

# Lecture 5: Document Summarization

# Pilsung Kang

# School of Industrial Management Engineering Korea University

# AGENDA

01 Word Cloud

---

02 Association Rules

---

03 Network Representation

---

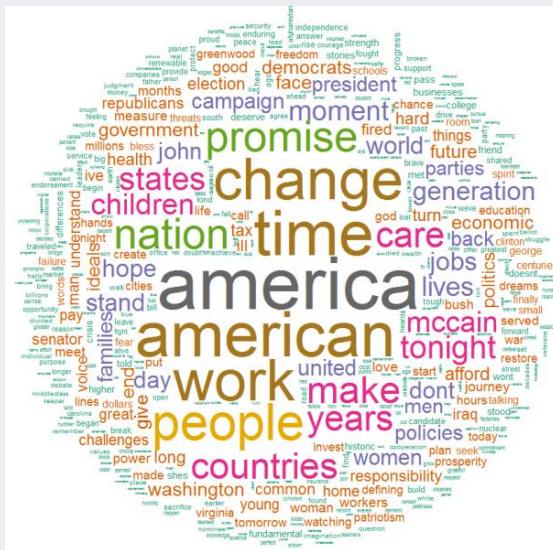
04 R Exercise

---

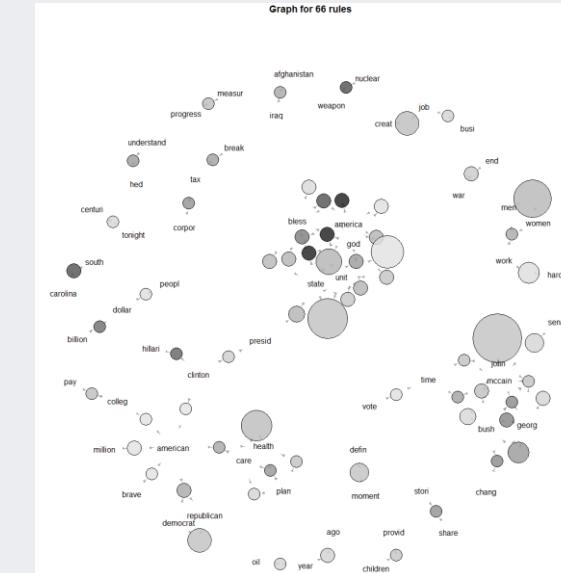
# Background

- How to visualize a large amount of text data at a glance?

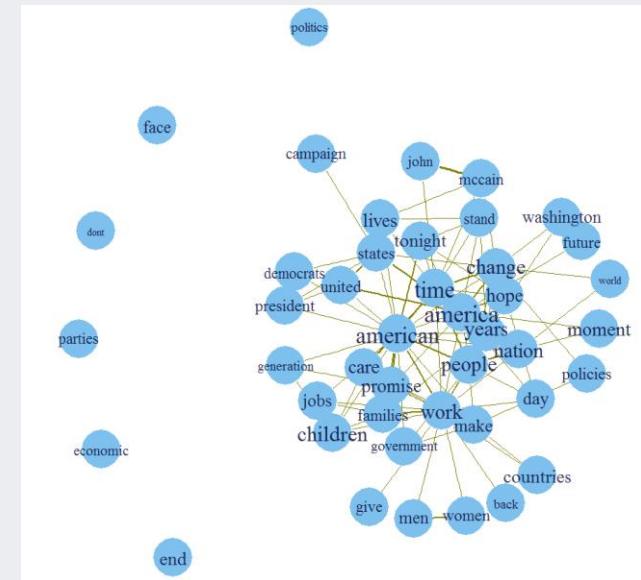
## Wordcloud



## Association Rules



## Keywords Network



	Frequency	Sentence/Phrase	Co-occurrence	Statistical significance
Wordcloud	○	✗	△	✗
Association Rules	○	○	○	○
Keyword Network	○	○	○	✗

# Wordcloud (Tagcloud)

- 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/>

# 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

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

# AGENDA

01 Word Cloud

---

02 Association Rules

---

03 Network Representation

---

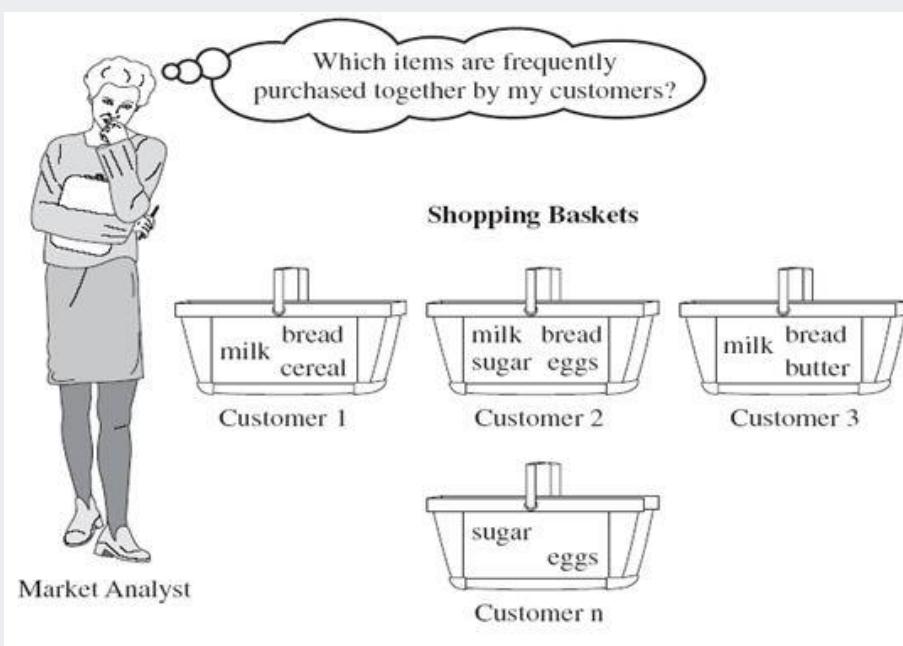
04 R Exercise

---

# Association rules

- Goal:

- ✓ Produce rules that define “what goes with what”
- ✓ “If X was purchased, then Y was also purchased” (in Market Basket Analysis; MBA)
- ✓ “If a word X is presented in a sentence (phrase), then a word Y is also presented in the same sentence (phrase)” (in Text Mining; TM)



> inspect(ares.pruned)		support	confidence	lift
1	{기적}	=> {한강}	0.04494382	1.0000000 22.250000
2	{기술}	=> {과학}	0.04494382	1.0000000 17.800000
3	{경제, 기술}	=> {창조}	0.03370787	1.0000000 11.125000
4	{경제, 과학}	=> {창조}	0.04494382	1.0000000 11.125000
5	{국민, 신뢰}	=> {정부}	0.03370787	1.0000000 11.125000
6	{융성}	=> {문화}	0.03370787	1.0000000 8.900000
7	{과학}	=> {창조}	0.04494382	0.8000000 8.900000
8	{민국}	=> {대한}	0.13483146	1.0000000 7.416667
9	{창조}	=> {경제}	0.08988764	1.0000000 6.846154
10	{오늘}	=> {민국}	0.04494382	0.8000000 5.933333
11	{오늘}	=> {대한}	0.04494382	0.8000000 5.933333
12	{과학}	=> {경제}	0.04494382	0.8000000 5.476923
13	{존경}	=> {여러분}	0.04494382	0.8000000 5.085714
14	{경제부흥, 국민}	=> {행복}	0.03370787	1.0000000 4.944444
15	{여러분, 히마}	=> {기적}	0.03370787	1.0000000 4.944444

# Association rules

- Features

- ✓ Rows are transactions (in MBA) or sentences/phrase (in TM)
- ✓ A Synthetic Example
  - 6 keywords in 10 sentences

Sentence	Word 1	Word 2	Word 3	Word 4
S <sub>1</sub>	Love	Movie	Football	
S <sub>2</sub>	Movie	Watch		
S <sub>3</sub>	Movie	Sleep		
S <sub>4</sub>	Love	Movie	Watch	
S <sub>5</sub>	Love	Sleep		
S <sub>6</sub>	Movie	Sleep		
S <sub>7</sub>	Movie	Watch		
S <sub>8</sub>	Love	Movie	Sleep	Football
S <sub>9</sub>	Love	Movie	Sleep	
S <sub>10</sub>	Party			

# Association rules

- Terminology
  - ✓ Antecedent – “IF” part
  - ✓ Consequent – “THEN” part
  - ✓ Item set – the items comprising the antecedent or consequent
  - ✓ Antecedent and consequent are disjoint (have no items in common)
- Generating rules
  - ✓ Many rules are possible (e.g., for sentence I)
    - If Love is presented, then Movie is also presented.
    - If Love and Movie are presented, then Football is also presented.
    - If Football is presented, then Love is also presented.
    - etc.

# Association rules: Performance measures

For the rule  $A \rightarrow B$

- Support

$$\text{Support}(A) = P(A) \text{ or } \text{Support}(A \rightarrow B) = P(A, B)$$

✓ Used to **find the frequent item sets**

- Confidence

$$\text{Confidence}(A \rightarrow B) = \frac{P(A, B)}{P(A)}$$

✓ Used to **generate meaningful rules**

- Lift

$$\text{Lift}(A \rightarrow B) = \frac{P(A, B)}{P(A) \times P(B)}$$

✓ Used to **evaluate the statistical significance of the generated rules**

- If lift = 1, then the antecedent and the consequents are statistically independent
- If lift > 1, then the rule is useful in finding consequent item sets

# Association rules

- How to generate effective association rules?
  - ✓ Ideally, create all possible combinations of items and see what rules are effective and what rules are not.
  - ✓ Computation time grows exponentially as the number of items increases.
- Brute-force approach
  - ✓ List all possible association rules
  - ✓ Compute the support and confidence for each rule
  - ✓ Prune rules that fail the minimum support and minimum confidence threshold
  - ✓ Computationally prohibitive!

# Association rules

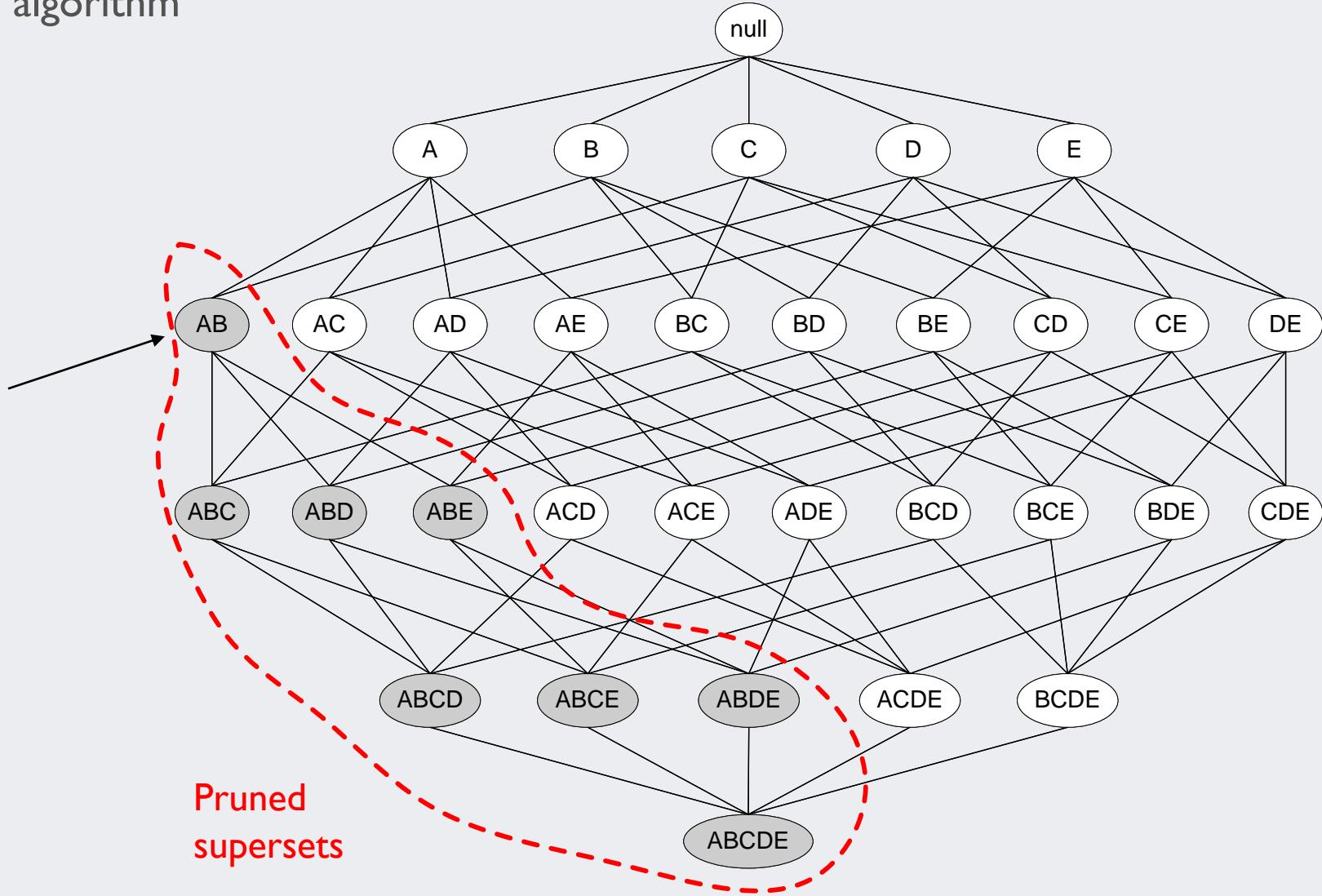
- A priori algorithm
  - ✓ Consider only “frequent item sets”
  - ✓ Support
    - Criterion for item set frequency  $P(A)$
    - #(%) of sentences that include the antecedent (both the antecedent and the consequent)
    - Support for the item set {Love, Movie} is 4 out of 10 sentences, or 40%
  - ✓ Support of an itemset never exceeds the support of its subsets, which is known as anti-monotone property of support.

# Association rules

- A Priori algorithm

Found to be  
Infrequent

Pruned  
supersets



# Association Rules: Generating Frequent Item Sets

- Set a minimum support criterion
  - Set the minimum support to 2 sentences or 20%

Sentence	Word 1	Word 2	Word 3	Word 4
S <sub>1</sub>	Love	Movie	Football	
S <sub>2</sub>	Movie	Watch		
S <sub>3</sub>	Movie	Sleep		
S <sub>4</sub>	Love	Movie	Watch	
S <sub>5</sub>	Love	Sleep		
S <sub>6</sub>	Movie	Sleep		
S <sub>7</sub>	Movie	Watch		
S <sub>8</sub>	Love	Movie	Sleep	Football
S <sub>9</sub>	Love	Movie	Sleep	
S <sub>10</sub>	Party			

# Association Rules: Generating Frequent Item Sets

2

- Generate the list of one-item sets that meets the support criterion

- Support {Movie} = 8/10 = 80%
- Support {Love} = 5/10 = 50%
- Support {Sleep} = 5/10 = 50%
- Support {Watch} = 3/10 = 30%
- Support {Football} = 2/10 = 20%
- Support {Party} = 1/10 = 10%

Party is removed because it does not meet the minimum support criterion

# Association Rules: Generating Frequent Item Sets

- Use the life of one-item sets to generate list of two-item sets that meet the support criterion

3

	Movie	Love	Sleep	Watch	Football
Movie		40%	40%	20%	20%
Love			30%	0%	20%
Sleep				0%	10%
Watch					0%
Football					

- Among the 10 possible item sets, six of them are still found to be frequent

# Association Rules: Generating Frequent Item Sets

4

- Use the list of two-item sets to generate the three-item sets.
- Continue up through k-item sets.

Set-size	Word 1	Word 2	Word 3	..	Word 6
1	Movie				
1	Love				
1	Sleep				
1	Watch				
1	Football				
2	Movie	Love			
2	Movie	Sleep			
2	Movie	Watch			
...	...	...			

# Association Rules

- A priori algorithm
  - ✓ Let  $k=1$
  - ✓ Generate frequent itemsets of length  $l$
  - ✓ Repeat until no new frequent itemsets are identified
    - Generate length  $(k+1)$  candidate itemsets from length  $k$  frequent itemsets
    - Prune candidate itemsets containing subsets of length  $k$  that are infrequent
    - Count the support of each candidate by scanning the DB
    - Eliminate candidates that are infrequent, leaving only those that are frequent

# Association Rules: Result

- Generated Rules

Rule: If all Antecedent items are purchased, then with Confidence percentage Consequent items will also be purchased.							
Row ID	Confidence %	Antecedent (A)	Consequent (C)	Support for A	Support for C	Support for A & C	Lift Ratio
6	100	Football	Love & Movie	2	4	2	2.5
2	100	Football	Love	2	5	2	2
4	100	Movie & Football	Love	2	5	2	2
3	100	Football	Movie	2	8	2	1.25
5	100	Love & Football	Movie	2	8	2	1.25
7	100	Watch	Movie	3	8	3	1.25
1	80	Love	Movie	5	8	4	1
8	80	Sleep	Movie	5	8	4	1

## ✓ Interpretation

- Support for A & C = 2 → There are two sentences that the words Football, Love, and Movie are presented together.
- 100% Confidence → If Football is presented in a sentence, then it is always that Love and Movie are also presented.
- Lift Ratio 2.5 → The association between the two item sets are 2.5 stronger than when they are assumed to be statistically independent.

# AGENDA

01 Word Cloud

---

02 Association Rules

---

03 Network Representation

---

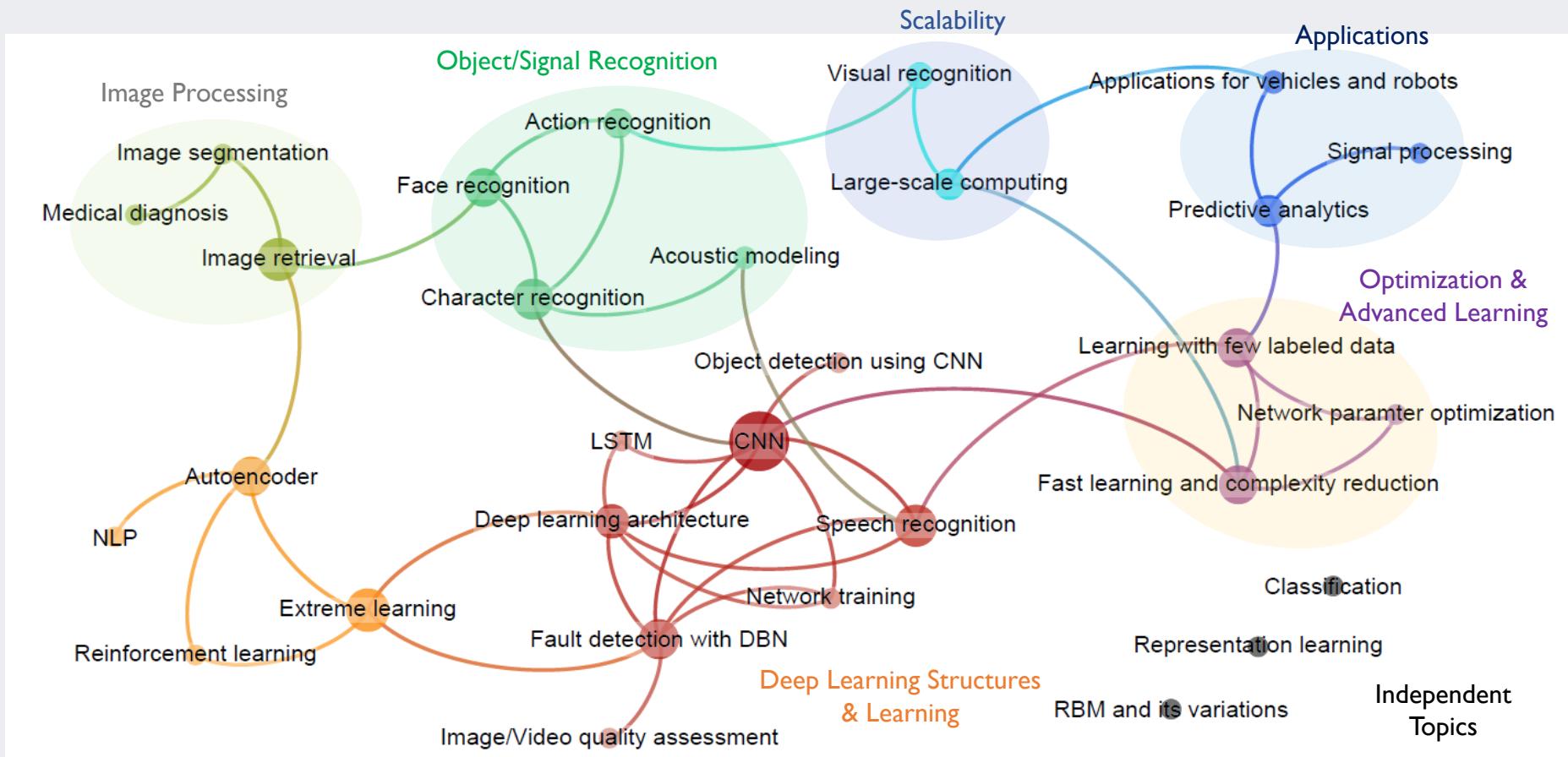
04 R Exercise

---

# Network Representation

Kim et al. (2016)

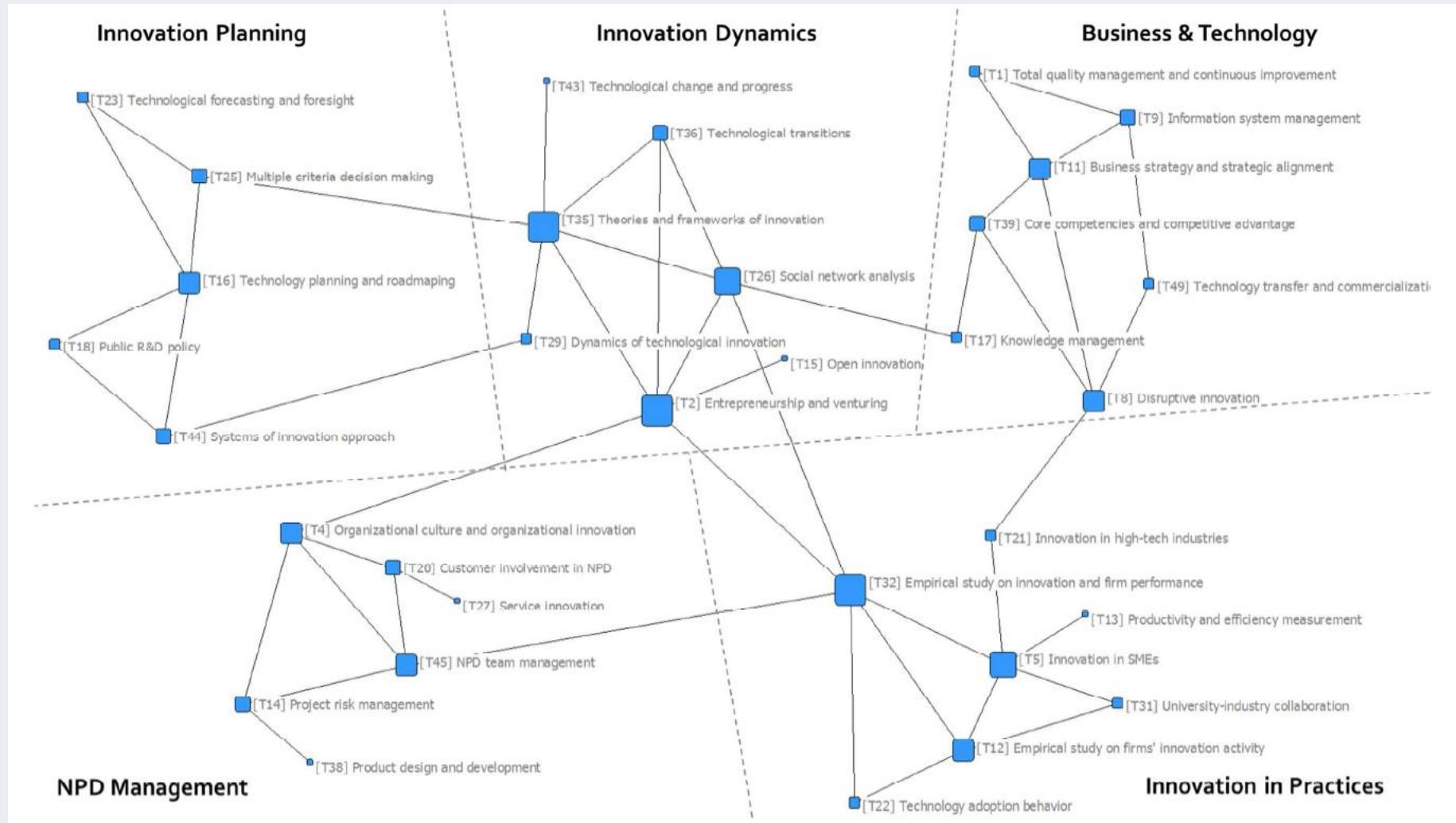
- Network of Deep Learning Topics Until Feb. 2016



# Network Representation

Lee and Kang (2017)

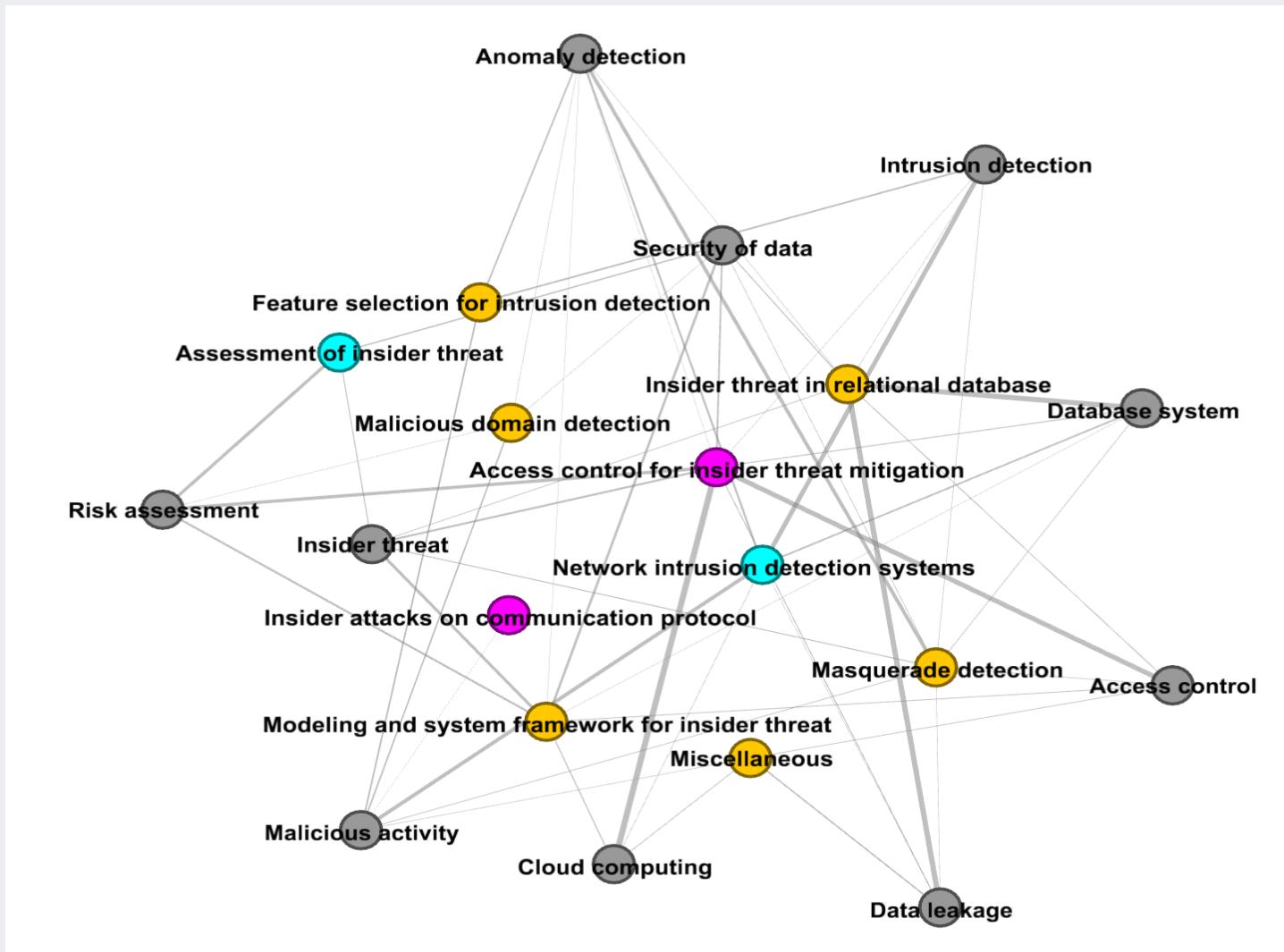
- Network of research topics for “technology and innovation management”



# Network Representation

Kim et al. (2016)

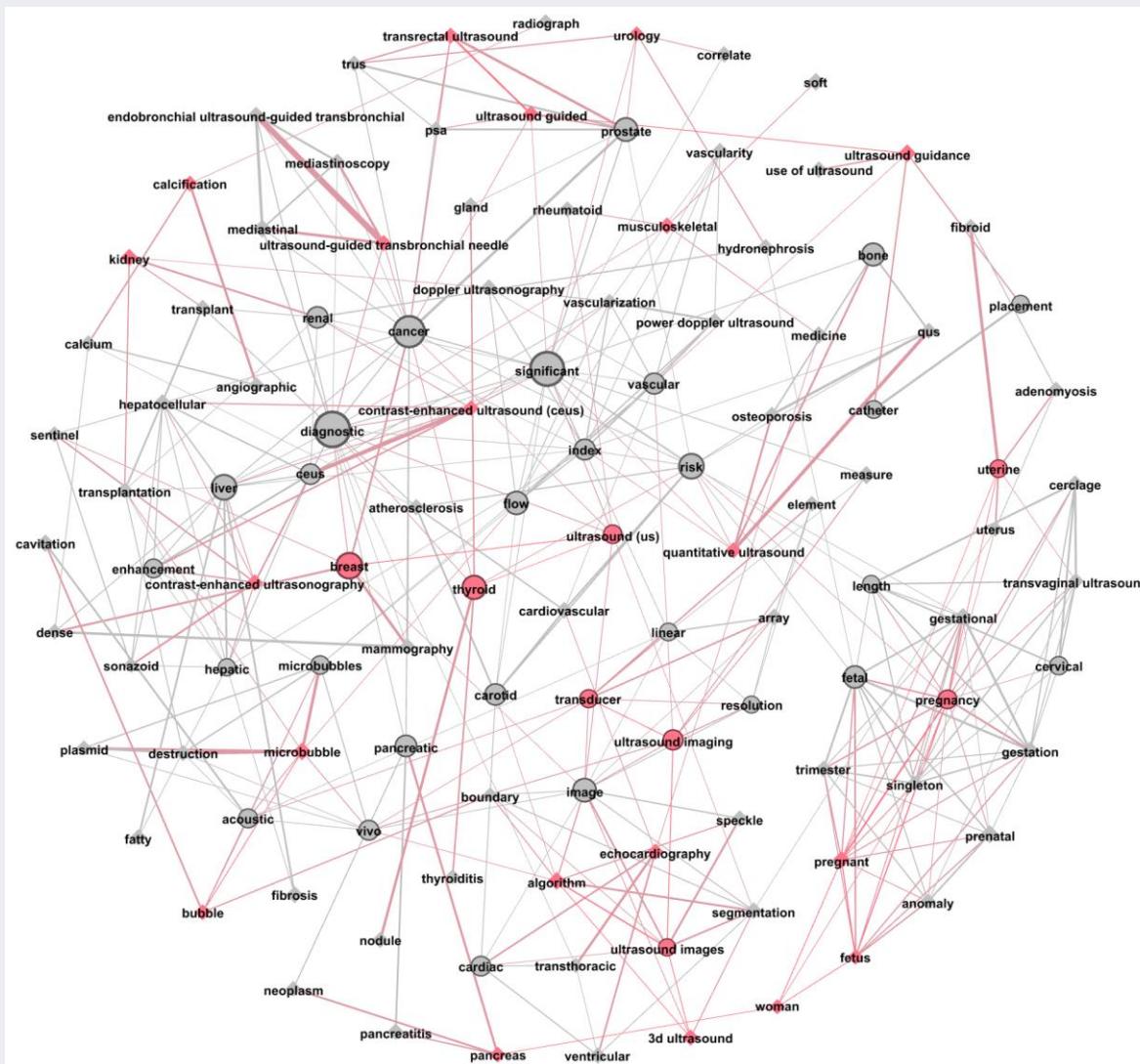
- Network of research topics for “Insider threats”



# Network Representation

Kim et al. (2017)

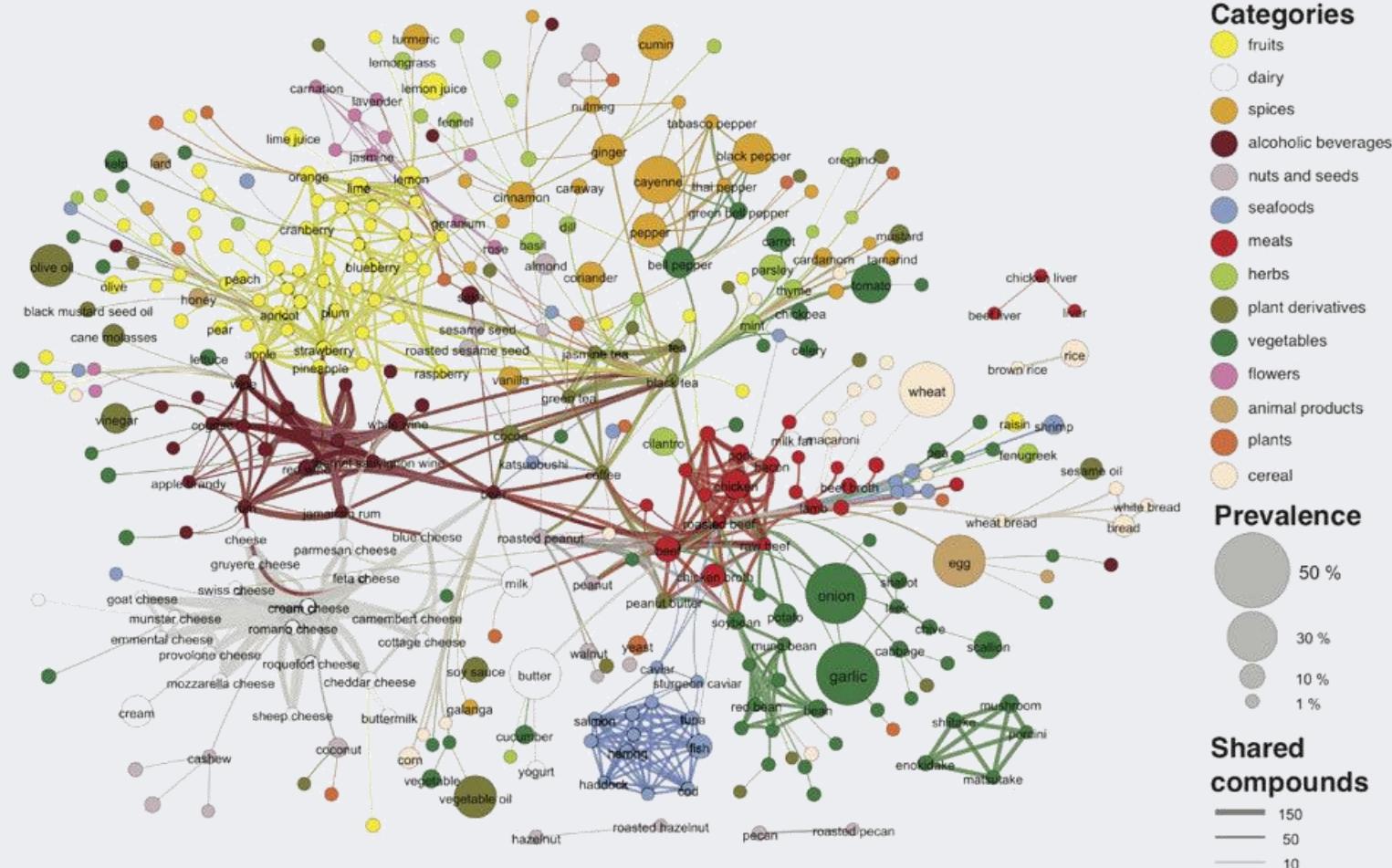
- Network of research topics for “Ultrasound and Ultrasonography”



# Network Representation

Ahn et al. (2011)

- Summarize a collection of text documents using a network
- ✓ The flavor network of the recipes



# Network Representation

Ahn et al. (2011)

- Summarize a collection of text documents using a network
  - ✓ The flavor network of the recipes

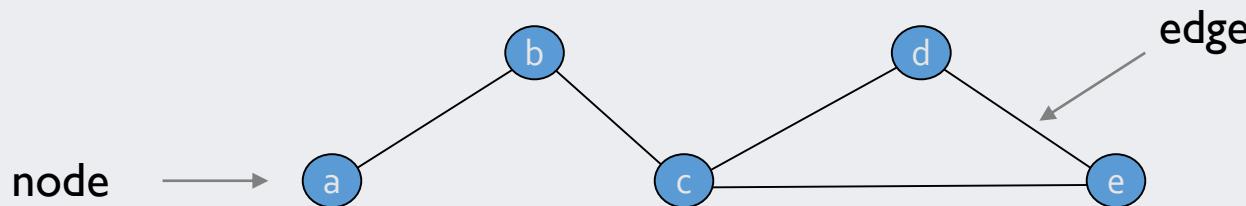
Table S2: Number of recipes and the detailed cuisines in each regional cuisine in the recipe dataset. Five groups have reasonably large size. We use all cuisine data when calculating the relative prevalence and flavor principles.

Cuisine set	Number of recipes	Cuisines included
North American	41525	American, Canada, Cajun, Creole, Southern soul food, Southwestern U.S.
Southern European	4180	Greek, Italian, Mediterranean, Spanish, Portuguese
Latin American	2917	Caribbean, Central American, South American, Mexican
Western European	2659	French, Austrian, Belgian, English, Scottish, Dutch, Swiss, German, Irish
East Asian	2512	Korean, Chinese, Japanese
Middle Eastern	645	Iranian, Jewish, Lebanese, Turkish
South Asian	621	Bangladeshi, Indian, Pakistani
Southeast Asian	457	Indonesian, Malaysian, Filipino, Thai, Vietnamese
Eastern European	381	Eastern European, Russian
African	352	Moroccan, East African, North African, South African, West African
Northern European	250	Scandinavian

# Network Representation

- What is a Network?

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

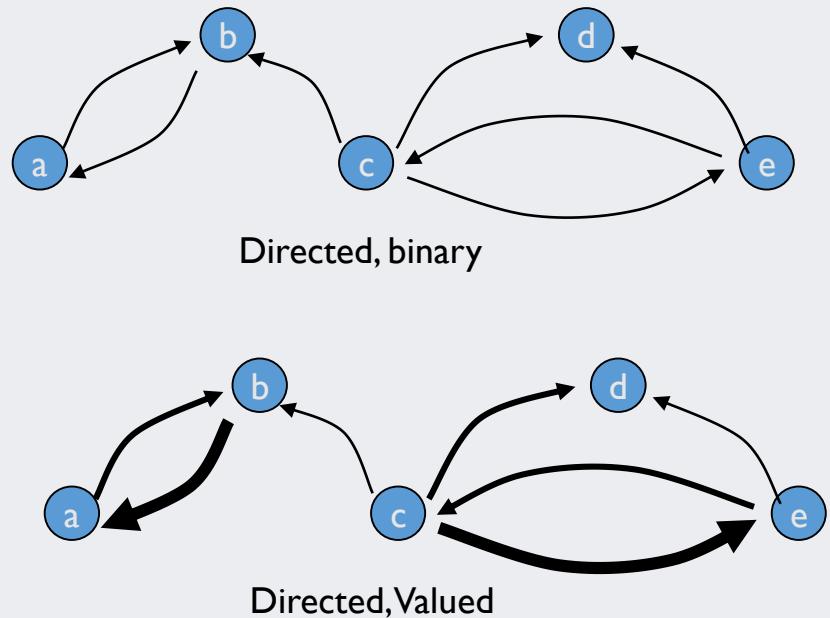
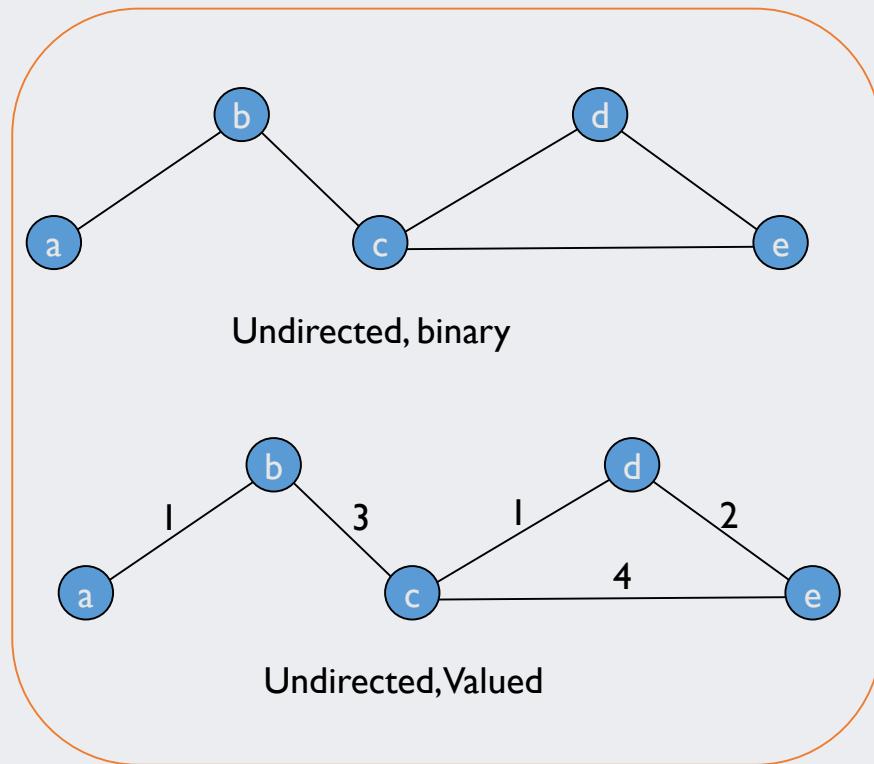


Points	Lines	Domain
Vertices	Edges, arcs	Math
Nodes	Links	Computer Science
Sites	Bonds	Physics
Actors	Ties, relations	Sociology
Keyword	Co-occurrence	Text Mining

# Network Representation: Connections

- Type of Connections

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

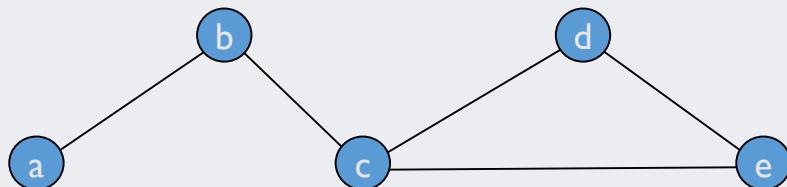


“Undirected networks are commonly used for Text Mining”

# Network Representation: Data Structure

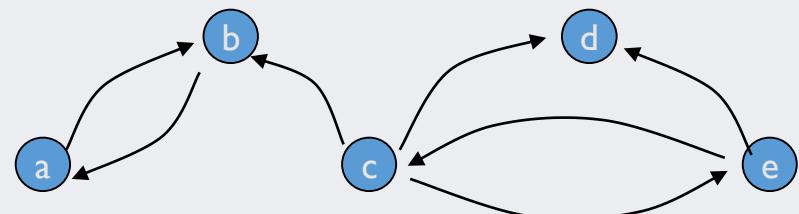
- Basic Data Structure

- ✓ From pictures to matrices



Undirected, binary

	a	b	c	d	e
a					
b					
c					
d					
e					



Directed, binary

	a	b	c	d	e
a					
b					
c					
d					
e					

# Network Representation: Data Structure

- Basic Data Structure

- ✓ From matrices to lists to save memory space

	a	b	c	d	e
a		1			
b	1		1		
c		1		1	1
d			1		1
e			1	1	

Adjacency List

a|b  
b|a c  
c|b d e  
d|c e  
e|c d

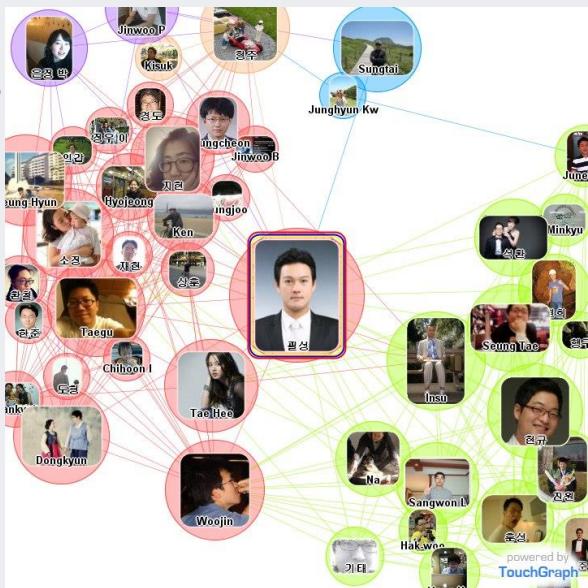
Arc List

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

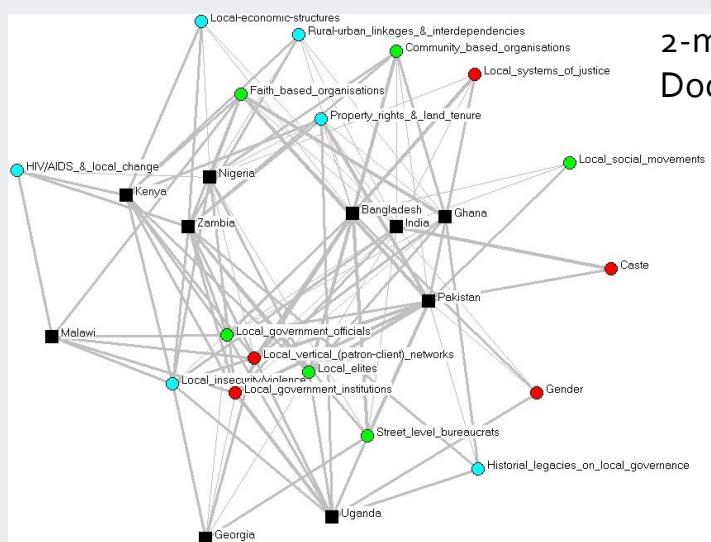
# Network Representation: Data Structure

- N-mode data
  - ✓ 1-mode data represent edges based on **direct** contact between actors in the network
  - ✓ 2-mode data represent **nodes from two separate classes**, where all ties are across classes
    - People as author on papers, events in the life history of people

1-mode  
Facebook friends



2-mode  
Documents & nation

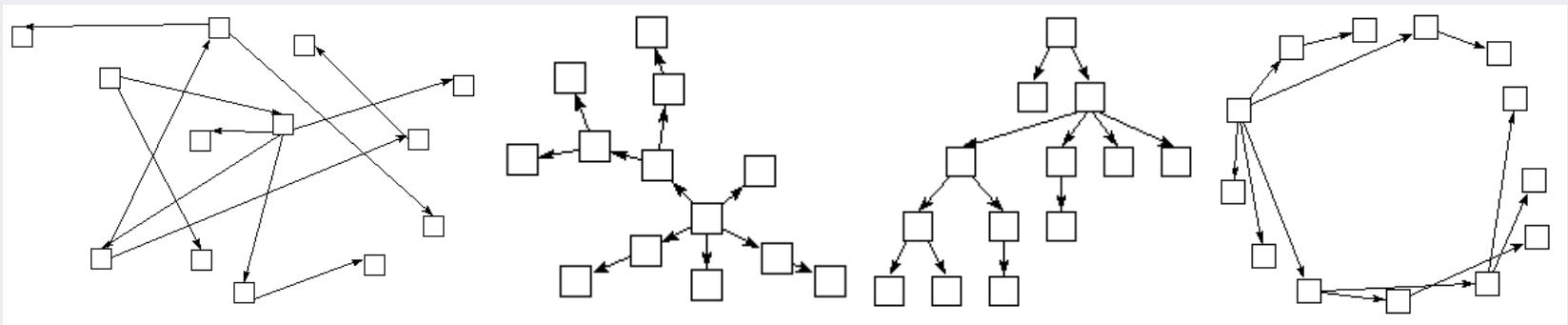


# Network Representation: Layout

- Graphical Representation

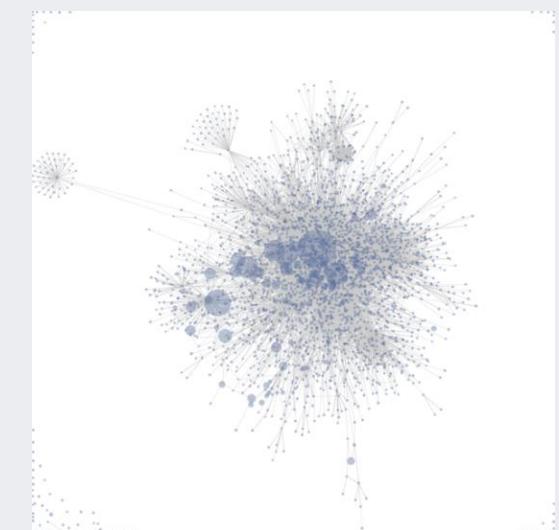
- ✓ No standard way to draw a sociogram

- 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



# Network Representation: Measuring Networks

- 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

# Network Representation: Measuring Networks

- Individual Node Level

- ✓ **Connectivity**

- refers to how actors in one part of the network are connected to actors in another part of the network

- **Reachability**

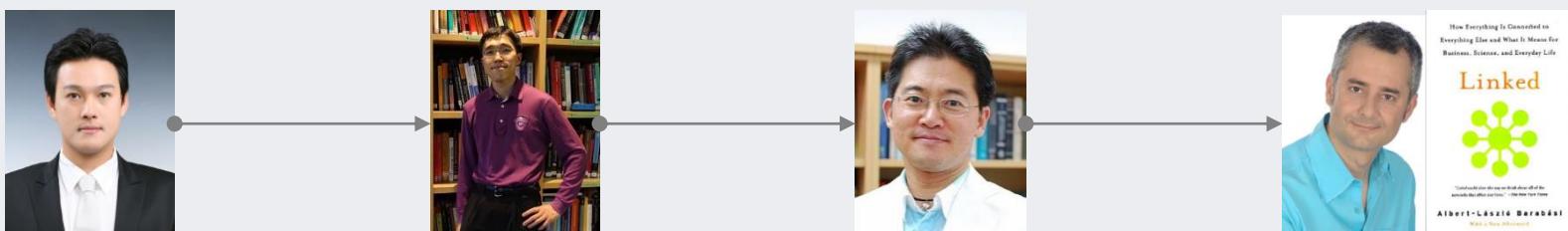
- Is it possible for actor i to reach actor j?
    - True only if there is a chain of contact from one actor to another

- **Distance**

- Given they can be reached, how many steps are they from each other?

- **Number of paths**

- How many different paths connect each pair?



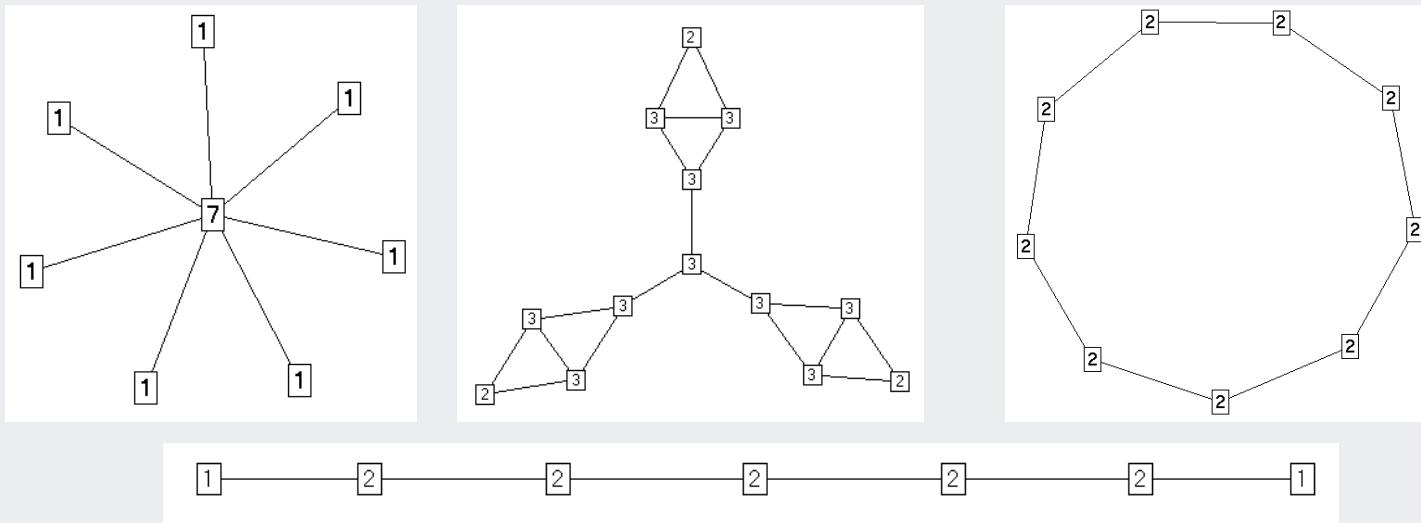
# Network Representation: Measuring Networks

- Individual Node Level

- ✓ **Centrality**

- Refers to **location**, identifying where an actor resides in a network
    - Commonly believed that actors in the “**center**” of the network are “**important**”
    - **Degree**
      - the number of edges connected to the actor

$$C_D = d(n_i) = X_{i+} = \sum_j X_{ij}$$



# Network Representation: Measuring Networks

- Individual Node Level

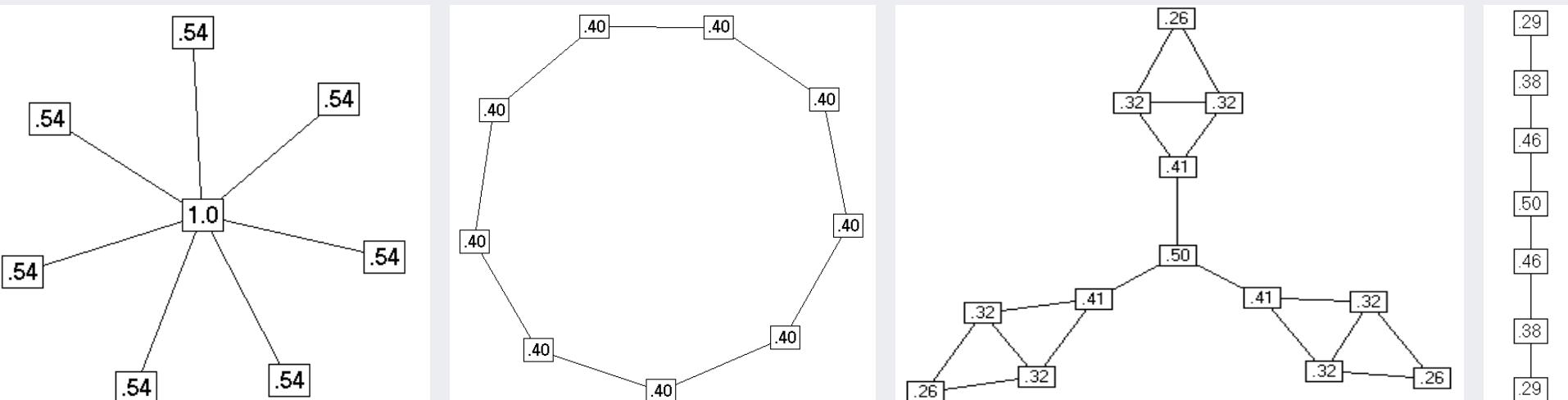
- ✓ Centrality

- Closeness

- An actor is considered important if it is relatively close to all other actors
    - based on the inverse of the distance of each actor to every other actor in the network, often normalized

$$C_c(n_i) = \left[ \sum_{j=1}^g d(n_i, n_j) \right]^{-1}$$

$$C'_c(n_i) = C_c(n_i)(g - 1)$$



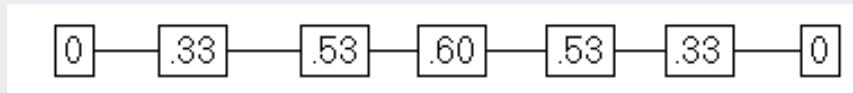
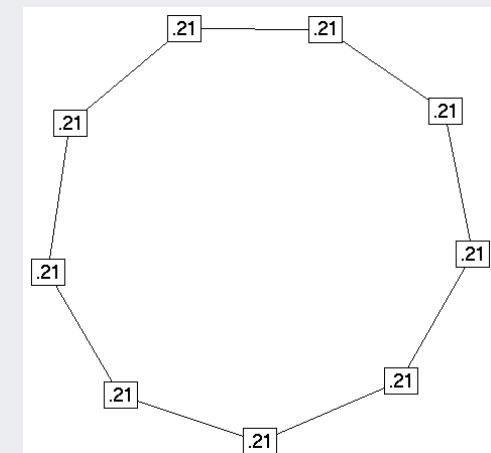
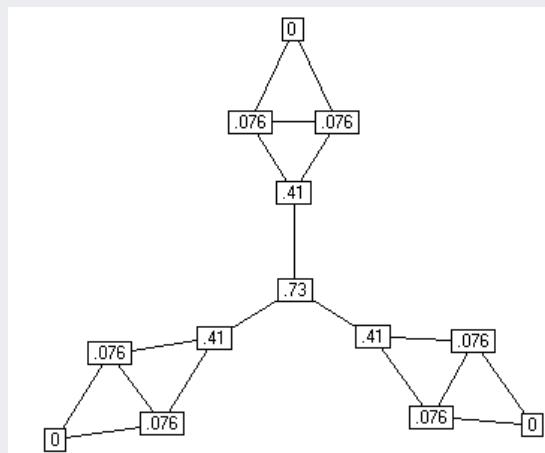
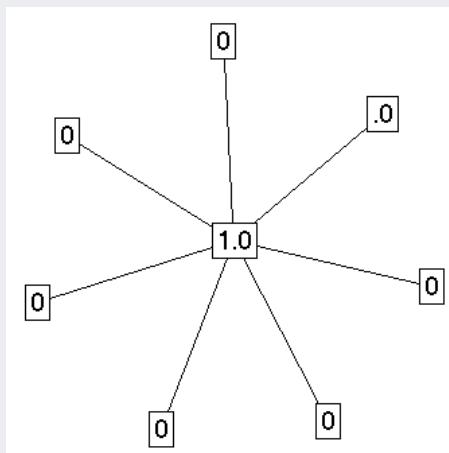
# Network Representation: Measuring Networks

- Individual Node Level

- ✓ Centrality

- Betweenness

- The number of shortest paths between i and k that actor j resides on

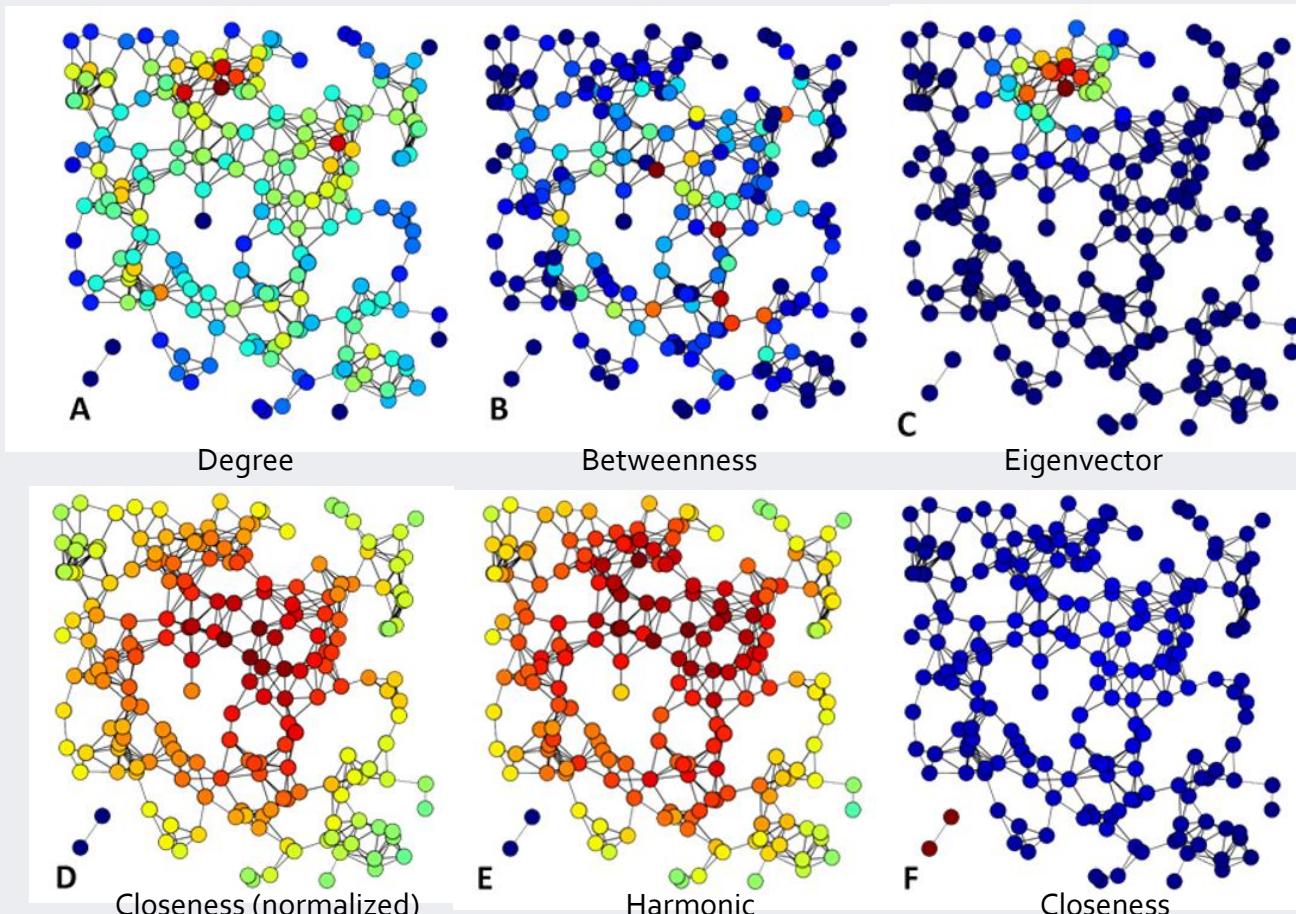


# Network Representation: Measuring Networks

- Individual Node Level

- ✓ Centrality

- Eigenvector centrality, Harmonic centrality

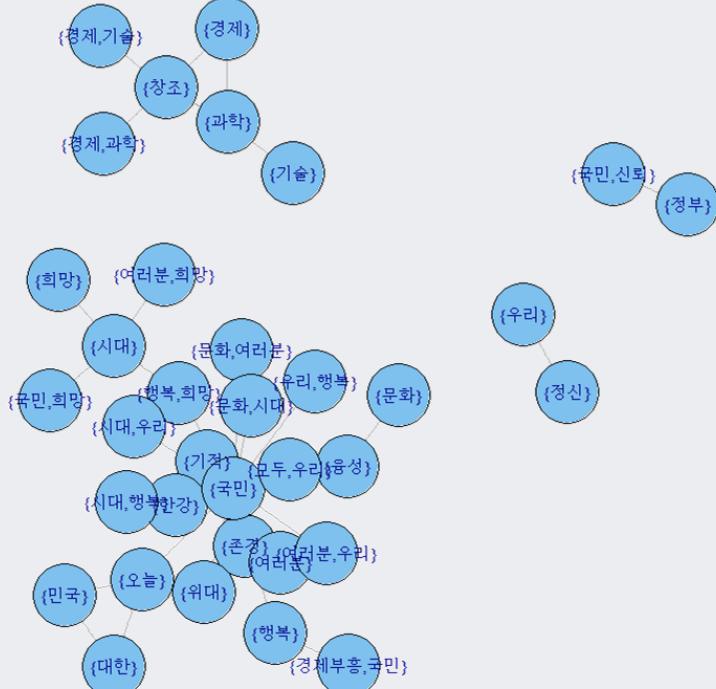


# Network Representation: Community Structure

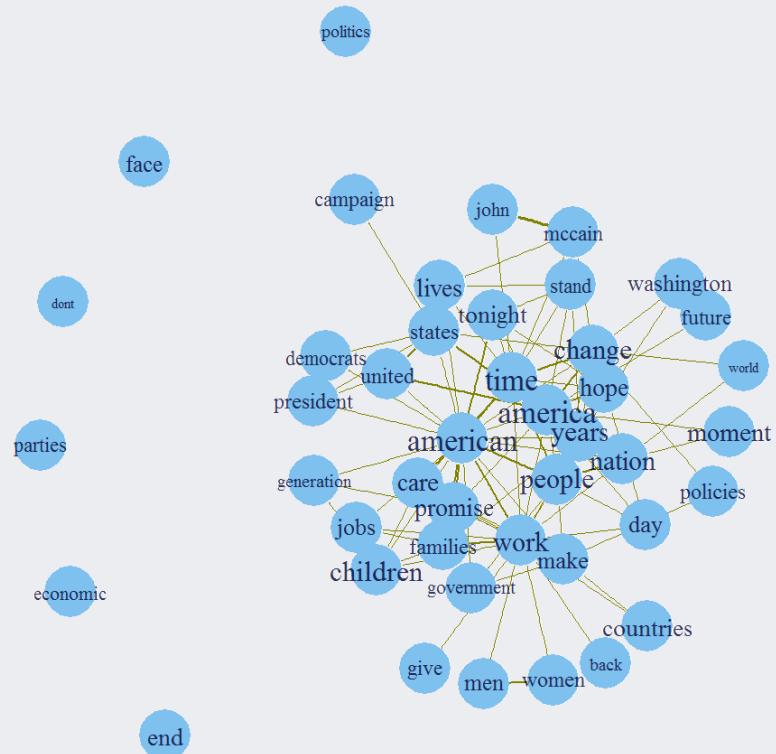
- Community Structure

- ✓ 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



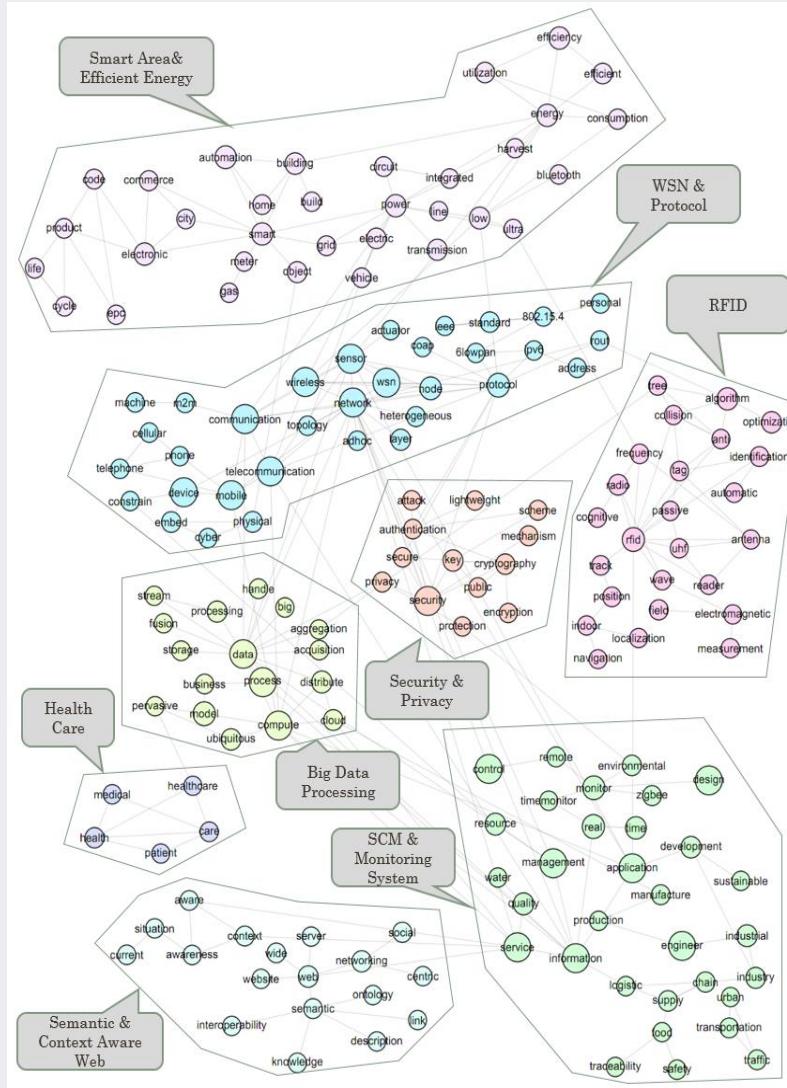
Difficult to identify the community structure



# Network Representation: Community Structure

Kim & Kang (2016)

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



# AGENDA

01 Word Cloud

---

02 Association Rules

---

03 Network Representation

---

04 R Exercise

---

# R Exercise: Word Cloud

- Dataset: Abstracts of the papers published since 2015
  - ✓ Journal 1: IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)
  - ✓ Journal 2: Journal of Banking and Finance (JoF)

The screenshot shows the IEEE Xplore Digital Library interface. At the top, there are links for 'BROWSE', 'MY SETTINGS', 'GET HELP', and 'WHAT CAN I ACCESS?'. Below this is a search bar with 'Enter Search Term' and a 'Search' button. Further down are buttons for 'Basic Search', 'Author Search', 'Publication Search', 'Advanced Search', and 'Other Search Options'. A banner for 'Global KU Frontier Spirit' is visible. The main content area displays the title 'Pattern Analysis and Machine Intelligence, IEEE Transactions on' and a 'Submit Your Manuscript' button. On the right, there are three performance metrics: Impact Factor (5.694), Eigenfactor (.04888), and Article Influence Score (3.155). A sidebar on the left provides information about the journal's aims and scope.

The screenshot shows the Journal of Banking & Finance website. At the top, there are social media icons and a navigation menu with 'Home', 'Journals', and 'Journal of Banking & Finance'. The main title 'Journal of Banking & Finance' is prominently displayed. Below it, there is a brief description: 'The Journal of Banking and Finance (JBF) publishes theoretical and empirical research papers spanning all the major research fields in finance and banking. The aim of the Journal of Banking and Finance is to provide an outlet for the increasing flow of scholarly research concerning financial institutions...'. There are also links for 'Supports Open Access', 'Editors: Carol Alexander, Geert Bekaert', and 'View Editorial Board'. A sidebar on the right lists various journal features and articles, such as 'Most Downloaded', 'Recent Articles', 'Most Cited', 'Open Access Articles', 'Risk and risk management in the credit card industry', 'Financial innovation: The bright and the dark sides', and 'Risk management, corporate governance, and bank performance in the financial crisis'.

# R Exercise: Word Cloud

- Construct Corpus

```
install.packages("tm", dependencies = TRUE)
install.packages("wordcloud", dependencies = TRUE)
install.packages("arules", dependencies = TRUE)
install.packages("arulesViz", dependencies = TRUE)
install.packages("igraph", dependencies = TRUE)
install.packages("stringi", dependencies = TRUE)

library(tm)
library(wordcloud)
library(arules)
library(arulesViz)
library(igraph)
library(stringi)

# Load the data
TPAMI <- read.csv("IEEE_TPAMI_2015_2017.csv", encoding = "UTF-8", stringsAsFactors = FALSE)
JoF <- read.csv("Journal of Banking and Finance_2015_2017.csv", encoding = "UTF-8", stringsAsFactors = FALSE)

# Construct the corpus for each journal with the abstracts
TPAMI.Corpus <- Corpus(VectorSource(TPAMI$Abstract))
JoF.Corpus <- Corpus(VectorSource(JoF$Abstract))
```

# R Exercise: Word Cloud

- Preprocessing

- ✓ To lower case, remove punctuations/numbers/stopwords, stemming

```
# Preprocessing
# 1: to lower case
TPAMI.Corpus <- tm_map(TPAMI.Corpus, content_transformer(stri_trans_tolower))
JoF.Corpus <- tm_map(JoF.Corpus, content_transformer(stri_trans_tolower))

# 2: remove punctuations
TPAMI.Corpus <- tm_map(TPAMI.Corpus, content_transformer(removePunctuation))
JoF.Corpus <- tm_map(JoF.Corpus, content_transformer(removePunctuation))

# 3. remove numbers
TPAMI.Corpus <- tm_map(TPAMI.Corpus, content_transformer(removeNumbers))
JoF.Corpus <- tm_map(JoF.Corpus, content_transformer(removeNumbers))

# 4. remove stopwords (SMART stopwords list)
myStopwords <- c(stopwords("SMART"))

TPAMI.Corpus <- tm_map(TPAMI.Corpus, removeWords, myStopwords)
JoF.Corpus <- tm_map(JoF.Corpus, removeWords, myStopwords)

# 5. Stemming
TPAMI.Corpus <- tm_map(TPAMI.Corpus, stemDocument)
JoF.Corpus <- tm_map(JoF.Corpus, stemDocument)

# 4. remove stopwords (with frequently used words)
myStopwords <- c(stopwords("SMART"), "financ", "american", "associ", "firm",
                  "model", "data", "algorithm", "method", "imag", "ieee", "propos",
                  "elsevi", "problem", "result", "find", "paper")

TPAMI.Corpus <- tm_map(TPAMI.Corpus, removeWords, myStopwords)
JoF.Corpus <- tm_map(JoF.Corpus, removeWords, myStopwords)
```

# R Exercise: Word Cloud

- Construct Term-Document Matrices

```
# Term-Document Matrix
TPAMI.TDM <- TermDocumentMatrix(TPAMI.Corpus, control = list(minWordLength = 1))
JoF.TDM <- TermDocumentMatrix(JoF.Corpus, control = list(minWordLength = 1))

as.matrix(TPAMI.TDM)[11:30,11:30]
as.matrix(JoF.TDM)[11:30,11:30]
```

Terms	Docs																												
	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30									
function	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0									
import	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	0	1	0	0									
induc	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
interfac	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
lab	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
lead	0	0	0	0	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0									
likelihood	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
line	0	0	0	0	0	0	0	0	2	0	0	1	0	0	0	0	0	0	0	0									
local	0	0	1	0	0	0	0	2	0	0	1	0	1	0	0	0	0	0	1	0									
locat	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
measur	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
multipl	0	0	0	0	0	0	1	0	0	0	0	0	0	0	3	0	0	0	0	0									
multiview	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0									
object	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	5	3									
occur	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
periscop	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
pixel	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0									
pool	0	0	0	0	0	0	0	9	0	0	0	0	0	0	0	0	0	0	0	0									
predat	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
probabilist	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0									

Terms	Docs																												
	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30									
elast	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
empir	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0								
entiti	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
environ	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
equiti	0	2	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0								
establish	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0								
estim	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	1	1	0	0	0									
fdi	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
financi	1	0	0	0	1	3	2	1	1	0	3	1	1	3	0	0	1	2	0	0									
host	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
influenc	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
inform	0	0	0	0	0	2	0	0	0	0	0	1	0	0	0	0	0	0	0	0									
insight	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
locat	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0									
measur	1	1	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	1	4	0									
multin	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
negat	0	0	0	1	0	3	0	0	0	0	0	3	1	0	0	0	0	0	0	0									
polici	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
probabl	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1									
provid	0	0	0	1	0	0	1	0	0	1	0	1	1	0	0	0	0	0	0	0									

# R Exercise: Word Cloud

- Frequently used words for each journal

```
# Frequently used words
findFreqTerms(TPAMI.TDM, lowfreq=100)
findFreqTerms(JoF.TDM, lowfreq=100)
```

```
> findFreqTerms(TPAMI.TDM, lowfreq=100)
[1] "approach"      "demonstr"     "estim"        "function"      "local"        "measur"
[7] "multipl"       "object"       "space"        "task"         "video"        "analysi"
[13] "applic"        "effect"       "framework"    "inform"        "optim"        "present"
[19] "shape"          "show"         "stateoftheart" "base"         "dataset"      "detect"
[25] "experiment"    "learn"        "process"      "track"        "comput"       "discrimin"
[31] "effici"         "evalu"        "exist"        "featur"       "formul"       "general"
[37] "human"          "larg"         "number"       "sampl"        "scene"        "set"
[43] "train"          "work"         "distribut"   "face"         "improv"      "match"
[49] "perform"        "structur"    "visual"       "classif"      "cluster"      "represent"
[55] "compar"         "experi"      "robust"       "challeng"    "recognit"    "segment"
[61] "label"          "graph"       "point"        ""             ""             ""
> findFreqTerms(JoF.TDM, lowfreq=100)
[1] "capit"          "countri"      "effect"       "empir"        "equiti"       "estim"        "financi"      "inform"       "measur"
[10] "polici"         "signific"    "bank"         "credit"      "import"       "larg"         "risk"         "examin"      "fund"
[19] "increas"        "manag"        "posit"       "return"      "level"        "loan"         "investig"    "market"      "relat"
[28] "show"           "stock"        "system"      "impact"      "structur"    "time"         "rate"         "suggest"    "higher"
[37] "investor"       "predict"     "price"       "trade"       "volatil"     "invest"      "perform"     "period"     "portfolio"
[46] "control"        "crisi"        "high"        "liquid"      "asset"        "lower"       "studi"        "evid"        "cost"
[55] "factor"         "option"      "debt"        "econom"    ""             ""             ""             ""             ""
```

# R Exercise: Word Cloud

- Construct Word Clouds

```
# Construct a Word Cloud with IEEE_TPAMI abstracts
TPAMI.wcmat <- as.matrix(TPAMI.TDM)

# calculate the frequency of words
TPAMI.word.freq <- sort(rowSums(TPAMI.wcmat), decreasing=TRUE)
TPAMI.keywords <- names(TPAMI.word.freq)
TPAMI.wcdat <- data.frame(word = TPAMI.keywords, freq = TPAMI.word.freq)

pal <- brewer.pal(8, "Dark2")
wordcloud(TPAMI.wcdat$word, TPAMI.wcdat$freq, min.freq=10, scale = c(5, 0.2),
          rot.per = 0.1, col=pal, random.order=F)

# Construct a Word Cloud with JoF abstracts
JoF.wcmat <- as.matrix(JoF.TDM)

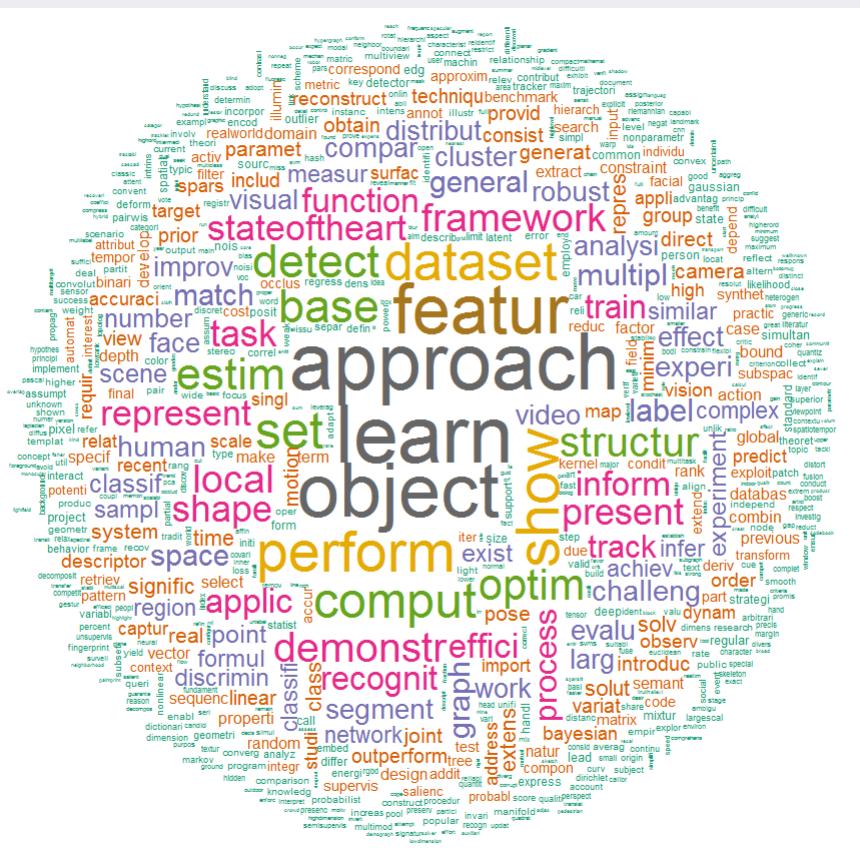
# calculate the frequency of words
JoF.word.freq <- sort(rowSums(JoF.wcmat), decreasing=TRUE)
JoF.keywords <- names(JoF.word.freq)
JoF.wcdat <- data.frame(word = JoF.keywords, freq = JoF.word.freq)

pal <- brewer.pal(8, "Dark2")
wordcloud(JoF.wcdat$word, JoF.wcdat$freq, min.freq=10, scale = c(5, 0.2),
          rot.per = 0.1, col=pal, random.order=F)
```

# R Exercise:Word Cloud

- Construct Word Clouds

TPAMI



JoF



# R Exercise: Association Rules

- find association rules for the words in each journal

```
# Association Rules for IEEE_TPAMI
TPAMI.tran <- as.matrix(t(TPAMI.TDM))
TPAMI.tran <- as(TPAMI.tran, "transactions")

TPAMI.rules <- apriori(TPAMI.tran, parameter=list(minlen=2, supp=0.04, conf=0.75))
inspect(TPAMI.rules)

TPAMI.rules.sorted <- sort(TPAMI.rules, by="lift")
subset.matrix <- is.subset(TPAMI.rules.sorted, TPAMI.rules.sorted)
subset.matrix[lower.tri(subset.matrix, diag=T)] <- NA
redundant <- colSums(subset.matrix, na.rm=T) >= 1
TPAMI.rules.pruned <- TPAMI.rules.sorted[!redundant]
inspect(TPAMI.rules.pruned)

# Plot the rules
plot(TPAMI.rules.pruned, method="graph")

# Association Rules for Journal of Banking and Finance
JoF.tran <- as.matrix(t(JoF.TDM))
JoF.tran <- as(JoF.tran, "transactions")

JoF.rules <- apriori(JoF.tran, parameter=list(minlen=2, supp=0.04, conf=0.7))
inspect(JoF.rules)

JoF.rules.sorted <- sort(JoF.rules, by="lift")
subset.matrix <- is.subset(JoF.rules.sorted, JoF.rules.sorted)
subset.matrix[lower.tri(subset.matrix, diag=T)] <- NA
redundant <- colSums(subset.matrix, na.rm=T) >= 1
JoF.rules.pruned <- JoF.rules.sorted[!redundant]
inspect(JoF.rules.pruned)

# Plot the rules
plot(JoF.rules.pruned, method="graph")
```

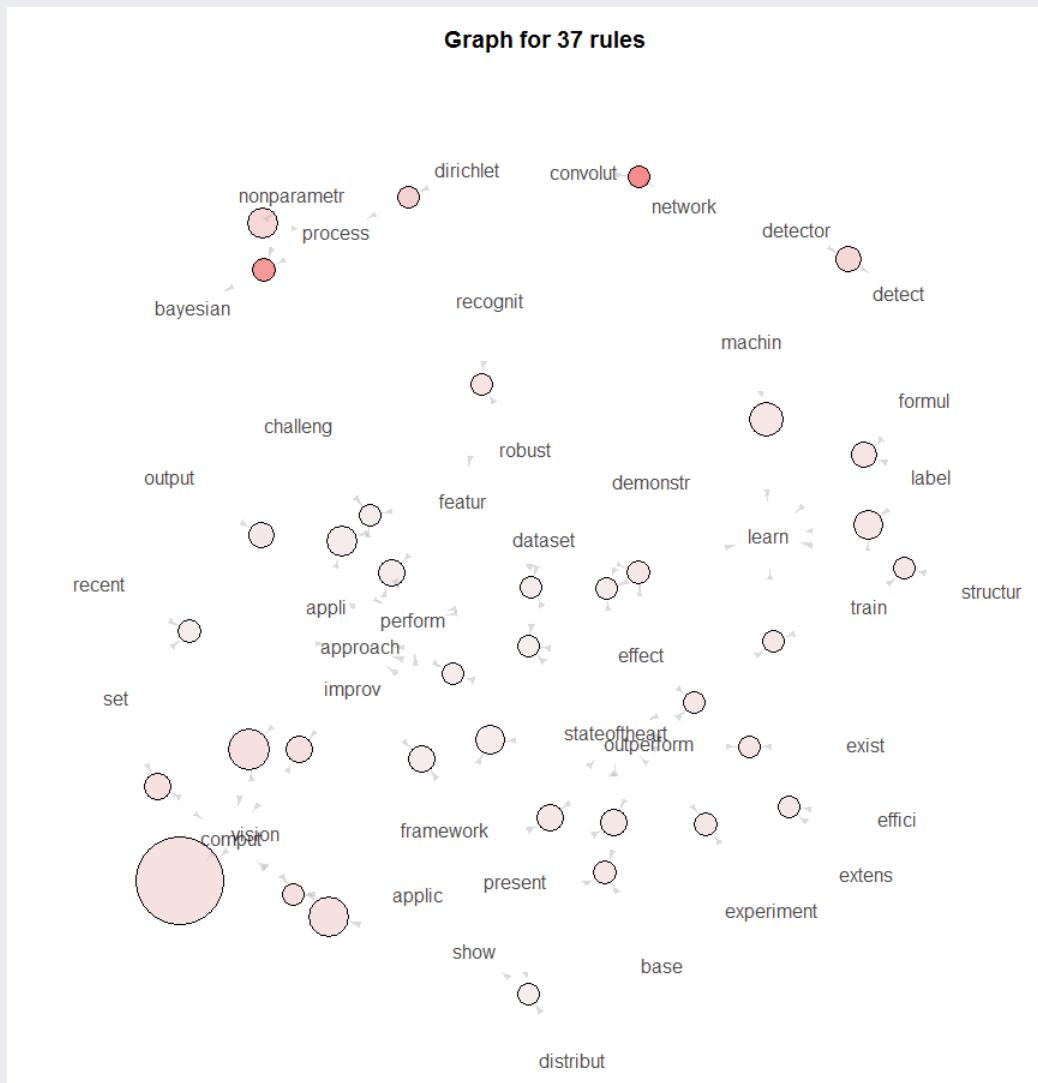
# R Exercise: Association Rules

- Generated rules from TPAMI Abstracts

lhs	rhs	support	confidence	lift
[1] {convolut}	=> {network}	0.04157044	0.8571429	8.247619
[2] {process,nonparametr}	=> {bayesian}	0.04387991	0.8260870	7.610546
[3] {dirichlet}	=> {process}	0.04157044	0.9473684	4.363942
[4] {detector}	=> {detect}	0.04618938	0.8000000	3.765217
[5] {nonparametr}	=> {process}	0.05311778	0.7931034	3.653338
[6] {set,vision}	=> {comput}	0.04849885	0.9545455	2.910691
[7] {applic,vision}	=> {comput}	0.04157044	0.9473684	2.888807
[8] {perform,vision}	=> {comput}	0.04849885	0.9130435	2.784140
[9] {approach,vision}	=> {comput}	0.06697460	0.9062500	2.763424
[10] {show,vision}	=> {comput}	0.06466513	0.9032258	2.754203
[11] {vision}	=> {comput}	0.12702079	0.9016393	2.749365
[12] {formul,label}	=> {learn}	0.04618938	0.8695652	2.614734
[13] {robust,recognit}	=> {featur}	0.04157044	0.7500000	2.557087
[14] {outperform,exist}	=> {stateoftheart}	0.04157044	0.8181818	2.393735
[15] {train,label}	=> {learn}	0.05080831	0.7857143	2.362599
[16] {train,structur}	=> {learn}	0.04157044	0.7826087	2.353261
[17] {machin}	=> {learn}	0.05773672	0.7812500	2.349175
[18] {outperform,show,base}	=> {stateoftheart}	0.04387991	0.7916667	2.316160
[19] {effect,outperform}	=> {stateoftheart}	0.04157044	0.7826087	2.289659
[20] {outperform,featur}	=> {learn}	0.04387991	0.7600000	2.285278
[21] {outperform,train}	=> {learn}	0.04157044	0.7500000	2.255208
[22] {outperform,experiment}	=> {stateoftheart}	0.04387991	0.7600000	2.223514
[23] {applic,outperform}	=> {stateoftheart}	0.04849885	0.7500000	2.194257
[24] {outperform,present}	=> {stateoftheart}	0.04849885	0.7500000	2.194257
[25] {effici,extens}	=> {stateoftheart}	0.04157044	0.7500000	2.194257
[26] {output}	=> {perform}	0.04618938	0.8333333	2.038606
[27] {featur,improv}	=> {perform}	0.04849885	0.7500000	1.834746
[28] {present,improv}	=> {perform}	0.04849885	0.7500000	1.834746
[29] {stateoftheart,dataset,robust}	=> {perform}	0.04157044	0.7500000	1.834746
[30] {approach,featur,challeng}	=> {perform}	0.04157044	0.7500000	1.834746
[31] {demonstr,stateoftheart,learn}	=> {perform}	0.04157044	0.7500000	1.834746
[32] {appli,featur}	=> {approach}	0.05311778	0.8214286	1.787330
[33] {stateoftheart,appli}	=> {approach}	0.04157044	0.8181818	1.780265
[34] {framework,outperform}	=> {approach}	0.05080831	0.7857143	1.709620
[35] {set,recent}	=> {approach}	0.04387991	0.7600000	1.653668
[36] {effect,stateoftheart,dataset}	=> {approach}	0.04157044	0.7500000	1.631910
[37] {present,distribut}	=> {show}	0.04157044	0.7500000	1.599754

# R Exercise: Association Rules

- Generated rules from TPAMI Abstracts



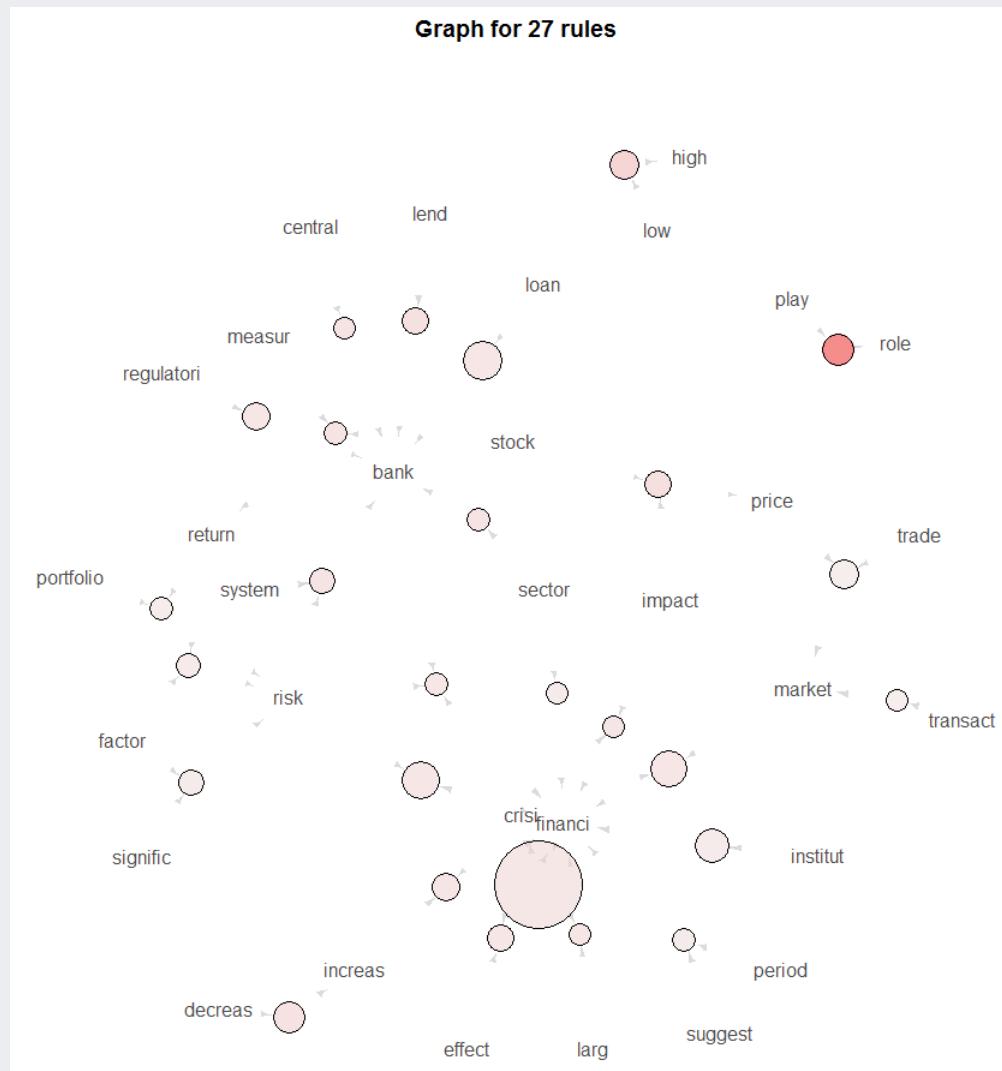
# R Exercise: Association Rules

- Generated rules from JoF Abstracts

	lhs	rhs	support	confidence	lift
[1]	{play}	=> {role}	0.05566219	0.9354839	8.550651
[2]	{low}	=> {high}	0.05182342	0.7105263	4.113158
[3]	{stock,impact}	=> {price}	0.04798464	0.7812500	3.083570
[4]	{lend}	=> {bank}	0.04798464	0.8333333	2.801075
[5]	{decreas}	=> {increas}	0.05566219	0.7073171	2.791759
[6]	{measur,stock}	=> {return}	0.04222649	0.7096774	2.641014
[7]	{central}	=> {bank}	0.04030710	0.7777778	2.614337
[8]	{risk,system}	=> {bank}	0.04606526	0.7741935	2.602289
[9]	{sector}	=> {bank}	0.04222649	0.7586207	2.549944
[10]	{larg,crisi}	=> {financi}	0.04030710	0.9545455	2.537338
[11]	{impact,crisi}	=> {financi}	0.04030710	0.9545455	2.537338
[12]	{increas,crisi}	=> {financi}	0.04990403	0.9285714	2.468294
[13]	{loan}	=> {bank}	0.06717850	0.7291667	2.450941
[14]	{market,crisi}	=> {financi}	0.06333973	0.9166667	2.436650
[15]	{bank,risk,crisi}	=> {financi}	0.04222649	0.9166667	2.436650
[16]	{regulatori}	=> {bank}	0.04990403	0.7222222	2.427599
[17]	{risk,crisi}	=> {financi}	0.06525912	0.8947368	2.378357
[18]	{effect,crisi}	=> {financi}	0.04798464	0.8928571	2.373360
[19]	{crisi}	=> {financi}	0.14203455	0.8604651	2.287257
[20]	{return,factor}	=> {risk}	0.04414587	0.7666667	1.987231
[21]	{sector}	=> {financi}	0.04030710	0.7241379	1.924877
[22]	{suggest,period}	=> {financi}	0.04222649	0.7096774	1.886438
[23]	{signific,factor}	=> {risk}	0.04606526	0.7272727	1.885120
[24]	{institut}	=> {financi}	0.05950096	0.7045455	1.872797
[25]	{return,portfolio}	=> {risk}	0.04222649	0.7096774	1.839512
[26]	{transact}	=> {market}	0.04030710	0.7500000	1.752242
[27]	{price,trade}	=> {market}	0.05182342	0.7105263	1.660019

# R Exercise: Association Rules

- Generated rules from JoF Abstracts



# R Exercise: Keyword Network

- Transform the Term-Frequency matrix into Term-Term co-occurrence matrix

```
# Section 3: Keyword Network -----
# For IEEE_TPAMI Abstracts
# Change it to a Boolean matrix
TPAMI.wcmat[TPAMI.wcmat >= 1] <- 1
# find the words that are used more than 30 times
freq.idx <- which(rowSums(TPAMI.wcmat) >= 30)
TPAMI.wcmat.freq <- TPAMI.wcmat[freq.idx,]

# Transform into a term-term adjacency matrix
TPAMI.ttmat <- TPAMI.wcmat.freq %*% t(TPAMI.wcmat.freq)

# inspect terms numbered 1 to 10
TPAMI.ttmat[1:10,1:10]
```

> TPAMI.ttmat[1:10,1:10]

Terms	approach	demonstr	estim	function	import	lead	local	measur	multipl	object
approach	199	85	52	35	26	16	43	26	41	57
demonstr	85	164	44	38	21	14	36	28	41	59
estim	52	44	106	19	13	8	22	13	23	34
function	35	38	19	84	13	8	16	18	22	26
import	26	21	13	13	61	6	11	12	11	22
lead	16	14	8	8	6	42	14	8	9	10
local	43	36	22	16	11	14	98	12	20	39
measur	26	28	13	18	12	8	12	69	17	27
multipl	41	41	23	22	11	9	20	17	90	37
object	57	59	34	26	22	10	39	27	37	131

# R Exercise: Keyword Network

- Build a graph & find communities
  - ✓ Undirected, simplification (remove loops), etc.

```
# Build a graph from the above matrix
TPAMI.graph <- graph.adjacency(TPAMI.ttm, weighted=T, mode = "undirected")

# remove loops
TPAMI.graph <- simplify(TPAMI.graph)

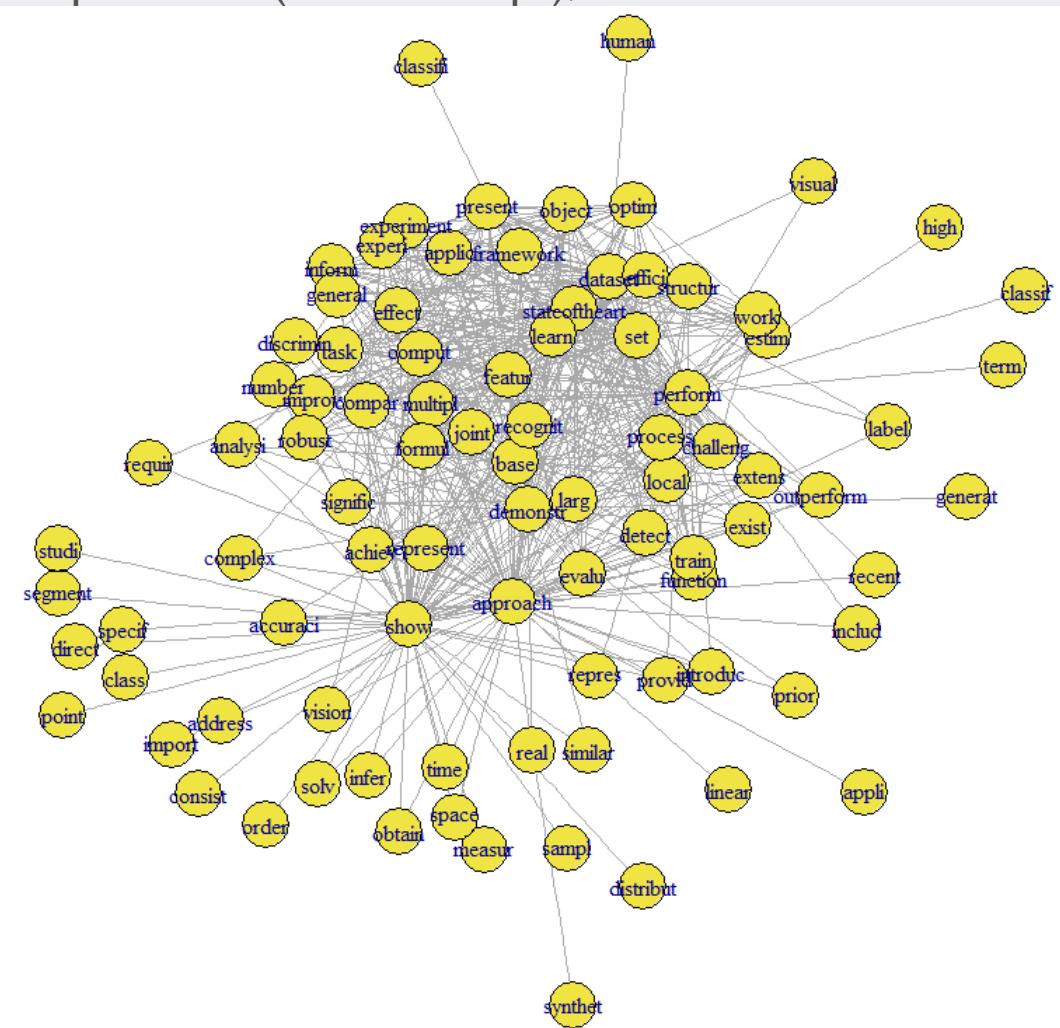
# set labels and degrees of vertices
V(TPAMI.graph)$label <- V(TPAMI.graph)$name
TPAMI.graph <- delete.edges(TPAMI.graph, which(E(TPAMI.graph)$weight <= 30))

TPAMI.graph <- delete.vertices(TPAMI.graph, which(degree(TPAMI.graph) == 0))
V(TPAMI.graph)$degree <- degree(TPAMI.graph)

# set seed to make the layout reproducible
set.seed(3952)
plot(TPAMI.graph, layout=layout.fruchterman.reingold)
plot(TPAMI.graph, layout=layout.kamada.kawai,
     vertex.size = 10, vertex.color = 4, vertex.label.cex = 1)
```

# R Exercise: Keyword Network

- Build a graph & find communities
    - ✓ Undirected, simplification (remove loops), etc.



# R Exercise: Keyword Network

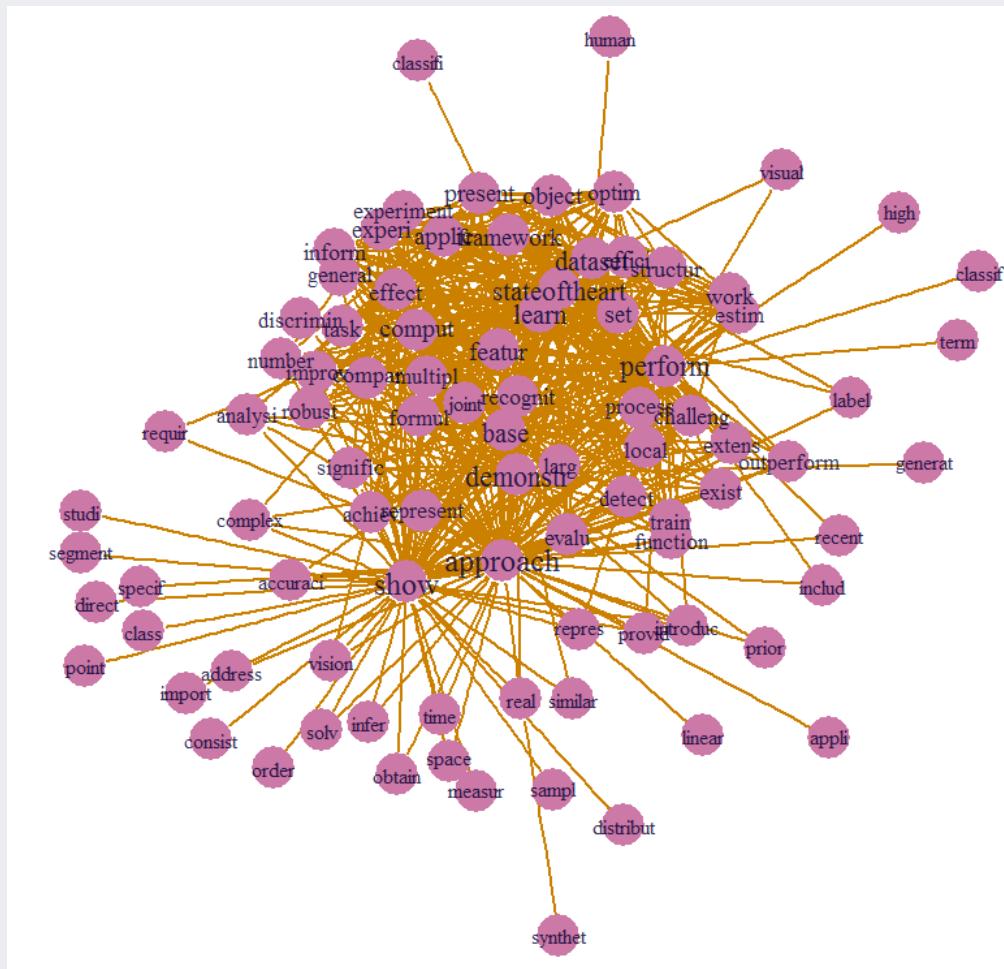
- Build a graph & find communities
  - ✓ Undirected, simplification (remove loops), etc.

```
# Make the network look better
V(TPAMI.graph)$label.cex <- 0.5*V(TPAMI.graph)$degree/max(V(TPAMI.graph)$degree)+1
V(TPAMI.graph)$label.color <- rgb(0, 0, 0.2, 0.8)
V(TPAMI.graph)$frame.color <- NA
egam <- 3*(log(E(TPAMI.graph)$weight+1))/max(log(E(TPAMI.graph)$weight+1))
E(TPAMI.graph)$color <- rgb(0.8, 0.5, 0)
E(TPAMI.graph)$width <- egam

# plot the graph in layout
plot(TPAMI.graph, layout=layout.kamada.kawai, vertex.size = 10, vertex.color = 7)
```

# R Exercise: Keyword Network

- Build a graph & find communities
    - ✓ Undirected, simplification (remove loops), etc.



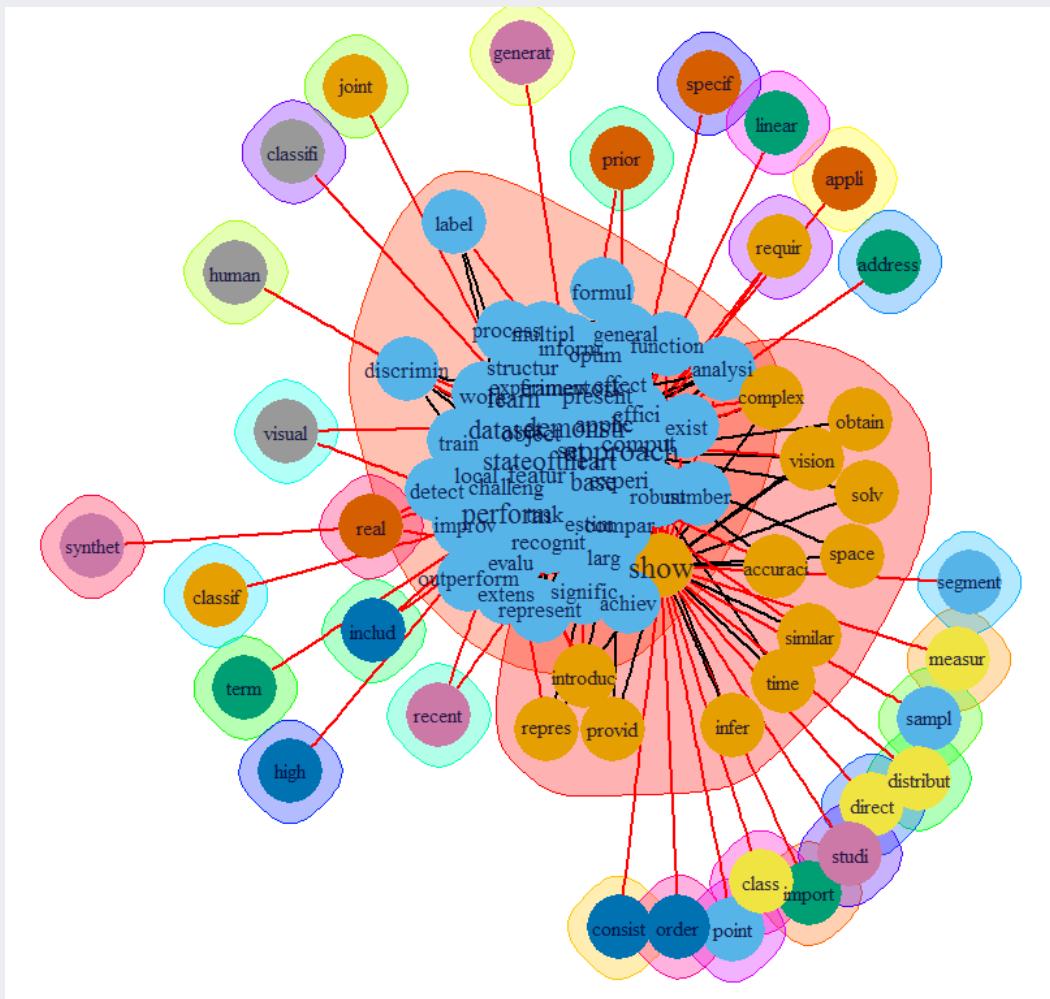
# R Exercise: Keyword Network

- Build a graph & find communities
  - ✓ Undirected, simplification (remove loops), etc.

```
# Plot the communities
TPAMI.community <- walktrap.community(TPAMI.graph)
modularity(TPAMI.community)
membership(TPAMI.community)
plot(TPAMI.community, TPAMI.graph)
```

# R Exercise: Keyword Network

- Build a graph & find communities
- ✓ Undirected, simplification (remove loops), etc.



# R Exercise: Keyword Network

- Transform the Term-Frequency matrix into Term-Term co-occurrence matrix

```
# For Journal of Banking and Finance Abstracts
# Change it to a Boolean matrix
JoF.wcmat[JoF.wcmat >= 1] <- 1
# find the words that are used more than 10 times
freq.idx <- which(rowSums(JoF.wcmat) >= 30)
JoF.wcmat.freq <- JoF.wcmat[freq.idx,]

# Transform into a term-term adjacency matrix
JoF.ttmat <- JoF.wcmat.freq %*% t(JoF.wcmat.freq)

# inspect terms numbered 5 to 10
JoF.ttmat[1:10,1:10]
```

```
> JoF.ttmat[1:10,1:10]
```

Terms	Terms									
	affect	analyz	approach	capit	corpor	countri	effect	empir	equiti	estim
affect	68	11	8	18	11	16	26	11	10	10
analyz	11	67	7	12	8	12	20	13	13	10
approach	8	7	51	14	4	7	11	11	9	12
capit	18	12	14	84	12	13	32	12	16	17
corpor	11	8	4	12	55	14	20	9	6	8
countri	16	12	7	13	14	66	30	11	10	7
effect	26	20	11	32	20	30	181	32	24	28
empir	11	13	11	12	9	11	32	88	12	16
equiti	10	13	9	16	6	10	24	12	69	8
estim	10	10	12	17	8	7	28	16	8	85

# R Exercise: Keyword Network

- Build a graph & find communities
  - ✓ Undirected, simplification (remove loops), etc.

```
# Build a graph from the above matrix
JoF.graph <- graph.adjacency(JoF.ttm, weighted=T, mode = "undirected")

# remove loops
JoF.graph <- simplify(JoF.graph)

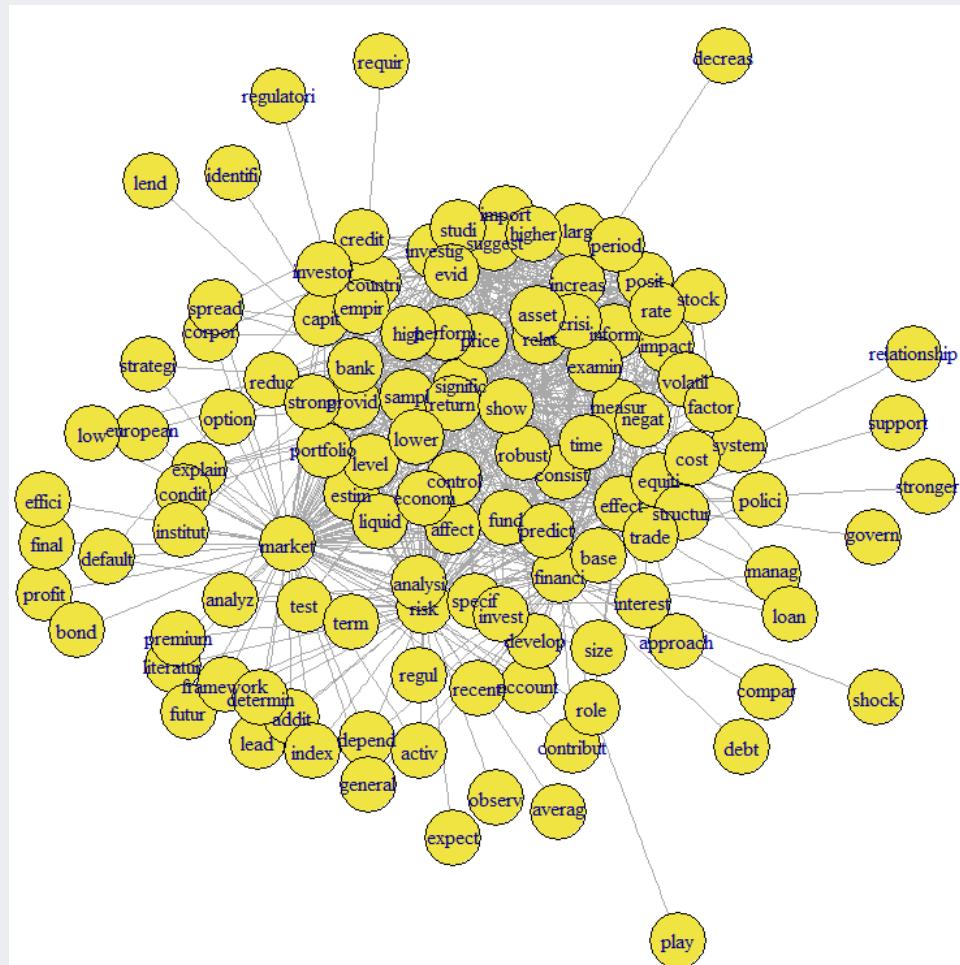
# set labels and degrees of vertices
V(JoF.graph)$label <- V(JoF.graph)$name
JoF.graph <- delete.edges(JoF.graph, which(E(JoF.graph)$weight <= 20))

JoF.graph <- delete.vertices(JoF.graph, which(degree(JoF.graph) == 0))
V(JoF.graph)$degree <- degree(JoF.graph)

# set seed to make the layout reproducible
set.seed(3952)
plot(JoF.graph, layout=layout.fruchterman.reingold)
plot(JoF.graph, layout=layout.kamada.kawai,
     vertex.size = 13, vertex.color = 4, vertex.label.cex = 1)
```

# R Exercise: Keyword Network

- Build a graph & find communities
  - ✓ Undirected, simplification (remove loops), etc.



# R Exercise: Keyword Network

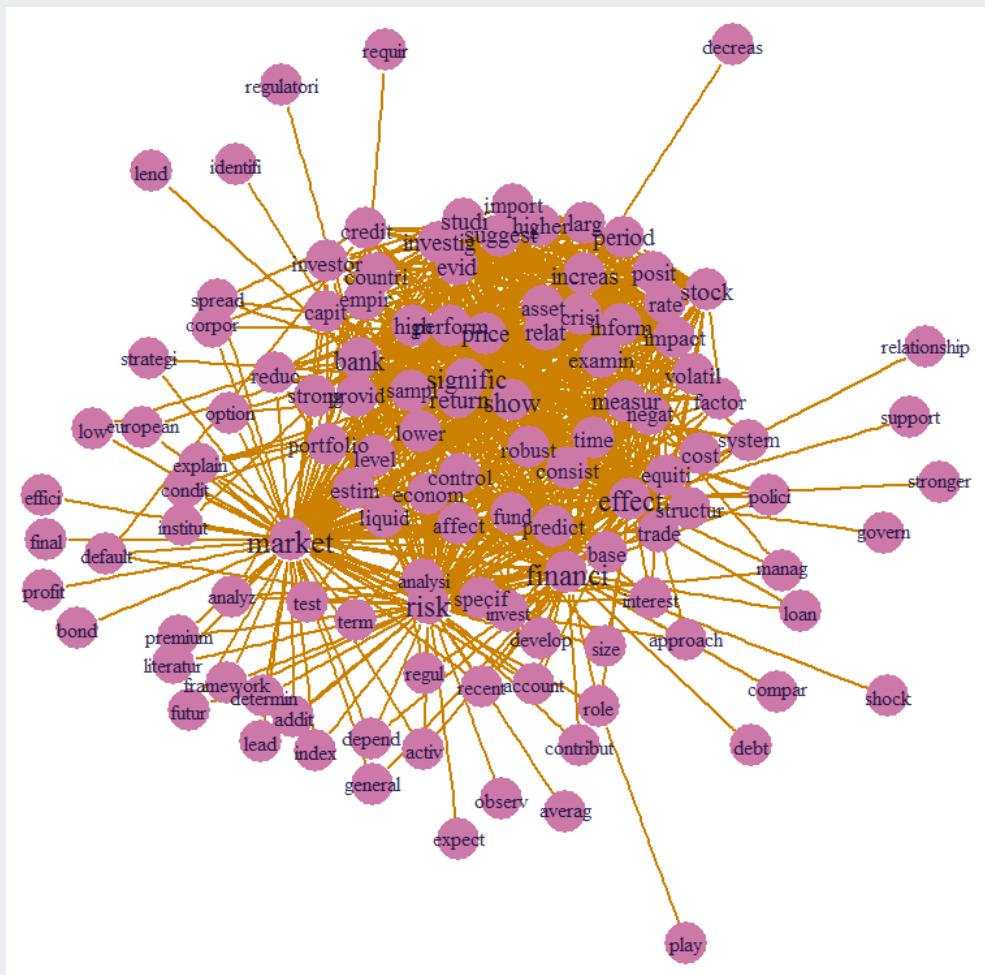
- Build a graph & find communities
  - ✓ Undirected, simplification (remove loops), etc.

```
# Make the network look better
V(JoF.graph)$label.cex <- 0.5*V(JoF.graph)$degree/max(V(JoF.graph)$degree)+1
V(JoF.graph)$label.color <- rgb(0, 0, 0.2, 0.8)
V(JoF.graph)$frame.color <- NA
egam <- 3*(log(E(JoF.graph)$weight+1))/max(log(E(JoF.graph)$weight+1))
E(JoF.graph)$color <- rgb(0.8, 0.5, 0)
E(JoF.graph)$width <- egam

# plot the graph in layout1
plot(JoF.graph, layout=layout.kamada.kawai, vertex.size = 10, vertex.color = 7)
```

# R Exercise: Keyword Network

- Build a graph & find communities
  - ✓ Undirected, simplification (remove loops), etc.



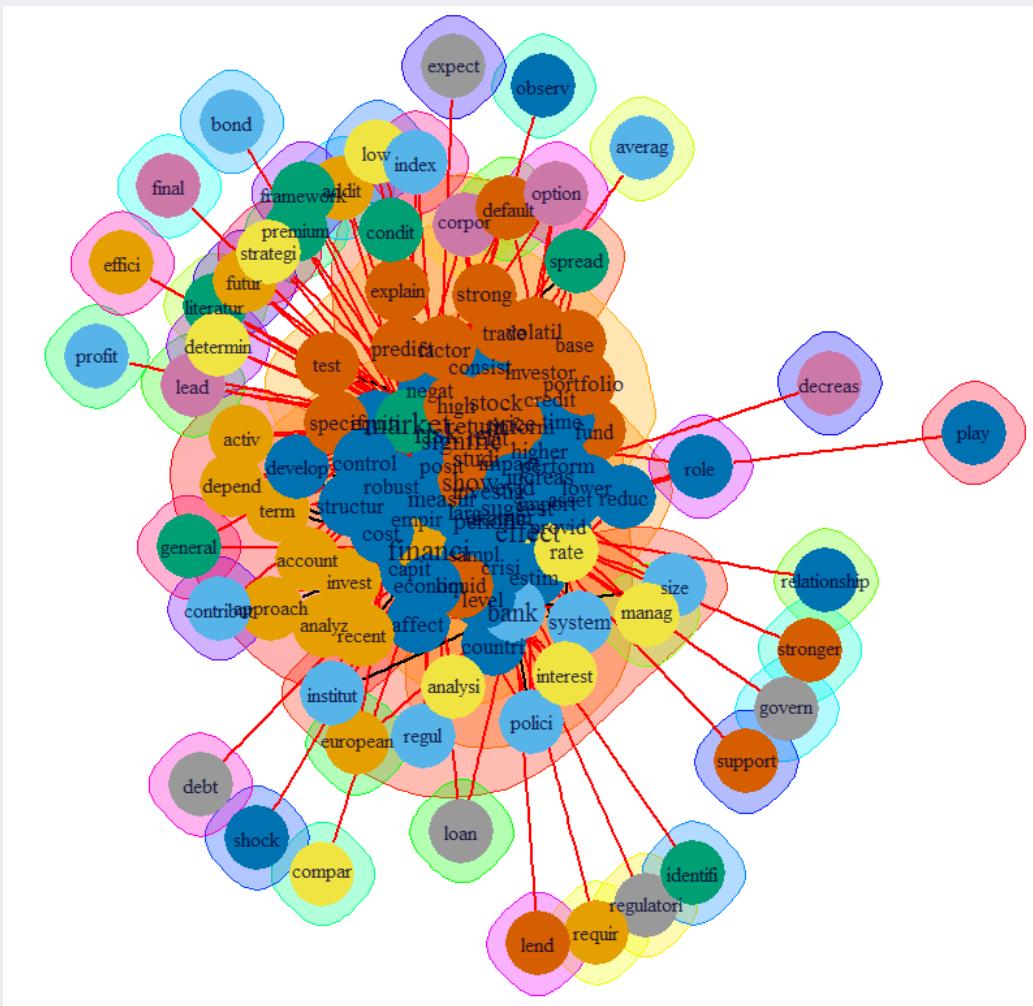
# R Exercise: Keyword Network

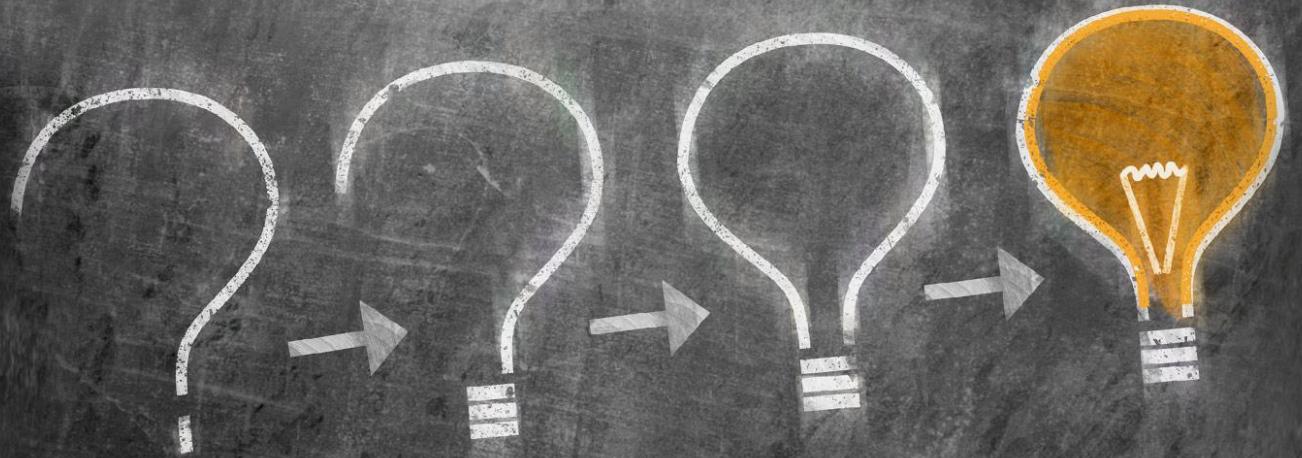
- Build a graph & find communities
  - ✓ Undirected, simplification (remove loops), etc.

```
# Plot the communities
JoF.community <- walktrap.community(JoF.graph)
modularity(JoF.community)
membership(JoF.community)
plot(JoF.community, JoF.graph)
```

# R Exercise: Keyword Network

- Build a graph & find communities
- ✓ Undirected, simplification (remove loops), etc.





# References

## Research Papers

- Ahn, Y. Y., Ahnert, S. E., Bagrow, J. P., & Barabási, A. L. (2011). Flavor network and the principles of food pairing. *Scientific reports*, 1.
- Kim, J. & Kang, P. (2016+). Analyzing international collaboration and identifying core topics for the “Internet of Things” based on network analysis and topic modeling, under review.
- Kim, J., Park, M., Kim, H., Cho, S., Kang, P., Lee, D., Yang, K., & Kim, K. (2016). 이상치 탐지 기법을 활용한 내부자 위협 탐지 방법론 개발. 대한산업공학회 추계학술대회, 서울.
- Kim, H., Park, M., & Kang, P. (2016). 토픽모델링과 사회연경망을 통한 딥러닝 연구동향 분석. 대한산업공학회 춘계공동학술대회, 제주.
- Lee, H. & Kang, P. (2017+). Identifying core topics in technology and innovation management studies: A topic model approach. *Journal of Technology Transfer*. Online available <https://link.springer.com/content/pdf/10.1007%2Fs10961-017-9561-4.pdf>

## Other Materials

- Kim, C., Kim, H., Cho, S., Kim, J., & Kang, P. (2017). 초음파 관련 임상연구 데이터의 텍스트마이닝 분석 플랫폼 개발, 산학과제 Granted by 삼성메디슨