

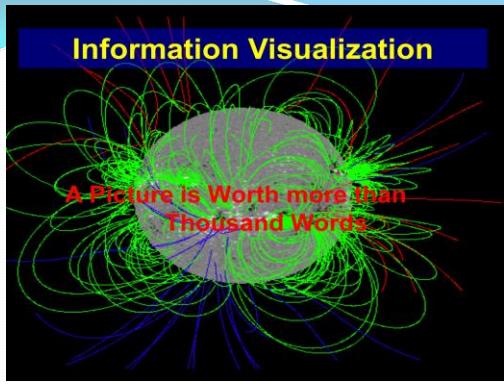
# Text Visualization and Knowledge Discovery for Text

Do Phuc, Ph. D, Assoc. Prof.  
Faculty of Information Systems,  
University of Information Technology, VNU-HCM



Assoc. Prof. Phuc Do, 2017

## Information Visualization



A Picture Is Worth more Than  
Thousands Words

Assoc. Prof. Phuc Do, 2017

# Text Visualization and Knowledge Discovery for Text

- Text is one of the greatest data and popular around us. Text contains our history and is a major mean to recording information and knowledge. With the Internet, text data has been generated with significantly speed. Currently, millions of websites are generating extraordinary amount of online text data everyday. The users of the social network like FaceBook, Twitter, ... produce billions of posting messages everyday.
- The explosion of the data makes a lot of challenges to discover and understand the knowledge in text.
- Text visualization techniques can be helpful for solving these problems. Various visualizations have been designed for showing the similarity of text documents, revealing and summarizing text content, and helping discover the big text data exploration.
- In this talk, we will show the techniques of text visualization and the application of text visualization in knowledge discovery for text. We also address the challenges for text visualization with big text data in big data era.

Assoc. Prof. Phuc Do, 2017

## Outline

- Data Visualization
- Examples
- Tree Visualization
- Graph Visualization
- Text Visualization

Assoc. Prof. Phuc Do, 2017

# Data Visualization

Assoc. Prof. Phuc Do, 2017

## Visualization & Information Visualization:

- Visualization:
  - “The use of computer-supported, interactive, visual representations of data to amplify cognition.”
  - Goal: discovery, decision making, explanation
- Information Visualization:
  - “The use of **computer-supported?**, interactive, visual representations of abstract data to amplify cognition.”

Assoc. Prof. Phuc Do, 2017

## Visualization Amplifies Cognition

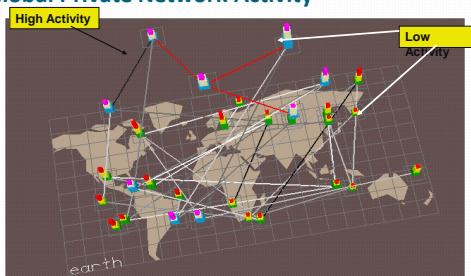
1. Increase memory and processing resources
2. Reduce information searching
3. Enhance detection of patterns
4. Enable perceptual inference operations
5. Use perception attention for monitoring
6. Encode information in a manipulative medium

Assoc Prof. Phuc Do, 2017

7

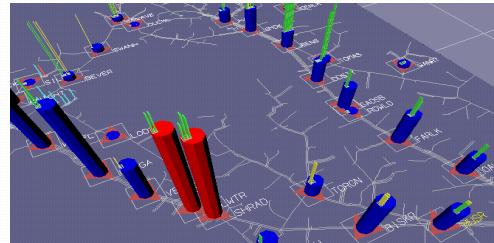
## Data Visualization Examples

### Global Private Network Activity



Phuc Do

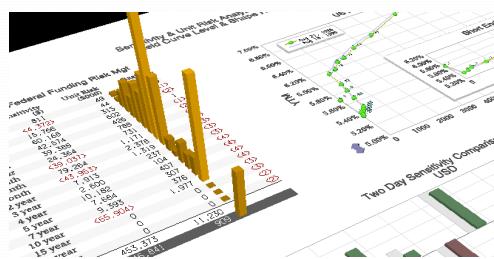
### Natural Gas Pipeline Analysis



**Height shows total flow through compressor stations.**

Phuc Do

### Risk Analysis Report



Phuc Do

### Telephone Polling Results

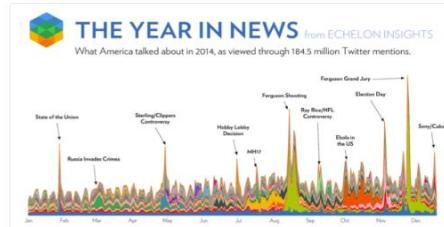


**Note:** On the "live" map, clicking on an area allows the user to drill down and see results for smaller areas.

Phuc Do

## 12) The Year in News

The best data visualizations communicate information and do so in an intuitive and beautiful way. Echelon Insights nailed it with this piece, which visualizes the most talked-about news stories of 2014 on Twitter. What do 184.5 million tweets look like? Rad spin art.



Phuc Do

## An example of Data Mining Visualization: Association rules

## List of Associations

Assoc Prof. Phuc Do, 201

16

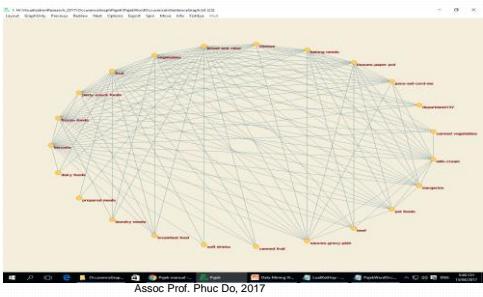
## Drawbacks

- Long list of rules
  - Many rules are redundant
  - Very hard to capture the important rules
  - How to access association rules effectively

## Graph of Association rule

- Association rules  $X \rightarrow Y$
  - Edge
    - Source: X
    - Target: Y
  - Set of Association rules
    - Graph  $(V,E)$
    - V: set of frequent itemsets
    - E: set of edges created from association rules

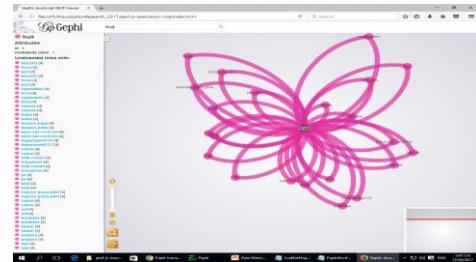
# Graph visualization by Pajek software



Assoc Prof. Phuc Do, 2017



# Graph Visualization by Gephi on the Web



Assoc Prof. Phuc Do, 201

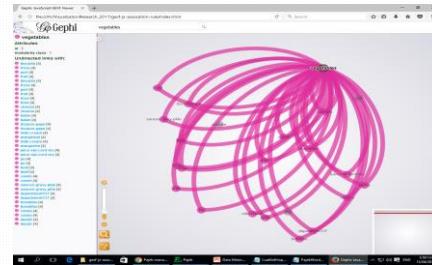
## Interactive mode in Graph visualization

- Very easy to capture the important association rules
  - Click on node to discover the co-occurrence items
  - Move around to discover the association rules

Assoc Prof. Phuc Do, 2017



## Discover the items co occurrence with vegetables



Assoc Prof. Phuc Do, 201

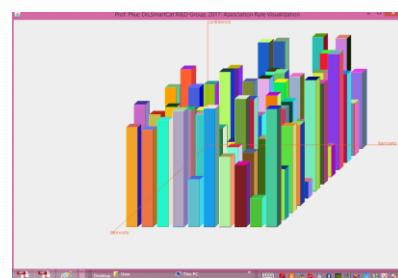
## Using cube to visualize AR

- Each rule is a cube
  - Click on Cube to visualize the association rule
  - The height of Cube: the confidence of AR.
  - Easy to discover the important rules based on the height of cube.

Assoc Prof Phuc Do, 2017



AR visualization, each AR is a cube



Assoc Prof Phuc Do 201

## Visualization of Relations

- Tree Visualization
- Graph Visualization

Assoc Prof. Phuc Do, 2017

25

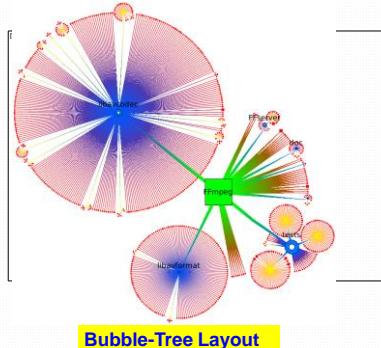
## Tree Visualization

- Treemap method: visualize the tree structure that use virtually every pixel of the display space to convey information
- Every subtree is represented by a rectangle, that is partitioned into smaller rectangles with correspond to its children.
- The position of the slicing lines determines the relative sizes of the child rectangles.
- For every child, repeat the slicing recursively, swapping the slicing direction from vertical to horizontal or conversely

Assoc Prof. Phuc Do, 2017

26

## Tree Visualization

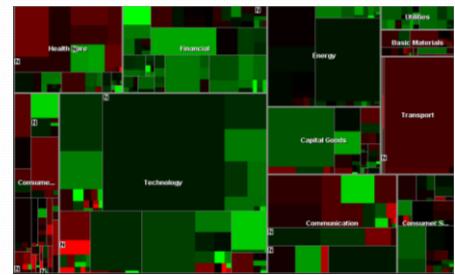


Bubble-Tree Layout

Assoc Prof. Phuc Do, 2017

27

## Tree Visualization



Assoc Prof. Phuc Do, 2017

28

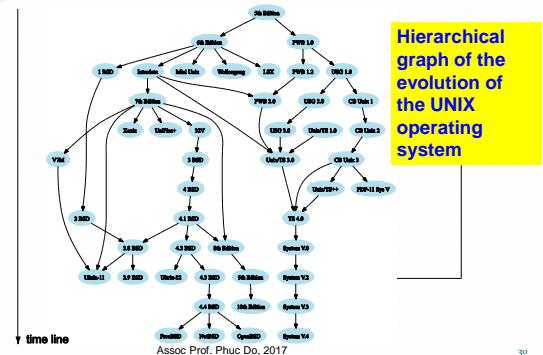
## Graph Visualization

- Graphs are the most general type of relational data
- Similar to a tree, it consists of nodes and edges
- Different from a tree, a child node may have multiple parent nodes
- A graph contains loops, or multiple paths between two nodes in the graph.

Assoc Prof. Phuc Do, 2017

29

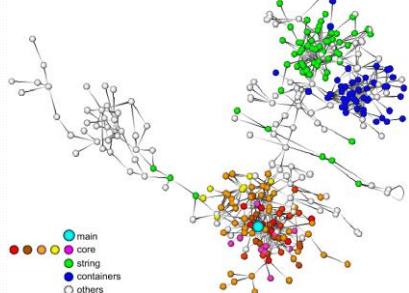
## Graph Visualization



Assoc Prof. Phuc Do, 2017

30

## Graph Visualization



Call Graph using a Force-Directed Layout

Assoc Prof. Phuc Do, 2017

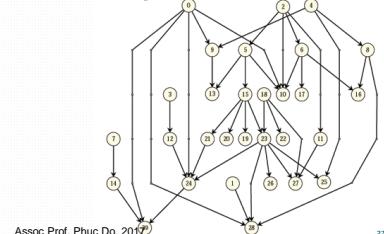
31

## Layout of Directed Graphs

- Layering (<http://www.csus.yk.ue/staff/NikolaNikolov/#phd>)

- Algorithm to calculate the co-ordinate of vertices

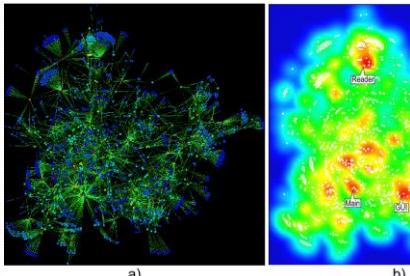
Crossings: 11



Assoc Prof. Phuc Do, 2017

32

## Graph Visualization



Force-Directed Layout

Assoc Prof. Phuc Do, 2017

Splatting;  
Dense Representation

33

## Text Visualization and Knowledge Discovery for Text

Assoc Prof. Phuc Do, 2017

34

## Why visualize text?

- **Understanding** – read a document
- **Summaries** – get the “key” topic of a document
- **Clustering** – group together similar contents
- **Correlate** – compare patterns in text to those in other data, e.g., correlate with social network

Assoc Prof. Phuc Do, 2017

35

## Summarize with Words

Assoc Prof. Phuc Do, 2017

36

## Words are (not) nominal?

## High dimensional

- High dimensional (10,000+)
  - Words have meanings and relations
    - Correlations: *Hong Kong, San Francisco, Bay Area*
    - Order: *April, February, January, June, March, May*
    - Hierarchy, antonyms & synonyms, entities, ... (ontology)

Assoc Prof. Phuc Do, 2017

37

## Tips: Tokenization and Stemming

- Tokenizer, Vntokenizer
    - "Thuyền ai thấp\_thoáng bên sông",
    - "Đưa câu mái dày chạnh\_lòng nước non"  
(Hue Fork Song)
  - Stemming
  - Name Entity Recognition

Assoc Prof. Phuc Do, 20

38

## Bag of Words and Vector Model

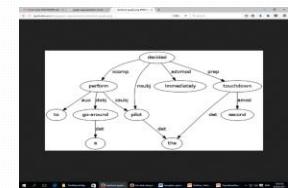
- Ignore ordering relationships within the text
  - A document  $\approx$  vector of term weights
    - Each dimension corresponds to a term (10,000+)
    - Each value represents the relevance
      - For example, simple term counts
  - Aggregate into a document-term matrix
    - Document vector space model
    - Frequency of word in document

Assoc Prof. Phuc Do, 2017

39

## Graph based text representation

- Frequent words: Vertices
  - Link: the co occurrence words in sentence.
  - Co occurrence graph
  - Preserve the order of words in sentence



40

## Analyzing document: word cloud

- **Strengths**

- Can help with initial idea about text

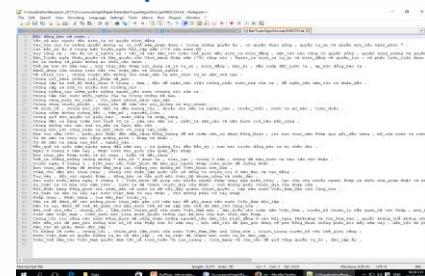
#### • Weaknesses

- Term frequency may not be meaningful
  - Does not show the structure of the text

Assoc Prof. Phuc Do, 2017

43

## **Proclamation of Independence of the Democratic Republic of Vietnam**



42



Word Cloud



Assoc Prof. Phuc Do, 2017 43



Word Cloud



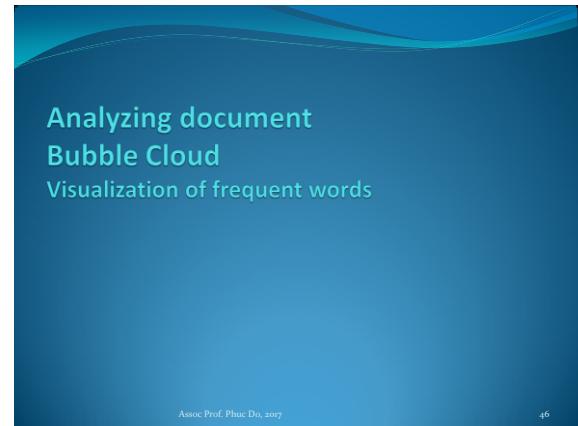
Assoc Prof. Phuc Do, 2017 44



Word Cloud



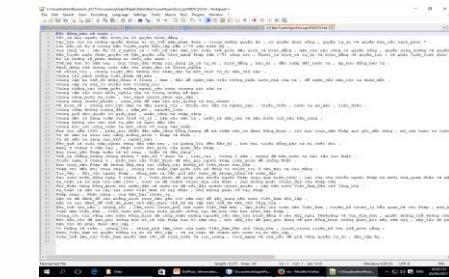
Assoc Prof. Phuc Do, 2017 45



Assoc Prof. Phuc Do, 2017 46



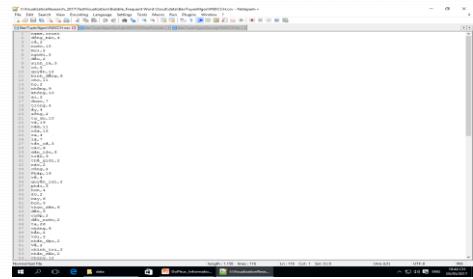
Proclamation of Independence of the Democratic Republic of Vietnam



Assoc Prof. Phuc Do, 2017 47



List of frequent words



Assoc Prof. Phuc Do, 2017 48

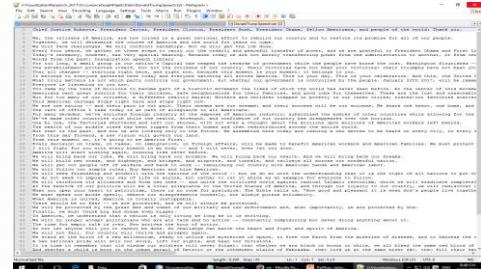
## Bubble Cloud



Assoc Prof. Phuc Do, 2017

49

## President Donald Trump inauguration speech



Assoc Prof. Phuc Do, 2017

50

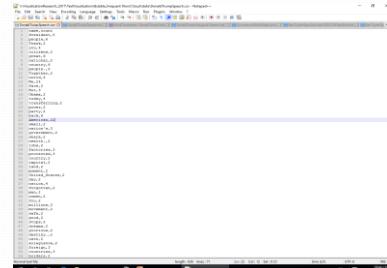
<https://qz.com/889611/the-most-frequently-used-words-in-donald-trumps-inauguration-speech-and-every-previous-one/>

- The most frequently used words in Donald Trump's inauguration speech—and every previous one
- After a bitter contest and controversial win, Donald J. Trump is now officially president of the United States. In his inaugural address, Trump continued the message of populism he maintained throughout his campaign. “We are transferring power from Washington, D.C.,” he said, “and giving it back to you, the people.” Trump also said he would invest in improving the country’s infrastructure, lamented the shuttering of factories, and referred to crime as “American carnage”.
- But the word Trump used most in his speech wasn’t “people,” “power,” or even “great.” It was “dream.” At one point, Trump said he would “bring back our dreams.” Later, when talking about the victims of the aforementioned carnage, he said “their dreams are our dreams”.
- We counted the words used in Trump’s speech, as well as the words in every inaugural address ever delivered, to find which were used the most. (This list filters out common articles like “the,” “and,” and “or,” as well as some terms that are inherently ubiquitous in inaugural speeches, like “government” and “America.”) The result gives us a sense of what each president, going back to 1789, focused on at the start of his term.

Assoc Prof. Phuc Do, 2017

51

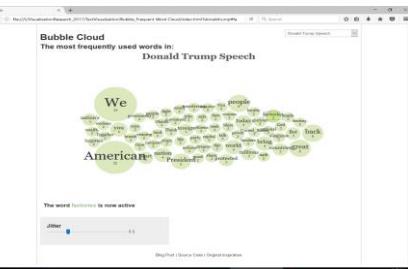
## List of frequent words



Assoc Prof. Phuc Do, 2017

52

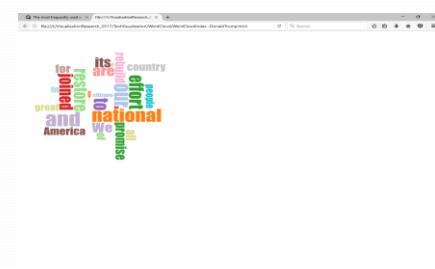
## Bubble Cloud



Assoc Prof. Phuc Do, 2017

53

## Word Cloud



Assoc Prof. Phuc Do, 2017

54

## Word Cloud



Assoc Prof. Phuc Do, 2017

55

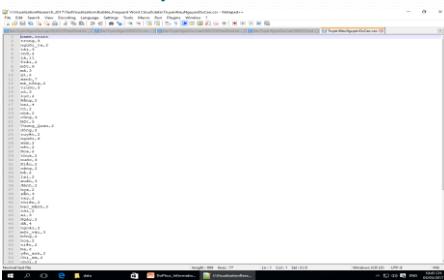
## Kieu Legend, Nguyen Du

- Trâm\_nâm trong cõi người ta .
- Chữ tài chữ mệnh khéo là ghê nhau .
- Trái qua mót cuộc bể dâu .
- Những điều trông thấy mà đau\_dồn lòng .
- Lá gi\_bé sắc\_tu\_phong .
- Trời xanh quen thói mà hông\_dành\_ghen .
- Cao thon lán giờ trước đèn ,
- Phòng\_tinh có lục\_còn truyền sứ\_xanh .
- Rảng nám\_Gia\_Tinh triều Minh ,
- Bón phumour phảng\_làng , hai kính\_vững\_vàng .
- Cỏ nhà viền ngoài họ\_Wrong ,
- Gia\_tu\_nghি\_cứng thường thường bức\_trung .
- Một trai con thứ rốt\_lòng .
- Vua Quang là chí , nỗi dòng\_nho\_gia .
- Dẫu\_hồng\_hai\_làm\_mưa .
- Thúy\_Kieu là chí , cõi là Thúy\_Vân .
- Mai cõi\_cách , tuyệt\_tinh\_thần .
- Một người mới ve : muối\_phân\_yen\_mười .
- Văn\_xem\_trang trong khác vòi ,
- Khuôn\_trang\_dày\_dần , nét\_ngài\_nó\_nang .

Assoc Prof. Phuc Do, 2017

56

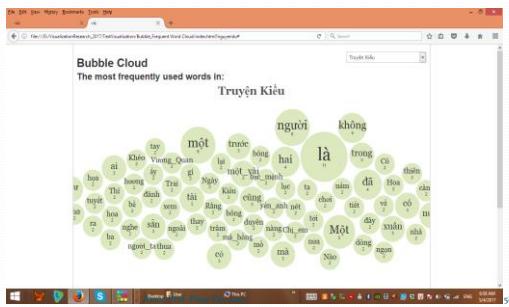
## List of Frequent words



Assoc Prof. Phuc Do, 2017

57

## Bubble Cloud



58

## Latent Dirichlet Allocation

Discover the implicit topics of a corpus

Assoc Prof. Phuc Do, 2017

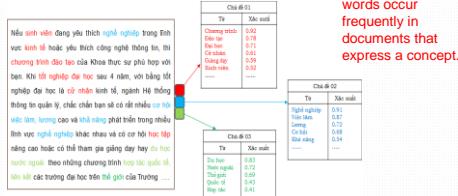
59

## Latent Dirichlet Allocation (LDA)

- Invented by Prof. David Blei, 2003
- Probabilistic Graphical Model (PGM)
- Generative probabilistic modeling
  - Treats data as observations
  - Contains hidden variables
  - Hidden variables reflect thematic structure of the collection.
- Infer hidden structure using posterior inference
  - Discovering topics in the collection.

60

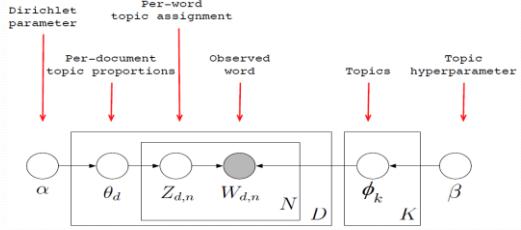
## Using LDA to discover the latent topics of messages



LDA is a generative probabilistic model of a corpus. The basic idea is that the documents are represented as random mixtures over latent topics, where a topic is characterized by a distribution over words.

61

## LDA as Graphical Model



$$P(\mathbf{W}, \mathbf{Z}, \theta, \phi; \alpha, \beta) = \prod_{i=1}^K P(\varphi_i; \beta) \prod_{j=1}^M P(\theta_j; \alpha) \prod_{d=1}^N P(Z_{d,j} | \theta_j) P(W_{d,j} | \varphi_{Z_{d,j}})$$

62

## Generating Process of LDA

- Choose  $N \sim Poisson(\xi)$
- Choose  $\theta \sim Dir(\alpha)$
- For each of the N words  $w_n$ 
  - Choose a topic  $z_n \sim Multinomial(\theta)$
  - Choose a word  $w_n$  from  $p(w_n | z_n, \beta)$ , a multinomial probability conditioned on the topic  $z_n$ .  $\beta$  is a  $k \times n$  matrix parameterized with the word probabilities.

$$[\beta]_{k \times V} \quad \beta_{ij} = p(w^j = 1 | z^i = 1)$$

## Gibbs Sampling for LDA

- Gibbs Sampling can be used to sample  $\phi$  and  $\theta$ . They are calculated by:

$$\hat{\phi}_k^{w_d} = \frac{N_k^{w_d} + \beta}{N_k^{(1)} + W\beta} \quad \hat{\theta}_k^d = \frac{N_k^d + \alpha}{N_k^{(1)} + K\alpha}$$

$\phi$  is a probability dist. of topic x word  
 $\theta$  is a probability dist. of doc x topic

65

## Inference and Parameter Estimation

- Parameters:
  - Corpus level:  $\alpha, \beta$ : sampled once in the corpus creating generative model.
  - Document level:  $\theta, z$ : sampled once to generate a document
- Inference: estimation of document-level parameters.  

$$p(\theta, z | w, \alpha, \beta) = \frac{p(\theta, z, w | \alpha, \beta)}{p(w | \alpha, \beta)}$$
- However, it is intractable to compute, which needs approximate inference.

## Latent topics without label

Topic 4	Topic 12	Topic 17			
doanh	0.03189	phòng	0.03925	công_ty	0.03245
bởi	0.01957	quản_lý	0.02322	viec_làm	0.01767
Bắc	0.01769	tín_chỉ	0.01933	khiển_lập	0.01485
hoạt_dong	0.01689	loc_phí	0.01402	nhân_vật	0.01093
doanh_vien	0.01433	giảng_dạy	0.01021	tham_quan	0.00952
sát_dóng	0.01277	thi	0.00967	hàng	0.00711
công_tác	0.00877	loc_lai	0.00833	người	0.00429
sát_hoc	0.00597	thời_khoi_hoc	0.00790	em_cu_trường	0.00641
quy định	0.00535	phòng_học	0.00718	hoat_dong	0.00500
phong_trào	0.00499	máy_chứa	0.00554	hoc_kì	0.00259

Topic 2	Topic 7	Topic 14			
chết_luong	0.04213	điều	0.06212	người_ngồi	0.04539
sinh_viên	0.02543	tô	0.05423	hợp_tác	0.03478
bết_tháo	0.01456	phép_qué	0.04562	đa_học	0.03941
tua_dầm	0.01213	tình_cảm	0.03412	nhân_người	0.03415
khát_nghem	0.01122	tâm_ý	0.02431	học_khi	0.02727
lỗi_nghi	0.00992	bản_bì	0.01321	quá_danh	0.01225
khô_hóe	0.00876	cam_giac	0.00872	tiêu_dùng	0.00912
aphelin_cửu	0.00612	thầy	0.00790	hoc_bóng	0.00823
phong_trào	0.00467	cô	0.00632	nhà_máu	0.00763
số_lưu	0.00421	qua_tặng	0.00452	đi_uguai	0.00649

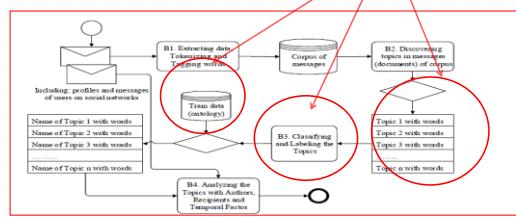
66

## Automatic labeling of topics

67

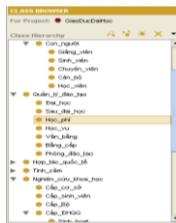
## Automatic labeling by classification

Latent topic: just a group of words, no topic label  
We assign label to this topic.



68

## A topic taxonomy of university



We build a topic taxonomy of university  
Then we collect 200-300 messages for each topic  
We use LDA to discover the words representing for this topic  
SVM is used for classification two sets of words.

69

70

Table 1. List of labeled topics and words distribution with probability

Topic 4 "social activity"	Topic 12 "management training"	Topic 17 "recruitment jobs"
union 0.03189	Office 0.03925	company 0.03245
associate 0.01957	management 0.02322	jobs 0.01767
Uncle Ho 0.01569	skills 0.01353	teacher 0.01426
activity 0.01689	school fee 0.01402	staff 0.01093
union member 0.01433	teaching 0.01021	visit 0.00952
content 0.01277	test 0.00967	Salary 0.00711
mission 0.00677	leaning 0.00833	people 0.00429
school year 0.00597	schedule 0.00790	environment 0.00641
regulation 0.00535	classroom 0.00718	activity 0.00541
movement 0.00499	projector 0.00554	learning 0.00359

## Visualize Themes in a Document Collection

## Topical Analysis

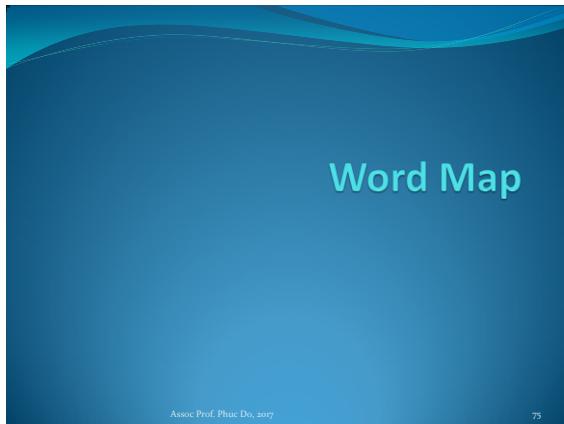
- Large document collections
  - Too large to manually read the source documents
  - Deeper analysis than the most common theme
- Statistical topic modeling
  - Analysis of word relationships
  - Extract *latent topics* belonging to the documents
- Latent Dirichlet Allocation (LDA)



## LDA: visualize topics



## LDA: visualize topics



## Word Map

Assoc Prof. Phuc Do, 2017

75



## WordTree (Wattenberg et al)



## Biển nhớ, Trịnh Công Sơn

- "Ngày\_mai em đi , biển nhớ tên em gọi về",
- "Ngày\_mai em đi , đồi\_núi nghiêng\_nghiêng đợi\_chờ","

Assoc Prof. Phuc Do, 2017

77



## Word Tree



## Graph based document summarization

Assoc Prof. Phuc Do, 2017 79

## Text Visualization Summary

- High Dimensionality
  - Where possible use **text to represent text...**
  - ... which terms are the most descriptive?
- Context & Semantics
  - Provide **relevant context** to aid understanding.
  - Show (or provide access to) the **source text**.
  - Modeling Abstraction
- Determine your **analysis task**.
  - Understand abstraction of your **language models**.
  - Match analysis task with appropriate tools and models.

Assoc Prof. Phuc Do, 2017 80

## Kieu legend of Nguyễn Du

Assoc Prof. Phuc Do, 2017 81

## Clustering the sentences

Assoc Prof. Phuc Do, 2017 82

## Connected Components

Assoc Prof. Phuc Do, 2017 83

## Discover the connected component

Assoc Prof. Phuc Do, 2017 84

## Doc. Similarity & Clustering

- In vector model, compute distance among docs
  - For TF.IDF, typically cosine distance
  - Similarity measure can be used to cluster
- Topic modeling approaches
  - Assume documents are a mixture of topics
  - Topics are (roughly) a set of co-occurring terms
  - Latent Semantic Analysis (LSA): reduce term matrix
  - Latent Dirichlet Allocation (LDA): statistical model

Assoc Prof. Phuc Do, 2017

85

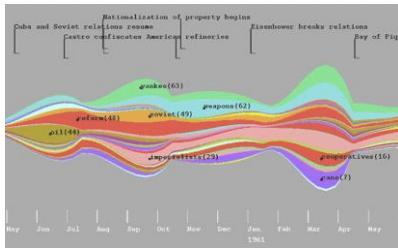
## Parallel Tag Clouds [Collins et al 09]



Assoc Prof. Phuc Do, 2017

86

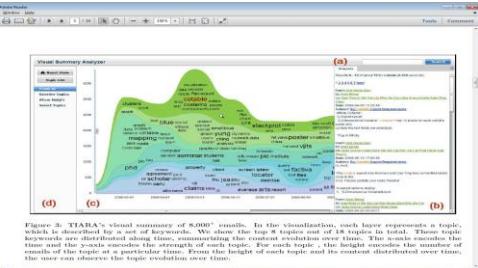
## ThemeRiver [Havre et al 99]



Assoc Prof. Phuc Do, 2017

87

## What topics they used in their email



Assoc Prof. Phuc Do, 2017

88

## Visualizing Sentiments and Emotions

Assoc Prof. Phuc Do, 2017

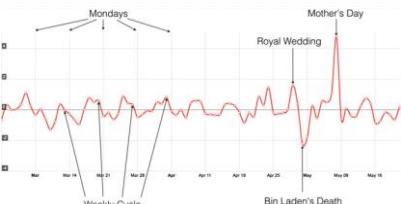
89

## Sentiment analysis

- Sentiment analysis, also known as opinion mining, is one of the most important text mining tasks and has been widely used for analyzing, for example, reviews or social media data for various of applications, including marketing and customer service.
- The result of sentiment analysis is negative, neutral, positive

Assoc Prof. Phuc Do, 2017

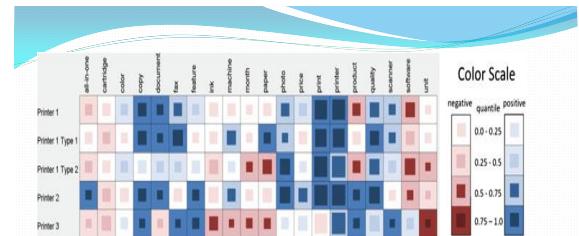
90



C. Nan and W. Cui, *Introduction to Text Visualization*, Atlantis Briefs in Artificial Intelligence 1, DOI 10.2991/978-94-6239-186-1, 6

Assoc Prof. Phuc Do, 2017

91



Summary report of printers: each row shows the attribute performances of a specific printer. Blue color represents comparatively positive user opinions and red color comparatively negative ones (see color scale). The size of an inner rectangle indicates the amount of customers that commented on an attribute. The larger the rectangle the more comments have been provided by the customers

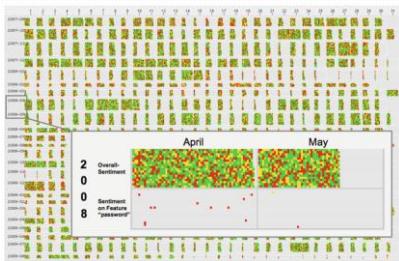
C. Nan and W. Cui, *Introduction to Text Visualization*,  
Atlantis Briefs in Artificial Intelligence 1, DOI 10.2991/978-94-6239-186-1, 6

Assoc Prof. Phuc Do, 2017

92

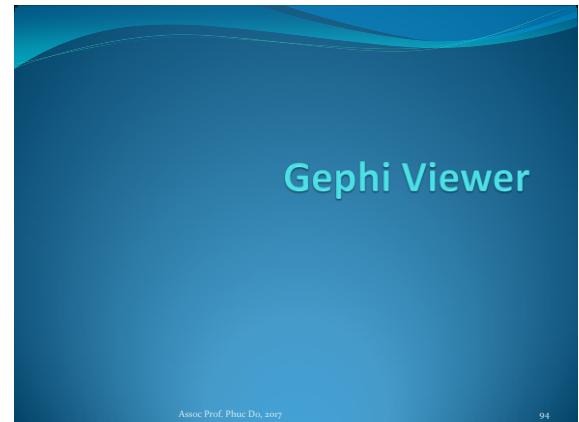


Pixel map calendar: each pixel indicates a document with the color encoding its overall sentiment, i.e., the average of all contained feature sentiments. Here, red indicates negative, green indicates positive, and yellow indicates neutral. In the background, the x-axis bins are days and y-axis bins are years with month.



Assoc Prof. Phuc Do, 2017

93



Assoc Prof. Phuc Do, 2017

94



Assoc Prof. Phuc Do, 2017

95



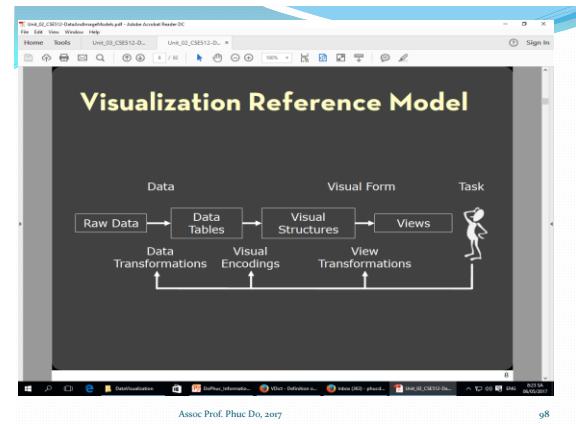
Assoc Prof. Phuc Do, 2017

96

# Visualization Design

Jeffrey Heer University of Washington, 2015

Assoc Prof. Phuc Do, 2017 97



```

graph LR
    subgraph Left [ ]
        direction TB
        T[task]
        D1[data  
physical type  
int, float, etc.  
abstract type  
nominal, ordinal, etc.]
        D2[domain  
metadata  
semantics  
conceptual model]
        D1 --- D2
    end
    subgraph Right [ ]
        direction TB
        PA[processing algorithms]
        M[mapping  
visual encoding  
visual metaphor]
        I[image  
visual channel  
retinal variables]
        PA --> M
        M --> I
    end
    T --> PA
    
```

Assoc Prof. Phuc Do, 2017 99

# Formalizing Design (Mackinlay 1986)

Assoc Prof. Phuc Do, 2017 100

## Visualization components

- Color
- Size
- Texture
- Proximity
- Annotation
- Interactivity
  - Selection / Filtering
  - Zoom
  - Animation

Assoc Prof. Phuc Do, 2017 101

## Visual Encoding Variables

- Position
- Size
- Value
- Texture
- Color
- Orientation
- Shape
- Others?

Assoc Prof. Phuc Do, 2017 102

## Choosing Visual Encodings

- **Challenge:**
- Assume 8 visual encodings and  $n$  data attributes.
- We would like to pick the “best” encoding among a combinatorial set of possibilities with size  $(n+1)^8$
- **Principle of Consistency:**
- The properties of the image (visual variables) should match the properties of the data.
- **Principle of Importance Ordering:**
- Encode the most important information in the most effective way.

Assoc Prof. Phuc Do, 2017

103

## Design Criteria (Mackinlay)

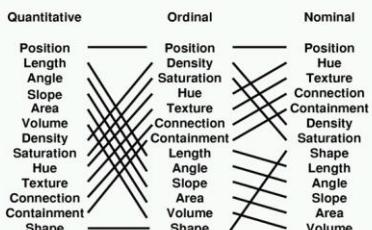
- **Expressiveness**
- A set of facts is expressible in a visual language if the sentences (i.e. the visualizations) in the language express all the facts in the set of data, and only the facts in the data.
- **Effectiveness**
- A visualization is more effective than another visualization if the information conveyed by one visualization is more readily perceived than the information in the other visualization.

Assoc Prof. Phuc Do, 2017

104

## Channel ranking varies by data type

spatial position best for all types



[Mackinlay, Automating the Design of Graphical Presentations of Relational Information, ACM TOC 5.2, 1986]

12

Assoc Prof. Phuc Do, 2017

105

## Visualization software resources

Assoc Prof. Phuc Do, 2017

106

## Visualization software

- Host language (C/C++/Java/Python) plus OpenGL
- Stat/math package with graphics
  - R
  - MATLAB
- Special-purpose info visualization software
  - Earth mapping, biological network visualization, etc.
- Browser-enabled graphics/info visualization packages
  - [Google Charts](#)
  - [Processing / Processing.js](#)
  - [D3](#)
  - Java + Flash

Assoc Prof. Phuc Do, 2017

107

Assoc Prof. Phuc Do, 2017

108

## Hands-on

- [HTML intro\\*](#)
- [Google charts](#)
- D3

## Resources

- Books

[Visual Complexity, Mapping Patterns of Information](#), Manuel Lima  
[The Visual Display of Quantitative Information](#), Edward Tufte  
[Information Visualization: Beyond the Horizon](#), Chaomei Chen  
[JavaScript: The Definitive Guide](#), David Flanagan  
[Getting Started with D3](#), Mike Dewar  
[Visualizing Data](#), Ben Fry  
[Interactive Data Visualization for the Web](#), Scott Murray

- Websites

[http://processing.org/](#)  
[http://d3js.org/](#), [https://github.com/mbostock/d3/wiki/API-Reference](#)  
[http://code.google.com/apis/ajax/playground/](#)  
[http://www.edwardtufte.com/tufte/](#)  
[http://www.visualcomplexity.com/](#)  
[http://www.webdesignerdepot.com/2009/06/50-great-examples-of-data-visualization/](#)

Assoc Prof. Phuc Do, 2017

109

## Conclusion

- Data Visualization , Text visualization are effective tools to visualize and discover knowledge from data and text.
- Big data is a motion to promote knowledge discovery and data visualization is a good tool to support big data discovery
- Research in big data analysis is concurrently process with data and text visualization

Assoc Prof. Phuc Do, 2017

110

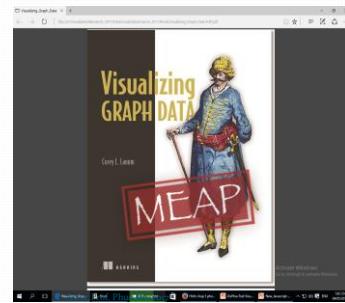
## Some Challenges in data and text Visualization

- Interactive Visualization of Very Large Graphs
- Develop novel interactive visualizations that are tightly coupled with automatic natural language processing methods to enable human users to explore their data.
- Develop Real-Time Visualization of Streaming Text
- How interactive visualization can assist investigative analysis for large corpus
- Optimizing the facilities of hardware and software for big data visualization and discovery.

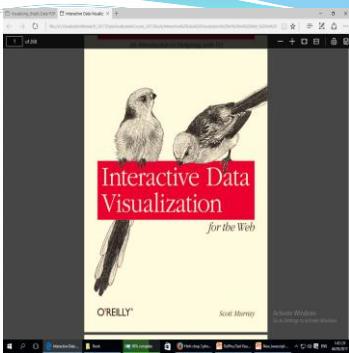
Assoc Prof. Phuc Do, 2017

111

## References

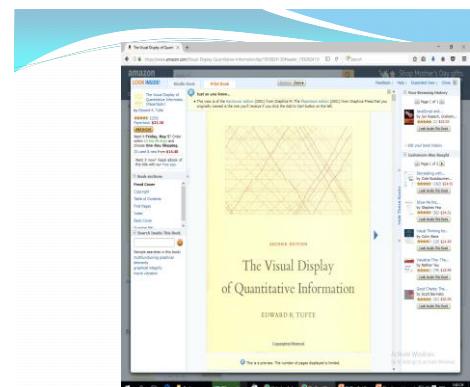


112



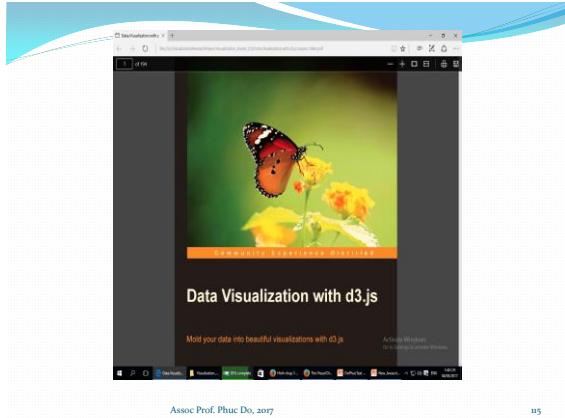
Assoc Prof. Phuc Do, 2017

113



Assoc Prof. Phuc Do, 2017

114



## Short Bio of Phuc Do

- Phuc Do is an Associate Professor of Information Systems at the University of Information Technology, Vietnam National University of Ho Chi Minh City.
- He currently works mostly on knowledge discovery technologies and their application in social network analysis by using topic modeling .
- Phuc Do has written about 80 research papers and books. He has been the Principal Investigator of several projects in Knowledge Discovery From Data, Text Analysis and Summarization, Topic Based Social Network Analysis, Citation Network Analysis.

Assoc Prof. Phuc Do, 2017

117

