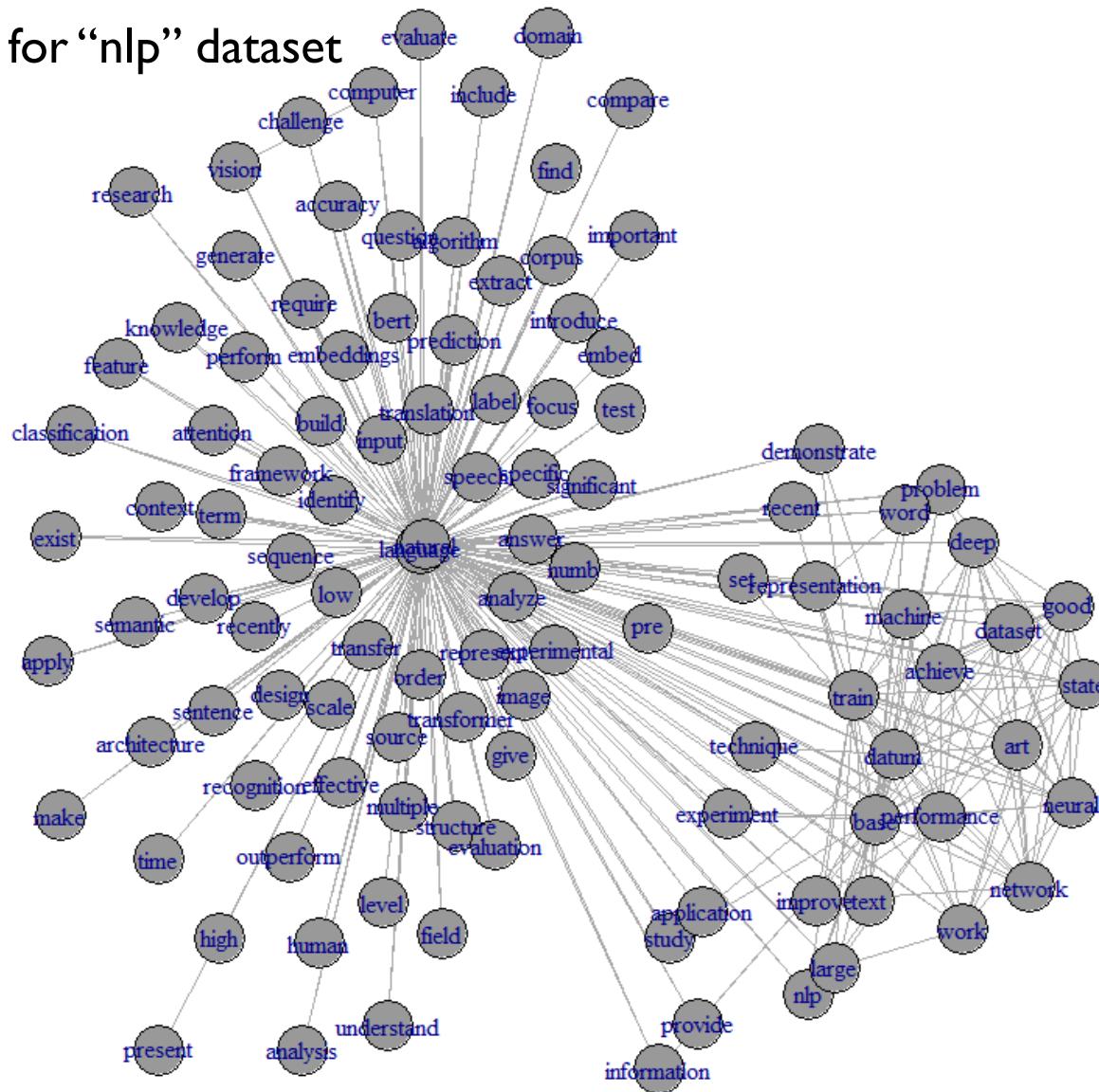
Text Summarization: Word Network

- Plot some networks for “nlp” dataset



Text Summarization: Word Network

- Plot some networks for “nlp” dataset



Text Summarization: Word Network

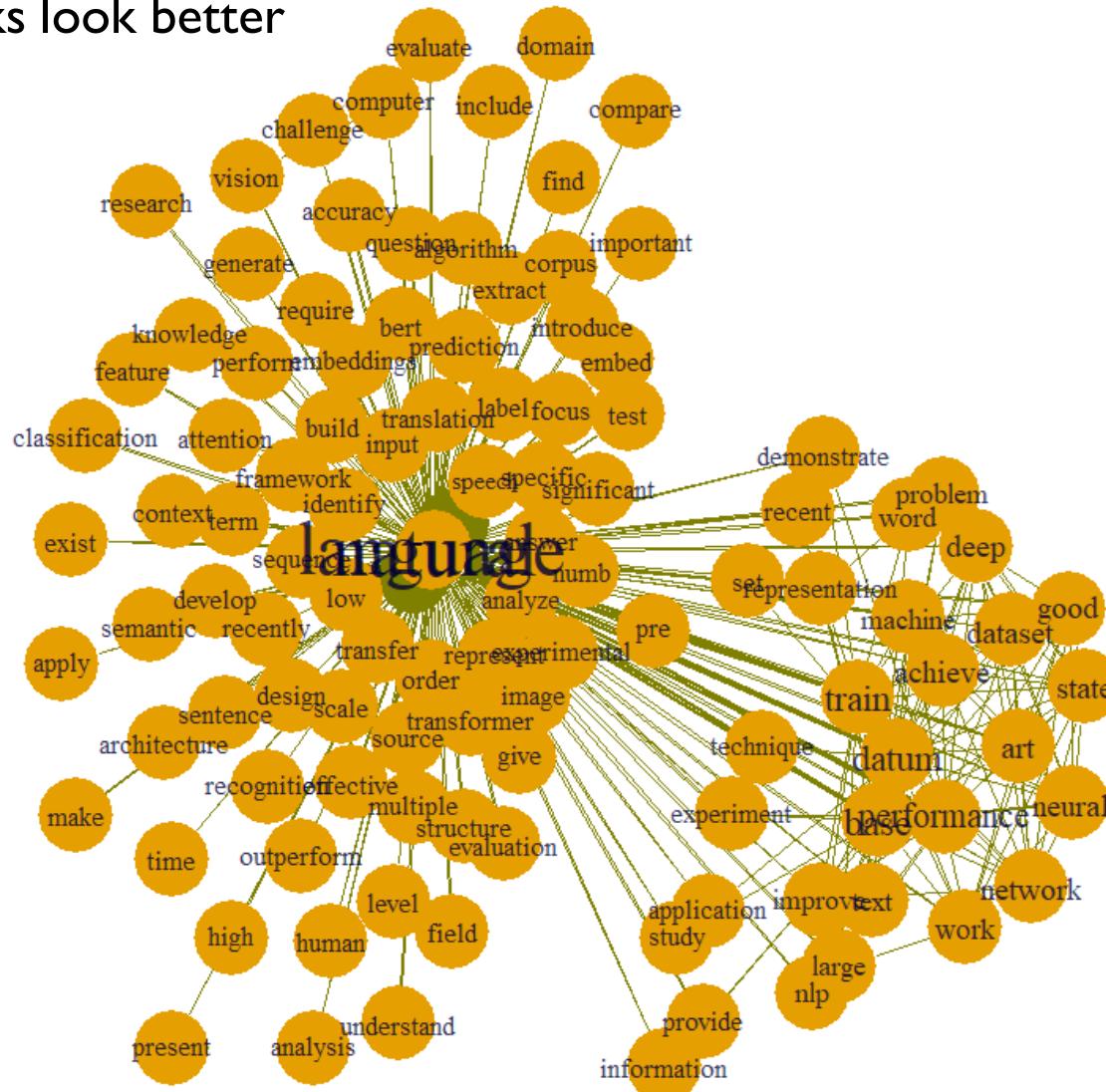
- Make networks look better

```
# Make the network look better
V(nlp_graph)$label.cex <- 2*V(nlp_graph)$degree/max(V(nlp_graph)$degree)+1
V(nlp_graph)$label.color <- rgb(0, 0, 0.2, 0.8)
V(nlp_graph)$frame.color <- NA
egam1 <- 300*(exp(E(nlp_graph)$weight/100))/sum(exp(E(nlp_graph)$weight/100))
E(nlp_graph)$color <- rgb(0.5, 0.5, 0)
E(nlp_graph)$width <- egam1

# plot the graph in layout1
plot(nlp_graph, layout=layout.kamada.kawai)
```

Text Summarization: Word Network

- Make networks look better



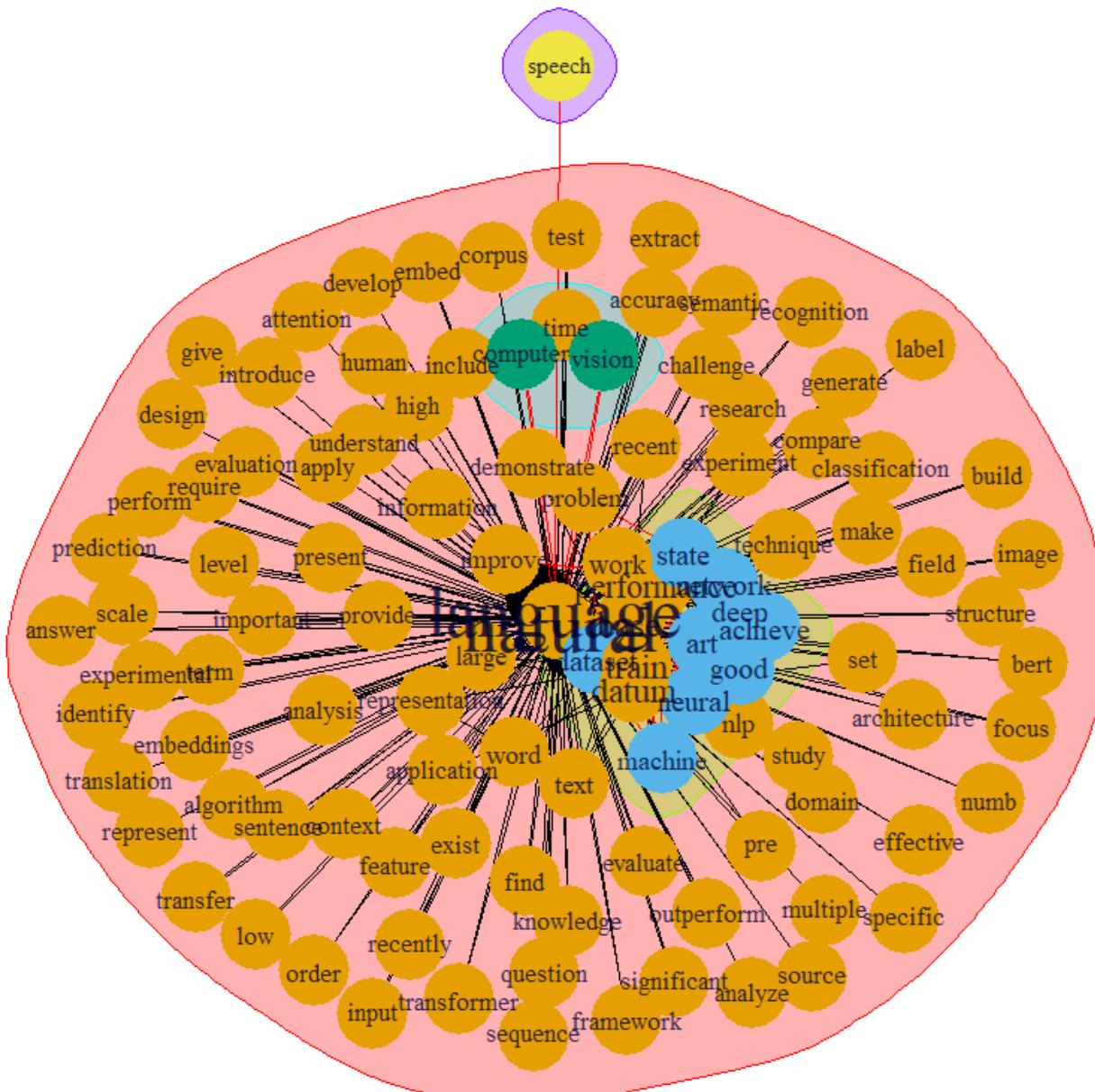
Text Summarization: Word Network

- Find community

```
# Find community
nlp_community <- walktrap.community(nlp_graph, steps = 3)
modularity(nlp_community)
membership(nlp_community)
plot(nlp_community, nlp_graph)
```

Text Summarization: Word Network

- Find community



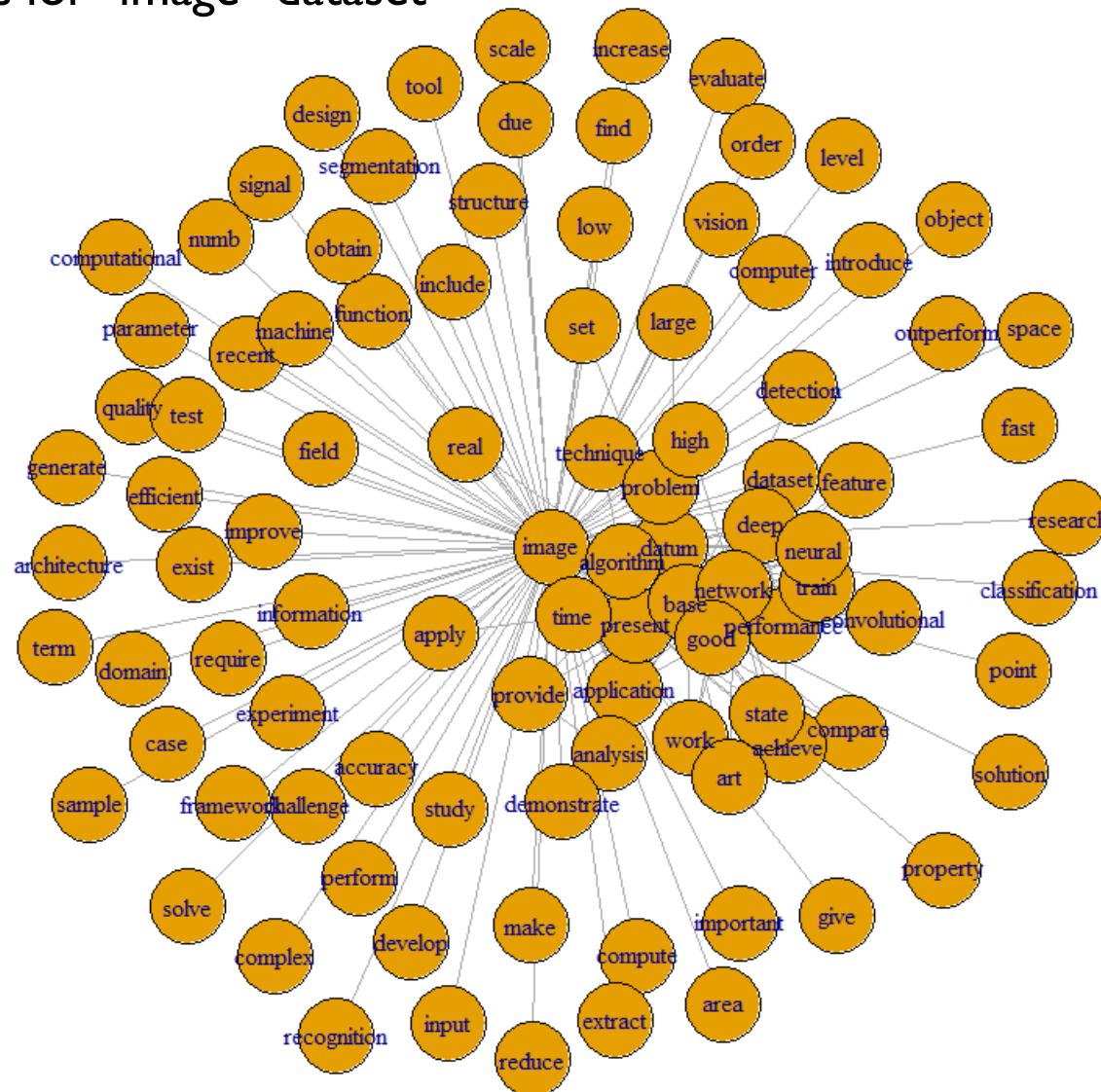
Text Summarization: Word Network

- Plot some networks for “image” dataset

```
# set seed to make the layout reproducible
set.seed(3952)
plot(image_graph, layout=layout.fruchterman.reingold)
plot(image_graph, layout=layout.kamada.kawai, vertex.size = 10,
      vertex.color = 8, vertex.label.cex = 1)
```

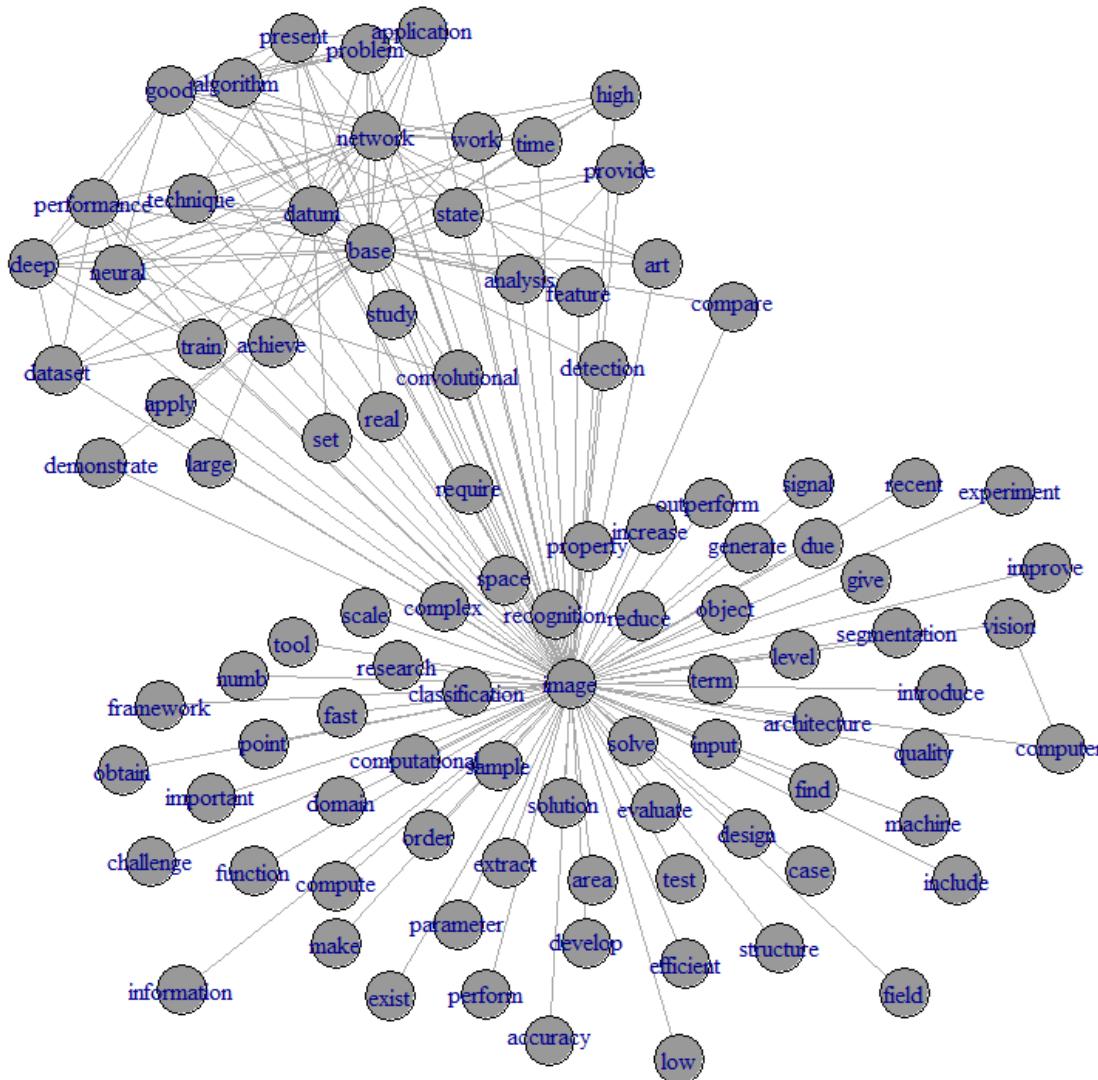
Text Summarization: Word Network

- Plot some networks for “image” dataset



Text Summarization: Word Network

- Plot some networks for “image” dataset



Text Summarization: Word Network

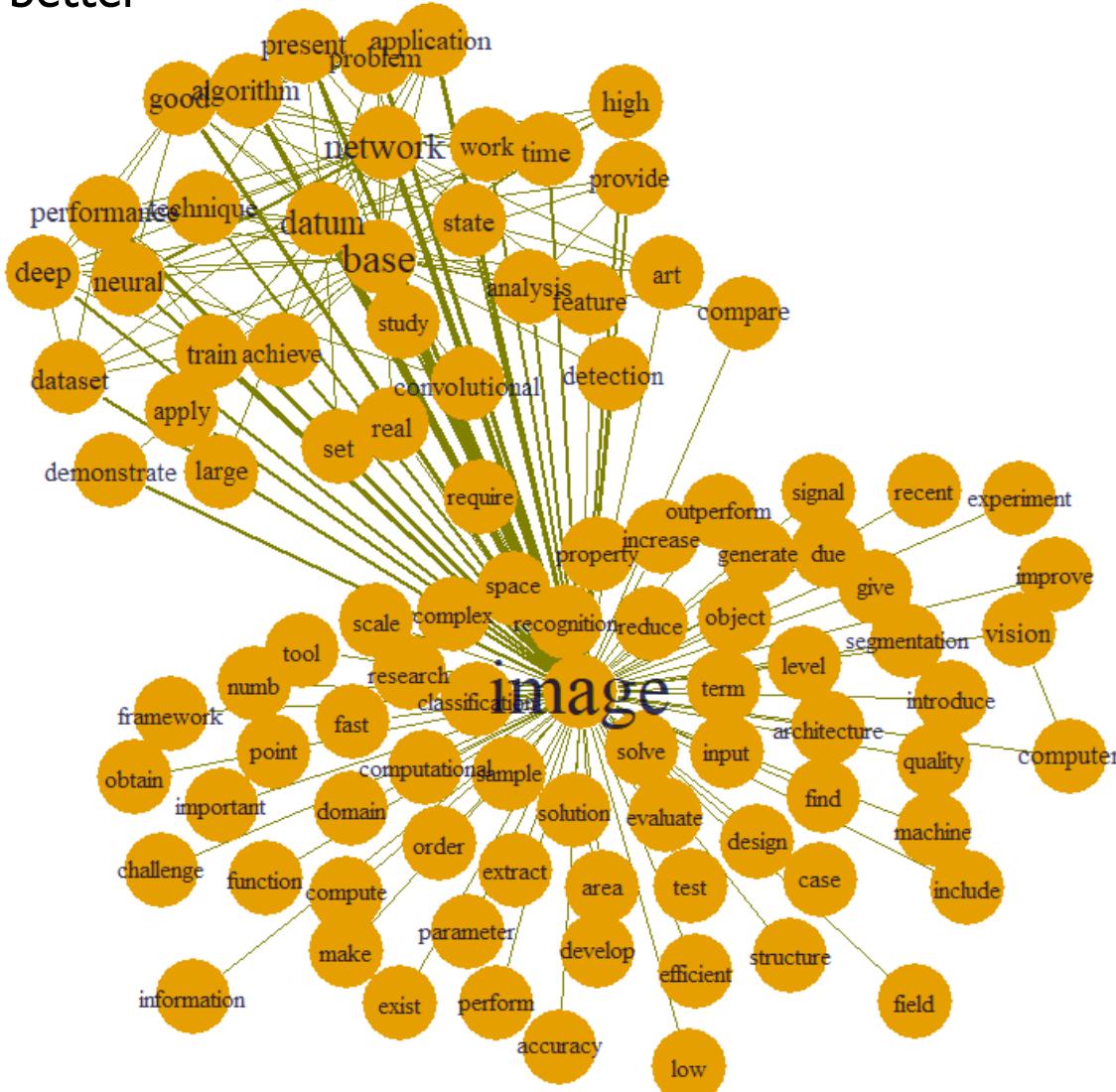
- Make networks look better

```
# Make the network look better
V(image_graph)$label.cex <-
  2*V(image_graph)$degree/max(V(image_graph)$degree)+1
V(image_graph)$label.color <- rgb(0, 0, 0.2, 0.8)
V(image_graph)$frame.color <- NA
egam1 <-
  300*(exp(E(image_graph)$weight/100))/sum(exp(E(image_graph)$weight/100))
E(image_graph)$color <- rgb(0.5, 0.5, 0)
E(image_graph)$width <- egam1

# plot the graph in layout1
plot(image_graph, layout=layout.kamada.kawai)
```

Text Summarization: Word Network

- Make networks look better



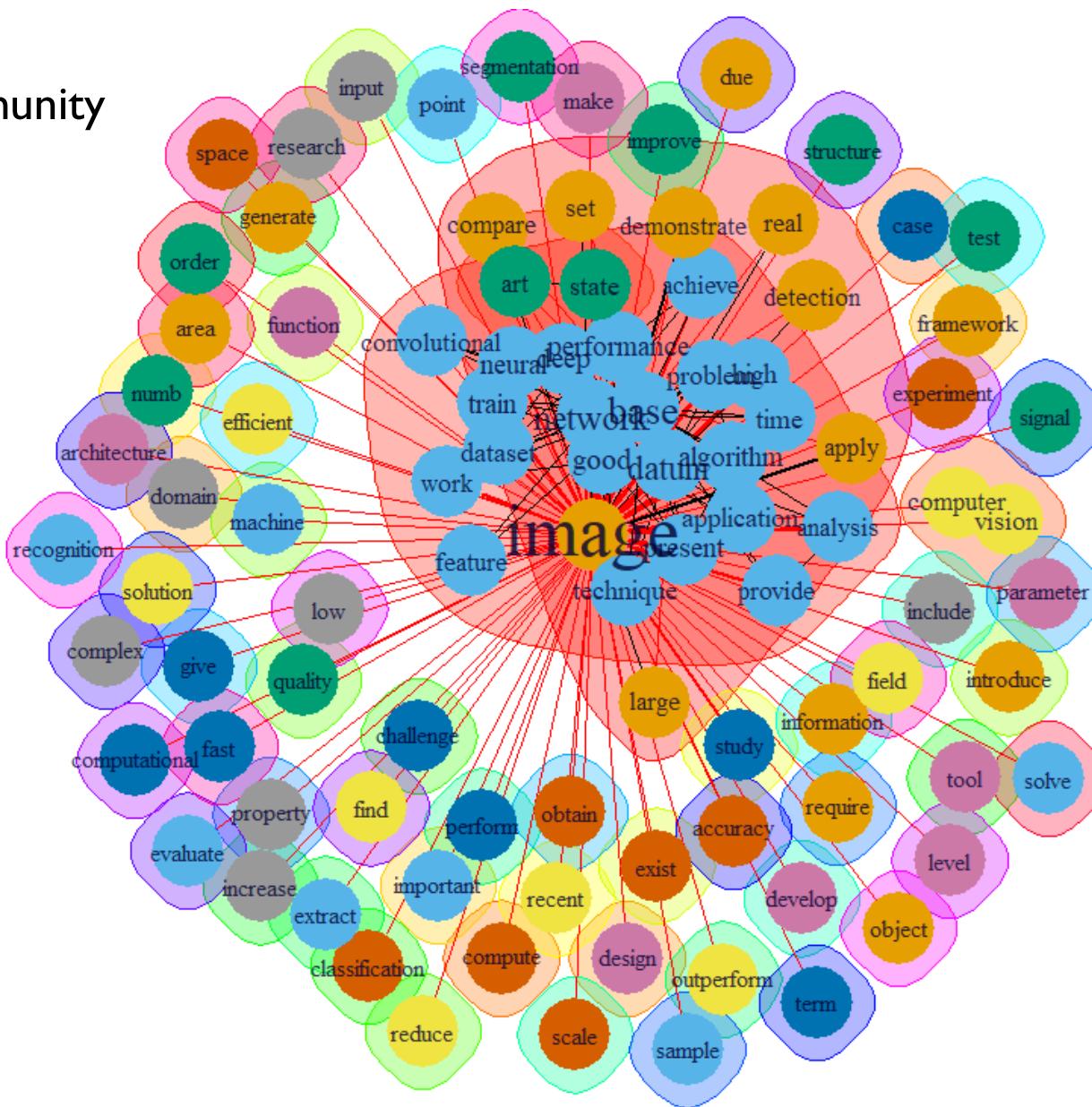
Text Summarization: Word Network

- Find community

```
# Find community
image_community <- walktrap.community(image_graph, steps = 3)
modularity(image_community)
membership(image_community)
plot(image_community, image_graph)
```

Text Summarization: Word Network

- Find community





ANY
questions?